

RIFT: GROUP-RELATIVE RL FINE-TUNING FOR REALISTIC AND CONTROLLABLE TRAFFIC SIMULATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Achieving both realism and controllability in closed-loop traffic simulation remains a key challenge in autonomous driving. Dataset-based methods reproduce realistic trajectories but suffer from *covariate shift* in closed-loop deployment, compounded by simplified dynamics models that further reduce reliability. Conversely, physics-based simulation methods enhance reliable and controllable closed-loop interactions but often lack expert demonstrations, compromising realism. To address these challenges, we introduce a dual-stage AV-centric simulation framework that conducts imitation learning pre-training in a data-driven simulator to capture trajectory-level realism and route-level controllability, followed by reinforcement learning fine-tuning in a physics-based simulator to enhance style-level controllability and mitigate covariate shift. In the fine-tuning stage, we propose *RIFT*, a novel group-relative RL fine-tuning strategy that evaluates all candidate modalities through group-relative formulation and employs a surrogate objective for stable optimization, enhancing style-level controllability and mitigating covariate shift while preserving the trajectory-level realism and route-level controllability inherited from IL pre-training. Extensive experiments demonstrate that *RIFT* improves realism and controllability in traffic simulation while simultaneously exposing the limitations of modern AV systems in closed-loop evaluation.

1 INTRODUCTION

Reliable closed-loop traffic simulation is critical for developing advanced autonomous vehicle (AV) systems, supporting training and evaluation Feng et al. (2023b); Ding et al. (2023). An ideal traffic simulation should possess two key properties: *realistic*, reflecting real-world driving behavior; *controllable*, enabling customizable traffic simulation according to user requirements.

To balance these two essential properties, existing traffic simulation methods adopt different trade-offs depending on the underlying platform, often favoring either realism or controllability, as illustrated in Figure 1. Methods based on data-driven simulators exploit real-world data to generate realistic trajectories by learning multimodal behavioral patterns through imitation learning (IL) Ngiam et al. (2021); Sun et al. (2022); Feng et al. (2023a); Mahjourian et al. (2024). In addition to realism, recent studies on data-driven simulators have pursued controllability by conditioning scenario generation on user-specified inputs—such as text conditions Zhang et al. (2024); Tan et al. (2023), goal conditions Tan et al. (2024); Rowe et al. (2024), or cost functions Zhong et al. (2023b); Jiang et al. (2023b); Zhong et al. (2023a)—producing scenarios that are both realistic and aligned with user requirements. However, their open-loop training paradigm introduces the *covariate shift* problem during closed-loop deployment, arising from the distribution mismatch between training and deployment states. Moreover, data-driven simulators often adopt simplified environment dynamics Gulino et al. (2023); Caesar et al. (2021), resulting in unrealistic interactions and state transitions that further degrade closed-loop reliability. In contrast, physics-based simulators provide fine-grained control over scenario configuration through physical engines, enabling high-fidelity closed-loop interactions. Nonetheless, the absence of expert demonstrations makes it challenging to reproduce realistic behavior. To mitigate this, several approaches employ reinforcement learning (RL) to directly acquire controllable behaviors through interaction with the simulator Ding et al. (2021); Hanselmann et al. (2022); Chen et al. (2024b); Zhang et al. (2023b), although often at the cost of realism. Other approaches enhance realism by injecting real-world traffic data into physics-based simulators Osiński et al. (2020); Li et al. (2023), but typically rely on log-replay or rule-based simulation, limiting

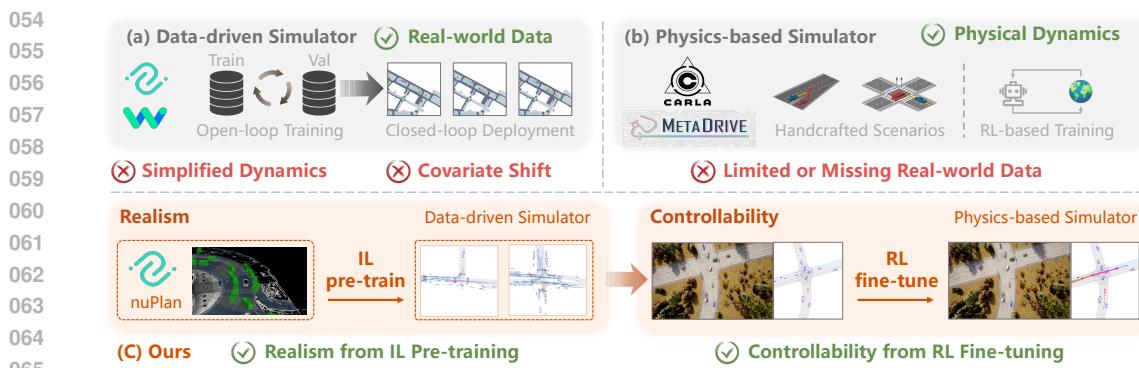


Figure 1: **Traffic Simulation across Different Platforms.** (a) Data-driven Simulator: employs imitation learning to replicate real-world driving behaviors, but suffers from covariate shift and simplified dynamics; (b) Physics-based Simulator: enables controllable scenario construction via high-fidelity closed-loop interaction, but lacks large-scale real-world data; (c) Our framework: combines IL pre-training in a data-driven simulator to ensure realism with RL fine-tuning in a physics-based simulator to enhance controllability.

controllability and interactivity. Despite recent advances, a fundamental trade-off persists between realism and controllability across both paradigms, making it challenging to achieve both simultaneously in interactive closed-loop scenarios.

Drawing inspiration from the widely adopted “pre-training and fine-tuning” paradigm in large language models (LLMs) Rafailov et al. (2023); Yu et al. (2025); Shao et al. (2024), we combine the strengths of two platforms. Specifically, we perform IL pre-training in a data-driven simulator to capture realism, followed by RL fine-tuning in a physics-based simulator to address covariate shift and enhance controllability.

Building on this insight, we propose a dual-stage AV-centric simulation framework (Figure 1) that unifies the strengths of data-driven and physics-based simulators through a “pre-training and fine-tuning” paradigm, balancing realism and controllability in traffic simulation. In Stage 1, we pre-train a planning model via IL to generate realistic and multimodal trajectories conditioned on given route-level reference lines. This stage achieves both trajectory-level realism, capturing realistic and multimodal behavior patterns, and route-level controllability, guaranteeing compliance with prescribed reference lines. In Stage 2, we identify critical background vehicles (CBVs) through route-level interaction analysis, focusing on those most likely to interact with the AV. For these CBVs, we leverage the IL pre-trained model from Stage 1, conditioned on their route-level reference lines, to automatically generate realistic and multimodal trajectories that remain route-level controllable. On top of these generated candidates, we introduce *RIFT*, a novel group-relative RL fine-tuning strategy that improves controllability over driving styles and mitigates covariate shift. Unlike prior methods Zhang et al. (2023a); Peng et al. (2024) that fine-tune only the best trajectory or action, *RIFT* evaluates all candidate modalities via group-relative formulation Shao et al. (2024) and employs a surrogate objective for stable optimization, enhancing style-level controllability and alleviating covariate shift while preserving the trajectory-level realism and route-level controllability established in Stage 1.

Our contributions can be summarized as:

- We propose a dual-stage AV-centric simulation framework that combines IL pre-training in a data-driven simulator and RL fine-tuning in a physics-based simulator, leveraging their complementary strengths to balance realism and controllability.
- We propose *RIFT*, a novel group-relative RL fine-tuning strategy that evaluates all candidate modalities through group-relative formulation and employs a surrogate objective for stable optimization, improving style-level controllability and alleviating covariate shift, while retaining the trajectory-level realism and route-level controllability inherited from IL pre-training.
- Extensive experiments demonstrate that *RIFT* enhances the realism and controllability of traffic simulation, effectively exposing the limitations of modern AV systems under closed-loop settings.

108 **2 RELATED WORK**

110 **Realistic Traffic Simulation.** A variety of generative architectures have been explored for realistic
 111 traffic simulation [Tan et al. \(2021\)](#); [Zhang et al. \(2023c\)](#); [Yang et al. \(2025\)](#), including conditional
 112 variational autoencoders [Suo et al. \(2021\)](#); [Rempe et al. \(2022\)](#); [Xu et al. \(2023\)](#) and diffusion-based
 113 models [Jiang et al. \(2024\)](#); [Chitta et al. \(2024\)](#); [Zhou et al. \(2025\)](#); [Tan et al. \(2025\)](#). However,
 114 maintaining long-term stability remains challenging due to the *covariate shift* between open-loop
 115 training and closed-loop deployment. Recent methods such as SMART [Wu et al. \(2024\)](#), GUMP [Hu et al. \(2024\)](#),
 116 Trajeglish [Phlion et al. \(2023\)](#), and MotionLM [Seff et al. \(2023\)](#) address this issue by
 117 formulating traffic simulation as a next-token prediction (NTP) task, leveraging discrete action spaces
 118 to improve closed-loop robustness. Despite these advances, most approaches remain confined to data-
 119 driven simulation platforms [Gulino et al. \(2023\)](#); [Caesar et al. \(2021\)](#); [Dauner et al. \(2024\)](#); [Montali et al. \(2023\)](#), which typically adopt simplified environment dynamics. Such oversimplifications limit
 120 the reliability of long-term closed-loop interactions, especially in complex and interactive scenarios.
 121

122 **Controllable Traffic Simulation.** Recent studies have introduced diverse conditioning mecha-
 123 nisms to generate traffic scenarios aligned with user preferences. CTG [Zhong et al. \(2023b\)](#) and
 124 MotionDiffuser [Jiang et al. \(2023b\)](#) employ diffusion models conditioned on cost-based signals.
 125 Language-conditioned methods, including CTG++ [Zhong et al. \(2023a\)](#), LCTGen [Tan et al. \(2023\)](#),
 126 and ProSim [Tan et al. \(2024\)](#), enable user specification through language prompts. Other strategies
 127 adopt guided sampling (SceneControl [Lu et al. \(2024\)](#)), retrieval-based generation (RealGen [Ding et al. \(2024\)](#)), or reward-driven causality modeling (CCDiff [Lin et al. \(2024\)](#)). Despite improving
 128 controllability, existing approaches remain confined to open-loop settings or simplified dynamics,
 129 and primarily target low-level control. High-level attributes such as driving style are underexplored,
 130 leaving the integration of realism and controllability in closed-loop simulation an open challenge.
 131

