

CarPlanner: Consistent Auto-regressive Trajectory Planning for Large-scale Reinforcement Learning in Autonomous Driving

Dongkun Zhang^{1,2} Jiaming Liang² Ke Guo² Sha Lu¹ Qi Wang²
 Rong Xiong^{1,✉} Zhenwei Miao^{2,†} Yue Wang¹
¹Zhejiang University ²Cainiao Network

Abstract

Trajectory planning is vital for autonomous driving, ensuring safe and efficient navigation in complex environments. While recent learning-based methods, particularly reinforcement learning (RL), have shown promise in specific scenarios, RL planners struggle with training inefficiencies and managing large-scale, real-world driving scenarios. In this paper, we introduce **CarPlanner**, a **Consistent auto-regressive Planner** that uses RL to generate multi-modal trajectories. The auto-regressive structure enables efficient large-scale RL training, while the incorporation of consistency ensures stable policy learning by maintaining coherent temporal consistency across time steps. Moreover, CarPlanner employs a generation-selection framework with an expert-guided reward function and an invariant-view module, simplifying RL training and enhancing policy performance. Extensive analysis demonstrates that our proposed RL framework effectively addresses the challenges of training efficiency and performance enhancement, positioning CarPlanner as a promising solution for trajectory planning in autonomous driving. To the best of our knowledge, we are the first to demonstrate that the RL-based planner can surpass both IL- and rule-based state-of-the-arts (SOTAs) on the challenging large-scale real-world dataset nuPlan. Our proposed CarPlanner surpasses RL-, IL-, and rule-based SOTA approaches within this demanding dataset.

1. Introduction

Trajectory planning [39] is essential in autonomous driving, utilizing outputs from perception and trajectory prediction modules to generate future poses for the ego vehicle. A controller tracks this planned trajectory, producing control commands for closed-loop driving. Recently, learning-based trajectory planning has garnered attention due to its potential to automate algorithm iteration, eliminate tedious rule design, and ensure safety and comfort in diverse real-

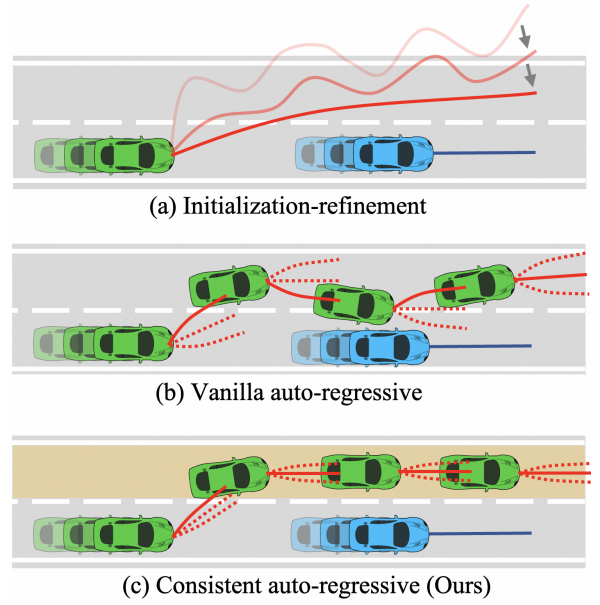


Figure 1. Frameworks for multi-step trajectory generation. (a) Initialization-refinement that generates an initial trajectory and refines it iteratively. (b) Vanilla auto-regressive models that decode subsequent poses sequentially. (c) Our consistent auto-regressive model that integrates time-consistent mode information.

world scenarios [39].

Most existing researches [3, 13, 19, 33] employ imitation learning (IL) to align planned trajectories with those of human experts. However, this approach suffers from distribution shift [32] and causal confusion [10]. Reinforcement learning (RL) offers a potential solution, addressing these challenges and providing richer supervision through reward functions. Although RL shows effectiveness in domains such as games [37], robotics [22], and language models [27], it still struggles with training inefficiencies and performance issues in the large-scale driving task. To the extent of our knowledge, no RL methods have yet achieved competitive results on large-scale open datasets such as nuPlan [2], which features diverse real-world scenarios.

Thus, this paper aims to tackle two key challenges in RL for trajectory planning: 1) training inefficiency and 2)

[†]Project lead. [✉]Corresponding author (rxiong@zju.edu.cn).

poor performance. Training inefficiency arises from the fact that RL typically operates in a model-free setting, necessitating an inefficient simulator running on a CPU to repeatedly roll out a policy for data collection. To overcome this challenge, we propose an efficient model-based approach utilizing neural networks as transition models. Our method is optimized for execution on hardware accelerators such as GPUs, rendering our time cost comparable to that of IL-based methods.

To apply RL to solve the trajectory planning problem, we formulate it as a multi-step sequential decision-making task utilizing a Markov Decision Process (MDP). Existing methods that generate the trajectory[†] in multiple steps generally fall into two categories: initialization-refinement [20, 23, 36, 48] and auto-regressive models [31, 35, 44, 49].

The first category, illustrated in Fig. 1 (a), involves generating an initial trajectory estimate and subsequently refining it through iterative applications of RL. However, recent studies, including Gen-Drive [21], suggest that it continues to lag behind SOTA IL and rule-based planners. One notable limitation of this approach is its neglect of the temporal causality inherent in the trajectory planning task. Additionally, the complexity of direct optimization over high-dimensional trajectory space can hinder the performance of RL algorithms. The second category consists of auto-regressive models, shown in Fig. 1 (b), which generate the poses of the ego vehicle recurrently using a single-step policy within a world model. In this category, ego poses at all time steps are consolidated to form the overall planned trajectory. As taking temporal causality into account, current auto-regressive models allow for interactive behaviors. However, a common limitation is their reliance on auto-regressively random sampling from action distributions to generate multi-modal trajectories. This vanilla auto-regressive procedure may compromise long-term consistency and unnecessarily expand the exploration space in RL, leading to poor performance.

To address the limitations of auto-regressive models, we introduce **CarPlanner**, a Consistent auto-regressive model designed for efficient, large-scale RL-based **Planner** training (see Fig. 1 (c)). The key insight of CarPlanner is its incorporation of consistent mode representation as conditions for the auto-regressive model. Specifically, we leverage a longitudinal-lateral decomposed mode representation, where the longitudinal mode is a scalar that captures average speeds, and the lateral mode encompasses all possible routes derived from the current state of the ego vehicle along with map information. This mode remains constant across time steps, providing stable and consistent guidance during policy sampling.

[†]In this paper, the term “trajectory” refers to the future poses of the ego vehicle or traffic agents. To avoid confusion, we use the term “state (action) sequence” to refer to the “trajectory” in the RL community.

Furthermore, we propose a universal reward function that suits large-scale and diverse scenarios, eliminating the need for scenario-specific reward designs. This function consists of an expert-guided and task-oriented term. The first term quantifies the displacement error between the ego-planned trajectory and the expert’s trajectory, which, along with the consistent mode representation, narrows down the policy’s exploration space. The second term incorporates common senses in driving tasks including the avoidance of collision and adherence to the drivable area. Additionally, we introduce an Invariant-View Module (IVM) to supply invariant-view input for policy, with the aim of providing time-agnostic policy input, easing the feature learning and embracing generalization. To achieve this, IVM preprocesses state and lateral mode by transforming agent, map, and route information into the ego’s current coordinate and by clipping information that is distant from the ego.

To our knowledge, we are the first to demonstrate that RL-based planner outperforms state-of-the-art (SOTA) IL and rule-based approaches on the challenging large-scale nuPlan dataset. In summary, the key contributions of this paper are highlighted as follows:

- We present **CarPlanner**, a consistent auto-regressive planner that trains an RL policy to generate consistent multi-modal trajectories.
- We introduce an expert-guided universal reward function and IVM to simplify RL training and improve policy generalization, leading to enhanced closed-loop performance.
- We conduct a rigorous analysis on the characteristics of IL and RL training, providing insights into their strengths and limitations, while highlighting the advantages of RL in tackling challenges such as distribution shift and causal confusion.
- Our framework showcases exceptional performance, surpassing all RL-, IL-, and rule-based SOTAs on the nuPlan benchmark. This underscores the potential of RL in navigating complex real-world driving scenarios.

2. Related Work

2.1. Imitation-based Planning

The use of IL [8, 17] to train planners based on human demonstrations has garnered significant interest recently. This approach leverages the driving expertise of experienced drivers who can safely and comfortably navigate a wide range of real-world scenarios, along with the added advantage of easily collectible driving data at scale [2, 11, 18]. Numerous studies [6, 20, 30, 33] have focused on developing innovative networks to enhance open-loop performance in this domain. However, the ultimate challenge of autonomous driving is achieving closed-loop operation, which is evaluated using driving-oriented met-

