# WcDT: World-centric Diffusion Transformer for Traffic Scene Generation

Chen Yang<sup>1\*</sup>, Yangfan He<sup>2\*</sup>, Aaron Xuxiang Tian<sup>3</sup>, Dong Chen<sup>4</sup>, Jianhui Wang<sup>5</sup>, Tianyu Shi<sup>6</sup>, Arsalan Heydarian<sup>7</sup>, Pei Liu<sup>8⊠</sup>

Abstract— In this paper, we introduce a novel approach for autonomous driving trajectory generation by harnessing the complementary strengths of diffusion probabilistic models (a.k.a., diffusion models) and transformers. Our proposed framework, termed the "World-centric Diffusion Transformer"(WcDT), optimizes the entire trajectory generation process, from feature extraction to model inference. To enhance the scene diversity and stochasticity, the historical trajectory data is first preprocessed into "Agent Move Statement" and encoded into latent space using Denoising Diffusion Probabilistic Models (DDPM) enhanced with Diffusion with Transformer (DiT) blocks. Then, the latent features, historical trajectories, HD map features, and historical traffic signal information are fused with various transformer-based encoders that is used to enhance the interaction of agents with other elements in the traffic scene. The encoded traffic scenes are then decoded by a trajectory decoder to generate multimodal future trajectories. Comprehensive experimental results show that the proposed approach exhibits superior performance in generating both realistic and diverse trajectories, showing its potential for integration into automatic driving simulation systems. Our code is available at https://github.com/yangchen1997/WcDT.

# I. INTRODUCTION

Autonomous driving is a transformative technology aimed at reducing driver fatigue and traffic congestion by enabling autonomous vehicle operation [1], [2], [3], [4]. Developing these algorithms involves iterative optimization for safety and performance [5], [6], [7], but real-world testing poses challenges due to time constraints, safety concerns, regulatory hurdles, and high costs [8], [9]. Simulators play a vital role in the cost-effective testing and evaluation of autonomous driving systems (ADS) by providing controllable environments [10], [11]. To be effective, they must realistically replicate traffic scenarios and driver behaviors [12]. Current simulators rely on replaying driving logs or

- <sup>1</sup>Department of Computer Science and Informatics, Cardiff University. yc19970530@gmail.com
- <sup>2</sup>College of Libera Arts, University of Minnesota Twin Cities. he000577@umn.edu
- <sup>3</sup>Independent researcher, USA. aarontian00@gmail.com
- <sup>4</sup>Agricultural & Biological Engineering, Mississippi State University. dc2528@msstate.edu
- <sup>5</sup>Information and Software Engineering, University of Electronic Science and Technology of China. jianhuiwang@std.uestc.edu.cn
- <sup>6</sup>Transportation Research Institute. University of Toronto.
- ty.shi@mail.utoronto.ca <sup>7</sup>Link Lab & Civil and Environmental Engineering, University of Virginia. heydarian@virginia.edu
- <sup>8</sup>Intelligent Transportation Thrust, The Kong Hong University of Science and Technology (Guangzhou). pliu061@connect.hkust-gz.edu.cn



Fig. 1. World-centric model and agent-centric model: (a) The conventional "Agent-centric model" is common in trajectory prediction, including on the Sim Agents leaderboard. (b) Our approach replaces complex coordinate transformations with position embeddings, enhancing efficiency in multiscenario, multi-agent trajectory generation.

heuristic controllers [13], [12], limiting diversity and unpredictability in real-world behavior, which impacts ADS validation [14]. Multimodal motion prediction approaches [15], [16], [17], [18], [19], [20] have shown promise in traffic scene generation but struggle to generate diverse actions for all agents using comprehensive global information [15]. Generative adversarial networks (GANs) and Variational Auto-Encoders (VAEs) have been applied to traffic scene generation but face limitations. These models often lack diversity, reflecting training data distributions [12], and GANs suffer from unstable adversarial training [14]. Additionally, they fail to capture agent trajectory smoothness, leading to unrealistic results [14], and typically focus on individual vehicle paths, neglecting all agents. Recently, diffusion models have emerged as a promising alternative for diverse traffic scenarios [21], [12], [14], [22], treating generation as an inverse diffusion process. However, these models often require agent-centric Cartesian coordinates [15] and generate only one trajectory per agent per inference. In this paper, we propose a novel framework for traffic scene generation tailored to autonomous driving, leveraging diffusion models and transformer-based encoder-decoder architectures. Our "World-centric Diffusion Transformer" (WcDT) framework optimizes trajectory generation from feature extraction to inference, enabling coherent and joint future movements for various agents in a single inference. Our contributions include:

• A new paradigm for simultaneous, consistent future movement generation for all agents in a single inference.