132 **Closed-Loop Fine-Tuning.** Covariate shift—the mismatch between open-loop training and closed-
 133 loop deployment—remains a key challenge for reliable long-term traffic simulation. To address
 134 this, recent work explores fine-tuning strategies in the closed-loop setting. Hybrid IL and RL meth-
 135 ods [Zhang et al. \(2023a\)](#); [Peng et al. \(2024\)](#); [Lu et al. \(2023\)](#) enhance robustness but typically fine-tune
 136 the entire model via RL, which often compromises realism due to the difficulty of designing human-
 137 aligned reward functions. Supervised fine-tuning approaches such as CAT-K [Zhang et al. \(2025\)](#)
 138 show strong performance but rely on expert demonstrations, limiting scalability. TrafficRLHF [Cao et al. \(2024\)](#) improves alignment through reinforcement learning with human feedback (RLHF), but
 139 demands costly human input and suffers from reward model instability. Moreover, most existing
 140 methods focus on optimizing the best action or trajectory, ignoring the inherent multimodality of
 141 traffic simulation, thus limiting behavioral diversity during fine-tuning.
 142

143 **3 BACKGROUND**

144 **3.1 TASK REDEFINITION**

147 Following the widely adopted paradigm for closed-loop training and evaluation in autonomous
 148 driving [Jia et al. \(2024\)](#); [Xu et al. \(2022\)](#), our simulation framework includes a single autonomous
 149 vehicle (AV) navigating a predefined global route, accompanied by multiple rule-based background
 150 vehicles (BVs), forming an AV-centric closed-loop simulation environment. These BVs either provide
 151 diverse interactive data for training or serve to evaluate the AV’s robustness. Building upon this setup,
 152 we identify a subset of critical background vehicles (CBVs) that are more likely to interact with the
 153 AV. For these CBVs, the rule-based control is replaced with a well-trained planning model, enabling
 154 the synthesis of realistic and controllable behaviors in interactive closed-loop scenarios.
 155

156 **3.2 CBV-CENTRIC REALISTIC TRAJECTORY GENERATION**

157 With recent advances in imitation learning, data-driven approaches have demonstrated strong per-
 158 formance in generating realistic, multimodal trajectories [Zheng et al. \(2025\)](#); [Huang et al. \(2023\)](#);
 159 [Hu et al. \(2023\)](#); [Jiang et al. \(2023a\)](#); [Sun et al. \(2024\)](#). In fully observable simulation environments,
 160 Pluto [Cheng et al. \(2024a\)](#) produces reliable, realistic, and multimodal trajectories by leveraging
 161 ground-truth states, while enabling route-level controllability through reference line encoding. These
 capabilities make Pluto a suitable choice for our planning model.
 162

CBV-Centric Scene Encoding. Following Cheng et al. (2024a), for each CBV in the scene, we extract its current feature F_{cbv} , the historical features of neighboring vehicles F_{neighbor} , and vectorized map features F_{map} . These features are encoded into $E_{\text{cbv}} \in \mathbb{R}^{1 \times D}$, $E_{\text{neighbor}} \in \mathbb{R}^{N_{\text{neighbor}} \times D}$, and $E_{\text{map}} \in \mathbb{R}^{N_{\text{map}} \times D}$, respectively, where N_{neighbor} and N_{map} denote the number of neighboring vehicles and map elements, and D is the embedding dimension. To model the interactions among these embeddings, we concatenate them and apply a global positional embedding (PE) to obtain the unified scene embedding $E_s \in \mathbb{R}^{(1+N_{\text{neighbor}}+N_{\text{map}}) \times D}$ as:

$$E_s = \text{concat}(E_{\text{cbv}}, E_{\text{neighbor}}, E_{\text{map}}) + \text{PE}. \quad (1)$$

This scene embedding E_s is then passed through N Transformer encoder blocks for feature aggregation, yielding the final CBV-centric scene embedding E_{enc} . Each encoder block follows the standard Transformer formulation. Specifically, the i -th block is defined as:

$$\begin{aligned} E_s^i &= E_s^{i-1} + \text{MHA}(\text{LayerNorm}(E_s^{i-1})), \\ E_s^i &= E_s^i + \text{FFN}(\text{LayerNorm}(E_s^i)), \end{aligned} \quad (2)$$

where MHA is the standard multi-head attention function, FFN is the feedforward network layer.

Multimodal Trajectory Decoding. To capture the multimodal nature of real-world driving behaviors, we adopt the longitudinal-lateral decoupling mechanism proposed in Cheng et al. (2024a). This approach leverages reference line information to construct high-level lateral queries $Q_{\text{lat}} \in \mathbb{R}^{N_{\text{ref}} \times D}$, and introduces learnable longitudinal queries $Q_{\text{lon}} \in \mathbb{R}^{N_{\text{lon}} \times D}$. These are concatenated and projected to form the multimodal navigation query $Q_{\text{nav}} \in \mathbb{R}^{N_{\text{ref}} \times N_{\text{lon}} \times D}$ as:

$$Q_{\text{nav}} = \text{Projection}(\text{concat}(Q_{\text{lat}}, Q_{\text{lon}})), \quad (3)$$

where N_{ref} and N_{lon} denote the number of reference lines and longitudinal anchors, respectively. The navigation query Q_{nav} and the scene embedding E_{enc} are then fed into N decoder blocks to model lateral, longitudinal, and cross-modal interactions. Each decoder block is structured as:

$$\begin{aligned} \hat{Q}_{\text{nav}}^{i-1} &= \text{SelfAttn}(\text{SelfAttn}(Q_{\text{nav}}^{i-1}, \text{dim} = 0), \text{dim} = 1), \\ Q_{\text{nav}}^i &= \text{CrossAttn}(\hat{Q}_{\text{nav}}^{i-1}, E_{\text{enc}}, E_{\text{enc}}). \end{aligned} \quad (4)$$

SelfAttn, CrossAttn denote multi-head self-attention and cross-attention, respectively. Given the decoder’s final output Q_{dec} , two MLP heads are applied to produce the CBV-centric multimodal trajectories $\mathcal{T} \in \mathbb{R}^{N_{\text{ref}} \times N_{\text{lon}} \times T \times 6}$ and their confidence scores $\mathcal{S} \in \mathbb{R}^{N_{\text{ref}} \times N_{\text{lon}}}$:

$$\mathcal{T} = \text{MLP}(Q_{\text{dec}}), \quad \mathcal{S} = \text{MLP}(Q_{\text{dec}}), \quad (5)$$

where T is the prediction horizon, and each trajectory point τ_t^i encodes $[p_x, p_y, \cos \theta, \sin \theta, v_x, v_y]$.

4 METHODOLOGY

Leveraging the IL pre-trained planning model described in Section 3.2, realistic and multimodal trajectories can be generated across diverse scenarios conditioned on reference lines. However, the open-loop training paradigm leaves the policy vulnerable to covariate shift, even with contrastive learning Halawa et al. (2022); Wang et al. (2023) or data augmentation Cheng et al. (2024a). To address this, we propose *RIFT*, a group-relative RL fine-tuning strategy that enhances style-level controllability and mitigates covariate shift while preserving the trajectory-level realism and route-level controllability from pre-training. The following sections detail *RIFT*’s implementation within the physics-based simulator.

4.1 ROUTE-LEVEL INTERACTION ANALYSIS

Following Feng et al. (2023b), we address the “curse of rarity” Liu & Feng (2024) by selectively intervening in a set of critical background vehicles (CBVs) at key moments, while keeping non-critical agents under rule-based control for efficiency. CBVs are identified via route-level interaction analysis between the AV’s predefined global route and the candidate routes of surrounding vehicles, selecting the vehicle with the highest interaction probability (details in Appendix B.2).

The corresponding route-level reference line is then used as a condition for the IL pre-trained planning model (Section 3.2) to synthesize realistic and multimodal trajectories. For each identified CBV, the model generates $N_{\text{ref}} \times N_{\text{lon}}$ candidate trajectories, from which the highest-scoring one is selected for closed-loop execution. This process promotes realistic route-level interactions with the AV and enables the construction of meaningful interactive scenarios.

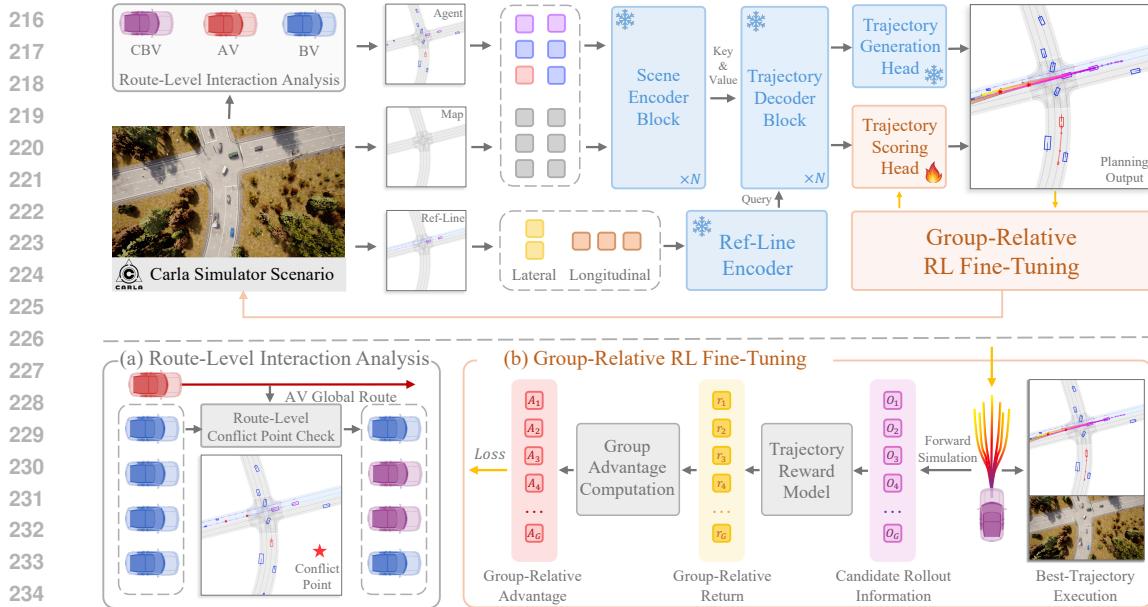


Figure 2: **Overview of the RIFT:** Building on the IL pre-trained model, *RIFT* performs route-level interaction analysis to identify critical background vehicles and the associated reference lines, enabling the generation of realistic and multimodal trajectories. To isolate style-level controllability from the trajectory-level realism and route-level controllability established during pre-training, only the scoring head is fine-tuned via *RIFT* while freezing other components. Specifically, *RIFT* computes group-relative advantages over all candidate rollouts, promoting alignment with user-preferred styles and mitigating covariate shift through RL fine-tuning.