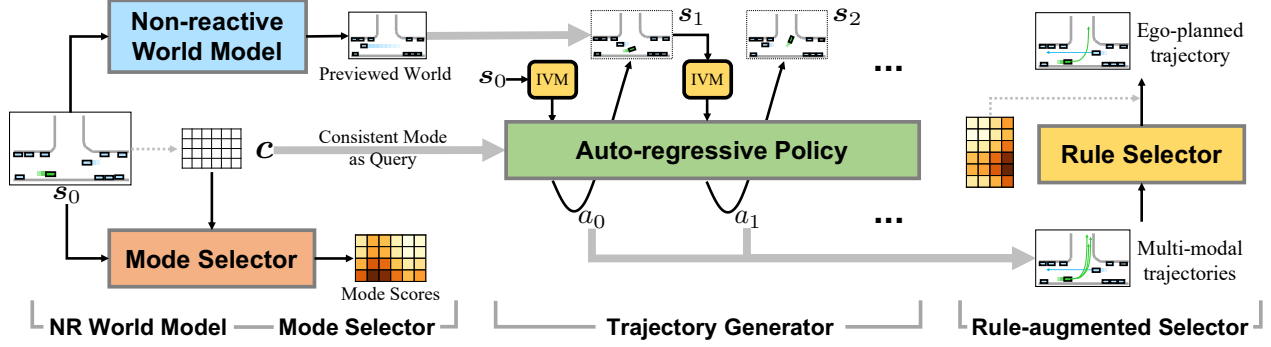


Figure 2. CarPlanner contains four parts. (1) The non-reactive world model takes initial state s_0 as input and predicts the future trajectories of traffic agents. (2) The mode selector outputs scores based on the initial state and the modes c . (3) The trajectory generator obeys an auto-regressive structure condition on the consistent mode and produces mode-aligned multi-modal trajectories. (4) The rule-augmented selector compensates the mode scores by safety, comfort, and progress metrics.

rics such as safety, adherence to traffic rules, comfort, and progress. This reveals a significant gap between the training and testing phases of planners. Moreover, IL is particularly vulnerable to issues such as distribution shift [32] and causal confusion [10]. The first issue results in suboptimal decisions when the system encounters scenarios that are not represented in the training data distribution. The second issue arises when networks inadvertently capture incorrect correlations and develop shortcuts based on input information, primarily due to the reliance on imitation loss from expert demonstrations. Despite efforts in several studies [1, 4, 5, 45] to address these challenges, the gap between training and testing remains substantial.

2.2. RL in Autonomous Driving

In the field of autonomous driving, RL has demonstrated effectiveness in addressing specific scenarios such as highway driving [24, 42], lane changes [16, 25], and unprotected left turns [24, 26, 43]. Most methods directly learn policies over the control space, which includes throttle, brake, and steering commands. Due to the high frequency of control command execution, the simulation process can be time-consuming, and exploration can be inconsistent [43]. Several works [43, 47] have proposed learning trajectory planners with actions defined as ego-planned trajectories, which temporally extend the exploration space and improve training efficiency. However, a trade-off exists between the trajectory horizon and training performance, as noted in ASAP-RL [43]. Increasing the trajectory horizon results in less reactive behaviors and a reduced amount of data, while a smaller trajectory horizon leads to challenges similar to those encountered in control space. Additionally, these methods typically employ a model-free setting, making them difficult to apply to the complex, diverse real-world scenarios found in large-scale driving datasets. In this paper, we propose adopting a model-based formulation that can facilitate RL training on large-scale datasets. Under this formulation, we aim to overcome the trajectory horizon

trade-off by using a world model, which can provide a preview of the world in which our policy can make multi-step decisions during testing.

3. Method

3.1. Preliminaries

MDP is used to model sequential decision problems, formulated as a tuple $\langle \mathcal{S}, \mathcal{A}, P, R, \rho_0, \gamma, T \rangle$. \mathcal{S} is the state space. \mathcal{A} is the action space. $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ [†] is the state transition probability. $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ denotes the reward function and is bounded. $\rho_0 \in \Delta(\mathcal{S})$ is the initial state distribution. T is the time horizon and γ is the discount factor of future rewards. The state-action sequence is defined as $\tau = (s_0, a_0, s_1, a_1, \dots, s_T)$, where $s_t \in \mathcal{S}$ and $a_t \in \mathcal{A}$ are the state and action at time step t . The objective of RL is to maximize the expected return:

$$\max_{\pi} \mathbb{E}_{s_t \sim P, a_t \sim \pi} \left[\sum_{t=0}^T \gamma^t R(s_t, a_t) \right]. \quad (1)$$

Vectorized state representation. State s_t contains map and agent information in vectorized representation [12]. Map information m includes the road network, traffic lights, etc, which are represented by polylines and polygons. Agent information includes the current and past poses of ego vehicle and other traffic agents, which are represented by polylines. The index of ego vehicle is 0 and the indices of traffic agents range from 1 to N . For each agent i , its history is denoted as $s_{t-H:t}^i, i \in \{0, 1, \dots, N\}$, where H is the history time horizon.

3.2. Problem Formulation

We model the trajectory planning task as a sequential decision process and decouple the auto-regressive models into policy and transition models. The key to connect trajectory planning and auto-regressive models is to define the

[†] $\Delta(\mathcal{X})$ denotes the set of probability distribution over set \mathcal{X} .

action as the next pose of ego vehicle, i.e., $a_t = s_{t+1}^0$. Therefore, after forwarding the auto-regressive model, the decoded pose is collected to be the ego-planned trajectory. Specifically, we can reduce the state-action sequence to the state sequence under this definition and vectorized representation:

$$\begin{aligned} & P(s_0, a_0, s_1, a_1, \dots, s_T) \\ &= P(m, s_{-H:0}^{0:N}, s_1^0, m, s_{1-H:1}^{0:N}, s_2^0, \dots, m, s_{T-H:T}^{0:N}) \quad (2) \\ &= P(m, s_{-H:0}^{0:N}, m, s_{1-H:1}^{0:N}, \dots, m, s_{T-H:T}^{0:N}) \\ &= P(s_0, s_1, \dots, s_T). \end{aligned}$$

The state sequence can be further formulated in an auto-regressive fashion and decomposed into policy and world model:

$$\begin{aligned} P(s_0, s_1, \dots, s_T) &= \rho_0(s_0) \prod_{t=0}^{T-1} P(s_{t+1}|s_t) \\ &= \rho_0(s_0) \prod_{t=0}^{T-1} P(s_{t+1}^0, s_{t+1}^{1:N}|s_t) \quad (3) \\ &= \rho_0(s_0) \prod_{t=0}^{T-1} \underbrace{\pi(a_t|s_t)}_{\text{Policy}} \underbrace{P(s_{t+1}^{1:N}|s_t)}_{\text{World Model}}. \end{aligned}$$

From Eq. (3), we can clearly identify the inherent problem associated with the typical auto-regressive approach: inconsistent behaviors across time steps arise from the policy distribution, which depends on random sampling from the action distribution.

To solve the above problem, we introduce consistent mode information c that remains unchanged across time steps into the auto-regressive fashion:

$$\begin{aligned} P(s_0, s_1, \dots, s_T) &= \int_c P(s_0, s_1, \dots, s_T, c) dc \\ &= \rho_0(s_0) \int_c P(c|s_0) P(s_1, \dots, s_T|c) dc \\ &= \rho_0(s_0) \prod_{t=0}^{T-1} \underbrace{P(s_{t+1}^{1:N}|s_t)}_{\text{World Model}} \int_c \underbrace{P(c|s_0)}_{\text{Mode Selector}} \prod_{t=0}^{T-1} \underbrace{\pi(a_t|s_t, c)}_{\text{Policy}} dc. \quad (4) \end{aligned}$$

Since we focus on the ego trajectory planning, the consistent mode c does not impact world model.

This consistent auto-regressive formulation defined in Eq. (4) reveals a generation-selection framework where the mode selector scores each mode based on the initial state s_0 and the trajectory generator generates multi-modal trajectories via sampling from the mode-conditioned policy.

Non-reactive world model. The world model formulated in Eq. (4) needs to be employed in every time step since it produces the poses of traffic agents at time step $t + 1$ based on current state s_t . In practice, this process is time-consuming and we do not observe a performance improvement by using this world model, therefore, we use trajectory predictors $P(s_{1:T}^{1:N}|s_0)$ as non-reactive world model

that produces all future poses of traffic agents in one shot given initial state s_0 .

3.3. Planner Architecture

The framework of our proposed **CarPlanner** is illustrated in Fig. 2, comprising four key components: 1) the non-reactive world model, 2) the mode selector, 3) the trajectory generator, and 4) the rule-augmented selector. Our planner operates within a generation-selection framework. Given an initial state s_0 and all possible N_{mode} modes, the trajectory selector evaluates and assigns scores to each mode. The trajectory generator then produces N_{mode} trajectories that correspond to their respective modes.

For trajectory generator, the initial state s_0 is replicated N_{mode} times, each associated with one of the N_{mode} modes, effectively creating N_{mode} parallel worlds. The policy is executed within these simulated environments. During the policy rollout, a trajectory predictor acts as the state transition model, generating future positions of traffic agents across all time horizons. The details are as follows:

3.3.1. Non-reactive World Model

This module takes the initial state s_0 as input and outputs the future trajectories of traffic agents. The initial state is processed by agent and map encoders, followed by a self-attention Transformer encoder [41] to fuse the agent and map features. The agent features are then decoded into future trajectories.

Agent and map encoders. The state s_0 contains both map and agent information. The map information m consists of $N_{m,1}$ polylines and $N_{m,2}$ polygons. The polylines describe lane centers and lane boundaries, with each polyline containing $3N_p$ points, where 3 corresponds to the lane center, the left boundary, and the right boundary. Each point is with dimension $D_m = 9$ and includes the following attributes: x, y, heading, speed limit, and category. When concatenated, the points of the left and right boundaries together with the center point yield a dimension of $N_{m,1} \times N_p \times 3D_m$. We leverage a PointNet [29] to extract features from the points of each polyline, resulting in a dimensionality of $N_{m,1} \times D$, where D represents the feature dimension. The polygons represent intersections, crosswalks, stop lines, etc, with each polygon containing N_p points. We utilize another PointNet to extract features from the points of each polygon, producing a dimension of $N_{m,2} \times D$. We then concatenate the features from both polylines and polygons to form the overall map features, resulting in a dimension of $N_m \times D$. The agent information A consists of N agents, where each agent maintains poses for the past H time steps. Each pose is with dimension $D_a = 10$ and includes the following attributes: x, y, heading, velocity, bounding box, time step, and category. Consequently, the agent information has a dimension of $N \times H \times D_a$. We apply another

PointNet to extract features from the poses of each agent, yielding an agent feature dimension of $N \times D$.

