Motion Forecasting Transformer with Global Intention Localization and Local Movement Refinement

Anonymous Author(s) Affiliation Address email

Abstract

Predicting multimodal future behavior of traffic participants is essential for robotic 1 vehicles to make safe decisions. Existing works explore to directly predict future 2 trajectories based on latent features or utilize dense goal candidates to identify 3 agent's destinations, where the former strategy converges slowly since all motion 4 modes are derived from the same feature while the latter strategy has efficiency 5 issue since its performance highly relies on the density of goal candidates. In 6 this paper, we propose the Motion TRansformer (MTR) framework that models 7 motion prediction as the joint optimization of global intention localization and 8 9 local movement refinement. Instead of using goal candidates, MTR incorporates 10 spatial intention priors by adopting a small set of learnable motion query pairs. Each motion query pair takes charge of trajectory prediction and refinement for a 11 specific motion mode, which stabilizes the training process and facilitates better 12 multimodal predictions. Experiments show that MTR achieves state-of-the-art 13 performance on both the marginal and joint motion prediction challenges, ranking 14 1^{st} on the leaderbaords of Waymo Open Motion Dataset. Code will be available. 15

16 1 Introduction

Motion forecasting is a fundamental task of modern autonomous driving systems. It has been receiving increasing attention in recent years [18, 43, 27, 52, 33] as it is crucial for robotic vehicles to understand driving scenes and make safe decisions. Motion forecasting requires to predict future behaviors of traffic participants by jointly considering the observed agent states and road maps, which is challenging due to inherently multimodal behaviors of the agent and complex scene environments.

To cover all potential future behaviors of the agent, existing approaches mainly fall into two different 22 lines: the goal-based methods and the direct-regression methods. The goal-based methods [18, 56] 23 adopt dense goal candidates to cover all possible destinations of the agent, predicting the probability 24 of each candidate being a real destination and then completing the full trajectory for each selected 25 candidate. Although these goal candidates alleviate the burden of model optimization by reducing 26 trajectory uncertainty, their density largely affects the performance of these methods: fewer candidates 27 will decrease the performance while more candidates will greatly increase computation and memory 28 cost. Instead of using goal candidates, the direct-regression methods [33, 44] directly predict a set 29 of trajectories based on the encoded agent feature, covering the agent's future behavior adaptively. 30 Despite the flexibility in predicting a broad range of agent behaviors, they generally converge slowly 31 as various motion modes are required to be regressed from the same agent feature without utilizing 32 any spatial priors. They also tend to predict the most frequent modes of training data since these 33 frequent modes dominate the optimization of the agent feature. In this paper, we present a unified 34 framework, namely Motion TRansformer (MTR), which takes the best of both types of methods. 35

³⁶ In our proposed MTR, we adopt a small set of novel motion query pairs to model motion prediction ³⁷ as the joint optimization of two tasks: The first global intention localization task aims to roughly ³⁸ identify agent's intention for achieving higher efficiency, while the second local movement refinement

task aims to adaptively refine each intention's predicted trajectory for achieving better accuracy. Our

⁴⁰ approach not only stabilizes the training process without depending on dense goal candidates but also

enables flexible and adaptive prediction by enabling local refinement for each motion mode.

Specifically, each motion query pair consists of two components, *i.e.*, a static intention query and a 42 dynamic searching query. The static intention queries are introduced for global intention localization, 43 where we formulate them based on a small set of spatially distributed intention points. Each static 44 intention query is the learnable positional embedding of an intention point for generating trajectory of 45 a specific motion mode, which not only stabilizes the training process by explicitly utilizing different 46 queries for different modes, but also eliminates the dependency on dense goal candidates by requiring 47 each query to take charge of a large region. The dynamic searching queries are utilized for local 48 movement refinement, where they are also initialized as the learnable embeddings of the intention 49 points but are responsible for retrieving fine-grained local features around each intention point. For 50 this purpose, the dynamic searching queries are dynamically updated according to the predicted 51 trajectories, which can adaptively gather latest trajectory features from a deformable local region for 52 iterative motion refinement. These two queries complement each other and have been empirically 53 demonstrated their great effectiveness in predicting multimodal future motion. Besides that, we also 54 propose a dense future prediction module. Existing works generally focus on modeling the agent 55 interaction over past trajectories while ignoring the future trajectories' interaction. To compensate for 56 such information, we adopt a simple auxiliary regression head to densely predict future trajectory 57 and velocity for each agent, which are encoded as additional future context features to benefit future 58 59 motion prediction of our interested agent. The experiments show that this simple auxiliary task works well and remarkably improves the performance of multimodal motion prediction. 60

Our contributions are three-fold: (1) We propose a novel motion decoder network with a new 61 concept of motion query pair, which adopts two types of queries to model motion prediction as joint 62 optimization of global intention localization and local movement refinement. It not only stabilizes 63 the training with mode-specific motion query pairs, but also enables adaptive motion refinement 64 by iteratively gathering fine-grained trajectory features. (2) We present an auxiliary dense future 65 prediction task to enable the future interactions between our interested agent and other agents. It 66 facilitates our framework to predict more scene-compliant trajectories for the interacting agents. (3) 67 By adopting these techniques, we propose MTR framework that explores transformer encoder-decoder 68 structure for multimodal motion prediction. Our approach achieves state-of-the-art performance 69 on both the marginal and joint motion prediction benchmarks of Waymo Open Motion Dataset 70 (WOMD) [13], outperforming previous best ensemble-free approaches with +8.48% mAP gains for 71 marginal motion prediction and +7.98% mAP gains for joint motion prediction. As of 19 May 2022, 72 our approach ranks 1st on both the marginal and joint motion prediction leaderboards of WOMD. 73

74 2 Related Work

Motion Prediction for Autonomous Driving. Recently, motion prediction has been extensively 75 studied due to the growing interest in autonomous driving, and it typically takes road map and 76 agent history states as input. To encode such scene context, early works [34, 29, 5, 11, 55, 3, 8] 77 typically rasterize them into an image so as to be processed with convolutional neural networks 78 (CNNs). LaneGCN [25] builds a lane graph toscalability capture map topology. VectorNet [15] is 79 widely adopted by recent works [18, 40, 33, 44] due to its efficiency and scalability, where both road 80 maps and agent trajectories are represented as polylines. We also adopt this vector representation, 81 but instead of building global graph of polylines, we propose to adopt transformer encoder on 82 local connected graph, which not only better maintains input locality structure but also is more 83 84 memory-efficient to enable larger map encoding for long-term motion prediction.