<sup>\*</sup> Equal contribution.

Corresponding author.

- A Diffusion-Transformer module that enhances scene diversity and stochasticity, integrating the world state efficiently.
- Benchmarking performance for realism and diversity in trajectory generation, validated on open traffic datasets.

## II. RELATED WORK

## A. Motion prediction-based methods

Recent developments in traffic scene generation utilize motion prediction methodologies to enhance realism in multimodal scenarios [15], [16], [17], [18], [19], [20]. For example, Multipath++ [16] advances its predecessor by integrating a context-aware fusion approach with Gaussian mixture models for more precise trajectory predictions. Similarly, Trafficsim [18] employs an implicit latent variable model for simulating multi-agent interactions. Transformer-based encoder-decoder architectures are also central to motion prediction [17], [19], [20]. The Scene Transformer [17] encodes interactions among agents using a global coordinate frame, enabling joint behavior prediction. The Motion Transformer (MTR) [19] optimizes both global intention and local movement, achieving top rankings in the Waymo Open Motion Dataset [23]. Leading the Waymo Open Sim Agents Challenge (WOSAC), the Multiverse Transformer (MVTA) introduces novel training and sampling methods along with a receding horizon prediction technique. However, while these approaches focus on local scene details, they often overlook the broader multimodal context. Our diffusion-based model addresses this by generating diverse actions for all agents in each inference, overcoming a significant limitation of traditional trajectory prediction methods.

#### B. Generative model-based methods

Generative adversarial networks (GANs) [24], [25] and Variational Auto-Encoders (VAEs) [26], [27] have been utilized for traffic scene generation. For example, [24] proposes a conditional generative neural system (CGNS) for probabilistic trajectory generation, while [27] develops a conditional VAE for multimodal, context-driven traffic scene generation. However, these methods often generate unrealistic trajectories due to their reliance on training data distribution and limited diversity [12]. Additionally, GANs can suffer from unstable training [14], and VAEs may be constrained by a simple Gaussian prior, limiting their expressiveness. Recently, diffusion models have emerged as a promising alternative to GANs and VAEs for generating realistic and diverse data [21], [12]. Notably, [21] applies a classifierguided diffusion approach to trajectory data with a probabilistic framework, while [12] introduces a conditional diffusion model for controllable traffic generation (CTG), allowing users to specify desired trajectory properties while maintaining realism and physical plausibility. However, these methods primarily focus on single-agent behaviors. Recent work on multi-agent trajectory generation using diffusion models includes SceneDM [14], which generates future motions of all agents, achieving state-of-the-art results on the Waymo Sim Agents Benchmark, and DJINN [22], which

produces traffic scenarios based on the joint states of all agents. A limitation of these models is that they predict individual agent trajectories per inference. In contrast, our approach integrates diffusion models with transformer-based encoder-decoder architectures to simultaneously generate joint, coherent future trajectories for all agents.

## III. METHODOLOGY

In this section, we present our novel WcDT framework for representing and generating complex traffic scenes. We first explain how traffic environments are modeled, followed by an introduction to the framework and its components.

## A. Traffic scene representation

Traffic environments are composed of multimodal data, such as road layouts, traffic signals, agent movements, and environmental conditions [28], [29]. To encode these elements in WcDT, we adopt a unified approach that captures both predicted and environmental (world) agents. Unlike existing methods that require transforming information to each agent's perspective [15], our approach simplifies this by:

- Using a unified Cartesian coordinates system for both predicted and world agents.
- Representing historical agent trajectories through movement statements instead of traditional coordinate vectors.

Key variables for simulating traffic scenarios in WcDT:

- $\mathcal{A}_{all}, \mathcal{A}_p, \mathcal{A}_w$ : Counts of all agents, predicted agents, and world agents, respectively.
- $\mathcal{T}_h$ ,  $\mathcal{T}_f$ : Historical and future time steps.
- $\mathscr{L}, \mathscr{P}$ : Lane lines and points within scenarios.
- $\mathscr{S}_{tl}$ : Traffic light states.
- $\mathscr{D}$ : Dimensionality of different traffic elements in a traffic scenario ( $\mathscr{D}_a$  represents the features of an agent,  $\mathscr{D}_t$  represents the features of a traffic light, and  $\mathscr{D}_m$  represents the features of a map element).