3.3.2. Mode Selector

This module takes s_0 and longitudinal-lateral decomposed mode information as input and outputs the probability of each mode.

Route-speed decomposed mode. To capture the longitudinal behaviors, we generate N_{lon} modes that represent the average speed of the trajectory associated with each mode. Each longitudinal mode $c_{\text{lon},j}$ is defined as a scalar value of $\frac{j}{N_{\text{lon}}}$, repeated across a dimension D . As a result, the dimensionality of the longitudinal modes is $N_{\text{lon}} \times D$. For lateral behaviors, we identify N_{lat} possible routes from the map using a graph search algorithm. These routes correspond to the lanes available for the ego vehicle. The dimensionality of these routes is $N_{\text{lat}} \times N_r \times D_m$. To extract meaningful representations, we employ another PointNet to aggregate the features of the N_r points along each route, producing a lateral mode with a dimension of $N_{\text{lat}} \times D$. To create a comprehensive mode representation c , we combine the lateral and longitudinal modes, resulting in a combined dimension of $N_{\text{lat}} \times N_{\text{lon}} \times 2D$. To align this mode information with other feature dimensions, we pass it through a linear layer, mapping it back to $N_{\text{lat}} \times N_{\text{lon}} \times D$. $N_{\text{mode}} = N_{\text{lat}}N_{\text{lon}}$.

Query-based Transformer decoder. This decoder is employed to fuse the mode features with map and agent features derived from s_0 . In this framework, the mode serves as the query, while the map and agent information act as the keys and values. The updated mode features are decoded through a multi-layer perceptron (MLP) to yield the scores for each mode, which are subsequently normalized using the softmax operator.

3.3.3. Trajectory Generator

This module operates in an auto-regressive manner, recurrently decoding the next pose of the ego vehicle, a_t , given the current state, s_t , and consistent mode information, c .

Invariant-view module (IVM). Before feeding the mode and state into the network, we preprocess them to eliminate time information. For the map and agent information in state s_t , we select the K nearest neighbors (KNN) [28] to the ego current pose and only feed these into the policy. K is set to the half of map and agent elements respectively. Regarding the routes that capture lateral behaviors, we filter out the segments where the point closest to the current pose of the ego vehicle is the starting point, retaining K_r points. In this case, K_r is set to a quarter of N_r points in one route. Finally, we transform the routes, agent, and map poses into the coordinate frame of the ego vehicle at the current time step t . We then subtract the historical time steps $t - H : t$ from the current time step t , yielding time steps in the range $-H : 0$.

Query-based Transformer decoder. We employ the same

backbone network architecture as the mode selector, but with different query dimensions. Due to the IVM and the fact that different modes yield distinct states, the map and agent information cannot be shared among modes. As a result, we fuse information for each individual mode. Specifically, the query dimension is $1 \times D$, while the dimensions of the keys and values are $(N + N_m) \times D$. The output feature dimension remains $1 \times D$. It is important to highlight that the Transformer decoder can process information from multiple modes in parallel, eliminating the need to handle each mode sequentially using a for loop.

Policy output. The mode feature is processed by two distinct heads: a policy head and a value head. Each head comprises its own MLP to produce the parameters for the action distribution and the corresponding value estimate. We employ a Gaussian distribution to model the action distribution, where actions are sampled from this distribution during training. In contrast, during inference, we utilize the mean of the distribution to determine the actions.

3.3.4. Rule-augmented Selector

This module first introduces a rule-based selector that takes the initial state s_0 , the multi-modal ego-planned trajectories, and the predicted future trajectories of agents as input. It calculates driving-oriented metrics such as safety, progress, comfort, and others. A comprehensive score is obtained by the weighted sum of the rule-based scores and the mode scores provided by the mode selector. The ego-planned trajectory with the highest score is selected as the output of the trajectory planner.

3.4. Training

We first train the non-reactive world model and freeze the weights during the training of the mode selector and trajectory generator. Instead of feeding all modes to the generator, we apply a winner-takes-all strategy, wherein a positive mode is assigned based on the ego ground-truth trajectory and serves as a condition for the trajectory generator.

Mode assignment. The positive lateral mode is determined by the endpoint of the ground-truth trajectory. The longitudinal distance from the starting position to this endpoint is divided into N_{lon} intervals, with the positive longitudinal mode corresponding to the relevant distance interval.

Loss function. For training the selector, we use cross-entropy loss that is the negative log-likelihood of the positive mode. The PPO loss consists of policy improvement loss, value estimation loss, and entropy loss for exploration. The full description can be found in the supplementary.

Reward function. To handle diverse scenarios, we use the negative displacement error (DE) between the ego future pose and the ground truth as a universal reward. We also introduce additional terms to improve trajectory quality: collision rate and drivable area compliance. If the future pose

collides or falls outside the drivable area, the reward is set to -1; otherwise, it is 0.

Mode dropout. To prevent over-reliance on mode or route information due to Transformers’ residual connections, we implement a mode dropout module during training that randomly masks the route to mitigate this issue.

Ego-history dropout. Previous works [1, 4, 5, 13] suggest that planners trained via IL may rely too heavily on past poses and neglect environmental state information. To counter this, we combine techniques from ChauffeurNet [1] and PlanTF [5] into an ego-history dropout module, randomly masking velocity and acceleration during IL training, but not during RL training. Further discussions can be found in Sec. 4.4.

4. Experiments

4.1. Experimental Setup

Dataset and simulator. We use nuPlan [2], a large-scale closed-loop platform for studying trajectory planning in autonomous driving, to evaluate the efficacy of our method. The nuPlan dataset contains driving log data over 1,500 hours collected by human expert drivers across 4 diverse cities. It includes complex, diverse scenarios such as lane follow and change, left and right turn, traversing intersections and bus stops, roundabouts, interaction with pedestrians, etc. As a closed-loop platform, nuPlan provides a simulator that uses scenarios from the dataset as initialization. During the simulation, traffic agents are taken over by log-replay (non-reactive) or an IDM [40] policy (reactive). The ego vehicle is taken over by user-provided planners. The simulator lasts for 15 seconds and runs at 10 Hz. At each timestamp, the simulator queries the planner to plan a trajectory, which is tracked by an LQR controller to generate control commands to drive the ego vehicle.

Benchmarks and metrics. We use two benchmarks: Test14-Random and Reduced-Val14 for comparing with other methods and analyzing the design choices within our method. The Test14-Random provided by PlanTF [5] contains 261 scenarios. The Reduced-Val14 provided by PDM [9] contains 318 scenarios.

We use the closed-loop score (CLS) provided by the official nuPlan devkit[†] to assess the performance of all methods. The CLS score comprehends different aspects such as safety (S-CR, S-TTC), drivable area compliance (S-Area), progress (S-PR), comfort, etc. Based on the different behavior types of traffic agents, CLS is detailed into CLS-NR (non-reactive) and CLS-R (reactive). To further analyze various components of our methods, we also use open-loop metrics, such as loss of trajectory generator and selector.

Implementation details. We follow PDM [9] to construct our training and validation splits. The size of the training

Type	Planners	CLS-NR	CLS-R
Rule	IDM [40]	70.39	72.42
	PDM-Closed [9]	90.05	91.64
IL	RasterModel [2]	69.66	67.54
	UrbanDriver [33]	63.27	61.02
	GC-PGP [14]	55.99	51.39
	PDM-Open [9]	52.80	57.23
	GameFormer [20]	80.80	79.31
	PlanTF [5]	86.48	80.59
	PEP [45]	91.45	89.74
	PLUTO [4]	<u>91.92</u>	<u>90.03</u>
RL	CarPlanner (Ours)	94.07	<u>91.1</u>

Table 1. Comparison with SOTAs in Test14-Random. Based on the type of trajectory generator, all methods are categorized into Rule, IL, and RL. The best result is in **bold** and the second best result is underlined.

Type	Planners	CLS-NR	S-CR	S-PR
Rule	PDM-Closed [9]	<u>91.21</u>	97.01	<u>92.68</u>
IL	GameFormer [20]	83.76	94.73	88.12
	PlanTF [5]	83.66	94.02	92.67
	Gen-Drive (Pretrain) [21]	85.12	93.65	86.64
RL	Gen-Drive (Finetune) [21]	87.53	95.72	89.94
	CarPlanner (Ours)	91.45	<u>96.38</u>	95.37

Table 2. Comparison with SOTAs in Reduced-Val14 with non-reactive traffic agents.

set is 176,218 where all available scenario types are used, with a number of 4,000 scenarios per type. The size of the validation set is 1,118 where 100 scenarios with 14 types are selected. We train all models with 50 epochs in 2 NVIDIA 3090 GPUs. The batch size is 64 per GPU. We use AdamW optimizer with an initial learning rate of 1e-4 and reduce the learning rate when the validation loss stops decreasing with a patience of 0 and decrease factor of 0.3. For RL training, we set the discount $\gamma = 0.1$ and the GAE parameter $\lambda = 0.9$. The weights of value, policy, and entropy loss are set to 3, 100, and 0.001, respectively. The number of longitudinal modes is set to 12 and a maximum number of lateral modes are set to 5.