4.2 GROUP-RELATIVE RL FINE-TUNING

Open-loop IL pre-training offers trajectory-level realism and route-level controllability; however, it inevitably suffers from covariate shift in closed-loop deployment, causing error accumulation and unrealistic long-term behaviors. Existing RL Schulman et al. (2017) and hybrid IL–RL methods Peng et al. (2024) partially mitigate covariate shift, but their optimization is restricted to the executed rollout, disregarding alternative candidates and degrading multimodality. More critically, covariate shift induces asymmetric degradation across model components: under the generation–selection paradigm, the generation head, conditioned on route-level priors, remains robust and consistently produces realistic multimodal candidates, whereas the scoring head, trained solely through imitation, is more vulnerable to distribution mismatch. These challenges motivate three key requirements for fine-tuning: (i) preserving multimodality, (ii) addressing asymmetric covariate shift, and (iii) ensuring stable policy improvement. We address these requirements through a unified framework that combines group-relative optimization, asymmetry-aware fine-tuning, and dual-clip stabilization.

To preserve multimodality, we adopt group-relative formulation Shao et al. (2024), which evaluates all candidate modalities within the group and assigns higher relative advantages to those better aligned with user-preferred styles. Considering closed-loop dynamics, we evaluate simulated rollouts rather than raw trajectories to mitigate plan–rollout deviation. Specifically, given $G = N_{\text{ref}} \times N_{\text{lon}}$ candidate trajectories $\mathcal{T} = \{\tau_i\}_{i=1}^G$ for a CBV at state s , we conduct forward simulation Dauner et al. (2023) (see Appendix B.6) to obtain rollouts $\tilde{\mathcal{T}} = \{\tilde{\tau}_i\}_{i=1}^G$. Each rollout is evaluated by a user-defined state-wise reward model StateWiseRM, yielding the corresponding discounted returns $\mathcal{R} = \{R_i\}_{i=1}^G$ from which we derive the group-relative advantages $\mathcal{A} = \{\hat{A}_i\}_{i=1}^G$ as follows:

$$R_i(s) = \sum_{t=0}^T \gamma^t [\text{StateWiseRM}(\tilde{\tau}_i^t, s)], \quad \hat{A}_i(s) = \frac{R_i(s) - \text{mean}(\mathcal{R})}{\sqrt{\text{Var}(\mathcal{R}) + \epsilon}}. \quad (6)$$

Here, \hat{A}_i quantifies the performance of each rollout relative to the group, promoting high-return rollouts without suppressing alternative modes.

In standard GRPO Shao et al. (2024), sampling from the old policy implicitly induces old-policy weighting. Extending this to our enumerated setting involves averaging terms weighted by $\pi_{\theta_{\text{old}}}$ in conjunction with the importance ratio $\rho_i(\theta) = \pi_\theta(\tau_i | s) / \pi_{\theta_{\text{old}}}(\tau_i | s)$, which yields a low-variance

270 estimate of the old-policy expectation over the enumerated support. The aggregated objective is:
 271

$$272 \quad \mathcal{J}(\theta) = \mathbb{E}_{s \sim \mathcal{D}} \left[\sum_{i=1}^G \pi_{\theta_{\text{old}}}(\tau_i | s) \min \left[\rho_i(\theta) \hat{A}_i, \text{clip}(\rho_i(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_i \right] \right] - \beta \mathbb{D}_{KL} [\pi_\theta || \pi_{ref}], \quad (7)$$

274 where π_{ref} denotes the IL pre-trained model. While exact over the enumerated support, this scheme
 275 overemphasizes frequent modes and under-represents rare but high-return ones, causing mode collapse
 276 and reduced diversity. To balance modality contributions, we adopt an equal-weight objective:

$$277 \quad \mathcal{J}(\theta) = \mathbb{E}_{s \sim \mathcal{D}} \left[\frac{1}{G} \sum_{i=1}^G \min \left[\rho_i(\theta) \hat{A}_i, \text{clip}(\rho_i(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_i \right] \right] - \beta \mathbb{D}_{KL} [\pi_\theta || \pi_{ref}]. \quad (8)$$

280 Under equal weighting, $\rho_i(\theta)$ regulates candidate updates rather than serving as a pure importance
 281 weight, removing old-policy bias and yielding balanced updates that preserve multimodality.

282 To address asymmetric covariate shift, we freeze the generation head to retain trajectory-level realism
 283 and fine-tune only the scoring head to enhance style-level controllability. In this setting, constraining
 284 the scoring head with the KL term to the IL pre-trained model would anchor learning to a biased
 285 reference, thereby hindering adaptation. We therefore remove the KL term, allowing the scoring head
 286 to adapt freely while leveraging the stable candidates provided by the frozen generation head.

287 Removing the KL term improves flexibility but raises stability concerns. Although the clipped-ratio
 288 mechanism in PPO constrains update magnitude, it proves insufficient in the group-relative setting.
 289 Specifically, when a rare trajectory under the old policy receives a higher probability from the current
 290 policy despite a negative advantage, the product $\rho_i(\theta) \hat{A}_i$ can become disproportionately large and
 291 destabilize learning. To address this, we incorporate the dual-clip surrogate from Dual-Clip PPO [Ye et al. \(2020\)](#); [Gao et al. \(2021\)](#), which lower-bounds clipped negative advantages. This establishes
 292 a trust-region-like constraint that guarantees bounded per-candidate updates (see [Theorem A.3](#)),
 293 thereby preventing extreme negative shifts while preserving responsiveness to user-preferred styles.
 294 The resulting surrogate objective, termed *RIFT*, is

$$295 \quad \mathcal{J}_{\text{RIFT}}(\theta) = \mathbb{E}_{s \sim \mathcal{D}} \left[\frac{1}{G} \sum_{i=1}^G \psi(\rho_i(\theta), \hat{A}_i) \right], \quad (9)$$

$$296 \quad \psi(\rho, \hat{A}) = \begin{cases} \min(\rho \hat{A}, \text{clip}(\rho, 1 - \epsilon, 1 + \epsilon) \hat{A}), & \hat{A} \geq 0, \\ \max(\min(\rho \hat{A}, \text{clip}(\rho, 1 - \epsilon, 1 + \epsilon) \hat{A}), c \hat{A}), & \hat{A} < 0 \end{cases} \quad (\epsilon > 0, c > 1).$$

300 This objective integrates multimodality preservation, asymmetry-aware fine-tuning, and stable optimization
 301 into a unified framework, enhancing style-level controllability and mitigating covariate shift
 302 while retaining trajectory-level realism and route-level controllability (analysis in [Appendix A](#)).
 303

304 5 EXPERIMENT

307 This section systematically addresses the following research questions: **Q1**: How does *RIFT* compare
 308 with representative baselines in terms of the realism and controllability of the generated traffic
 309 scenarios? **Q2**: How can the generated traffic scenario be effectively utilized to support downstream
 310 autonomous driving tasks? **Q3**: How do the components of *RIFT* contribute to overall performance,
 311 and to what extent is style-level controllability preserved under varying user-specified driving styles?

312 5.1 EXPERIMENT SETUPS

314 Under the dual-stage AV-centric simulation framework, we adopt Pluto [Cheng et al. \(2024a\)](#) as our
 315 planning model for its well-established performance and open-source implementation. To ensure
 316 fair comparison, we use the official IL pre-trained checkpoint provided by Pluto, trained on the
 317 nuPlan dataset [Caesar et al. \(2021\)](#). Simulations are conducted in CARLA [Dosovitskiy et al. \(2017\)](#),
 318 leveraging Bench2Drive [Jia et al. \(2024\)](#) to support AV-centric closed-loop simulation and evaluation.
 319 Implementation details, training protocols, and evaluation settings are described in [Appendix B](#).

320 **Baseline.** To systematically evaluate the effectiveness of *RIFT* in traffic simulation, we compare it
 321 against the following baselines, with implementation details provided in [Appendix B.5](#).
 322

- 323 • **Pure RL/IL:** Methods trained solely with RL or IL, without fine-tuning, including *Pluto* [Cheng et al. \(2024a\)](#), as well as *FREA*, *FPPO-RS*, and *PPO*, all from [Chen et al. \(2024b\)](#).

Table 1: **Comparison in Controllability and Realism.** Metrics are evaluated under the PDM-Lite [Beißwenger \(2024\)](#) AV setting across three random seeds, with the **best** and the **second-best** results highlighted accordingly.

Method	Type	Kinematic Metrics			Interaction Metrics			Map Metrics	
		S-SW ↑	S-WD ↓	A-SW ↑	CPK ↓	RP ↑	2D-TTC ↑	ACT ↑	ORR ↓
Pluto	IL	0.88 ± 0.01	5.81 ± 0.06	0.90 ± 0.01	5.06 ± 2.69	564.14 ± 114.41	2.50 ± 1.48	2.44 ± 1.39	0.24 ± 0.15
PPO	RL	0.95 ± 0.01	4.45 ± 0.15	0.89 ± 0.02	13.95 ± 2.34	409.51 ± 30.38	2.59 ± 1.60	2.52 ± 1.57	9.17 ± 2.39
FREA	RL	0.93 ± 0.01	5.10 ± 0.14	0.93 ± 0.01	30.42 ± 5.28	292.81 ± 68.54	2.71 ± 1.40	2.67 ± 1.41	9.01 ± 2.09
FPPO-RS	RL	0.87 ± 0.01	5.80 ± 0.11	0.80 ± 0.03	21.39 ± 3.23	356.79 ± 26.19	2.55 ± 1.69	2.53 ± 1.68	8.60 ± 0.25
SFT-Pluto	SFT	0.88 ± 0.02	6.01 ± 0.19	0.87 ± 0.02	6.33 ± 2.23	780.48 ± 41.05	2.20 ± 1.64	2.12 ± 1.51	0.06 ± 0.07
RS-Pluto	SFT+RLFT	0.93 ± 0.00	5.40 ± 0.15	0.92 ± 0.01	4.11 ± 3.90	819.40 ± 74.07	2.27 ± 1.45	2.23 ± 1.43	1.05 ± 0.31
RTR-Pluto	SFT+RLFT	0.85 ± 0.00	6.24 ± 0.16	0.81 ± 0.03	6.98 ± 2.59	481.60 ± 70.19	2.55 ± 1.60	2.47 ± 1.51	0.08 ± 0.09
PPO-Pluto	RLFT	0.95 ± 0.01	4.96 ± 0.31	0.90 ± 0.02	6.89 ± 3.19	683.57 ± 38.12	2.66 ± 1.50	2.60 ± 1.43	0.07 ± 0.13
REINFORCE-Pluto	RLFT	0.92 ± 0.01	5.63 ± 0.19	0.90 ± 0.02	6.98 ± 0.86	813.70 ± 24.76	2.39 ± 1.64	2.30 ± 1.55	1.37 ± 1.13
GRPO-Pluto	RLFT	0.94 ± 0.04	4.96 ± 0.89	0.96 ± 0.00	7.24 ± 4.04	892.65 ± 65.27	2.65 ± 1.44	2.61 ± 1.48	0.10 ± 0.08
RIFT-Pluto (ours)	RLFT	0.97 ± 0.01	4.46 ± 0.43	0.93 ± 0.01	6.83 ± 2.62	995.33 ± 84.62	2.74 ± 1.30	2.71 ± 1.32	0.36 ± 0.20