For different traffic objects, we represent them as follows:

- Agent move statement and features: To mitigate the impact of varying agent positions on historical and future trajectories, we introduce absolute states for past and prospective agent states. For agent *i* at time step *t*, the state  $s_t^i$  is defined as  $s_t^i = [(x_t x_{t-1}), (y_t y_{t-1}), (\theta_t \theta_{t-1}), (v_t v_{t-1})]$ , where  $x_t, y_t, \theta_t$ , and  $v_t$  represent longitudinal position, lateral position, heading angle, and velocity, respectively. The feature space for each agent is  $[\mathscr{A}, \mathscr{T}_h 1, \mathscr{D}_a]$ .
- *Traffic light feature:* The traffic light dataset for each scenario, denoted as  $[\mathscr{S}_{tl}, \mathscr{T}_h, \mathscr{D}_t]$ , contains the positions and operational statuses of signals over historical intervals. For any traffic signal point  $s_{tl} \in \mathscr{S}_{tl}$ , this information is represented using a one-hot encoding of signal states and spatial positions at each historical moment.



Fig. 2. Overview of WcDT, which consists of the following modules: (a) Agent action generation and agent to agent cross attention blocks; (b) The traffic scene encoder extracts temporal and spatial features in the traffic scene, including: other agents, traffic signals, HD maps; (c) The multimodal trajectory decoder is used to generate possible future actions for all predicted agents.

• *Map feature:* The map features, denoted as  $[1, \mathcal{L}, \mathcal{P}, \mathcal{D}_m]$ , encompass key lane details in a traffic scenario, including positions and types. Each lane line  $l_t \in \mathcal{L}$  at the current time step is represented positions of all points along the lane and using one-hot encoding to specify the its type.

Figure 2 shows an overview of our proposed WcDT framework for traffic scene generation, including three major components: action diffusion, scene encoder, and trajectory decoder, which are detailed in the following subsections.

## B. Action Diffusion

To enhance trajectory diversity in WcDT, we encode agent actions into latent space to increase variability. These latent features are then input into the scene encoder. We use Denoising Diffusion Probabilistic Models (DDPM) [30] for action encoding. Although DDPM traditionally employs U-Net architectures, recent research [31] demonstrates that transformers can achieve comparable performance without U-Net's inductive biases. Consequently, we replace U-Net with Diffusion Transformers (DiTs) to improve performance and ensure diverse agent trajectories. Figure 3 shows the architecture of conditional DiT blocks for encoding "latent action features". The network takes random noise, time steps, and historical trajectories as inputs and produces latent action features for the scene encoder. The DDPM loss function guides the network to generate latent features consistent with agent kinematics [30], thus enhancing trajectory variability. The loss function for the DiT module is as follows:

$$\mathscr{L}_{diff} = ||\varepsilon - \varepsilon_{\theta} (\sqrt{\bar{\alpha}_{t}} x_{0} + \sqrt{1 - \bar{\alpha}_{t}} \varepsilon, t)||^{2}, \qquad (1)$$

where  $\bar{\alpha}_t$  are hyperparameters for diffusion model training,  $\varepsilon_{\theta}$  represents the diffusion model with DiT blocks, and  $\varepsilon$  is Gaussian noise.

#### C. Scene Encoder

In traffic scenes, agents like vehicles, pedestrians, and bicycles, along with map features and traffic signals, are present. To generate diverse trajectories, we use embedding blocks of different sizes and layers. These blocks encode agents' characteristics, bypassing the need for agent-specific coordinate transformations. The Pose-Embedding encodes positional data  $p_i$  into a 1D matrix, while Feature-Embedding



Fig. 3. Overview of the developed conditional DiT blocks and illustration of DiT block processing, where action latent features are integrated with map, object, and traffic light using multi-head attention for traffic scene generation.

translates attributes like height, width, and type into another matrix. For agent *i*:

$$E_p = \phi_p[x_i, y_i], \quad E_f = \phi_f[f_w, f_h, f_{type}], \tag{2}$$

where  $\phi_p$  and  $\phi_f$  represent linear transformations. The final agent embedding  $E_A$  is:

$$E_A = \text{ReLU}(\text{LayerNorm}(\text{Concat}(E_p, E_f))).$$
 (3)

To represent traffic scenarios, the encoding process integrates world agents, maps, and traffic lights. Features are processed into embeddings, refined via neural network blocks, including multi-head self-attention for detailed analysis and crossattention for feature relationships. Attention layers, replacing CNNs, dynamically capture long-range dependencies. The self-attention encoding for world agents is:

$$q_{A_p} = W^{Q \times h} E_p, \quad k_{A_p} = W^{K \times h} E_p, \quad v_{A_p} = W^{V \times h} E_p, \quad (4)$$

where  $W^Q$ ,  $W^K$ ,  $W^V$  are learnable parameters, and attention is calculated as:

Cross-attention for encoding map and traffic light features follows a similar process:

$$q_{A_p} = W^{Q \times h} E_{A_p}, \quad k_m = W^{K \times h} E_m, \quad v_m = W^{V \times h} E_m, \quad (6)$$

$$\propto_m = \operatorname{Softmax}\left(\frac{q_{A_p}^I}{\sqrt{d_k}}k_M\right), \quad \operatorname{Cross}_m = \propto_{A_p} v_M.$$
 (7)

**Spatial and Temporal Fusion.** We propose Temporal-Spatial Fusion Attention layers to capture the dynamic nature of traffic scenarios by integrating multimodal data. The agent's features are augmented with latent action features (from Eq. 3) and processed through self-attention layers to identify key temporal-spatial insights, ensuring an accurate understanding of traffic dynamics.

# D. Trajectory Decoder

The trajectory decoder translates fused traffic features into agents' future trajectories using GRU and MLP blocks. Drawing inspiration from [15], we employ a multimodal output mechanism to handle agents with varied behaviors. To reduce the influence of different initial positions, the decoder outputs agents' move statements and their likelihoods. The trajectory for model  $\mathcal{M}$  is computed as shown in Fig. 4:

$$\operatorname{Traj}_{a}^{m} = \operatorname{Pos}_{a} + \sum_{i=t}^{T_{f}} [\Delta x_{m}, \Delta y_{m}, \Delta \theta_{m}],$$
(8)

where  $\text{Pos}_a$  is the agent's current position, and  $\text{Traj}_a^m$  denotes future trajectory points, computed by adding the displacement  $[\Delta x_m, \Delta y_m, \Delta \theta_m]$  for each time step. The speed is:

$$\text{Speed}_{a}^{m} = \frac{[\Delta x_{m}, \Delta y_{m}]}{\Delta t},$$
(9)

calculated from the displacement over time  $\Delta t$ . This kinematic approach outputs the trajectory as  $[x_a^m, y_a^m, \theta_a^m, v_a^m]$ .

## E. Loss functions

Our model aims to ensure generated trajectories adhere to scene constraints while maintaining diversity. The trajectory with the lowest loss from the multimodal set is selected, and its deviation from the ground truth is measured using the Huber loss [32]:

$$\mathscr{L}_{reg} = \text{Huber}(\text{Traj}_p, \text{Traj}_{gt}), \tag{10}$$

where  $\operatorname{Traj}_p$  and  $\operatorname{Traj}_{gt}$  are the predicted and ground truth trajectories, respectively. We also introduce a classification loss to identify the modality closest to the ground truth, where the modality with the smallest AED is used as the classification target:

$$\mathscr{L}_{cls} = -\sum_{i=1}^{M} y_i log(p_i), \qquad (11)$$

The total loss is a combination of diffusion, regression, and classification losses:

$$\mathscr{L}_{total} = \mathscr{L}_{diff} + \mathscr{L}_{reg} + \mathscr{L}_{cls}.$$
 (12)

Here,  $\mathscr{L}_{diff}$  is the standard diffusion model loss, computed as the L2 loss between predicted and original noise.



Fig. 4. Illustration of the trajectory generation process.

# IV. EXPERIMENTS AND RESULTS

#### A. Experimental Setup

**Dataset.** We use the Waymo Motion Prediction dataset [23], containing 576,012 driving scenarios. The data is divided into 486,995 training, 44,097 validation, and 44,920 testing scenarios. Each scenario lasts 9 seconds, sampled at 10 Hz, with only the first second of the testing scenarios available for generating future trajectories for the next 8 seconds. **Metrics.** We employ established evaluation metrics [33], [35], [18] and Sim Agents Challenge metrics [33] to assess the realism and diversity of generated trajectories. These metrics cover kinematic, object interaction, and map-based aspects. We minimize the negative log-likelihood (NLL):

$$NLL^* = -\frac{1}{|\mathcal{D}|} \sum_{i=0}^{|\mathcal{D}|} Logq^{world}(o_{\geq t,i}|o_{< t,i}), \qquad (13)$$

where  $o_{<t,i}$  represents historical observations, and  $o_{\ge t,i}$  denotes future observations.