4.2. Comparison with SOTAs

SOTAs. We categorize the methods into Rule, IL, and RL based on the type of trajectory generator. (1) PDM [9] wins the nuPlan challenge 2023, its IL-based and rule-based variants are denoted as PDM-Closed and PDM-Open, respectively. PDM-Closed follows the generation-selection framework where IDM is used to generate multiple candidate trajectories and rule-based selector considering safety, progress, and comfort is used to select the best trajectory. (2) PLUTO [4] also obeys the generation-selection framework and uses contrastive IL to incorporate various data augmentation techniques and trains the generator. (3) Gen-Drive [21] is a concurrent work that follows a pretrain-finetune pipeline where IL is used to pretrain a diffusion-based planner and RL is used to finetune the denoising process based on a reward model trained by AI preference.

[†]<https://github.com/motional/nuplan-devkit>

Design Choices				Closed-loop metrics (\uparrow)					Open-loop metrics (\downarrow)	
Reward DE	Reward Quality	Coord Trans	KNN	CLS-NR	S-CR	S-Area	S-PR	S-Comfort	Loss Selector	Loss Generator
\times	\checkmark	\checkmark	\checkmark	31.79	95.74	98.45	33.10	48.84	1.03	30.3
\checkmark	\times	\checkmark	\checkmark	90.44	97.49	96.91	93.33	90.73	0.99	1221.6
\checkmark	\checkmark	\times	\checkmark	90.78	96.92	98.46	91.37	94.23	1.00	2130.7
\checkmark	\checkmark	\checkmark	\times	92.73	98.07	98.46	94.69	93.44	1.03	2083.6
\checkmark	\checkmark	\checkmark	\checkmark	94.07	99.22	99.22	95.06	91.09	1.03	1624.5

Table 3. Ablation studies on the design choices in RL training. Experimental results are based on the Test14-random non-reactive benchmark.

Results. We compare our method with SOTAs in Test14-Random and Reduced-Val14 benchmark as shown in Tab. 1 and Tab. 2. Overall, our CarPlanner demonstrates superior performance, particularly in non-reactive environments.

In the non-reactive setting, our method achieves the highest scores across all metrics, with an improvement of 4.02 and 2.15 compared to PDM-Closed and PLUTO, establishing the potential of RL and the superior performance of our proposed framework. Moreover, CarPlanner reveals substantial improvement in the progress metric S-PR compared to PDM-Closed in Tab. 2 and comparable collision metric S-CR, indicating the ability of our method to improving driving efficiency while maintaining safe driving. Importantly, we do not apply any techniques commonly used in IL such as data augmentation [4, 5] and ego-history masking [13], underscoring the intrinsic capability of our approach to solving the closed-loop task.

In the reactive setting, while our method performs well, it falls slightly short of PDM-Closed. This discrepancy arises because our model was trained exclusively in non-reactive settings and has not interacted with the IDM policy used by reactive settings; as a result, our model is less robust to disturbances generated by reactive agents during testing.

4.3. Ablation Studies

We investigate the effects of different design choices in RL training. The results are shown in Tab. 3.

Influence of reward items. The results demonstrate that the DE and quality rewards are complementary. When using the DE reward only, the planner tends to generate static trajectories and achieves a low progress metric. This occurs because the ego vehicle begins in a safe, drivable state, but moving forward is at risk of collisions or leaving the drivable area. On the other hand, compared to using DE reward only, incorporating the quality reward significantly improves closed-loop metrics. For instance, the S-CR metric rises from 97.49 to 99.22, and the S-Area metric rises from 96.91 to 99.22. These improvements indicate that the quality reward encourages safe and comfortable behaviors.

Effectiveness of IVM. The results show that the coordinate transformation and KNN techniques in IVM notably improve closed-loop metrics and generator loss. For instance, with the coordinate transformation technique, the

overall closed-loop score increases from 90.78 to 94.07, and S-PR rises from 91.37 to 95.06. These improvements are attributed to the enhanced accuracy of value estimation in RL, leading to generalized driving in closed-loop.

4.4. Extension to IL

In addition to designing for RL training, we also extend the CarPlanner to incorporate IL. We conduct rigorous analysis to compare the effects of various design choices in IL and RL training, as summarized in Tab. 4. Our findings indicate that while mode dropout and scorer side task contribute to both IL and RL training, ego-history dropout and backbone sharing, often effective in IL, are less suitable for RL.

Ego-history dropout. This technique effectively addresses the causal confusion issue in IL by preventing over-reliance on ego history, leading to more generalized planning behavior. Our experiments confirm that ego-history dropout positively impacts IL training, as it improves performance across closed-loop metrics like S-CR and S-Area. However, in RL training, we observe a negative impact on advantage estimation due to ego-history dropout, which significantly affects the accuracy of value estimation, leading to closed-loop performance degradation. This indicates that RL training naturally overcomes the causal confusion issue of IL, since explicit task-oriented metrics can be incorporated into reward functions, highlighting the potential of RL to pushing the boundaries of learning-based planning.

Backbone sharing. This choice, often used in IL-based multi-modal planners, promotes feature sharing across tasks to improve generalization. While backbone sharing helps IL by balancing losses across trajectory generator and selector, we find it adversely affects RL training. Specifically, backbone sharing leads to higher losses for both the trajectory generator and selector in RL, indicating that gradients from each task interfere. The divergent objectives in RL for trajectory generation and selection tasks seem to conflict, reducing overall policy performance. Consequently, we avoid backbone sharing in our RL framework to maintain task-specific gradient flow and improve policy quality.

4.5. Qualitative Results

We provide qualitative results as shown in Fig. 3. In this scenario, ego vehicle is required to execute a right

	Design Choices				Closed-loop metrics (\uparrow)					Open-loop metrics (\downarrow)	
Loss Type	Mode Dropout	Scorer Side Task	Ego-history Dropout	Backbone Sharing	CLS-NR	S-CR	S-Area	S-PR	S-Comfort	Loss Selector	Loss Generator
IL	\times	\times	\times	\times	90.82	97.29	98.45	92.15	94.57	<u>1.04</u>	147.5
	\checkmark	\times	\times	\times	91.21	96.54	98.46	91.44	96.92	1.07	<u>153.0</u>
	\checkmark	\checkmark	\times	\times	91.51	96.91	98.46	95.30	96.91	1.04	162.3
	\checkmark	\checkmark	\checkmark	\times	92.72	98.06	98.84	94.88	95.35	<u>1.04</u>	167.5
	\checkmark	\checkmark	\checkmark	\checkmark	93.41	<u>98.85</u>	98.85	93.87	96.15	<u>1.04</u>	174.3
RL	\times	\times	\times	\times	91.67	98.84	98.84	91.69	90.73	<u>1.04</u>	1812.6
	\checkmark	\times	\times	\times	93.46	98.07	99.61	94.26	92.28	1.09	2254.6
	\checkmark	\checkmark	\times	\times	94.07	99.22	99.22	95.06	91.09	1.03	1624.5
	\checkmark	\checkmark	\checkmark	\times	89.51	97.27	98.44	90.93	83.20	1.05	5424.3
	\checkmark	\checkmark	\checkmark	\checkmark	88.66	95.54	98.84	92.82	86.05	1.21	1928.1

Table 4. Effect of different components on IL and RL Loss using our CarPlanner. Experimental results are based on the Test14-random non-reactive benchmark.

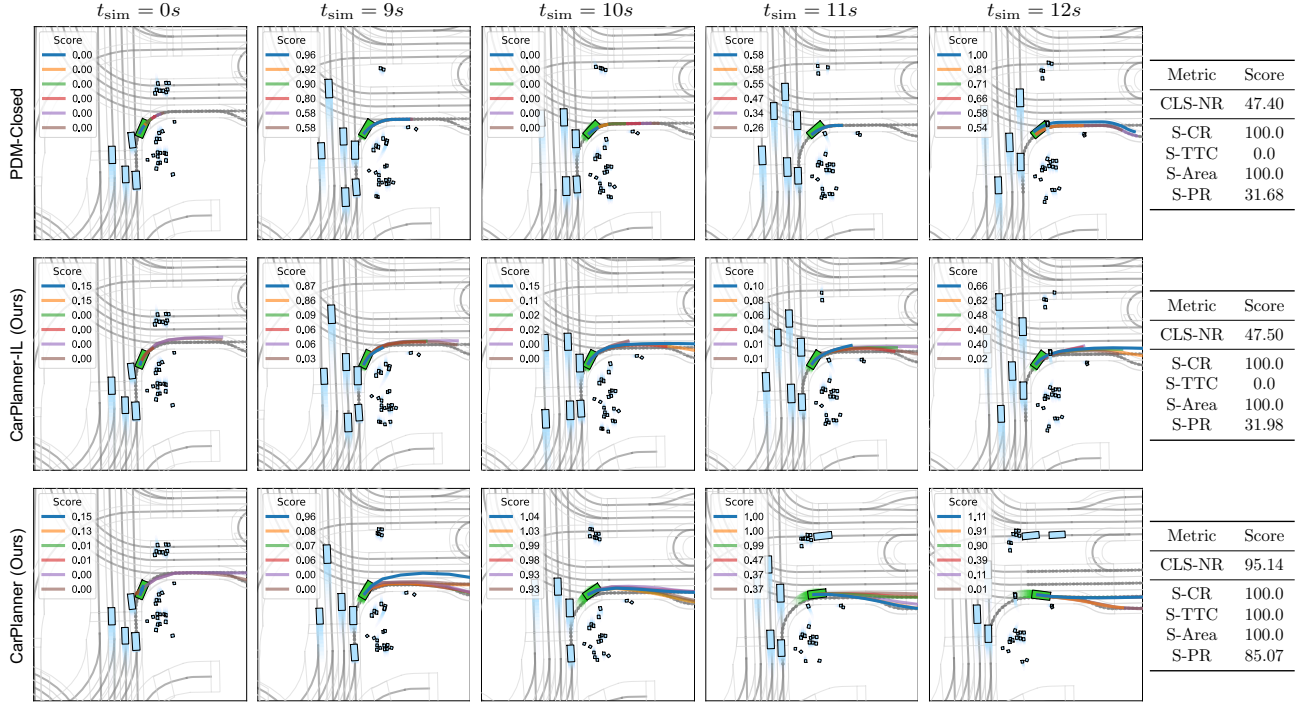


Figure 3. Qualitative comparison of PDM-Closed and our method in non-reactive environments. The scenario is annotated as waiting_for_pedestrian_to_cross. In each frame shot, ego vehicle is marked as green. Traffic agents are marked as sky blue. Lineplot with blue is the ego planned trajectory.