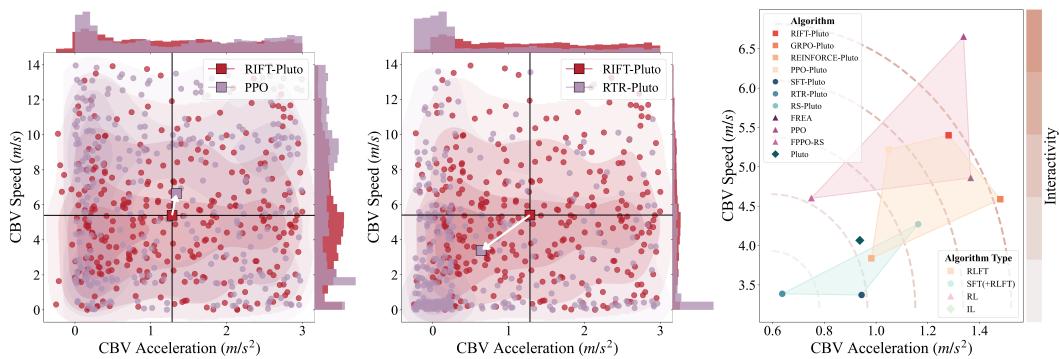


Figure 3: **Speed and Acceleration Distribution.** RL-based methods tend to be interactive but unnatural, whereas supervised methods are overly conservative. *RIFT* strikes a balance, yielding higher interactivity with realistic distributional profiles, reducing hesitation while maintaining safe interactions.

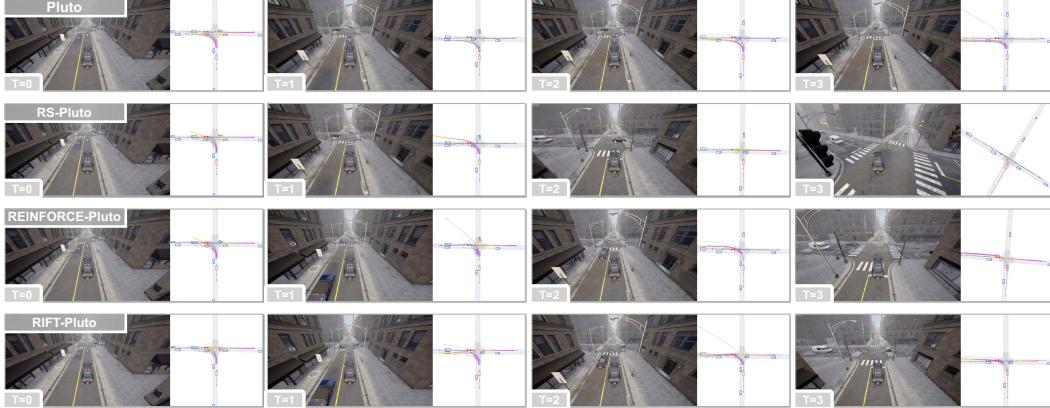


Figure 4: Temporal comparisons illustrating *RIFT*'s superior performance over other baselines under AV-centric closed-loop simulation. CBV is marked in purple, AV (PDM-Lite) is in red, and BVs are in blue.

- **RLFT/SFT:** Methods that fine-tune the pre-trained Pluto model using either RL or supervised objectives, including *PPO-Pluto* Schulman et al. (2017), *REINFORCE-Pluto* Sutton et al. (1999), *GRPO-Pluto* Shao et al. (2024), and *SFT-Pluto*.
- **Hybrid:** Methods that combine RL and supervised fine-tuning, including *RTR-Pluto* Zhang et al. (2023a) and *RS-Pluto* Peng et al. (2024).

All methods are fine-tuned on the scoring head to ensure fair comparisons, while isolating style-level controllability from trajectory-level realism and route-level controllability, as confirmed by the ablation studies in Section 5.4. Following the realism standards of the Sim Agent Challenge in WOSAC Montali et al. (2023), we adopt a normal style reward for all RL-based baselines, with details in Appendix B.7. Results under an aggressive style reward are reported in Section 5.4.

Metrics. Building on the WOSAC evaluation framework, we categorize our evaluation metrics into three groups: *kinematic metrics*, *interaction metrics*, and *map metrics*. Kinematic metrics capture distributional motion properties (S-SW, S-WD, A-SW), as in [Chen et al. \(2024a\)](#), with the absence of ground-truth trajectories in CARLA precluding displacement-based measures (e.g., ADE, FDE). Interaction metrics evaluate agent interactions through collision frequency (Collision Per Kilometer, CPK), driving efficiency (Route Progress, RP), and safety-critical measures, including 2D-TTC [Guo et al. \(2023\)](#) and ACT [Venthuruthiyil & Chunchu \(2022\)](#). Map metrics evaluate adherence to road geometry through the Off-Road Rate (ORR). Collectively, these metrics comprehensively evaluate realism and controllability in closed-loop simulation; detailed definitions are in [Appendix B.8](#).

5.2 REALISTIC AND CONTROLLABLE TRAFFIC SCENARIO GENERATION (Q1)

Main Results. To address **Q1**, we evaluate the controllability and realism of the generated scenario across CBV methods, with results summarized in [Table 1](#). *RIFT* consistently outperforms all baselines in both aspects across most settings. While supervised learning methods achieve slightly lower CPK and ORR, this improvement is primarily due to their inherently conservative behavior, derived from the expert PDM-Lite [Beißwenger \(2024\)](#), which prioritizes safety by avoiding risky maneuvers.

This conservative tendency is further highlighted in [Figure 3](#), where supervised policies exhibit significantly lower speed and acceleration profiles. In contrast, *RIFT* strikes a more favorable balance between safety and interactivity. It achieves superior safety performance, as reflected by higher 2D-TTC and ACT scores, while avoiding the overly cautious behaviors typical of supervised approaches. As shown in [Figure 3](#), *RIFT* demonstrates higher average speed and acceleration, indicating more interactive behavior, while maintaining realistic motion profiles.

Qualitative Results. To further demonstrate the effectiveness of *RIFT*, we compare closed-loop simulations against representative baselines, as shown in [Figure 4](#). Baseline methods often suffer from unstable or low-quality trajectory selection in closed-loop settings, whereas *RIFT* consistently selects smooth, high-quality trajectories with superior temporal consistency. Further qualitative examples are presented in [Appendix D.3](#).

5.3 GENERATED TRAFFIC SCENARIOS FOR CLOSED-LOOP AV EVALUATION (Q2)

To address **Q2**, we assess the suitability of traffic scenarios generated by different CBV methods for closed-loop AV evaluation. Following KING [Hanselmann et al. \(2022\)](#), we adopt PDM-Lite [Beißwenger \(2024\)](#)—a rule-based planner with privileged access—as a reference to evaluate two key scenario properties: feasibility, measured by Driving Score (DS), and naturalness, captured by our proposed Blocked Rate (BR). A high DS indicates that the AV can reliably complete the scenario, while a low BR reflects realistic interactions without excessive obstruction from surrounding vehicles. Together, DS and BR offer a principled basis for evaluating scenario quality.

To further assess the ability of each scenario to reveal weaknesses in learning-based planners, we compare PlanT [Renz et al. \(2022\)](#), UniAD [Hu et al. \(2023\)](#), and VAD [Jiang et al. \(2023a\)](#) with PDM-Lite. As these models are sensitive to subtle or adversarial interactions, informative scenarios should induce noticeable performance drops. As shown in [Table 2](#), traffic generated by *RIFT* achieves the highest DS and lowest BR under PDM-Lite, while also causing the largest degradation across all learning-based planners. These results confirm that *RIFT* generates interactive and feasible scenarios that effectively expose limitations of modern AV systems. See [Appendix C](#) for detailed results.

5.4 ABLATION STUDY (Q3)

Building on the design choices introduced in [Section 4.2](#), we systematically ablate five components of *RIFT*: weighting scheme (Old-Weight vs. Equal-Weight), fine-tuning module (Scoring Head vs. All Head), KL regularization (w/ KL vs. w/o KL), policy clipping (Dual-Clip vs. PPO-Clip), and style preference (Normal vs. Aggressive). All experiments share identical settings, and results are reported in [Table 3](#).

Equal-Weight vs. Old-Weight. Replacing old-policy weighting with equal weighting eliminates the likelihood bias toward frequent modes and enables balanced updates across all candidates. This leads to improved exploitation of high-return rollouts and better multimodality preservation.

432 Table 2: **Comparison of AV Evaluation across CBV Methods.** Each metric is evaluated across three random
 433 seeds, with the **best** and the **second-best** results highlighted accordingly.