85 Given the encoded scene context features, existing works explore various strategies to model multimodal future motion. Early works [1, 19, 37, 41, 38] propose to generate a set of trajectory samples 86 to approximate the output distribution. Some other works [9, 20, 31, 35, 39] parameterize multi-87 modal predictions with Gaussian Mixture Models (GMMs) to generate compact distribution. HOME 88 series [17, 16] generate trajectories with sampling on a predicted heatmap. IntentNet [7] considers 89 intention prediction as a classification with 8 high level actions, while [27] proposes a region-based 90 training strategy. Goal-based methods [56, 38, 14, 28] are another kinds of models where they first 91 estimate several goal points of the agents and then complete full trajectory for each goal. 92



Figure 1: The architecture of MTR framework. (a) indicates the dense future prediction module, which predicts a single trajectory for each agent (*e.g.*, drawn as yellow dashed curves in the above of (a)). (b) indicates the dynamic map collection module, which collects map elements along each predicted trajectory (*e.g.*, drawn as the shadow region along each trajectory in the above part of (b)) to provide trajectory-specific feature for motion decoder network. (c) indicates the motion decoder network, where \mathcal{K} is the number of motion query pairs, T is the number of future frames, D is hidden feature dimension and N is the number of transformer decoder layers. The predicted trajectories, motion query pairs, and query content features are the outputs from last decoder layer and will be taken as input to next decoder layer. For the first decoder layer, both two components of motion query pair are initialized as predefined intention points, the predicted trajectories are replaced with the intention points for initial map collection, and query content features are initialized as zeros.

Recently, the large-scale Waymo Open Motion Dataset (WOMD) [13] is proposed for long-term 93 motion prediction. To address this challenge, DenseTNT [18] adopts a goal-based strategy to classify 94 endpoint of trajectory from dense goal points. Other works directly predict the future trajectories 95 based on the encoded agent features [33] or latent anchor embedding [44]. However, the goal-based 96 strategy has the efficiency concern due to a large number of goal candidates, while the direct-97 regression strategy converges slowly as the predictions of various motion modes are regressed from 98 the same agent feature. In contrast, our approach adopts a small set of learnable motion query pairs, 99 which not only eliminate the large number of goal candidates but also alleviate the optimization 100 burden by utilizing mode-specific motion query pairs for predicting different motion modes. 101

Some very recent works [42, 22, 21] also achieve top performance on WOMD by exploring Mix-and-Match block [42], a variant of MultiPath++ [22] or heterogeneous graph [21]. However, they generally focus on exploring various structures for encoding scene context, while how to design a better motion decoder for multimodal motion prediction is still underexplored. In contrast, our approach focuses on addressing this challenge with a novel transformer-based motion decoder network.

Transformer. Transformer [45] has been widely applied in natural language processing [10, 2] and computer vision [12, 47, 4, 46, 53]. Our approach is inspired by DETR [4] and its follow-up works [58, 30, 26, 54], especially DAB-DETR [26], where the object query is considered as the positional embedding of a spatial anchor box. Motivated by their great success in object detection, we introduce a novel concept of motion query pair to model multimodal motion prediction with prior intention points, where each motion query pair takes charge of predicting a specific motion mode and also enables iterative motion refinement by combining with transformer decoders.

114 **3** Motion TRansformer (MTR)

We propose Motion TRansformer (MTR), which adopts a novel transformer encoder-decoder structure with iterative motion refinement for predicting multimodal future motion. The overall structure is illustrated in Figure 1. In Sec. 3.1, we introduce our encoder network for scene context modeling. In Sec. 3.2, we present motion decoder network with a novel concept of motion query pair for predicting multimodal trajectories. Finally, in Sec. 3.3, we introduce the optimization process of our framework.

120 3.1 Transformer Encoder for Scene Context Modeling

The future behaviors of the agents highly depend on the agents' interaction and road map. To encode such scene context, existing approaches have explored various strategies by building global interacting graph [15, 18] or summarizing map features to agent-wise features [33, 44]. We argue that the locality structure is important for encoding scene context, especially for the road map. Hence, we propose a transformer encoder network with local self-attention to better maintain such structure information.

Input representation. We follow the vectorized representation [15] to organize both input trajectories 126 and road map as polylines. For the motion prediction of a interested agent, we adopt the agent-centric 127 strategy [56, 18, 44] that normalizes all inputs to the coordinate system centered at this agent. Then, 128 a simple polyline encoder is adopted to encode each polyline as an input token feature for the 129 transformer encoder. Specifically, we denote the history state of N_a agents as $A_{in} \in \mathbb{N}^{N_a \times t \times C_a}$, 130 where t is the number of history frames, C_a is the number of state information (e.g., location, heading 131 angle and velocity), and we pad zeros at the positions of missing frames for trajectories that have less than t frames. The road map is denoted as $M_{in} \in \mathbb{R}^{N_m \times n \times C_m}$, where N_m is the number of map 132 133 polylines, n is the number of points in each polyline and C_m is the number of attributes of each point 134 (e.g., location and road type). Both of them are encoded by a PointNet-like [36] polyline encoder as: 135

$$A_{\rm p} = \phi \left({\rm MLP}(A_{\rm in}) \right), \quad M_{\rm p} = \phi \left({\rm MLP}(M_{\rm in}) \right), \tag{1}$$

where $MLP(\cdot)$ is a multilayer perceptron network, and ϕ is max-pooling to summarize each polyline features as agent features $A_p \in \mathbb{R}^{N_a \times D}$ and map features $M_p \in \mathbb{R}^{N_m \times D}$ with feature dimension D.