turn while navigating around pedestrians. In this case, Our method demonstrates a smooth, efficient performance. From $t_{\text{sim}} = 0s$ to $t_{\text{sim}} = 9s$, all methods wait for the pedestrians to cross the road. At $t_{\text{sim}} = 10s$, an unexpected pedestrian goes back and prepares to re-cross the road. PDM-Closed is unaware of this situation and takes an emergency stop but still intersects with this pedestrian. In contrast, our IL variant displays an awareness of the pedestrian’s movements and consequently conducts a braking maneuver. However, it still remains close to the pedestrian. Our RL method avoids this hazard by starting up early up to $t_{\text{sim}} = 9s$ and achieves the highest progress and safety metrics.

5. Conclusion

In this paper, we introduce CarPlanner, a consistent auto-regressive planner aiming at large-scale RL training. Thanks to the proposed framework, we train an RL-based planner that outperforms existing RL-, IL-, and rule-based SOTAs. Furthermore, we provide analysis indicating the characteristics of IL and RL, highlighting the potential of RL to take a further step toward learning-based planning.

Limitations and future work. RL needs delicate design and is prone to input representation. Besides, RL can overfit its training environment and suffer from performance drop in unseen environments. Future work aims to develop robust RL algorithms to tackle these limitations.

References

- [1] Mayank Bansal, Alex Krizhevsky, and Abhijit Ogale. Chauffeurnet: Learning to drive by imitating the best and synthesizing the worst. In *Proc. of Robotics: Science and Systems (RSS)*, 2019. 3, 6
- [2] Holger Caesar, Juraj Kabzan, Kok Seang Tan, Whye Kit Fong, Eric Wolff, Alex Lang, Luke Fletcher, Oscar Beijbom, and Sammy Omari. nuplan: A closed-loop ml-based planning benchmark for autonomous vehicles. *arXiv preprint arXiv:2106.11810*, 2021. 1, 2, 6
- [3] Sergio Casas, Abbas Sadat, and Raquel Urtasun. Mp3: A unified model to map, perceive, predict and plan. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 14403–14412, 2021. 1
- [4] Jie Cheng, Yingbing Chen, and Qifeng Chen. Pluto: Pushing the limit of imitation learning-based planning for autonomous driving. *arXiv preprint arXiv:2404.14327*, 2024. 3, 6, 7
- [5] Jie Cheng, Yingbing Chen, Xiaodong Mei, Bowen Yang, Bo Li, and Ming Liu. Rethinking imitation-based planners for autonomous driving. In *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*, pages 14123–14130. IEEE, 2024. 3, 6, 7
- [6] Kashyap Chitta, Aditya Prakash, Bernhard Jaeger, Zehao Yu, Katrin Renz, and Andreas Geiger. Transfuser: Imitation with transformer-based sensor fusion for autonomous driving. *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, 2022. 2
- [7] Ignasi Clavera, Yao Fu, and Pieter Abbeel. Model-augmented actor-critic: Backpropagating through paths. In *Proc. of the Int. Conf. on Learning Representations (ICLR)*, 2020. 4
- [8] Felipe Codevilla, Matthias Müller, Antonio López, Vladlen Koltun, and Alexey Dosovitskiy. End-to-end driving via conditional imitation learning. In *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*, pages 4693–4700, 2018. 2
- [9] Daniel Dauner, Marcel Hallgarten, Andreas Geiger, and Kashyap Chitta. Parting with misconceptions about learning-based vehicle motion planning. In *Proc. of the Conf. on Robot Learning (CoRL)*, 2023. 6
- [10] Pim De Haan, Dinesh Jayaraman, and Sergey Levine. Causal confusion in imitation learning. *Proc. of the Advances in Neural Information Processing Systems (NIPS)*, 32, 2019. 1, 3
- [11] 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 *Proc. of the IEEE/CVF Intl. Conf. on Computer Vision (ICCV)*, pages 9710–9719, 2021. 2
- [12] 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 *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 11525–11533, 2020. 3
- [13] Ke Guo, Wei Jing, Junbo Chen, and Jia Pan. CCIL: Context-conditioned imitation learning for urban driving. In *Proc. of Robotics: Science and Systems (RSS)*, 2023. 1, 6, 7
- [14] Marcel Hallgarten, Martin Stoll, and Andreas Zell. From prediction to planning with goal conditioned lane graph traversals. *arXiv preprint arXiv:2302.07753*, 2023. 6
- [15] Nicklas A Hansen, Hao Su, and Xiaolong Wang. Temporal difference learning for model predictive control. In *Proc. of the Intl. Conf. on Machine Learning (ICML)*, pages 8387–8406. PMLR, 2022. 4
- [16] Xiangkun He, Haohan Yang, Zhongxu Hu, and Chen Lv. Robust lane change decision making for autonomous vehicles: An observation adversarial reinforcement learning approach. *IEEE Transactions on Intelligent Vehicles*, 8(1):184–193, 2022. 3
- [17] Jonathan Ho, Jayesh Gupta, and Stefano Ermon. Model-free imitation learning with policy optimization. In *Proc. of the Intl. Conf. on Machine Learning (ICML)*, pages 2760–2769, 2016. 2
- [18] John Houston, Guido Zuidhof, Luca Bergamini, Yawei Ye, Long Chen, Ashesh Jain, Sammy Omari, Vladimir Iglovikov, and Peter Ondruska. One thousand and one hours: Self-driving motion prediction dataset. In *Proc. of the Conf. on Robot Learning (CoRL)*, pages 409–418, 2021. 2
- [19] Yihan Hu, Jiazhi Yang, Li Chen, Keyu Li, Chonghao Sima, Xizhou Zhu, Siqi Chai, Senyao Du, Tianwei Lin, Wenhai Wang, et al. Planning-oriented autonomous driving. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 17853–17862, 2023. 1
- [20] Zhiyu Huang, Haochen Liu, and Chen Lv. Gameformer: Game-theoretic modeling and learning of transformer-based interactive prediction and planning for autonomous driving. In *Proc. of the IEEE/CVF Intl. Conf. on Computer Vision (ICCV)*, pages 3903–3913, 2023. 2, 6, 3
- [21] Zhiyu Huang, Xinshuo Weng, Maximilian Igl, Yuxiao Chen, Yulong Cao, Boris Ivanovic, Marco Pavone, and Chen Lv. Gen-drive: Enhancing diffusion generative driving policies with reward modeling and reinforcement learning fine-tuning. *arXiv preprint arXiv:2410.05582*, 2024. 2, 6
- [22] Julian Ibarz, Jie Tan, Chelsea Finn, Mrinal Kalakrishnan, Peter Pastor, and Sergey Levine. How to train your robot with deep reinforcement learning: lessons we have learned. *Intl. Journal of Robotics Research (IJRR)*, 40(4-5):698–721, 2021. 1
- [23] Chiyu Jiang, Andre Cornman, Cheolho Park, Benjamin Sapp, Yin Zhou, Dragomir Anguelov, et al. Motiondiffuser: Controllable multi-agent motion prediction using diffusion. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 9644–9653, 2023. 2
- [24] Edouard Leurent and Jean Mercat. Social attention for autonomous decision-making in dense traffic. *arXiv preprint arXiv:1911.12250*, 2019. 3
- [25] Guofa Li, Yifan Yang, Shen Li, Xingda Qu, Nengchao Lyu, and Shengbo Eben Li. Decision making of autonomous vehicles in lane change scenarios: Deep reinforcement learning approaches with risk awareness. *Transportation research part C: emerging technologies*, 134:103452, 2022. 3