Method	PDM-Lite		PlanT		UniAD		VAD	
	DS ↑	BR ↑	DS	ΔDS ↓	DS	ΔDS ↓	DS	ΔDS ↓
Pluto	77.84 ± 2.20	23.33 ± 5.77	42.52 ± 4.72	-35.32	73.73 ± 1.24	-4.11	66.87 ± 2.11	-10.97
PPO	76.26 ± 0.12	30.00 ± 0.00	36.39 ± 1.11	-39.87	69.79 ± 1.41	-6.47	67.64 ± 1.27	-8.62
FREA	83.53 ± 0.13	20.00 ± 0.00	39.61 ± 1.34	-43.92	69.29 ± 5.22	-14.24	67.57 ± 5.37	-15.96
FPPO-RS	83.52 ± 0.09	20.00 ± 0.00	38.85 ± 4.91	-44.67	75.13 ± 5.18	-8.39	69.15 ± 2.79	-14.37
SFT-Pluto	86.09 ± 2.04	13.33 ± 5.77	39.41 ± 4.97	-47.28	77.49 ± 5.93	-9.20	68.89 ± 0.87	-17.80
RS-Pluto	89.32 ± 1.41	13.33 ± 5.77	42.05 ± 4.08	-47.27	80.62 ± 0.78	-8.70	69.48 ± 5.02	-19.84
RTR-Pluto	87.64 ± 1.56	10.00 ± 0.00	40.08 ± 2.38	-47.56	77.69 ± 2.82	-9.95	66.27 ± 4.53	-21.37
PPO-Pluto	85.63 ± 2.02	16.67 ± 5.77	41.86 ± 2.78	-43.77	77.14 ± 3.36	-8.49	68.62 ± 3.16	-17.01
REINFORCE-Pluto	92.17 ± 3.45	10.00 ± 10.00	45.25 ± 1.75	-46.92	79.89 ± 1.97	-12.28	70.28 ± 3.58	-21.89
GRPO-Pluto	89.86 ± 2.10	6.67 ± 5.77	47.24 ± 5.67	-42.62	81.02 ± 0.64	-8.84	72.55 ± 0.74	-17.31
RIFT-Pluto (ours)	94.78 ± 1.37	0.00 ± 0.00	44.28 ± 3.15	-50.50	73.79 ± 6.53	-20.99	68.24 ± 3.23	-26.54

447 Table 3: **Ablation Study on RIFT.** Evaluation under PDM-Lite AV setting with three random seeds.

Method	Kinematic Metrics			Interaction Metrics			Map Metrics	
	S-SW ↑	S-WD ↓	A-SW ↑	CPK ↓	RP ↑	2D-TTC ↑	ACT ↑	ORR ↓
w/ Old-Weight	0.82 (-0.15)	6.24 (+1.78)	0.85 (-0.08)	7.51 (+0.68)	574.51 (-420.82)	2.70 (-0.04)	2.68 (-0.03)	0.00 (-0.36)
w/ All-Head	0.96 (-0.01)	4.70 (+0.24)	0.94 (+0.01)	7.84 (+0.01)	827.12 (-168.21)	2.83 (+0.09)	2.76 (+0.05)	0.43 (+0.07)
w/ KL	0.93 (-0.04)	5.33 (+0.87)	0.90 (-0.03)	7.05 (+0.22)	815.06 (-180.27)	2.76 (+0.02)	2.73 (+0.02)	0.38 (+0.02)
w/ PPO-Clip	0.91 (-0.06)	5.92 (+1.46)	0.94 (+0.01)	2.03 (-4.80)	655.39 (-339.94)	2.57 (-0.17)	2.54 (-0.17)	0.04 (-0.32)
w/ Aggressive	0.97 (+0.00)	3.89 (-0.57)	0.94 (+0.01)	8.41 (+1.58)	1053.76 (+58.43)	2.93 (+0.19)	2.88 (+0.17)	0.91 (+0.55)
RIFT-Pluto (ours)	0.97	4.46	0.93	6.83	995.33	2.74	2.71	0.36

456 **Scoring Head vs. All Head.** Freezing the generation head is crucial for retaining trajectory-level
 457 realism and route-level controllability. Fine-tuning all heads (*w/ All Head*) disrupts the pre-trained
 458 generation head and slightly degrades realism metrics, whereas fine-tuning only the scoring head
 459 achieves better controllability without compromising realism.

460 **w/ KL vs. w/o KL.** Anchoring the scoring head to the IL pre-trained reference via KL regularization
 461 (*w/ KL*) constrains adaptation to a biased reference under asymmetric covariate shift. Removing
 462 this term improves controllability while maintaining realism, confirming that free adaptation of the
 463 scoring head yields more effective policy improvement.

464 **Dual-Clip vs. PPO-Clip.** Replacing dual-clip with standard PPO clipping (*w/ PPO-Clip*) results in
 465 overly conservative behaviors and reduced efficiency, as extreme negative updates can dominate and
 466 suppress positive learning signals. Dual-clip bounds such updates while preserving responsiveness to
 467 high-return rollouts, producing more realistic and efficient behavior.

468 **Normal vs. Aggressive.** Adopting a more aggressive reward that emphasizes efficiency increases
 469 route progress but also raises collision and off-road rates, illustrating the efficiency–safety trade-off.
 470 These results demonstrate that *RIFT* supports flexible style shaping while maintaining stability and
 471 multimodality. Additional qualitative insights on controllability are provided in Appendix D.1.

474 6 CONCLUSION

475 In this work, we propose a dual-stage AV-centric simulation framework that conducts IL pre-training
 476 in a data-driven simulator to capture trajectory-level realism and route-level controllability, followed
 477 by RL fine-tuning in a physics-based simulator to address covariate shift and enhance style-level
 478 controllability. During fine-tuning, we introduce *RIFT*, a novel group-relative RL fine-tuning strategy
 479 that evaluates all candidate modalities using the group-relative formulation combined with a surrogate
 480 objective for optimization, thereby enhancing style-level controllability and mitigating covariate
 481 shift, while preserving the trajectory-level realism and route-level controllability established in IL
 482 pre-training. Extensive experiments demonstrate that *RIFT* generates scenarios with superior realism
 483 and controllability, effectively revealing the limitations of modern AV systems and further bridging
 484 the gap between traffic simulation and reliable closed-loop evaluation. Due to space limits, limitations
 485 and future directions are in Appendix E.2, and experimental reproducibility details are in Appendix B.

486 REFERENCES
487

- 488 Jens Beißwenger. PDM-Lite: A rule-based planner for carla leaderboard 2.0. <https://github.com/OpenDriveLab/DriveLM/blob/DriveLM-CARLA/docs/report.pdf>, 2024.
489 Accessed: 2025-04-09. 7, 8, 19, 20, 21, 22, 23
- 490 Holger Caesar, Juraj Kabzan, Kok Seang Tan, Whye Kit Fong, Eric Wolff, Alex Lang, Luke Fletcher,
491 Oscar Beijbom, and Sammy Omari. nuplan: A closed-loop ml-based planning benchmark for
492 autonomous vehicles. *arXiv preprint arXiv:2106.11810*, 2021. 1, 3, 6, 19
- 493 Yulong Cao, Boris Ivanovic, Chaowei Xiao, and Marco Pavone. Reinforcement learning with human
494 feedback for realistic traffic simulation. In *2024 IEEE International Conference on Robotics and
495 Automation (ICRA)*, pp. 14428–14434. IEEE, 2024. 3
- 496 Di Chen, Meixin Zhu, Hao Yang, Xuesong Wang, and Yinhai Wang. Data-driven traffic simulation:
497 A comprehensive review. *IEEE Transactions on Intelligent Vehicles*, 2024a. 8, 21, 22
- 498 Keyu Chen, Yuheng Lei, Hao Cheng, Haoran Wu, Wenchao Sun, and Sifa Zheng. FREA: Feasibility-
499 guided generation of safety-critical scenarios with reasonable adversariality. In *8th Annual
500 Conference on Robot Learning*, 2024b. URL <https://openreview.net/forum?id=3bcujpPikC>. 1, 6, 19
- 501 Jie Cheng, Yingbing Chen, and Qifeng Chen. Pluto: Pushing the limit of imitation learning-based
502 planning for autonomous driving. *arXiv preprint arXiv:2404.14327*, 2024a. 3, 4, 6, 19
- 503 Jie Cheng, Yingbing Chen, Xiaodong Mei, Bowen Yang, Bo Li, and Ming Liu. Rethinking imitation-
504 based planners for autonomous driving. In *2024 IEEE International Conference on Robotics and
505 Automation (ICRA)*, pp. 14123–14130. IEEE, 2024b. 20
- 506 Kashyap Chitta, Daniel Dauner, and Andreas Geiger. Sledge: Synthesizing driving environments
507 with generative models and rule-based traffic. In *European Conference on Computer Vision*, pp.
508 57–74. Springer, 2024. 3
- 509 Marco Cusumano-Towner, David Hafner, Alex Hertzberg, Brody Huval, Aleksei Petrenko, Eugene
510 Vinitsky, Erik Wijmans, Taylor Killian, Stuart Bowers, Ozan Sener, et al. Robust autonomy
511 emerges from self-play. *arXiv preprint arXiv:2502.03349*, 2025. 20
- 512 Daniel Dauner, Marcel Hallgarten, Andreas Geiger, and Kashyap Chitta. Parting with misconceptions
513 about learning-based vehicle motion planning. In *Conference on Robot Learning*, pp. 1268–1281.
514 PMLR, 2023. 5
- 515 Daniel Dauner, Marcel Hallgarten, Tianyu Li, Xinshuo Weng, Zhiyu Huang, Zetong Yang, Hongyang
516 Li, Igor Gilitschenski, Boris Ivanovic, Marco Pavone, et al. Navsim: Data-driven non-reactive
517 autonomous vehicle simulation and benchmarking. *Advances in Neural Information Processing
518 Systems*, 37:28706–28719, 2024. 3
- 519 Wenhao Ding, Baiming Chen, Bo Li, Kim Ji Eun, and Ding Zhao. Multimodal safety-critical
520 scenarios generation for decision-making algorithms evaluation. *IEEE Robotics and Automation
521 Letters*, 6(2):1551–1558, 2021. 1
- 522 Wenhao Ding, Chejian Xu, Mansur Arief, Haohong Lin, Bo Li, and Ding Zhao. A survey on
523 safety-critical driving scenario generation—a methodological perspective. *IEEE Transactions on
524 Intelligent Transportation Systems*, 2023. 1
- 525 Wenhao Ding, Yulong Cao, Ding Zhao, Chaowei Xiao, and Marco Pavone. Realgen: Retrieval
526 augmented generation for controllable traffic scenarios. In *European Conference on Computer
527 Vision*, pp. 93–110. Springer, 2024. 3
- 528 Alexey Dosovitskiy, German Ros, Felipe Codella, Antonio Lopez, and Vladlen Koltun. Carla: An
529 open urban driving simulator. In *Conference on robot learning*, pp. 1–16. PMLR, 2017. 6, 18, 19
- 530 Lan Feng, Quanyi Li, Zhenghao Peng, Shuhan Tan, and Bolei Zhou. Trafficgen: Learning to generate
531 diverse and realistic traffic scenarios. In *2023 IEEE International Conference on Robotics and
532 Automation (ICRA)*, pp. 3567–3575. IEEE, 2023a. 1