**Implementation Details.** We train our model for 128 epochs using two NVIDIA A100 GPUs, with Adam optimizer [36]. We set the batch size to 128, the initial learning rate to  $2 \times 10^{-4}$ , and apply a cosine annealing scheduler [37] for learning rate adjustment. The architecture includes 2 DiT blocks, 4 Other Agent Former blocks, 4 Map Former blocks, and 2 Traffic Light Former blocks, with the Trajectory Decoder using 2 MLP blocks. Both Multi-Head Self-Attention and Cross-Attention mechanisms are configured with 8 attention heads. We test two model variants: WcDT-64 (64 hidden units) and WcDT-128 (128 hidden units). The origin for all scenarios is set at the current self-driving vehicle's location.

#### B. Comparison with state-of-the-art methods

We compare our proposed WcDT model against several state-of-the-art benchmarks submitted to the Sim Agent Challenge<sup>1</sup> [33], including Random Agent, Constant Velocity [33], MTR+++ [34], WayFormer [28], MULTI-PATH++ [16], MVTA [20], and MVTE [20]. MVTE, MVTA, and MTR+++ show advanced capabilities in generating

<sup>&</sup>lt;sup>1</sup>Sim Agent Leaderboard as of 02-04-2024: https://waymo.com/ open/challenges/2023/sim-agents/

#### TABLE I

The Sim Agents Leaderboard results evaluate methods using 10 similarity metrics (kinematic, object interaction, and map-based) and distance error metrics (ADE and MinADE), where higher similarity values indicate better performance, and lower ADE/MinADE values represent more accurate trajectory predictions.

Method	Linear Speed	Linear Accel	Ang Speed	Ang Accel	Dist to Obj	Collision	TTC	Dist to Road Edge	Offroad	Composite Metric	ADE	MinADE
Random Agent [33]	0.002	0.044	0.074	0.120	0.000	0.006	0.734	0.178	0.325	0.163	50.740	50.707
Constant Velocity [33]	0.074	0.058	0.019	0.035	0.208	0.202	0.737	0.454	0.325	0.238	7.924	7.924
MTR+++ [34]	0.412	0.107	0.484	0.437	0.346	0.414	0.797	0.654	0.577	0.470	2.129	1.682
WayFormer [28]	0.408	0.127	0.473	0.437	0.358	0.403	0.810	0.645	0.589	0.472	2.588	1.694
MULTIPATH++ [16]	0.432	0.230	0.515	0.452	0.344	0.420	0.813	0.639	0.583	0.489	5.308	2.052
MVTA [20]	0.437	0.220	0.533	0.481	0.373	0.436	0.830	0.654	0.629	0.509	3.938	1.870
MVTE [20]	0.443	0.222	0.535	0.481	0.382	0.451	0.832	0.664	0.641	0.517	3.873	1.677
WcDT (Ours)	0.515	0.370	0.543	0.508	0.548	0.629	0.846	0.738	0.608	0.743	2.045	1.472

realistic and feasible motion trajectories for autonomous vehicles. Table I summarizes the evaluation results. The Random Agent method [33], which generates random trajectories, performs the worst with a composite score of 0.163. Constant Velocity [33], which predicts based on the last known heading and speed, improves slightly, scoring 0.238 in the composite metric. Our WcDT model achieves the highest composite metric of 0.743, indicating strong performance and outperforming MVTE [20] in specific metrics such as Linear Speed, Linear Acceleration, Angle Speed, Distance to Object, and Distance to Road Edge, while also achieving a better MinADE score. This highlights WcDT's ability to generate precise and contextually appropriate trajectory predictions, showing its strength in interpreting dynamic traffic environments.

## C. Ablation studies on diffusion model

We evaluate the effect of the "latent action features" encoding, comparing random noise inputs, the Unet network, and our custom DiT block. As shown in Table II, the DiT module consistently achieves the lowest ADE and MinADE scores, along with the highest composite score. This demonstrates that the diffusion model with the DiT block significantly improves action diversity while maintaining realistic trajectories.

#### TABLE II

Ablation Study on Diffusion Model: Evaluating the diffusion model's contribution and comparing impacts of Dit Blocks

Random Noise	Unet	Dit Blocks	ADE↓	minADE↓	Composite Metric ↑
✓			4.843	2.715	0.326
	$\checkmark$		4.163	1.907	0.480
		$\checkmark$	2.045	1.472	0.743

#### TABLE III

ABLATION STUDY ON SCENE-ENCODER COMPONENTS: WE ASSESS THE IMPORTANCE OF EACH MODULE BY ADDING OR REMOVING IT FROM THE SCENE ENCODER AND EVALUATING PERFORMANCE USING ADE AND MINADE.