Scene context encoding with local transformer encoder. The local structure of scene context is important for motion prediction. For example, the relation of two parallel lanes is important for modelling the motion of changing lanes, but adopting attention on global connected graph equally considers relation of all lanes. In contrast, we introduce such prior knowledge to context encoder by adopting local attention, which better maintains the locality structure and are more memory-efficient. Specifically, the attention module of *j*-th transformer encoder layer can be formulated as:

$$G^{j} = \text{MultiHeadAttn}\left(\text{query}=G^{j-1} + \text{PE}_{G^{j-1}}, \text{ key}=\kappa(G^{j-1}) + \text{PE}_{\kappa(G^{j-1})}, \text{ value}=\kappa(G^{j-1})\right), \quad (2)$$

where MultiHeadAttn (\cdot, \cdot, \cdot) is the multi-head attention layer [45], $G^0 = [A_p, M_p] \in \mathbb{N}^{(N_a + N_m) \times D}$ concatenating the features of agents and map, and $\kappa(\cdot)$ denotes *k*-nearest neighbor algorithm to find *k* closest polylines for each query polyline. PE denotes sinusoidal position encoding of input tokens, where we utilize the latest position for each agent and utilize polyline center for each map polyline. Thanks to such local self-attention, our framework can encode a much larger area of scene context. The encoder network finally generates both agent features $A_{past} \in \mathbb{R}^{N_a \times D}$ and map features $M \in \mathbb{R}^{N_m \times D}$, which are considered as the scene context inputs of the following decoder network.

Dense future prediction for future interactions. Interactions with other agents heavily affect behaviors of our interested agent, and the pioneer works propose to model the multi-agent interactions with hub-host based network [59], dynamic relational reasoning [24], social spatial-temporal network [51], etc. However, most existing works generally focus on learning such interactions over past trajectories while ignoring the interactions of future trajectories. Therefore, considering that the encoded features *A* have already learned rich context information of all agents, we propose to densely predict both future trajectories and velocities of all agents by adopting a simple regression head on *A*:

$$S_{1:T} = \mathsf{MLP}(A_{\mathsf{past}}),\tag{3}$$

where $S_i \in \mathbb{R}^{N_a \times 4}$ includes future position and velocity of each agent at time step i, and T is the 158 number of future frames to be predicted. The predicted trajectories $S_{1:T}$ are encoded by adopting 159 the same polyline encoder as Eq. (1) to encode the agents' future states as features $A_{\text{future}} \in \mathbb{R}^{N_a \times D}$ 160 which are then utilized to enhance the above features A by using a feature concatenation and three 161 MLP layers as $A = MLP([A_{past}, A_{future}])$. This auxiliary task provides additional future context 162 information to the decoder network, facilitating the model to predict more scene-compliant future 163 trajectories for the interested agent. The experiments in Table 3 demonstrates that this simple and 164 light-weight auxiliary task can effectively improve the performance of multimodal motion prediction. 165

166 **3.2 Transformer Decoder with Motion Query Pair**

Given the scene context features, a transformer-based motion decoder network is adopted for multimodal motion prediction, where we propose *motion query pair* to model motion prediction as the joint optimization of global intention localization and local movement refinement. Each motion query pair contains two types of queries, *i.e.*, static intention query and dynamic searching query, for conducting global intention localization and local movement refinement respectively. As shown ¹⁷² in Figure 2, our motion decoder network contains stacked transformer decoder layers for iteratively ¹⁷³ refining the predicted trajectories with motion query pairs. Next, we illustrate the detailed structure.

Global intention localization aims to localize agent's potential motion intentions in an efficient and effective manner. We propose *static intention query* to narrow down the uncertainty of future trajectory by utilizing different intention queries for different motion modes. Specifically, we generate

 \mathcal{K} representative intention points $I \in \mathbb{R}^{\mathcal{K} \times 2}$ 177 by adopting k-means clustering algorithm on 178 the endpoints of ground-truth (GT) trajecto-179 ries, where each intention point represents an 180 implicit motion mode that considers both mo-181 tion direction and velocity. We model each 182 static intention query as the learnable posi-183 tional embedding of the intention point as: 184

$$Q_I = \mathrm{MLP}\left(\mathrm{PE}(I)\right),\tag{4}$$

where $PE(\cdot)$ is the sinusoidal position encod-185 ing, and $Q_I \in \mathbb{R}^{\mathcal{K} \times D}$. Notably, each in-186 tention query takes charge of predicting tra-187 jectories for a specific motion mode, which 188 stabilizes the training process and facilitates 189 190 predicting multimodal trajectories since each 191 motion mode has their own learnable embedding. Thanks to their learnable and adaptive 192 properties, we only need a small number of 193 queries (e.g., 64 queries in our setting) for 194 efficient intention localization, instead of us-195 ing densely-placed goal candidates [56, 18] 196 to cover the destinations of the agents. 197



Figure 2: The network structure of our motion decoder network with motion query pair.

Local movement refinement aims to complement with global intention localization by iteratively gathering fine-grained trajectory features for refining the trajectories. We propose *dynamic searching query* to adaptively probe trajectory features for each motion mode. Each dynamic searching query is also the position embedding of a spatial point, which is initialized with its corresponding intention point but will be dynamically updated according to the predicted trajectory in each decoder layer. Specifically, given the predicted future trajectories $Y_{1:T}^j = \{Y_i^j \in \mathbb{R}^{K \times 2} \mid i = 1, \dots, T\}$ in *j*-th decoder layer, the dynamic searching query of (j + 1)-th decoder layer is updated as follows:

$$Q_S^{j+1} = \mathsf{MLP}\left(\mathsf{PE}(Y_T^j)\right). \tag{5}$$

As shown in Figure 3, for each motion query pair, we propose a *dynamic map collection* module to extract fine-grained trajectory features by querying map features from a trajectory-aligned local region, which is implemented by collecting L polylines whose centers are closest to the predicted trajectory. As the agent's behavior largely depends on road maps, this local movement refinement strategy enables to continually focus on latest local context information for iterative motion refinement.