- [26] Weiwei Liu, Wei Jing, Ke Guo, Gang Xu, Yong Liu, et al. Traco: Learning virtual traffic coordinator for cooperation with multi-agent reinforcement learning. In *Conference on Robot Learning*, pages 2465–2477. PMLR, 2023. 3
- [27] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Proc. of the Advances in Neural Information Processing Systems (NIPS)*, 35:27730–27744, 2022. 1
- [28] Leif E Peterson. K-nearest neighbor. *Scholarpedia*, 4(2): 1883, 2009. 5
- [29] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 652–660, 2017. 4
- [30] Katrin Renz, Kashyap Chitta, Otniel-Bogdan Mercea, A Sophia Koepke, Zeynep Akata, and Andreas Geiger. Plant: Explainable planning transformers via object-level representations. In *Proc. of the Conf. on Robot Learning (CoRL)*, 2022. 2
- [31] Nicholas Rhinehart, Rowan McAllister, Kris Kitani, and Sergey Levine. Precog: Prediction conditioned on goals in visual multi-agent settings. In *Proc. of the IEEE/CVF Intl. Conf. on Computer Vision (ICCV)*, pages 2821–2830, 2019. 2
- [32] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 627–635, 2011. 1, 3
- [33] Oliver Scheel, Luca Bergamini, Maciej Wolczyk, Błażej Osiński, and Peter Ondruska. Urban driver: Learning to drive from real-world demonstrations using policy gradients. In *Proc. of the Conf. on Robot Learning (CoRL)*, pages 718–728, 2022. 1, 2, 6
- [34] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015. 1
- [35] Ari Seff, Brian Cera, Dian Chen, Mason Ng, Aurick Zhou, Nigamaa Nayakanti, Khaled S Refaat, Rami Al-Rfou, and Benjamin Sapp. Motionlm: Multi-agent motion forecasting as language modeling. In *Proc. of the IEEE/CVF Intl. Conf. on Computer Vision (ICCV)*, pages 8579–8590, 2023. 2
- [36] Shaoshuai Shi, Li Jiang, Dengxin Dai, and Bernt Schiele. Motion transformer with global intention localization and local movement refinement. *Proc. of the Advances in Neural Information Processing Systems (NIPS)*, 35:6531–6543, 2022. 2
- [37] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016. 1
- [38] Simon Suo, Sebastian Regalado, Sergio Casas, and Raquel Urtasun. Trafficsim: Learning to simulate realistic multi-agent behaviors. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 10400–10409, 2021. 3
- [39] Ardi Tampuu, Tambet Matiisen, Maksym Semikin, Dmytro Fishman, and Naveed Muhammad. A survey of end-to-end driving: Architectures and training methods. *IEEE Transactions on Neural Networks and Learning Systems*, 33(4): 1364–1384, 2020. 1
- [40] Martin Treiber, Ansgar Hennecke, and Dirk Helbing. Congested traffic states in empirical observations and microscopic simulations. *Physical review E*, 62(2):1805, 2000. 6
- [41] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Proc. of the Advances in Neural Information Processing Systems (NIPS)*, 30, 2017. 4
- [42] Huanjie Wang, Shihua Yuan, Mengyu Guo, Xueyuan Li, and Wei Lan. A deep reinforcement learning-based approach for autonomous driving in highway on-ramp merge. *Proceedings of the Institution of Mechanical engineers, Part D: Journal of Automobile engineering*, 235(10-11):2726–2739, 2021. 3
- [43] Letian Wang, Jie Liu, Hao Shao, Wenshuo Wang, Ruobing Chen, Yu Liu, and Steven L Waslander. Efficient Reinforcement Learning for Autonomous Driving with Parameterized Skills and Priors. In *Proc. of Robotics: Science and Systems (RSS)*, Daegu, Republic of Korea, 2023. 3
- [44] Wei Wu, Xiaoxin Feng, Ziyang Gao, and Yuheng Kan. Smart: Scalable multi-agent real-time simulation via next-token prediction. *arXiv preprint arXiv:2405.15677*, 2024. 2
- [45] Dongkun Zhang, Jiaming Liang, Sha Lu, Ke Guo, Qi Wang, Rong Xiong, Zhenwei Miao, and Yue Wang. Pep: Policy-embedded trajectory planning for autonomous driving. *IEEE Robotics and Automation Letters*, 2024. 3, 6
- [46] Zhejun Zhang, Alexander Liniger, Christos Sakaridis, Fisher Yu, and Luc V Gool. Real-time motion prediction via heterogeneous polyline transformer with relative pose encoding. *Proc. of the Advances in Neural Information Processing Systems (NIPS)*, 36, 2024. 3
- [47] Tong Zhou, Letian Wang, Ruobing Chen, Wenshuo Wang, and Yu Liu. Accelerating reinforcement learning for autonomous driving using task-agnostic and ego-centric motion skills. In *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, pages 11289–11296. IEEE, 2023. 3
- [48] Yang Zhou, Hao Shao, Letian Wang, Steven L Waslander, Hongsheng Li, and Yu Liu. Smartrefine: A scenario-adaptive refinement framework for efficient motion prediction. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 15281–15290, 2024. 2
- [49] Zikang Zhou, Haibo Hu, Xinhong Chen, Jianping Wang, Nan Guan, Kui Wu, Yung-Hui Li, Yu-Kai Huang, and Chun Jason Xue. Behaviorgpt: Smart agent simulation for autonomous driving with next-patch prediction. *arXiv preprint arXiv:2405.17372*, 2024. 2

CarPlanner: Consistent Auto-regressive Trajectory Planning for Large-scale Reinforcement Learning in Autonomous Driving

Supplementary Material

6. Training Procedure

Algorithm 1 outlines the training process for the CarPlanner framework. The procedure involves two primary steps: (1) training the non-reactive world model, and (2) training the mode selector and the trajectory generator. The definitions of the loss functions are given in the following.

Loss of non-reactive world model. The non-reactive world model β is trained to simulate agent trajectories based on the initial state \mathbf{s}_0 . For each data sample $(\mathbf{s}_0, s_{1:T}^{1:N, \text{gt}}) \in \mathcal{D}$, the model predicts trajectories $s_{1:T}^{1:N} = \beta(\mathbf{s}_0)$, and the training objective minimizes the L1 loss:

$$L_{\text{wm}} = \frac{1}{T} \sum_{t=1}^T \sum_{n=1}^N \|s_t^n - s_t^{n, \text{gt}}\|_1. \quad (5)$$

Mode selector loss. The mode selector assigns scores $\sigma = f_{\text{selector}}(\mathbf{s}_0, \mathbf{c})$ to each candidate mode c_i , and the ground-truth positive mode c^* is used to compute the cross-entropy loss:

$$L_{\text{selector}} = - \sum_{i=1}^{N_{\text{mode}}} \mathbb{I}(c_i = c^*) \log \sigma_i, \quad (6)$$

where M is the number of candidate modes, and \mathbb{I} is the indicator function.

Generator loss with RL. The PPO loss consists of three parts: policy, value, and entropy loss.

The policy loss is defined as:

$$\begin{aligned} & \text{PolicyLoss}(a_{0:T-1}, d_{0:T-1, \text{new}}, d_{0:T-1}, A_{0:T-1}) \\ &= - \frac{1}{T} \sum_{t=0}^{T-1} \min(r_t A_t, \text{clip}(r_t, 1 - \epsilon, 1 + \epsilon) A_t), \end{aligned} \quad (7)$$

where the ratio r_t is given by $r_t = \frac{\text{Prob}(a_t, d_{t, \text{new}})}{\text{Prob}(a_t, d_t)}$, $d_{t, \text{new}}$ and d_t are the policy distributions (mean and standard deviation of Gaussian distribution) at time step t induced by π and π_{old} respectively, the function $\text{Prob}(a, d)$ calculates the probability of a given action a under a distribution d , and A_t is the advantage estimated using GAE [34].

The value and entropy loss are defined as:

$$\text{ValueLoss}(V_{0:T-1, \text{new}}, \hat{R}_{0:T-1}) = \frac{1}{T} \sum_{t=0}^{T-1} \|V_{t, \text{new}} - \hat{R}_t\|_2^2, \quad (8)$$

$$\text{Entropy}(d_{0:T-1, \text{new}}) = \frac{1}{T} \sum_{t=0}^{T-1} \mathcal{H}(d_{t, \text{new}}), \quad (9)$$

Algorithm 1 Training Procedure of CarPlanner

```

1: Input: Dataset  $\mathcal{D}$  containing initial states  $\mathbf{s}_0$  and ground-truth
   trajectories  $s_{1:T}^{0:N, \text{gt}}$ , longitudinal modes  $\mathbf{c}_{\text{lon}}$ , discount factor  $\gamma$ ,
   GAE parameter  $\lambda$ , update interval  $I$ .
2: Require: Non-reactive world model  $\beta$ , mode selector  $f_{\text{selector}}$ ,
   policy  $\pi$ , policy old  $\pi_{\text{old}}$ .
3: Step 1: Training World Model
4: for  $(\mathbf{s}_0, s_{1:T}^{1:N, \text{gt}}) \in \mathcal{D}$  do
5:   Simulate agent trajectories  $s_{1:T}^{1:N} \leftarrow \beta(\mathbf{s}_0)$ 
6:   Calculate loss  $L_{\text{wm}} \leftarrow \text{L1Loss}(s_{1:T}^{1:N}, s_{1:T}^{1:N, \text{gt}})$ 
7:   Backpropagate and update  $\beta$  using  $L_{\text{wm}}$ 
8: end for
9: Step 2: Training Selector and Generator
10: Initialize training_step  $\leftarrow 0$ 
11: Initialize policy old  $\pi_{\text{old}} \leftarrow \pi$ 
12: for  $(\mathbf{s}_0, s_{1:T}^{0, \text{gt}}) \in \mathcal{D}$  do
13:   Non-Reactive World Model:
14:   Simulate agent trajectories  $s_{1:T}^{1:N} \leftarrow \beta(\mathbf{s}_0)$ 
15:   Mode Assignment:
16:   Determine  $\mathbf{c}_{\text{lat}}$  based on  $\mathbf{s}_0$ 
17:   Concatenate  $\mathbf{c}_{\text{lat}}$  and  $\mathbf{c}_{\text{lon}}$  to get  $\mathbf{c}$ 
18:   Determine positive mode  $c^*$  based on  $s_{1:T}^{0, \text{gt}}$  and  $\mathbf{c}$ 
19:   Mode Selector Loss:
20:   Compute scores  $\sigma \leftarrow f_{\text{selector}}(\mathbf{s}_0, \mathbf{c})$ 
21:    $L_{\text{selector}} \leftarrow \text{CrossEntropyLoss}(\sigma, c^*)$ 
22:   Generator Loss:
23:   if Reinforcement Learning (RL) Training then
24:     Use  $\pi_{\text{old}}$ ,  $\mathbf{s}_0$ ,  $c^*$ , and  $s_{1:T}^{1:N}$  to collect rollout data
      $(\mathbf{s}_{0:T-1}, a_{0:T-1}, d_{0:T-1}, V_{0:T-1}, R_{0:T-1})$ 
25:     Compute advantage  $A_{0:T-1}$  and return  $\hat{R}_{0:T-1}$  using
     GAE [34]:  $A_{0:T-1}, \hat{R}_{0:T-1} \leftarrow \text{GAE}(R_{0:T-1}, V_{0:T-1}, \gamma, \lambda)$ 
26:     Compute policy distribution and value estimates:
      $(d_{0:T-1, \text{new}}, V_{0:T-1, \text{new}}) \leftarrow \pi(\mathbf{s}_{0:T-1}, a_{0:T-1}, c^*)$ 
27:      $L_{\text{generator}} \leftarrow \text{ValueLoss}(V_{0:T-1, \text{new}}, \hat{R}_{0:T-1}) +$ 
      $\text{PolicyLoss}(d_{0:T-1, \text{new}}, d_{0:T-1}, A_{0:T-1}) -$ 
      $\text{Entropy}(d_{0:T-1, \text{new}})$ 
28:   else if Imitation Learning (IL) Training then
29:     Use  $\pi$ ,  $\mathbf{s}_0$ ,  $c^*$ , and  $s_{1:T}^{1:N}$  to collect action sequence
      $a_{0:T-1}$ 
30:     Stack action sequence as ego-planned trajectory
      $s_{1:T}^0 \leftarrow \text{Stack}(a_{0:T-1})$ 
31:      $L_{\text{generator}} \leftarrow \text{L1Loss}(s_{1:T}^0, s_{1:T}^{0, \text{gt}})$ 
32:   end if
33:   Overall Loss:
34:    $L \leftarrow L_{\text{selector}} + L_{\text{generator}}$ 
35:   Backpropagate and update  $f_{\text{selector}}$ ,  $\pi$  using  $L$ 
36:   Policy Update:
37:   Increment training_step  $\leftarrow$  training_step + 1
38:   if training_step %  $I == 0$  then
39:     Update  $\pi_{\text{old}} \leftarrow \pi$ 
40:   end if
41: end for