Spatial &Temporal Attention	Other Agent Former	HD Map Former	Traffic Light Former	ADE↓	minADE↓
	√	~	√	3.035	1.973
$\checkmark$		$\checkmark$	$\checkmark$	3.490	1.883
$\checkmark$	$\checkmark$		$\checkmark$	2.960	2.130
$\checkmark$	$\checkmark$	$\checkmark$		2.593	1.865
✓	$\checkmark$	$\checkmark$	$\checkmark$	2.045	1.472

#### D. Ablation studies on traffic scene encoder

As shown in Table III, the HD Map Former is crucial for accurate trajectory generation, with its removal leading to the worst ADE and MinADE scores. The Spatial and Temporal Attention blocks significantly enhance the encoder's understanding of traffic, with their absence resulting in sub-optimal ADE (3.035) and MinADE (1.973) scores. The Traffic Light and Other Agent Formers further boost accuracy, and the full encoder setup delivers the best results, demonstrating the effectiveness of our approach.

# E. Ablation studies on trajectory decoder

We evaluate the trajectory decoder by assessing the contributions of GRU and MLP blocks, focusing on how different configurations impact trajectory prediction. As shown in Table IV, MLP layers play a crucial role, with their absence resulting in the highest ADE and MinADE scores (3.759 and 2.475). This highlights the importance of MLP blocks in refining the decoder's predictive capabilities. GRU blocks are also essential, helping the model leverage historical data effectively. The best results are achieved when both MLP and GRU blocks are used, demonstrating their combined impact on enhancing trajectory generation.

TABLE IV

Ablation Study on the Trajectory Decoder: We assessed the

IMPACT OF VARIOUS COMPONENTS AND THE DIMENSIONALITY OF HIDDEN UNITS IN THE TRAJECTORY DECODER ON ADE AND MINADE PERFORMANCE.

MLP Blocks	<b>GRU Blocks</b>	Dimension	ADE↓	minADE↓
	$\checkmark$	64	4.174	2.532
	$\checkmark$	128	3.759	2.475
$\checkmark$		64	3.612	1.805
$\checkmark$		128	3.857	1.780
$\checkmark$	$\checkmark$	64	2.857	1.735
$\checkmark$	$\checkmark$	128	2.045	1.472

#### F. Ablation studies on network structures

We examine the effects of varying the number of modalities and attention heads on trajectory generation. Table V shows that WcDT-128 consistently outperforms WcDT-64, indicating that more attention layers enhance prediction accuracy. Additionally, the multi-modality configuration in WcDT-64 yields better results than the single-modality setup, as it enables the model to process more information, improving



Fig. 5. Visualization results of the ground truth and WTSGM-generated trajectories in cruise scenarios.

its understanding of the driving environment and leading to more precise predictions.

TABLE V IMPACT OF MULTIMODAL TRAJECTORY DECODERS AND DIMENSION OF ATTENTION BLOCKS ON SCENE GENERATION PERFORMANCE

Method	Multimodal	Attention Block Heads	ADE↓	minADE↓
WcDT-64	1	8	3.758	3.758
WcDT-64	10	8	3.475	2.548
WcDT-64	10	16	3.729	2.470
WcDT-64	30	16	3.647	1.962
WcDT-128	10	8	2.948	1.781
WcDT-128	30	16	2.045	1.472

# G. Visualization Results

Figure 5 illustrates the generated trajectories for randomly sampled Waymo dataset scenarios. The input features include map elements (black dotted lines) and the initial 1-second trajectories of various agents (dots), with each agent's trajectory represented by a unique color. The intensity of the colors deepens over time, reflecting the temporal progression of each agent's movement. The left side demonstrates lane-changing maneuvers, highlighting the model's ability to predict diverse and accurate trajectories in dynamic driving conditions. On the right, the figure showcases the model's performance in more complex intersection scenarios, further underscoring its robustness and precision in handling challenging traffic environments.

## V. CONCLUSION

This paper introduced a novel traffic scene generation framework that optimizes trajectory generation through Diffusion with Transformer (DiT) blocks. The model effectively fuses latent features, historical trajectories, HD maps, and traffic signal data using transformer-based encoders with attention mechanisms. A key contribution is the multimodal trajectory decoder, which generates a wide range of future trajectories, enhancing the diversity and realism of the generated traffic scenes. Experimental results show that our approach sets a new standard for realism and diversity in traffic scene generation. Future work will focus on improving robustness for more complex urban scenarios and handling more agents.