Attention module with motion query pair. In each decoder layer, static intention query is utilized to propagate information among different motion intentions, while dynamic searching query is utilized to aggregate trajectory-specific features from scene context features. Specifically, we utilize static intention query as the position embedding of self-attention module as follows:

$$C_{\text{sa}}^{j} = \text{MultiHeadAttn}(\text{query}=C^{j-1} + Q_{I}, \text{ key}=C^{j-1} + Q_{I}, \text{ value}=Q_{I}),$$
(6)

where $C^{j-1} \in \mathbb{R}^{K \times D}$ is query content features from (j-1)-th decoder layer, C^0 is initialized to zeros, and $C_{sa}^j \in \mathbb{R}^{K \times D}$ is the updated query content. Next, we utilize dynamic searching query as query position embedding of cross attention to probe trajectory-specific features from the outputs of encoder. Inspired by [30, 26], we concatenate content features and position embedding for both query and key to decouple their contributions to the attention weights. Two cross-attention modules are adopted separately for aggregating features from both agent features A and map features M as:

$$C_{A}^{j} = \text{MultiHeadAttn}(\text{query}=[C_{\text{sa}}^{j}, Q_{S}^{j}], \text{ key}=[A, \text{PE}_{A}], \text{ value}=A),$$

$$C_{M}^{j} = \text{MultiHeadAttn}(\text{query}=[C_{\text{sa}}^{j}, Q_{S}^{j}], \text{ key}=[\alpha(M), \text{PE}_{\alpha(M)}], \text{ value}=\alpha(M)), \quad (7)$$

$$C^{j} = \text{MLP}([C_{A}^{j}, C_{M}^{j}])$$



Figure 3: The illustration of dynamic map collection module for iterative motion refinement.

where $[\cdot, \cdot]$ indicates feature concatenation, $\alpha(M)$ is the aforementioned dynamic map collection module to collect *L* trajectory-aligned map features for motion refinement. Note that for simplicity, in Eq. (6) and (7), we omit the residual connection and feed-forward network in transformer layer [45].

Finally, $C^j \in \mathbb{R}^{K \times D}$ is the updated query content features for each motion query pair in *j*-th layer.

Multimodal motion prediction with Gaussian Mixture Model. For each decoder layer, we append a prediction head to C^{j} for generating future trajectories. As the behaviors of the agents are highly multimodal, we follow [9, 44] to represent the distribution of predicted trajectories with Gaussian Mixture Model (GMM) at each time step. Specifically, for each future time step $i \in \{1, \dots, T\}$, we predict the probability p and parameters $(\mu_x, \mu_y, \sigma_x, \sigma_y, \rho)$ of each Gaussian component as follows

$$Z_{1:T}^{j} = \mathrm{MLP}(C^{j}), \tag{8}$$

where $Z_i^j \in \mathbb{R}^{\mathcal{K} \times 6}$ includes \mathcal{K} Gaussian components $\mathcal{N}_{1:\mathcal{K}}(\mu_x, \sigma_x; \mu_y, \sigma_y; \rho)$ with probability distribution $p_{1:\mathcal{K}}$. The predicted distribution of agent's position at time step *i* can be formulated as:

$$P_i^j(o) = \sum_{k=1}^{\mathcal{K}} p_k \cdot \mathcal{N}_k(o_x - \mu_x, \sigma_x; o_y - \mu_y, \sigma_y; \rho).$$
(9)

where $P_i^j(o)$ is the occurrence probability of the agent at spatial position $o \in \mathbb{R}^2$. The predicted trajectories $Y_{1:T}^j$ can be generated by simply extracting the predicted centers of Gaussian components.

233 **3.3 Training Losses**

Our model is trained end-to-end with two training losses. The first auxiliary loss is L1 regression loss 234 to optimize the outputs of Eq. (3). For the second Gaussian regression loss, we adopt negative log-235 likelihood loss according to Eq. (9) to maximum the likelihood of ground-truth trajectory. Inspired 236 by [9, 44], we adopt a hard-assignment strategy that selects one closest motion query pair as positive 237 Gaussian component for optimization, where the selection is implemented by calculating the distance 238 between each intention point and the endpoint of GT trajectory. The Gaussian regression loss is 239 adopted in each decoder layer, and the final loss is the sum of the auxiliary regression loss and all the 240 Gaussian regression loss with equal loss weights. Please refer to appendix for more loss details. 241

242 **4** Experiments

243 4.1 Experimental Setup

Dataset and metrics. We evaluate our approach on the large-scale Waymo Open Motion Dataset 244 (WOMD) [13], which mines interesting interactions from real-world traffic scenes and is currently 245 the most diverse interactive motion dataset. There are two tasks in WOMD with separate evaluation 246 metrics: (1) The marginal motion prediction challenge that independently evaluates the predicted 247 motion of each agent (up to 8 agents per scene). (2) The joint motion prediction challenge that needs 248 to predict the joint future positions of 2 interacting agents for evaluation. Both of them provide 1 249 second of history data and aim to predict 6 marginal or joint trajectories of the agents for 8 seconds 250 into the future. There are totally 487k training scenes, and about 44k validation scenes and 44k testing 251 scenes for each challenge. We utilize the official evaluation tool to calculate the evaluation metrics, 252 where the mAP and miss rate are the most important ones as in the official leaderboard [49, 48]. 253

Table 1: Performance comparison of marginal motion prediction on the validation and test set of Waymo Open Motion Dataset. †: The results are shown in *italic* for reference since their performance is achieved with model ensemble techniques. We only evaluate our default setting MTR on the test set by submitting to official test server due to the limitation of submission times of WOMD.

	Method	Reference	minADE \downarrow	minFDE \downarrow	Miss Rate ↓	mAP ↑
Test	MotionCNN [23]	CVPRw 2021	0.7400	1.4936	0.2091	0.2136
	ReCoAt [57]	CVPRw 2021	0.7703	1.6668	0.2437	0.2711
	DenseTNT [18]	ICCV 2021	1.0387	1.5514	0.1573	0.3281
	SceneTransformer [33]	ICLR 2022	0.6117	1.2116	0.1564	0.2788
	MTR (Ours)	-	0.6050	1.2207	0.1351	0.4129
	[†] MultiPath++ [44]	ICRA 2022	0.5557	1.1577	0.1340	0.4092
	[†] MTR-Advanced-ens (Ours)	-	0.5640	1.1344	0.1160	0.4492
Val	MTR (Ours)	-	0.6046	1.2251	0.1366	0.4164
	MTR-e2e (Ours)	-	0.5160	1.0404	0.1234	0.3245
	[†] MTR-ens (Ours)	-	0.5686	1.1534	0.1240	0.4323
	[†] MTR-Advanced-ens (Ours)	-	0.5597	1.1299	0.1167	0.4551

Table 2: Performance comparison of joint motion prediction on the interactive validation and test set of Waymo Open Motion Dataset.