```

Parameter	Value
Feature dimension D	256
Static point dimension D_m	9
Agent pose dimension D_a	10
Activation	ReLU
Number of layers	3
Number of attention heads	8
Dropout	0.1
discount factor γ	0.1
GAE parameter λ	0.9
Clip range ϵ	0.2
Update interval I	8
Weight of selector loss	1
Weight of value loss	3
Weight of policy loss	100
Weight of entropy loss	0.001
Weight of IL loss	1

Table 5. Hyperparameters of model architecture, PPO-related parameters, and loss weights.

where $V_{t,\text{new}}$ and R_t are the predicted and actual returns, and \mathcal{H} represents the entropy of the policy distribution d .

Generator loss with IL. In IL, the generator minimizes the trajectory error between the ego-planned trajectory $s_{1:T}^0$ and the ground-truth trajectory $s_{1:T}^{0,\text{gt}}$. The loss is defined as:

$$L_{\text{generator}} = \frac{1}{T} \sum_{t=1}^T \left\| s_t^0 - s_t^{0,\text{gt}} \right\|_1. \quad (10)$$

7. Implementation Details

The hyperparameters of model architecture, PPO-related parameters, and loss weights are summarized in Tab. 5. Note that the magnitude of value, policy, and entropy loss are 10^3 , 10^0 , and 10^{-3} respectively. The trajectory generator generates trajectories with a time horizon of 8 seconds in 1-second interval, corresponding to time horizon $T = 8$. During testing, these trajectories are interpolated to a 0.1-second interval. The ratio of scores generated by the rule and mode selectors is set at 1 : 0.3. In cases where no ego candidate trajectory satisfies the safety criteria evaluated by the rule selector, an emergency stop is triggered. For the Test14-Random benchmark, a replanning frequency of 10Hz is employed, adhering to the official nuPlan simulation configuration. In contrast, for the Reduced-Val14 benchmark, a replanning frequency of 1Hz is used to ensure a fair comparison with Gen-Drive [21].

Model Type	Design Choices		Closed-loop metrics (\uparrow)				
	Random Sample	Guide Reward	CLS-NR	S-CR	S-Area	S-PR	S-Comfort
Vanilla	✓	Progress	67.56 \pm 0.38	90.97 \pm 0.78	94.64 \pm 1.72	72.17 \pm 0.21	64.21 \pm 1.29
		DE	86.89 \pm 0.28	97.34 \pm 0.37	96.36 \pm 0.18	89.90 \pm 0.11	94.03 \pm 0.65
Consistent	✗	FE	88.14	96.86	98.43	91.39	73.73
	✗	DE	94.07	99.22	99.22	95.06	91.09

Table 6. Comparison of vanilla and consistent auto-regressive frameworks with different guide reward design. Experimental results are based on the Test14-Random non-reactive benchmark.

World Model	Closed-loop metrics (\uparrow)				
	CLS-NR	S-CR	S-Area	S-PR	S-Comfort
Reactive	91.03	96.92	99.23	91.28	90.00
Non-reactive	94.07	99.22	99.22	95.06	91.09

Table 7. Comparison of the usage of reactive and non-reactive world models. Experimental results are based on the Test14-Random non-reactive benchmark.

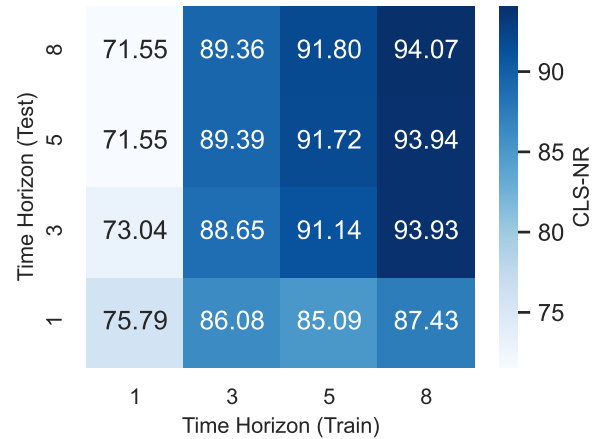


Figure 4. Performance of different training time horizons under different testing time horizons. The value in each cell is the CLS-NR metric on the Test14-Random non-reactive benchmark.

8. Ablation Study on RL Training

In this part, we examine the performance of vanilla and consistent auto-regressive frameworks, the use of reactive and non-reactive world model in RL training, and the impact of varying the time horizon.

Vanilla vs. consistent auto-regressive framework. The results are shown in Tab. 6. The consistent auto-regressive framework generates multi-modal trajectories by conditioning on mode representations. In contrast, the vanilla framework relies on random sampling from the action Gaussian distribution to produce multi-modal trajectories. To ensure comparability in the number of modes generated by both frameworks, we sample 60 trajectories in parallel for the vanilla framework. Given that random sampling introduces variability, we average the results across 3 random seeds. For the consistent framework, we use displacement

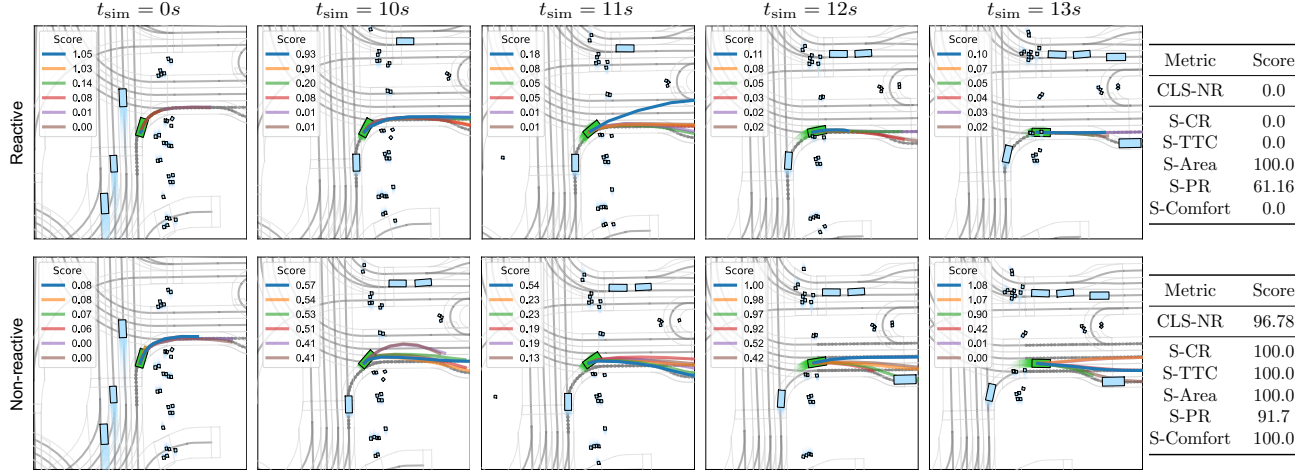


Figure 5. Qualitative comparison of using reactive and non-reactive world model in non-reactive environments. The scenario is annotated as `waiting_for_pedestrian_to_cross`. In each frame shot, ego vehicle is marked as green. Traffic agents are marked as sky blue. Lineplot with blue is the ego planned trajectory.

error (DE) and final error (FE) as guide functions to assist the policy in generating mode-aligned trajectories. For the vanilla framework, DE is compared against a progress reward, which encourages longitudinal movement along the route while discouraging excessive lateral deviations that move the vehicle too far from any possible route. Overall, our proposed consistent framework outperforms the vanilla framework in terms of closed-loop performance, highlighting the benefits of incorporating consistency. Additionally, we find that DE serves as an effective guide function for policy training, further enhancing closed-loop performance.