#### References

- Y. Huang, J. Du, Z. Yang, Z. Zhou, L. Zhang, and H. Chen, "A survey on trajectory-prediction methods for autonomous driving," *IEEE Transactions on Intelligent Vehicles*, vol. 7, no. 3, pp. 652–674, 2022.
- [2] S. Teng, X. Hu, P. Deng, B. Li, Y. Li, Y. Ai, D. Yang, L. Li, Z. Xuanyuan, F. Zhu *et al.*, "Motion planning for autonomous driving: The state of the art and future perspectives," *IEEE Transactions on Intelligent Vehicles*, 2023.
- [3] B. Paden, M. Čáp, S. Z. Yong, D. Yershov, and E. Frazzoli, "A survey of motion planning and control techniques for self-driving urban vehicles," *IEEE Transactions on intelligent vehicles*, vol. 1, no. 1, pp. 33–55, 2016.
- [4] D. Chen, K. Zhang, Y. Wang, X. Yin, Z. Li, and D. Filev, "Communication-efficient decentralized multi-agent reinforcement learning for cooperative adaptive cruise control," *IEEE Transactions* on *Intelligent Vehicles*, 2024.
- [5] L. Chen, Y. Li, C. Huang, B. Li, Y. Xing, D. Tian, L. Li, Z. Hu, X. Na, Z. Li *et al.*, "Milestones in autonomous driving and intelligent vehicles: Survey of surveys," *IEEE Transactions on Intelligent Vehicles*, vol. 8, no. 2, pp. 1046–1056, 2022.
- [6] W. Liu, M. Hua, Z. Deng, Z. Meng, Y. Huang, C. Hu, S. Song, L. Gao, C. Liu, B. Shuai *et al.*, "A systematic survey of control techniques and applications in connected and automated vehicles," *IEEE Internet of Things Journal*, 2023.
- [7] S. Ge, Y. Xie, K. Liu, Z. Ding, E. Hu, L. Chen, and F.-Y. Wang, "The use of intelligent vehicles and artificial intelligence in mining operations: Ethics, responsibility, and sustainability," *IEEE Transactions on Intelligent Vehicles*, vol. 8, no. 2, pp. 1021–1024, 2023.
- [8] M. O'Kelly, A. Sinha, H. Namkoong, R. Tedrake, and J. C. Duchi, "Scalable end-to-end autonomous vehicle testing via rare-event simulation," *Advances in neural information processing systems*, vol. 31, 2018.