-	Method	Reference	minADE \downarrow	minFDE \downarrow	Miss Rate ↓	mAP↑
Test	Waymo LSTM baseline [13]	ICCV 2021	1.9056	5.0278	0.7750	0.0524
	HeatIRm4 [32]	CVPRw 2021	1.4197	3.2595	0.7224	0.0844
	AIR^2 [50]	CVPRw 2021	1.3165	2.7138	0.6230	0.0963
	SceneTransformer [33]	ICLR 2022	0.9774	2.1892	0.4942	0.1192
	M2I [40]	CVPR 2022	1.3506	2.8325	0.5538	0.1239
	MTR (Ours)	-	0.9181	2.0633	0.4411	0.2037
Val	MTR (Ours)	-	0.9132	2.0536	0.4372	0.1992

Implementation details. For the context encoding, we stack 6 transformer encoder layers. The 254 road map is represented as multiple polylines, where each polyline contains up to 20 map points 255 (about 10m in WOMD). We select $N_m = 768$ nearest map polylines around the interested agent. 256 The number of neighbors in encoder's local self-attention is set to 16. The hidden feature dimension 257 is set as D = 256. For the decoder modules, we stack 6 decoder layers. L is set to 128 to collect the 258 closest map polylines from context encoder for motion refinement. By default, we utilize 64 motion 259 query pairs where their intention points are generated by conducting k-means clustering algorithm on 260 the training set. To generate 6 future trajectories for evaluation, we use non-maximum suppression 261 (NMS) to select top 6 predictions from 64 predicted trajectories by calculating the distances between 262 their endpoints, and the distance threshold is set as 2.5m. Please refer to Appendix for more details. 263

Training details. Our model is trained in an end-to-end manner by AdamW optimizer with a learning
rate of 0.0001 and batch size of 80 scenes. We train the model for 60 epochs with 8 GPUs (NVDIA
RTX 8000), and the learning rate is decayed by a factor of 0.5 every 5 epochs from epoch 30. The
weight decay is set as 0.01 and we do not use any data augmentation.

MTR-e2e for end-to-end motion prediction. We also propose an end-to-end variant of MTR, called MTR-e2e, where only 6 motion query pairs are adopted so as to remove NMS post processing. In the training process, instead of using static intention points for target assignment as in MTR, MTR-e2e selects positive mixture component by calculating the distances between its 6 predicted trajectories and the GT trajectory, since 6 intention points are too sparse to well cover all potential future motions.

273 4.2 Main Results

Performance comparison for marginal motion prediction. Table 1 shows our main results for marginal motion prediction, our MTR outperforms previous ensemble-free approaches [18, 33] with remarkable margins, increasing the mAP by +8.48% and decreasing the miss rate from 15.64\% to 13.51\%. In particular, our single-model results of MTR also achieve better mAP than the latest work MultiPath++ [44], where it uses a novel model ensemble strategy that boosts its performance.

Table 1 also shows the comparison of MTR variants. MTR-e2e achieves better minADE and minFDE by removing NMS post-processing, while MTR achieves better mAP since it learns explicit meaning of each motion query pair that produces more confident intention predictions. We also propose a simple model ensemble strategy to merge the predictions of MTR and MTR-e2e and utilize NMS to remove redundant predictions (denoted as MTR-ens), and it takes the best of both models and achieves much better mAP. By adopting such ensemble strategy to 7 variants of our framework (*e.g.*,

Table 3: Effects of different components in MTR framework. All models share the same encoder network. "latent learnable embedding" indicates using 6 latent learnable embeddings as queries of decoder network, and "iterative refinement" indicates using 6 stacked decoders for motion refinement.

<i><i><i>αι ι τ ι</i></i></i>	• •	* / * *	~ ~					
Global Intention	Iterative	Local Movement	Dense Future		minEDE	Miss Data	AD A	
Localization	Refinement	Refinement	efinement Prediction min/		IIIIIIFDE↓	wiiss Kate \downarrow	mar	
Latent learnable embedding	×	×	×	0.6829	1.4841	0.2128	0.2633	
Static intention query	×	×	×	0.7036	1.4651	0.1845	0.3059	
Static intention query	\checkmark	×	×	0.6919	1.4217	0.1776	0.3171	
Static intention query	\checkmark	\checkmark	×	0.6833	1.4059	0.1756	0.3234	
Static intention query	\checkmark	×	\checkmark	0.6735	1.3847	0.1706	0.3284	
Static intention query	\checkmark	\checkmark	\checkmark	0.6697	1.3712	0.1668	0.3437	



Table 4: Effects of local self-attention in transformer encoder. "#polyline" is the number of input map polylines used for context encoding, and a large number of polylines indicate that there is a larger map context around the interested agent. "OOM" indicates running out of memory.

Attention	#Polyline	$\min ADE \downarrow$	$minFDE \downarrow$	$MR\downarrow$	$mAP \uparrow$
Global	256	0.683	1.4031	0.1717	0.3295
Global	512	0.6783	1.4018	0.1716	0.3280
Global	768	OOM	OOM	OOM	OOM
Local	256	0.6724	1.3835	0.1683	0.3372
Local	512	0.6707	1.3749	0.1670	0.3392
Local	768	0.6697	1.3712	0.1668	0.3437
Local	1024	0.6757	1.3782	0.1663	0.3452

Figure 4: MTR framework with different number of motion query pairs, and two different colored lines demonstrate different strategies for selecting the positive mixture component during training process.

more decoder layers, different number of queries, larger hidden dimension), our advanced ensemble results (denoted as MTR-Advanced-ens) achieve best performance on the test set leaderboard.

Performance comparison for joint motion prediction. To evaluate our approach for joint motion 287 prediction, we combine the marginal predictions of two interacting agents into joint prediction as 288 in [6, 13, 40], where we take the top 6 joint predictions from 36 combinations of these two agents. 289 The confidence of each combination is the product of marginal probabilities. Table 2 shows that 290 our approach outperforms state-of-the-arts [33, 40] with large margins on all metrics. Particularly, 291 our MTR boosts the mAP from 12.39% to 20.37% and decreases the miss rate from 49.42% to 292 44.11%. The remarkable performance gains demonstrate the effectiveness of MTR for predicting 293 scene-consistent future trajectories. Besides that, we also provide some qualitative results in Figure 5 294 to show our predictions in complicated interacting scenarios. 295

As of May 19, 2022, our MTR ranks 1^{st} on the motion prediction leaderboard of WOMD for both two challenges [49, 48]. The significant improvements manifest the effectiveness of MTR framework.