- 540 Shuo Feng, Haowei Sun, Xintao Yan, Haojie Zhu, Zhengxia Zou, Shengyin Shen, and Henry X Liu.
 541 Dense reinforcement learning for safety validation of autonomous vehicles. *Nature*, 615(7953):
 542 620–627, 2023b. 1, 4
- 543 Yiming Gao, Bei Shi, Xueying Du, Liang Wang, Guangwei Chen, Zhenjie Lian, Fuham Qiu, Guoan
 544 Han, Weixuan Wang, Deheng Ye, et al. Learning diverse policies in moba games via macro-goals.
 545 *Advances in Neural Information Processing Systems*, 34:16171–16182, 2021. 6
- 546 Cole Gulino, Justin Fu, Wenjie Luo, George Tucker, Eli Bronstein, Yiren Lu, Jean Harb, Xinlei
 547 Pan, Yan Wang, Xiangyu Chen, et al. Waymax: An accelerated, data-driven simulator for large-
 548 scale autonomous driving research. *Advances in Neural Information Processing Systems*, 36:
 549 7730–7742, 2023. 1, 3
- 550 Hongyu Guo, Kun Xie, and Mehdi Keyvan-Ekbatani. Modeling driver’s evasive behavior during
 551 safety-critical lane changes: Two-dimensional time-to-collision and deep reinforcement learning.
 552 *Accident Analysis & Prevention*, 186:107063, 2023. 8, 22
- 553 Marah Halawa, Olaf Hellwich, and Pia Bideau. Action-based contrastive learning for trajectory
 554 prediction. In *European conference on computer vision*, pp. 143–159. Springer, 2022. 4
- 555 Niklas Hanselmann, Katrin Renz, Kashyap Chitta, Apratim Bhattacharyya, and Andreas Geiger.
 556 King: Generating safety-critical driving scenarios for robust imitation via kinematics gradients. In
 557 *European Conference on Computer Vision*, pp. 335–352. Springer, 2022. 1, 8, 22
- 558 Peter E Hart, Nils J Nilsson, and Bertram Raphael. A formal basis for the heuristic determination
 559 of minimum cost paths. *IEEE Transactions on Systems Science and Cybernetics*, 4(2):100–107,
 560 1968. doi: 10.1109/TSSC.1968.300136. 18
- 561 Yihan Hu, Jiazhi Yang, Li Chen, Keyu Li, Chonghao Sima, Xizhou Zhu, Siqi Chai, Senyao Du, Tian-
 562 wei Lin, Wenhui Wang, Lewei Lu, Xiaosong Jia, Qiang Liu, Jifeng Dai, Yu Qiao, and Hongyang Li.
 563 Planning-oriented autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer
 Vision and Pattern Recognition*, 2023. 3, 8, 22
- 564 Yihan Hu, Siqi Chai, Zhenning Yang, Jingyu Qian, Kun Li, Wenxin Shao, Haichao Zhang, Wei Xu,
 565 and Qiang Liu. Solving motion planning tasks with a scalable generative model. In *European
 566 Conference on Computer Vision*, pp. 386–404. Springer, 2024. 3
- 567 Zherui Huang, Xing Gao, Guanjie Zheng, Licheng Wen, Xuemeng Yang, and Xiao Sun. Safety-critical
 568 traffic simulation with adversarial transfer of driving intentions. *arXiv preprint arXiv:2503.05180*,
 569 2025. 21
- 570 Zhiyu Huang, Haochen Liu, and Chen Lv. Gameformer: Game-theoretic modeling and learning of
 571 transformer-based interactive prediction and planning for autonomous driving. In *Proceedings
 572 of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 3903–3913, October
 573 2023. 3
- 574 Xiaosong Jia, Zhenjie Yang, Qifeng Li, Zhiyuan Zhang, and Junchi Yan. Bench2drive: To-
 575 wards multi-ability benchmarking of closed-loop end-to-end autonomous driving. *arXiv preprint
 576 arXiv:2406.03877*, 2024. 3, 6, 17
- 577 Bo Jiang, Shaoyu Chen, Qing Xu, Bencheng Liao, Jiajie Chen, Helong Zhou, Qian Zhang, Wenyu Liu,
 578 Chang Huang, and Xinggang Wang. Vad: Vectorized scene representation for efficient autonomous
 579 driving. *ICCV*, 2023a. 3, 8, 22
- 580 Chiyu Jiang, Andre Cornman, Cheolho Park, Benjamin Sapp, Yin Zhou, Dragomir Anguelov, et al.
 581 Motiondiffuser: Controllable multi-agent motion prediction using diffusion. In *Proceedings of the
 582 IEEE/CVF conference on computer vision and pattern recognition*, pp. 9644–9653, 2023b. 1, 3
- 583 Chiyu Max Jiang, Yijing Bai, Andre Cornman, Christopher Davis, Xiukun Huang, Hong Jeon,
 584 Sakshum Kulshrestha, John Wheatley Lambert, Shuangyu Li, Xuanyu Zhou, Carlos Fuentes,
 585 Chang Yuan, Mingxing Tan, Yin Zhou, and Dragomir Anguelov. Scenediffuser: Efficient and
 586 controllable driving simulation initialization and rollout. In *The Thirty-eighth Annual Conference
 587 on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=a4qT29Levh>. 3

- 594 Quanyi Li, Zhenghao Peng, Lan Feng, Zhizheng Liu, Chenda Duan, Wenjie Mo, and Bolei Zhou.
 595 Scenarionet: Open-source platform for large-scale traffic scenario simulation and modeling.
 596 *Advances in Neural Information Processing Systems*, 2023. 1
 597
- 598 Haohong Lin, Xin Huang, Tung Phan-Minh, David S Hayden, Huan Zhang, Ding Zhao, Siddhartha
 599 Srinivasa, Eric M Wolff, and Hongge Chen. Causal composition diffusion model for closed-loop
 600 traffic generation. *arXiv preprint arXiv:2412.17920*, 2024. 3, 22
- 601 Henry X Liu and Shuo Feng. Curse of rarity for autonomous vehicles. *nature communications*, 15
 602 (1):4808, 2024. 4
- 603 Jack Lu, Kelvin Wong, Chris Zhang, Simon Suo, and Raquel Urtasun. Scenecontrol: Diffusion
 604 for controllable traffic scene generation. In *2024 IEEE International Conference on Robotics and*
 605 *Automation (ICRA)*, pp. 16908–16914. IEEE, 2024. 3
- 606 Yiren Lu, Justin Fu, George Tucker, Xinlei Pan, Eli Bronstein, Rebecca Roelofs, Benjamin Sapp,
 607 Brandyn White, Aleksandra Faust, Shimon Whiteson, et al. Imitation is not enough: Robustifying
 608 imitation with reinforcement learning for challenging driving scenarios. In *2023 IEEE/RSJ*
 609 *International Conference on Intelligent Robots and Systems (IROS)*, pp. 7553–7560. IEEE, 2023.
 610 3
- 611 Reza Mahjourian, Rongbing Mu, Valerii Likhosherstov, Paul Mougin, Xiukun Huang, Joao Messias,
 612 and Shimon Whiteson. Unigen: Unified modeling of initial agent states and trajectories for
 613 generating autonomous driving scenarios. In *2024 IEEE International Conference on Robotics*
 614 *and Automation (ICRA)*, pp. 16367–16373. IEEE, 2024. 1
- 615 Nico Montali, John Lambert, Paul Mougin, Alex Kuefler, Nicholas Rhinehart, Michelle Li, Cole
 616 Gulino, Tristan Emrich, Zoey Yang, Shimon Whiteson, et al. The waymo open sim agents
 617 challenge. *Advances in Neural Information Processing Systems*, 36:59151–59171, 2023. 3, 7, 21,
 618 22
- 619 Jiquan Ngiam, Benjamin Caine, Vijay Vasudevan, Zhengdong Zhang, Hao-Tien Lewis Chiang, Jeffrey
 620 Ling, Rebecca Roelofs, Alex Bewley, Chenxi Liu, Ashish Venugopal, et al. Scene transformer: A
 621 unified architecture for predicting multiple agent trajectories. *arXiv preprint arXiv:2106.08417*,
 622 2021. 1
- 623 Błażej Osiński, Piotr Miłoś, Adam Jakubowski, Paweł Zięcina, Michał Martyniak, Christopher Galias,
 624 Antonia Breuer, Silviu Homoceanu, and Henryk Michalewski. Carla real traffic scenarios—novel
 625 training ground and benchmark for autonomous driving. *arXiv preprint arXiv:2012.11329*, 2020.
 626 1
- 627 Zhenghao Peng, Wenjie Luo, Yiren Lu, Tianyi Shen, Cole Gulino, Ari Seff, and Justin Fu. Improving
 628 agent behaviors with rl fine-tuning for autonomous driving. In *European Conference on Computer*
 629 *Vision*, pp. 165–181. Springer, 2024. 2, 3, 5, 7, 20
- 630 Jonah Philion, Xue Bin Peng, and Sanja Fidler. Trajeglish: Traffic modeling as next-token prediction.
 631 In *The Twelfth International Conference on Learning Representations*, 2023. 3
- 632 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
 633 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances*
 634 *in Neural Information Processing Systems*, 36:53728–53741, 2023. 2
- 635 Davis Rempe, Jonah Philion, Leonidas J Guibas, Sanja Fidler, and Or Litany. Generating useful
 636 accident-prone driving scenarios via a learned traffic prior. In *Proceedings of the IEEE/CVF*
 637 *Conference on Computer Vision and Pattern Recognition*, pp. 17305–17315, 2022. 3
- 638 Katrin Renz, Kashyap Chitta, Otniel-Bogdan Mercea, A. Sophia Koepke, Zeynep Akata, and Andreas
 639 Geiger. Plant: Explainable planning transformers via object-level representations. In *Conference*
 640 *on Robotic Learning (CoRL)*, 2022. 8, 22
- 641 Luke Rowe, Roger Girgis, Anthony Gosselin, Bruno Carrez, Florian Golemo, Felix Heide, Liam
 642 Paull, and Christopher Pal. CtRL-sim: Reactive and controllable driving agents with offline
 643 reinforcement learning. In *8th Annual Conference on Robot Learning*, 2024. URL <https://openreview.net/forum?id=MfIUKzihC8.1>