- [9] X. Hu, S. Li, T. Huang, B. Tang, R. Huai, and L. Chen, "How simulation helps autonomous driving: A survey of sim2real, digital twins, and parallel intelligence," *IEEE Transactions on Intelligent Vehicles*, 2023.
- [10] S. Grollius, M. Ligges, J. Ruskowski, and A. Grabmaier, "Concept of an automotive lidar target simulator for direct time-of-flight lidar," *IEEE Transactions on Intelligent Vehicles*, 2021.
- [11] E. Weiss and J. C. Gerdes, "High speed emulation in a vehicle-in-theloop driving simulator," *IEEE Transactions on Intelligent Vehicles*, vol. 8, no. 2, pp. 1826–1836, 2022.
- [12] Z. Zhong, D. Rempe, D. Xu, Y. Chen, S. Veer, T. Che, B. Ray, and M. Pavone, "Guided conditional diffusion for controllable traffic simulation," in 2023 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2023, pp. 3560–3566.
- [13] M. Treiber, A. Hennecke, and D. Helbing, "Congested traffic states in empirical observations and microscopic simulations," *Physical review E*, vol. 62, no. 2, p. 1805, 2000.
- [14] Z. Guo, X. Gao, J. Zhou, X. Cai, and B. Shi, "Scenedm: Scene-level multi-agent trajectory generation with consistent diffusion models," arXiv preprint arXiv:2311.15736, 2023.
- [15] Y. Chai, B. Sapp, M. Bansal, and D. Anguelov, "Multipath: Multiple probabilistic anchor trajectory hypotheses for behavior prediction," arXiv preprint arXiv:1910.05449, 2019.
- [16] B. Varadarajan, A. Hefny, A. Srivastava, K. S. Refaat, N. Nayakanti, A. Cornman, K. Chen, B. Douillard, C. P. Lam, D. Anguelov *et al.*, "Multipath++: Efficient information fusion and trajectory aggregation for behavior prediction," in 2022 International Conference on Robotics and Automation (ICRA). IEEE, 2022, pp. 7814–7821.
- [17] J. Ngiam, B. Caine, V. Vasudevan, Z. Zhang, H.-T. L. Chiang, J. Ling, R. Roelofs, A. Bewley, C. Liu, A. Venugopal *et al.*, "Scene transformer: A unified architecture for predicting multiple agent trajectories," *arXiv preprint arXiv:2106.08417*, 2021.
- [18] S. Suo, S. Regalado, S. Casas, and R. Urtasun, "Trafficsim: Learning to simulate realistic multi-agent behaviors," *arXiv preprint* arXiv:2101.06557, 2021.
- [19] S. Shi, L. Jiang, D. Dai, and B. Schiele, "Motion transformer with global intention localization and local movement refinement," Advances in Neural Information Processing Systems, 2022.
- [20] Y. Wang, T. Zhao, and F. Yi, "Multiverse transformer: 1st place solution for waymo open sim agents challenge 2023," arXiv preprint arXiv:2306.11868, 2023.
- [21] M. Janner, Y. Du, J. B. Tenenbaum, and S. Levine, "Planning with diffusion for flexible behavior synthesis," arXiv preprint arXiv:2205.09991, 2022.
- [22] M. Niedoba, J. W. Lavington, Y. Liu, V. Lioutas, J. Sefas, X. Liang, D. Green, S. Dabiri, B. Zwartsenberg, A. Scibior *et al.*, "A diffusion-model of joint interactive navigation," *arXiv preprint* arXiv:2309.12508, 2023.
- [23] Z. Sun, J. Wang, Y. Chen, J. Xu, X. Zhang, Y. Li, Y. Zhang, Z. Liu, J. Guo, T. Huang *et al.*, "Large scale interactive motion forecasting for autonomous driving: The waymo open motion dataset," in *Proceedings* of the IEEE/CVF International Conference on Computer Vision, 2021, pp. 11 271–11 281.
- [24] J. Li, H. Ma, and M. Tomizuka, "Conditional generative neural system for probabilistic trajectory prediction," in 2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). IEEE, 2019, pp. 6150–6156.
- [25] R. Bhattacharyya, B. Wulfe, D. J. Phillips, A. Kuefler, J. Morton, R. Senanayake, and M. J. Kochenderfer, "Modeling human driving behavior through generative adversarial imitation learning," *IEEE Transactions on Intelligent Transportation Systems*, vol. 24, no. 3, pp. 2874–2887, 2022.

- [26] W. Ding, W. Wang, and D. Zhao, "Multi-vehicle trajectories generation for vehicle-to-vehicle encounters," in 2019 IEEE International Conference on Robotics and Automation (ICRA), 2019.
- [27] G. Oh and H. Peng, "Cvae-h: Conditionalizing variational autoencoders via hypernetworks and trajectory forecasting for autonomous driving," arXiv preprint arXiv:2201.09874, 2022.
- [28] N. Nayakanti, R. Al-Rfou, A. Zhou, K. Goel, K. S. Refaat, and B. Sapp, "Wayformer: Motion forecasting via simple & efficient attention networks," in 2023 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2023, pp. 2980–2987.
- [29] Y. Gao, Z. Chen, J. Wang, X. Zhang, and Y. Zhang, "Vectornet: Encoding hd maps and agent dynamics from vectorized representation," *arXiv preprint arXiv:2006.05262*, 2020.
- [30] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Advances in neural information processing systems*, vol. 33, pp. 6840–6851, 2020.
- [31] W. Peebles and S. Xie, "Scalable diffusion models with transformers," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 4195–4205.
- [32] J. H. Friedman, "Greedy function approximation: A gradient boosting machine," Annals of Statistics, vol. 29, no. 5, pp. 1189–1232, 2001.
- [33] N. Montali, J. Lambert, P. Mougin, A. Kuefler, N. Rhinehart, M. Li, C. Gulino, T. Emrich, Z. Yang, S. Whiteson *et al.*, "The waymo open sim agents challenge," *arXiv preprint arXiv:2305.12032*, 2023.
- [34] S. Shi, L. Jiang, D. Dai, and B. Schiele, "Mtr++: Multi-agent motion prediction with symmetric scene modeling and guided intention querying," *arXiv preprint arXiv:2306.17770*, 2023.
- [35] J. Gil, L. Martín, C. Montes, and A. Ortega, "A fast procedure for computing the distance between complex objects in three-dimensional space," *Computer graphics forum*, vol. 10, no. 4, pp. 331–340, 1991.
- [36] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *Computer Science*, 2014.
- [37] I. Loshchilov and F. Hutter, "Sgdr: Stochastic gradient descent with warm restarts," arXiv e-prints, vol. abs/1608.03983, 2016.