298 4.3 Ablation Study

We study the effectiveness of each component in MTR. For efficiently conducting ablation experiments, we uniformly sampled 20% frames (about 97*k* scenes) from the WOMD training set according to their default order, and we empirically find that it has similar distribution with the full training set. All models are evaluated with marginal motion prediction metric on the validation set of WOMD.

Effects of the motion decoder network. We study the effectiveness of each component in our 303 decoder network, including global intention localization, iterative refinement and local movement 304 refinement. Table 3 shows that all components contributes remarkably to the final performance 305 in terms of the official ranking metric mAP. Especially, our proposed static intention queries with 306 intention points achieves much better mAP (*i.e.*, +4.26%) than the latent learnable embeddings 307 thanks to its mode-specific querying strategy, and both the iterative refinement and local movement 308 refinement strategy continually improve the mAP from 30.59% to 32.34% by aggregating more 309 fine-grained trajectory features for motion refinement. 310

Effects of dense future prediction. Table 3 shows that our proposed dense future prediction module significantly improves the quality of predicted trajectories (*e.g.*, +1.78% mAP), which verifies that future interactions of the agents' trajectories are important for motion prediction and our proposed strategy can learn such interactions to predict more reliable trajectories.







(a) V2 is passing the intersection to turn left with high speed. Our model predicts multimodal crosswalk while **V1** is on the right-turn lane behaviors for V1: turn left or make a U-turn. In any case, V1 is predicted to yield for V2.

(b) P2 is passing the road through the to turn right. Both V1 and V3 are predicted to yield for P2.

(c) Our model predicts multimodal behaviors for V1: go straight and turn right, since it still has a distance to the intersection. V2 is predicted to yield for V1 when turning left, since V1 is moving fast towards the intersection.

Figure 5: Qualitative results of MTR framework on WOMD. There are two interested agents in each scene (green rectangle), where our model predicts 6 multimodal future trajectories for each of them. For other agents (blue rectangle), a single trajectory is predicted by dense future prediction module. We use gradient color to visualize the trajectory waypoints at different future time step, and trajectory confidence is visualized by setting different transparent. Abbreviation: Vehicle (V), Pedestrian (P).

Effects of local attention for context encoding. Table 4 shows that by taking the same number of 315 map polylines as input, local self-attention in transformer encoder achieves better performance than 316 global attention (*i.e.*, +0.77% mAP for 256 polylines and +1.12% mAP for 512 polylines), which 317 verifies that the input local structure is important for motion prediction and introducing such prior 318 knowledge with local attention can benefit the performance. More importantly, local attention is more 319 memory-efficient and the performance keeps growing when improving the number of map polylines 320 from 256 to 1,024, while global attention will run out of memory due to its quadratic complexity. 321

Effects of the number of motion query pairs with different training strategies. As mentioned 322 323 before, during training process, MTR and MTR-e2e adopt two different strategies for assigning 324 positive mixture component, where MTR depends on static intention points (denoted as α) while MTR-e2e utilizes predicted trajectories (denoted as β). Figure 4 investigates the effects of the number 325 of motion query pairs under these two strategies, where we have the following observations: (1) 326 When increasing the number of motion query pairs, strategy α achieves much better mAP and miss 327 rate than strategy β . Because intention query points can ensure more stable training process since 328 each intention query points is responsible to a specific motion mode. In contrast, strategy β depends 329 on unstable predictions and the positive component may randomly switch among all components, 330 so a large number of motion query pairs are hard to be optimized with strategy β . (2) The explicit 331 meaning of each intention query point also illustrates the reason that strategy α consistently achieves 332 much better mAP than strategy β , since it can predict trajectories with more confident scores to 333 benefit mAP metric. (3) From another side, when decreasing the number of motion query pairs, the 334 miss rate of strategy α greatly increases, since a limit number of intention query points can not well 335 cover all potential motions of agents. Conversely, strategy β works well for a small number of motion 336 query pairs since its queries are not in charge of specific region and can globally adapt to any region. 337

5 Conclusion 338

In this paper, we present MTR, a novel framework for multimodal motion prediction. The motion 339 query pair is defined to model motion prediction as the joint optimization of global intention localiza-340 tion and local movement refinement. The global intention localization adopts a small set of learnable 341 static intention queries to efficiently capture agent's motion intentions, while the local movement re-342 finement conducts iterative motion refinement by continually probing fine-grained trajectory features. 343 The experiments on both marginal and joint motion prediction challenges of large-scale WOMD 344 dataset show that our approach achieves state-of-the-art performance. 345

Limitations. The proposed framework adopts an agent-centric strategy to predict multimodal 346 future trajectories for one interested agent, which leads to redundant context encoding if there are 347 multiple interested agents in the same scene. Although the dense future prediction module partially 348 compensates for this limitation, it can only predict a single future trajectory for each agent. Hence, how 349 to develop a joint motion prediction framework that can simultaneously predict multimodal motion for 350 multiple agents is one important future work. Besides, the rule-based NMS post-processing can result 351 in suboptimal predictions for minADE and minFDE metrics, and how to develop a learning-based 352 353 module to produce a required number of future trajectories (e.g., 6 trajectory) from full multimodal predictions (e.g., 64 predictions) is also worth exploring for a more robust framework. 354