- 648 John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional
 649 continuous control using generalized advantage estimation. [arXiv preprint arXiv:1506.02438](#), 2015.
 650 [21](#)
- 651 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
 652 optimization algorithms. [arXiv preprint arXiv:1707.06347](#), 2017. [5](#), [7](#), [19](#)
- 653 Ari Seff, Brian Cera, Dian Chen, Mason Ng, Aurick Zhou, Nigamaa Nayakanti, Khaled S Refaat,
 654 Rami Al-Rfou, and Benjamin Sapp. Motionlm: Multi-agent motion forecasting as language
 655 modeling. In [Proceedings of the IEEE/CVF International Conference on Computer Vision](#), pp.
 656 8579–8590, 2023. [3](#)
- 657 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Huawei Zhang,
 658 Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical
 659 reasoning in open language models. [arXiv preprint arXiv:2402.03300](#), 2024. [2](#), [5](#), [7](#), [19](#)
- 660 Samuel Sanford Shapiro and Martin B Wilk. An analysis of variance test for normality (complete
 661 samples). [Biometrika](#), 52(3-4):591–611, 1965. [21](#)
- 662 Qiao Sun, Xin Huang, Brian C Williams, and Hang Zhao. Intersim: Interactive traffic simulation via
 663 explicit relation modeling. In [2022 IEEE/RSJ International Conference on Intelligent Robots and](#)
 664 [Systems \(IROS\)](#), pp. 11416–11423. IEEE, 2022. [1](#)
- 665 Wenchao Sun, Xuewu Lin, Yining Shi, Chuang Zhang, Haoran Wu, and Sifa Zheng. Sparsedrive:
 666 End-to-end autonomous driving via sparse scene representation. [arXiv preprint arXiv:2405.19620](#),
 667 2024. [3](#)
- 668 Simon Suo, Sebastian Regalado, Sergio Casas, and Raquel Urtasun. Trafficsim: Learning to simulate
 669 realistic multi-agent behaviors. In [Proceedings of the IEEE/CVF Conference on Computer Vision](#)
 670 and Pattern Recognition, pp. 10400–10409, 2021. [3](#)
- 671 Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for
 672 reinforcement learning with function approximation. [Advances in neural information processing](#)
 673 [systems](#), 12, 1999. [7](#), [19](#)
- 674 Shuhan Tan, Kelvin Wong, Shenlong Wang, Sivabalan Manivasagam, Mengye Ren, and Raquel
 675 Urtasun. Scenegen: Learning to generate realistic traffic scenes. In [Proceedings of the IEEE/CVF](#)
 676 [Conference on Computer Vision and Pattern Recognition](#), pp. 892–901, 2021. [3](#)
- 677 Shuhan Tan, Boris Ivanovic, Xinshuo Weng, Marco Pavone, and Philipp Krahenbuehl. Language
 678 conditioned traffic generation. [arXiv preprint arXiv:2307.07947](#), 2023. [1](#), [3](#)
- 679 Shuhan Tan, Boris Ivanovic, Yuxiao Chen, Boyi Li, Xinshuo Weng, Yulong Cao, Philipp Krähenbühl,
 680 and Marco Pavone. Promptable closed-loop traffic simulation. In [8th Annual Conference on Robot](#)
 681 [Learning](#), 2024. [1](#), [3](#)
- 682 Shuhan Tan, John Lambert, Hong Jeon, Sakshum Kulshrestha, Yijing Bai, Jing Luo, Dragomir
 683 Anguelov, Mingxing Tan, and Chiyu Max Jiang. Scenediffuser++: City-scale traffic simulation via
 684 a generative world model. In [Proceedings of the IEEE/CVF Conference on Computer Vision and](#)
 685 [Pattern Recognition \(CVPR\)](#), pp. 1570–1580, June 2025. [3](#)
- 686 CARLA Team. CARLA Autonomous Driving Leaderboard. <https://leaderboard.carla.org/>, 2025. Accessed: 2025-04-09. [17](#)
- 687 Leonid Nisonovich Vaserstein. Markov processes over denumerable products of spaces, describing
 688 large systems of automata. [Problemy Peredachi Informatsii](#), 5(3):64–72, 1969. [21](#)
- 689 Suvin P Venturuthiyil and Mallikarjuna Chunchu. Anticipated collision time (act): A
 690 two-dimensional surrogate safety indicator for trajectory-based proactive safety assessment.
 691 [Transportation research part C: emerging technologies](#), 139:103655, 2022. [8](#), [22](#)
- 692 Yuning Wang, Pu Zhang, Lei Bai, and Jianru Xue. Fend: A future enhanced distribution-aware
 693 contrastive learning framework for long-tail trajectory prediction. In [Proceedings of the IEEE/CVF](#)
 694 [conference on computer vision and pattern recognition](#), pp. 1400–1409, 2023. [4](#)

- Wei Wu, Xiaoxin Feng, Ziyan Gao, and Yuheng Kan. Smart: Scalable multi-agent real-time motion generation via next-token prediction. *Advances in Neural Information Processing Systems*, 37: 114048–114071, 2024. 3
- Chejian Xu, Wenhao Ding, Weijie Lyu, Zuxin Liu, Shuai Wang, Yihan He, Hanjiang Hu, Ding Zhao, and Bo Li. Safebench: A benchmarking platform for safety evaluation of autonomous vehicles. *Advances in Neural Information Processing Systems*, 35:25667–25682, 2022. 3
- Danfei Xu, Yuxiao Chen, Boris Ivanovic, and Marco Pavone. Bits: Bi-level imitation for traffic simulation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 2929–2936. IEEE, 2023. 3
- Xintao Yan, Zhengxia Zou, Shuo Feng, Haojie Zhu, Haowei Sun, and Henry X Liu. Learning naturalistic driving environment with statistical realism. *Nature communications*, 14(1):2037, 2023. 21
- Xiuyu Yang, Shuhan Tan, and Philipp Krähenbühl. Long-term traffic simulation with interleaved autoregressive motion and scenario generation. *arXiv preprint arXiv:2506.17213*, 2025. 3
- Deheng Ye, Zhao Liu, Mingfei Sun, Bei Shi, Peilin Zhao, Hao Wu, Hongsheng Yu, Shaojie Yang, Xipeng Wu, Qingwei Guo, et al. Mastering complex control in moba games with deep reinforcement learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 6672–6679, 2020. 6
- Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, et al. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*, 2025. 2
- Chris Zhang, James Tu, Lunjun Zhang, Kelvin Wong, Simon Suo, and Raquel Urtasun. Learning realistic traffic agents in closed-loop. In *7th Annual Conference on Robot Learning*, 2023a. URL <https://openreview.net/forum?id=yobahDU4HPP>. 2, 3, 7, 19, 21
- Jiawei Zhang, Chejian Xu, and Bo Li. Chatscene: Knowledge-enabled safety-critical scenario generation for autonomous vehicles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15459–15469, 2024. 1
- Linrui Zhang, Zhenghao Peng, Quanyi Li, and Bolei Zhou. Cat: Closed-loop adversarial training for safe end-to-end driving. In *7th Annual Conference on Robot Learning*, 2023b. 1
- Zhejun Zhang, Alexander Liniger, Dengxin Dai, Fisher Yu, and Luc Van Gool. Trafficbots: Towards world models for autonomous driving simulation and motion prediction. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 1522–1529. IEEE, 2023c. 3
- Zhejun Zhang, Peter Karkus, Maximilian Igl, Wenhao Ding, Yuxiao Chen, Boris Ivanovic, and Marco Pavone. Closed-loop supervised fine-tuning of tokenized traffic models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2025. 3
- Yinan Zheng, Ruiming Liang, Kexin ZHENG, Jinliang Zheng, Liyuan Mao, Jianxiong Li, Weihao Gu, Rui Ai, Shengbo Eben Li, Xianyuan Zhan, and Jingjing Liu. Diffusion-based planning for autonomous driving with flexible guidance. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=wM2sfVgMDH>. 3
- Ziyuan Zhong, Davis Rempe, Yuxiao Chen, Boris Ivanovic, Yulong Cao, Danfei Xu, Marco Pavone, and Baishakhi Ray. Language-guided traffic simulation via scene-level diffusion. *arXiv preprint arXiv:2306.06344*, 2023a. 1, 3
- Ziyuan Zhong, Davis Rempe, Danfei Xu, Yuxiao Chen, Sushant Veer, Tong Che, Baishakhi Ray, and Marco Pavone. Guided conditional diffusion for controllable traffic simulation. In *2023 IEEE international conference on robotics and automation (ICRA)*, pp. 3560–3566. IEEE, 2023b. 1, 3
- Yunsong Zhou, Naisheng Ye, William Ljungbergh, Tianyu Li, Jiazhi Yang, Zetong Yang, Hongzi Zhu, Christoffer Petersson, and Hongyang Li. Decoupled diffusion sparks adaptive scene generation. *arXiv preprint arXiv:2504.10485*, 2025. 3

756	A Theoretical Analysis	16
757	A.1 Setting	16
758	A.2 Listwise View and Diversity Pressure	16
759	A.3 Clipping as Stability Control	16
760	A.4 Smoothness w.r.t. Policy Divergence	17
761	A.5 Variance and Enumeration	17
762	A.6 Convergence of Stochastic Ascent	17
763	A.7 Why RIFT Preserves Multimodality	17
764		
765	B Experimental Details	17
766	B.1 Experiment Framework	17
767	B.2 Route-level Analysis for CBV Identification	18
768	B.3 Algorithm Framework	18
769	B.4 Training Details	18
770	B.5 Baselines Detailed Description	19
771	B.6 Forward Simulation	20
772	B.7 State-Wise Reward Model Setup	20
773	B.8 Controllability and Realism Metrics	21
774		
775	C AV Evaluation Details	22
776	C.1 AV Methods Implementation	22
777	C.2 AV Evaluation Metrics	22
778	C.3 End-to-End AV Visualization	23
779		
780	D Additional Results	24
781	D.1 Detailed Qualitative Results of Style-Level Controllability	24
782	D.2 Detailed Analysis in Driving Comfort	24
783	D.3 Visualization of the AV-Centric Closed-Loop Simulation	26
784		
785	E Discussion and Broader Implications	26
786	E.1 Use of Large Language Models (LLMs)	26
787	E.2 Limitations and Future Work	26
788	E.3 Social Impact	26
789		
790		
791		
792		
793		
794		
795		
796		
797		
798		
799		
800		
801		
802		
803		
804		
805		
806		
807		
808		
809		

810 A THEORETICAL ANALYSIS

811 A.1 SETTING

814 For each $s \sim \mathcal{D}$, a frozen trajectory generation head yields $\mathcal{C}(s) = \{\tau_i\}_{i=1}^G$. The trajectory score head
 815 defines $\pi_\theta(\tau_i | s)$ on $\mathcal{C}(s)$. Finite-horizon simulation provides returns

$$816 \quad R_i(s) = \sum_{t=0}^T \gamma^t \text{StateWiseRM}(\tilde{\tau}_i^t, s). \quad (10)$$

819 Uniform (within-group) moments:

$$820 \quad \mu_{\text{uni}}(s) = \frac{1}{G} \sum_{j=1}^G R_j(s), \quad \sigma_{\text{uni}}^2(s) = \frac{1}{G} \sum_{j=1}^G (R_j(s) - \mu_{\text{uni}}(s))^2. \quad (11)$$

823 Uniform, centered advantages:

$$824 \quad \hat{A}_i(s) = \frac{R_i(s) - \mu_{\text{uni}}(s)}{\sqrt{\sigma_{\text{uni}}^2(s) + \varepsilon}}, \quad \frac{1}{G} \sum_{i=1}^G \hat{A}_i(s) = 0. \quad (12)$$

827 Let $\rho_i(\theta) = \pi_\theta(\tau_i | s)/\pi_{\theta_{\text{old}}}(\tau_i | s)$. Define the *RIFT* surrogate

$$828 \quad \mathcal{J}_{\text{RIFT}}(\theta) = \mathbb{E}_s \left[\frac{1}{G} \sum_{i=1}^G \psi(\rho_i(\theta), \hat{A}_i(s)) \right], \quad (13)$$

831 with dual-clip kernel

$$832 \quad \psi(\rho, \hat{A}) = \begin{cases} \min(\rho \hat{A}, \text{clip}(\rho, 1 - \epsilon, 1 + \epsilon) \hat{A}), & \hat{A} \geq 0, \\ \max(\min(\rho \hat{A}, \text{clip}(\rho, 1 - \epsilon, 1 + \epsilon) \hat{A}), c \hat{A}), & \hat{A} < 0, \end{cases} \quad (\epsilon > 0, c > 1). \quad (14)$$

835 **Assumptions.** (A) Support floor: there exists $\pi_{\min} > 0$ such that $\pi_{\theta_{\text{old}}}(\tau_i | s) \geq \pi_{\min}$ for all
 836 (s, i) , and $\pi_\theta > 0 \Rightarrow \pi_{\theta_{\text{old}}} > 0$ on $\mathcal{C}(s)$. (B) Boundedness: $|\hat{A}_i(s)| \leq A_{\max}$. (C) Regularity:
 837 $\log \pi_\theta(\tau_i | s)$ is L -Lipschitz and C^2 on compact Θ .