Reactive vs. non-reactive world model. We compare the performance of the CarPlanner framework when trained with reactive and non-reactive world models. The reactive world model shares a similar architecture with the autoregressive planner for the ego vehicle, utilizing relative pose encoding [46] as the backbone network to extract features of traffic agents and predict their subsequent poses. The training loss and hyperparameters are consistent with those used for the non-reactive world model. As shown in Tab. 7, except for the S-Area metric, using non-reactive world model outperforms the reactive world model in our current implementation. The primary difference lies in the assumptions about traffic agents: the reactive world model assumes that the ego vehicle can negotiate with traffic agents and share the same priority, whereas in the non-reactive model, traffic agents do not respond to the ego vehicle, effectively assigning them higher priority. A representative example is presented in Fig. 5. When trained with the reactive world model, the planner assumes pedestrians will yield to the vehicle, leading it to attempt to move forward. However, at $t_{\text{sim}} = 12s$, the planner collides with pedestrians, triggering an emergency brake, which negatively impacts safety, progress, and comfort metrics. Although the performance

of using reactive world model is not satisfied currently, it is a more realistic assumption and we will further investigate this in future work.

Time horizon. We evaluate the CarPlanner framework by training it with different time horizons, including 1, 3, 5, and 8 seconds, and testing the planners in each time horizon. The results in Fig. 4 confirm that increasing the time horizon has a positive effect on the performance for both training and testing. A special case is when the training time horizon is set to 1, all tested time horizons exhibit poor performance, highlighting the importance of multi-step learning in RL. Additionally, the observation that increasing the training time horizon enhances closed-loop performance suggests the potential for further improvements by extending the time horizon beyond 8 seconds. However, due to current limitations in data preparation, which is designed for horizons up to 8 seconds, expanding the time horizon would not provide map information or ground-truth trajectories, hindering further analysis. Consequently, we leave this exploration for future work.

9. Comparison with Differentiable Loss

In typical IL setting, the supervision signal provided to the trajectory generator is the displacement error (DE) between the ego-planned trajectory and the ground-truth trajectory. Several works [4, 20, 38] propose to convert non-differentiable metrics, such as avoiding collision (Col) and adherence to drivable area (Area), into differentiable loss functions that can directly backpropagate to the generator. In contrast, CarPlanner leverages an RL framework, which introduces surrogate objectives to indirectly optimize these non-differentiable metrics.

In this part, we compare these two approaches which

	Supervision Signals			Closed-loop metrics (\uparrow)					Open-loop metrics (\downarrow)	
Loss Type	DE	Col	Area	CLS-NR	S-CR	S-Area	S-PR	S-Comfort	Col Mean [Min, Max]	Area Mean [Min, Max]
IL	✓	✗	✗	93.41	98.85	<u>98.85</u>	93.87	96.15	0.17 [0.00 , 0.47]	0.09 [0.00 , 0.40]
	✓	✓	✗	<u>93.67</u>	99.23	<u>98.85</u>	<u>94.63</u>	94.23	0.16 [0.00 , 0.43]	<u>0.07</u> [0.00 , 0.27]
	✓	✗	✓	93.12	98.46	98.84	92.88	94.21	<u>0.15</u> [0.00 , 0.44]	0.08 [0.00 , 0.30]
	✓	✓	✓	93.32	98.46	98.46	94.05	<u>95.77</u>	<u>0.15</u> [0.00 , 0.43]	0.09 [0.00 , 0.39]
RL	✓	✗	✗	90.44	97.49	96.91	93.33	90.73	0.17 [0.00 , 0.49]	0.14 [0.00 , 0.51]
	✓	✓	✓	94.07	<u>99.22</u>	99.22	95.06	91.09	0.12 [0.00 , 0.39]	0.05 [0.00 , 0.22]

Table 8. Comparison with different loss types and supervision signals. Closed-loop results are based on the Test14-Random non-reactive benchmark. Open-loop results are on validation set.

	Design Choices						Closed-loop metrics (\uparrow)				
Loss Type	Model Type	Mode Type	Mode Dropout	Scorer Side Task	Ego-history Dropout	Backbone Sharing	CLS-NR	S-CR	S-Area	S-PR	S-Comfort
IL	Vanilla	-	-	-	✓	-	86.48	97.09	97.29	88.05	94.19
	Consistent	Lon	✗	✓	✓	✓	88.79	96.67	96.08	89.63	94.90
	Consistent	Lon-Lat	✓	✓	✓	✓	93.41	98.85	98.85	93.87	96.15
RL	Vanilla	-	-	-	✗	-	85.56	97.27	95.70	89.17	93.36
	Consistent	Lon	✗	✓	✗	✗	90.57	97.30	97.68	92.20	94.59
	Consistent	Lon-Lat	✓	✓	✗	✗	94.07	99.22	99.22	95.06	91.09

Table 9. Effect of different mode representations. Experimental results are based on the Test14-Random non-reactive benchmark.

provide rich supervision signals to the trajectory generator. The results are summarized in Tab. 8. In IL training, the Col and Area metrics are converted into differentiable loss functions, whereas in RL training, Col and Area are treated as reward functions, contributing to the quality reward as described in the main paper. It is important to note that the implementations for differentiable loss functions and reward functions are identical, except that gradient flow is enabled for differentiable loss functions. The open-loop metrics compute the Col and Area values across all candidate multi-modal trajectories, with the Mean, Min, and Max referring to the mean, minimum, and maximum values of the Col and Area metrics within the candidate trajectory set.

Our findings suggest that incorporating Col loss benefits the open-loop Col metric and improves the closed-loop S-CR metrics, thereby enhancing closed-loop performance. However, incorporating Area loss results in better open-loop Area metrics but deteriorates closed-loop performance. Compared to differentiable loss functions, RL with Col and Area as quality rewards yields the trajectory set with the highest overall quality, as evidenced by smaller Mean and Max metrics in open-loop metrics. This improvement can be attributed to RL’s ability to optimize the reward-to-go using surrogate objectives that account for future rewards, while differentiable loss functions are limited to timewise-aligned optimization in our current implementation. This distinction is illustrated in Fig. 6: in (a), the loss at time step t is directly computed from s_t^0 , meaning that during

backward propagation, the loss at time step t cannot influence the optimization of prior time steps. In (b), however, the non-differentiable reward is aggregated into a return (reward-to-go), which serves as a reference for computing the loss at time step t . Through this process, the reward at time step t can influence the trajectory at earlier time steps t' ($t' < t$). In the future, we aim to combine the advantages of differentiable loss which can provide low-variance gradients, and RL which can provide long-term foresight, by model-based RL optimization techniques [7, 15].

10. Effect of Mode Representation

In this part, we examine the impact of mode representations on performance. The results are presented in Tab. 9. For both the vanilla and consistent frameworks, we disable the use of random sampling to focus solely on mode-aligned trajectories. As a result, the vanilla framework can only generate single-modal trajectories, leading to the lowest performance. In the consistent framework, we explore two types of mode representations: Lon and Lon-Lat. The Lon representation assigns modes based on longitudinal movements along the route, whereas the Lon-Lat representation decomposes modes by both longitudinal and lateral movements. Aligned with the main paper, we use ego-history dropout and backbone sharing only for IL training. For the Lon representation, we close mode dropout since it does not rely on any map or agent representation in initial state.

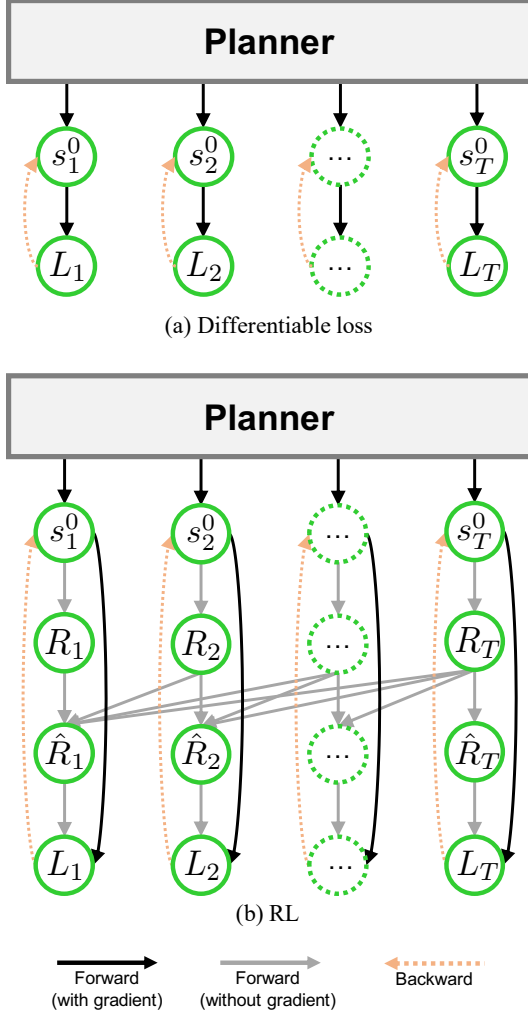


Figure 6. The computational graph of differentiable loss (a) and RL (b) framework for optimizing same metrics such as displacement errors, collision avoidance, and adherence to drivable area.

The results indicate that introducing consistency provides greater benefits to RL training, with the Lon-Lat representation proving to be more effective than the Lon representation. This suggests that decomposing mode representations into both longitudinal and lateral components enhances the model's ability by providing more explicit mode information.