355 References

- [1] Alexandre Alahi, Kratarth Goel, Vignesh Ramanathan, Alexandre Robicquet, Li Fei-Fei, and Silvio
 Savarese. Social lstm: Human trajectory prediction in crowded spaces. In *CVPR*, 2016.
- [2] Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. In *ICLR*, 2022.
- [3] Yuriy Biktairov, Maxim Stebelev, Irina Rudenko, Oleh Shliazhko, and Boris Yangel. Prank: motion
 prediction based on ranking. In *NeurIPS*, 2020.
- [4] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey
 Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020.
- [5] Sergio Casas, Cole Gulino, Renjie Liao, and Raquel Urtasun. Spagnn: Spatially-aware graph neural
 networks for relational behavior forecasting from sensor data. In *ICRA*, 2020.
- [6] Sergio Casas, Cole Gulino, Simon Suo, Katie Luo, Renjie Liao, and Raquel Urtasun. Implicit latent
 variable model for scene-consistent motion forecasting. In *ECCV*, 2020.
- [7] Sergio Casas, Wenjie Luo, and Raquel Urtasun. Intentnet: Learning to predict intention from raw sensor
 data. In *CoRL*, 2018.
- [8] Sergio Casas, Abbas Sadat, and Raquel Urtasun. Mp3: A unified model to map, perceive, predict and plan.
 In *CVPR*, 2021.
- [9] Yuning Chai, Benjamin Sapp, Mayank Bansal, and Dragomir Anguelov. Multipath: Multiple probabilistic anchor trajectory hypotheses for behavior prediction. In *CoRL*, 2019.
- Io] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *arXiv preprint arXiv:1810.04805*, 2018.
- [11] Nemanja Djuric, Vladan Radosavljevic, Henggang Cui, Thi Nguyen, Fang-Chieh Chou, Tsung-Han Lin,
 Nitin Singh, and Jeff Schneider. Uncertainty-aware short-term motion prediction of traffic actors for
 autonomous driving. In WACV, 2020.
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and
 Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*,
 2021.
- [13] Scott Ettinger, Shuyang Cheng, Benjamin Caine, Chenxi Liu, Hang Zhao, Sabeek Pradhan, Yuning Chai,
 Ben Sapp, Charles R Qi, Yin Zhou, et al. Large scale interactive motion forecasting for autonomous
 driving: The waymo open motion dataset. In *ICCV*, 2021.
- [14] Liangji Fang, Qinhong Jiang, Jianping Shi, and Bolei Zhou. Tpnet: Trajectory proposal network for motion
 prediction. In *CVPR*, 2020.
- [15] 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 *CVPR*, 2020.
- [16] Thomas Gilles, Stefano Sabatini, Dzmitry Tsishkou, Bogdan Stanciulescu, and Fabien Moutarde. Gohome:
 Graph-oriented heatmap output for future motion estimation. In *arXiv preprint arXiv:2109.01827*, 2021.
- [17] Thomas Gilles, Stefano Sabatini, Dzmitry Tsishkou, Bogdan Stanciulescu, and Fabien Moutarde. Home:
 Heatmap output for future motion estimation. In *ITSC*, 2021.
- [18] Junru Gu, Chen Sun, and Hang Zhao. Densetnt: End-to-end trajectory prediction from dense goal sets. In
 ICCV, 2021.
- [19] Agrim Gupta, Justin Johnson, Li Fei-Fei, Silvio Savarese, and Alexandre Alahi. Social gan: Socially
 acceptable trajectories with generative adversarial networks. In *CVPR*, 2018.
- ³⁹⁷ [20] Joey Hong, Benjamin Sapp, and James Philbin. Rules of the road: Predicting driving behavior with a ³⁹⁸ convolutional model of semantic interactions. In *CVPR*, 2019.
- Xiaosong Jia, Penghao Wu, Li Chen, Hongyang Li, Yu Liu, and Junchi Yan. Hdgt: Heterogeneous
 driving graph transformer for multi-agent trajectory prediction via scene encoding. In *arXiv preprint arXiv:2205.09753*, 2022.
- 402 [22] Stepan Konev. Mpa: Multipath++ based architecture for motion prediction. In *arXiv preprint* 403 *arXiv:2206.10041*, 2022.

- [23] Stepan Konev, Kirill Brodt, and Artsiom Sanakoyeu. Motionenn: A strong baseline for motion prediction in autonomous driving. In *Workshop on Autonomous Driving, CVPR*, 2021.
- [24] Jiachen Li, Fan Yang, Masayoshi Tomizuka, and Chiho Choi. Evolvegraph: Multi-agent trajectory
 prediction with dynamic relational reasoning. In *NeurIPS*, 2020.
- [25] Ming Liang, Bin Yang, Rui Hu, Yun Chen, Renjie Liao, Song Feng, and Raquel Urtasun. Learning lane
 graph representations for motion forecasting. In *ECCV*, 2020.
- [26] Shilong Liu, Feng Li, Hao Zhang, Xiao Yang, Xianbiao Qi, Hang Su, Jun Zhu, and Lei Zhang. DAB-DETR:
 Dynamic anchor boxes are better queries for DETR. In *ICLR*, 2022.
- 412 [27] Yicheng Liu, Jinghuai Zhang, Liangji Fang, Qinhong Jiang, and Bolei Zhou. Multimodal motion prediction
 413 with stacked transformers. In *CVPR*, 2021.
- [28] Karttikeya Mangalam, Harshayu Girase, Shreyas Agarwal, Kuan-Hui Lee, Ehsan Adeli, Jitendra Malik,
 and Adrien Gaidon. It is not the journey but the destination: Endpoint conditioned trajectory prediction. In
 ECCV, 2020.
- Francesco Marchetti, Federico Becattini, Lorenzo Seidenari, and Alberto Del Bimbo. Mantra: Memory
 augmented networks for multiple trajectory prediction. In *CVPR*, 2020.
- [30] Depu Meng, Xiaokang Chen, Zejia Fan, Gang Zeng, Houqiang Li, Yuhui Yuan, Lei Sun, and Jingdong
 Wang. Conditional detr for fast training convergence. In *ICCV*, 2021.
- [31] 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 *ICRA*,
 2020.
- 424 [32] Xiaoyu Mo, Zhiyu Huang, and Chen Lv. Multi-modal interactive agent trajectory prediction using 425 heterogeneous edge-enhanced graph attention network. In *Workshop on Autonomous Driving, CVPR*, 2021.
- [33] Jiquan Ngiam, Vijay Vasudevan, Benjamin Caine, Zhengdong Zhang, Hao-Tien Lewis Chiang, Jeffrey
 Ling, Rebecca Roelofs, Alex Bewley, Chenxi Liu, Ashish Venugopal, et al. Scene transformer: A unified
 architecture for predicting future trajectories of multiple agents. In *ICLR*, 2022.
- [34] Seong Hyeon Park, Gyubok Lee, Jimin Seo, Manoj Bhat, Minseok Kang, Jonathan Francis, Ashwin
 Jadhav, Paul Pu Liang, and Louis-Philippe Morency. Diverse and admissible trajectory forecasting through
 multimodal context understanding. In *ECCV*, 2020.
- [35] Tung Phan-Minh, Elena Corina Grigore, Freddy A Boulton, Oscar Beijbom, and Eric M Wolff. Covernet:
 Multimodal behavior prediction using trajectory sets. In *CVPR*, 2020.
- [36] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d
 classification and segmentation. In *CVPR*, 2017.
- [37] Nicholas Rhinehart, Kris M Kitani, and Paul Vernaza. R2p2: A reparameterized pushforward policy for
 diverse, precise generative path forecasting. In *ECCV*, 2018.
- [38] Nicholas Rhinehart, Rowan McAllister, Kris Kitani, and Sergey Levine. Precog: Prediction conditioned on
 goals in visual multi-agent settings. In *ICCV*, 2019.
- [39] Tim Salzmann, Boris Ivanovic, Punarjay Chakravarty, and Marco Pavone. Trajectron++: Dynamically feasible trajectory forecasting with heterogeneous data. In *ECCV*, 2020.
- [40] Qiao Sun, Xin Huang, Junru Gu, Brian C Williams, and Hang Zhao. M2i: From factored marginal
 trajectory prediction to interactive prediction. In *CVPR*, 2022.
- [41] Charlie Tang and Russ R Salakhutdinov. Multiple futures prediction. In *NeurIPS*, 2019.
- [42] Xiaocheng Tang, Soheil Sadeghi Eshkevari, Haoyu Chen, Weidan Wu, Wei Qian, and Xiaoming
 Wang. Golfer: Trajectory prediction with masked goal conditioning mnm network. In *arXiv preprint arXiv:2207.00738*, 2022.
- [43] Ekaterina Tolstaya, Reza Mahjourian, Carlton Downey, Balakrishnan Vadarajan, Benjamin Sapp, and
 Dragomir Anguelov. Identifying driver interactions via conditional behavior prediction. In *ICRA*, 2021.
- [44] Balakrishnan Varadarajan, Ahmed Hefny, Avikalp Srivastava, Khaled S Refaat, Nigamaa Nayakanti, Andre
 Cornman, Kan Chen, Bertrand Douillard, Chi Pang Lam, Dragomir Anguelov, et al. Multipath++: Efficient
 information fusion and trajectory aggregation for behavior prediction. In *ICRA*, 2022.

- [45] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
 Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- [46] Jiashun Wang, Huazhe Xu, Medhini Narasimhan, and Xiaolong Wang. Multi-person 3d motion prediction
 with multi-range transformers. In *NeurIPS*, 2021.
- [47] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *CVPR*,
 2018.
- [48] Waymo. Waymo open dataset interaction prediction challenge 2021. https://waymo.com/open/
 challenges/2021/interaction-prediction/, 2021. Accessed: 2022-05-19.
- [49] Waymo. Waymo open dataset motion prediction challenge 2021. https://waymo.com/open/
 challenges/2021/motion-prediction, 2021. Accessed: 2022-05-19.
- [50] David Wu and Yunnan Wu. Air2 for interaction prediction. In *Workshop on Autonomous Driving, CVPR*,
 2021.
- Yi Xu, Dongchun Ren, Mingxia Li, Yuehai Chen, Mingyu Fan, and Huaxia Xia. Tra2tra: Trajectory-to trajectory prediction with a global social spatial-temporal attentive neural network. In *RA-L*, 2021.
- 467 [52] Maosheng Ye, Tongyi Cao, and Qifeng Chen. Tpcn: Temporal point cloud networks for motion forecasting.
 468 In *CVPR*, 2021.
- [53] Fangao Zeng, Bin Dong, Tiancai Wang, Xiangyu Zhang, and Yichen Wei. Motr: End-to-end multiple-object
 tracking with transformer. In *arXiv preprint arXiv:2105.03247*, 2021.
- [54] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M Ni, and Heung-Yeung Shum.
 Dino: Detr with improved denoising anchor boxes for end-to-end object detection. In *arXiv preprint arXiv:2203.03605*, 2022.
- Yifan Zhang, Jinghuai Zhang, Jindi Zhang, Jianping Wang, Kejie Lu, and Jeff Hong. A novel learning
 framework for sampling-based motion planning in autonomous driving. In *AAAI*, 2020.
- [56] Hang Zhao, Jiyang Gao, Tian Lan, Chen Sun, Benjamin Sapp, Balakrishnan Varadarajan, Yue Shen,
 Yi Shen, Yuning Chai, Cordelia Schmid, et al. Tnt: Target-driven trajectory prediction. In *CoRL*, 2020.
- [57] Chen Lv Zhiyu Huang, Xiaoyu Mo. Recoat: A deep learning framework with attention mechanism for
 multi-modal motion prediction. In *Workshop on Autonomous Driving, CVPR*, 2021.
- [58] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable
 transformers for end-to-end object detection. In *ICLR*, 2021.
- [59] Yanliang Zhu, Deheng Qian, Dongchun Ren, and Huaxia Xia. Starnet: Pedestrian trajectory prediction
 using deep neural network in star topology. In *IROS*, 2019.

484 Checklist

486

487

488

489 490

491

492

494

495

- 485 1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes] See Section 5.
 - (c) Did you discuss any potential negative societal impacts of your work? [No] Our work is only for academic research purpose.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 493 2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [N/A]
 - (b) Did you include complete proofs of all theoretical results? [N/A]
- 496 3. If you ran experiments...
- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] Our implementation details are posted on Section 4.1 and the supplemental materials. The data is public dataset and our code will be available.

501 502 503	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See Section 4.1 for training details, and the data splits of ablation experiments are mentioned in Section 4.3.	
504 505	(c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [No] We fix the random seed for all our experiments.	
506 507	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See Section 4.1.	
508	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets	
509	(a) If your work uses existing assets, did you cite the creators? [Yes] See our reference.	
510	(b) Did you mention the license of the assets? [No] They are publicly available for	
511	academic use.	
512	(c) Did you include any new assets either in the supplemental material or as a URL? [No]	
513 514	(d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [Yes] They are publicly available for academic use.	
515 516 517	(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No] We only use the public Waymo Open Motion Dataset.	
518	5. If you used crowdsourcing or conducted research with human subjects	
519 520	 (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] 	
521 522	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]	
523 524	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]	