839 A.2 LISTWISE VIEW AND DIVERSITY PRESSURE

841 Consider the unclipped uniform surrogate

$$843 \quad L_{\text{RIFT}}(\theta) = \mathbb{E}_s \left[\frac{1}{G} \sum_{i=1}^G \rho_i(\theta) \hat{A}_i(s) \right] = \mathbb{E}_s \left[\frac{1}{G} \sum_{i=1}^G \frac{\hat{A}_i(s)}{\pi_{\theta_{\text{old}}}(\tau_i | s)} \pi_\theta(\tau_i | s) \right]. \quad (15)$$

845 **Proposition A.1** (Pairwise ascent and diversity). Fix s and shift an infinitesimal mass δ from j to i in
 846 $\pi_\theta(\cdot | s)$. Then $\delta L_{\text{RIFT}}(\theta) = \frac{\delta}{G} \left(\frac{\hat{A}_i(s)}{\pi_{\theta_{\text{old}}}(\tau_i | s)} - \frac{\hat{A}_j(s)}{\pi_{\theta_{\text{old}}}(\tau_j | s)} \right)$. Hence ascent moves mass toward larger
 847 \hat{A}/π_{old} , amplifying underrepresented high-quality candidates when π_{old} is peaky.

849 **Corollary A.2** (Top-1 Fisher consistency under uniform reference). If $\pi_{\theta_{\text{old}}}$ is uniform on $\mathcal{C}(s)$ and
 850 $i^*(s) = \arg \max_i \hat{A}_i(s)$ is unique, any global maximizer of L_{RIFT} concentrates $\pi_\theta(\cdot | s)$ on $i^*(s)$.

852 A.3 CLIPPING AS STABILITY CONTROL

854 Clipping is a pointwise pessimistic transform: for any $x = \rho \hat{A}$,

$$855 \quad \min(\rho \hat{A}, \text{clip}(\rho) \hat{A}) \leq \rho \hat{A}.$$

856 Summed over mixed signs, there is no global monotone lower bound for L_{RIFT} ; instead, clipping
 857 serves to bound the value and the gradient.

858 **Lemma A.3** (Bounded values and gradients). If $|\hat{A}| \leq A_{\max}$, then for all (s, i) : (i) Value bounds:
 859 $\psi \in [0, (1 + \epsilon)\hat{A}]$ for $\hat{A} \geq 0$, and $\psi \in [c\hat{A}, 0]$ for $\hat{A} < 0$. (ii) Gradient bounds:

$$861 \quad \left| \frac{\partial \psi}{\partial \log \pi_\theta} \right| \leq \begin{cases} (1 + \epsilon)|\hat{A}|, & \hat{A} \geq 0, \\ c|\hat{A}|, & \hat{A} < 0, \end{cases}$$

863 and on the negative branch when the dual-clip is active ($\psi = c\hat{A}$) the derivative is 0.

864 A.4 SMOOTHNESS W.R.T. POLICY DIVERGENCE
865

866 Write $w_i(s) = \hat{A}_i(s)/\pi_{\theta_{\text{old}}}(\tau_i | s)$ and note $|w_i| \leq A_{\max}/\pi_{\min}$ under Assumption A with $\pi_{\min} =$
867 $\inf_{s,i} \pi_{\theta_{\text{old}}}(i | s) > 0$ (label-smoothing in practice). Then the unclipped surrogate is linear in π_θ :

$$868 L_{\text{RIFT}}(\theta) - L_{\text{RIFT}}(\theta') = \mathbb{E}_s \left[\frac{1}{G} \sum_i w_i(s) (\pi_\theta(i | s) - \pi_{\theta'}(i | s)) \right].$$

871 **Lemma A.4** (Lipschitz continuity via KL). *For any θ, θ' ,*

$$873 |L_{\text{RIFT}}(\theta) - L_{\text{RIFT}}(\theta')| \leq \frac{A_{\max}}{\pi_{\min}} \sqrt{2 \mathbb{E}_s [\text{KL}(\pi_\theta(\cdot | s) \| \pi_{\theta'}(\cdot | s))]}.$$

876 *Proof.* By Hölder and Pinsker: $|\sum_i w_i \Delta \pi| \leq \|w\|_\infty \|\Delta \pi\|_1 \leq (A_{\max}/\pi_{\min}) \sqrt{2 \text{KL}(\pi_\theta \| \pi_{\theta'})}$, then
877 average over s . \square

878 **Lemma A.5** (Lipschitz continuity of clipped surrogate). *Because $\partial \psi / \partial \pi_\theta(i | s)$ is bounded by
879 A_{\max}/π_{\min} whenever the active branch is differentiable,*

$$881 |\mathcal{J}_{\text{RIFT}}(\theta) - \mathcal{J}_{\text{RIFT}}(\theta')| \leq \frac{A_{\max}}{\pi_{\min}} \sqrt{2 \mathbb{E}_s [\text{KL}(\pi_\theta(\cdot | s) \| \pi_{\theta'}(\cdot | s))]}.$$

884 A.5 VARIANCE AND ENUMERATION

885 Let $f_i(s; \theta) = \psi(\rho_i(\theta), \hat{A}_i(s))$. Exact enumeration yields $\text{Var}\left(\frac{1}{G} \sum_i f_i | s\right) = 0$ (assuming $\tilde{\tau}_i$ and
886 their evaluations are fixed during the update; otherwise, environment randomness still induces nonzero
887 variance), while sampling i i.i.d. within the group gives conditional variance $\text{Var}(f_i | s)/N$ for N
888 samples.

890 A.6 CONVERGENCE OF STOCHASTIC ASCENT

892 **Theorem A.6** (Convergence to a stationary point). *Under Assumptions A–C, with step sizes $\eta_k > 0$,
893 $\sum_k \eta_k = \infty$, $\sum_k \eta_k^2 < \infty$, and unbiased bounded-variance stochastic subgradients, the iterates of
894 stochastic subgradient ascent on $\mathcal{J}_{\text{RIFT}}$ satisfy*

$$895 \liminf_{k \rightarrow \infty} \mathbb{E}[\text{dist}(0, \partial^C \mathcal{J}_{\text{RIFT}}(\theta_k))] = 0,$$

897 where ∂^C denotes the Clarke generalized gradient.

899 *Sketch.* By Lemma A.3, generalized gradients are uniformly bounded; regularity of $\log \pi_\theta$ on compact
900 Θ implies Lipschitz continuity. Robbins–Monro / Kushner–Yin results for non-smooth stochastic
901 approximation apply. \square

903 A.7 WHY RIFT PRESERVES MULTIMODALITY

905 By Proposition A.1, ascent compares \hat{A}/π_{old} : under peaky π_{old} , underrepresented high- \hat{A} candidates
906 receive stronger positive updates, preserving and enhancing diversity. In the special case π_{old} is
907 (approximately) uniform, RIFT reduces to a listwise ranking ascent that directly promotes larger \hat{A} .

909 B EXPERIMENTAL DETAILS

911 B.1 EXPERIMENT FRAMEWORK

913 Our framework for reliable AV-centric closed-loop simulation is developed upon well-established
914 traffic simulation platforms, notably the CARLA Leaderboard Team (2025) and Bench2Drive Jia et al.
915 (2024), which serve as standard benchmarks in autonomous driving research. Traditionally, these
916 platforms use predefined scenarios along the AV’s global route to evaluate the multi-dimensional
917 performance of AV methods. In contrast, we replace these static scenarios with dynamically generated
traffic flows by randomly spawning background vehicles around the AV’s global path and simulating

918 their behavior using rule-based driving policies, as described in Section 3.1. Through the CBV
 919 identification mechanism outlined in Appendix B.2, we naturally introduce interactions between the
 920 AV and CBVs, thereby generating continuous, interactive scenarios over time. This framework serves
 921 as the foundation for both the training and evaluation processes in this paper.
 922

923 B.2 ROUTE-LEVEL ANALYSIS FOR CBV IDENTIFICATION 924

925 Identifying Critical Background Vehicles (CBVs) is essential to our AV-centric closed-loop simulation.
 926 Let \mathcal{V}_{AV} denote the autonomous vehicle (AV), and $\mathcal{V}_{BV} = \{\mathcal{V}_i\}_{i=1}^N$ represent the set of background
 927 vehicles in the environment. The AV navigates along a predefined global route $\mathcal{P} = \{p_k\}_{k=1}^M$,
 928 where each p_k corresponds to a waypoint along the route. The goal of CBV identification is to
 929 select background vehicles that are likely to share the AV’s destination and have similar estimated
 930 travel distance, thereby facilitating route-level interactions between the AV and CBVs. The primary
 931 criterion for identifying CBVs is the relative *distance-to-goal* difference between the AV and each
 932 background vehicle. This is mathematically expressed as:
 933

$$934 \quad | \hat{D}_{\text{global}}(p_k, \mathcal{V}_i) - \hat{D}_{\text{global}}(p_k, \mathcal{V}_{AV}) | < \delta, \quad (16)$$

937 where, $\hat{D}_{\text{global}}(p_k, \mathcal{V}_i)$ and $\hat{D}_{\text{global}}(p_k, \mathcal{V}_{AV})$ denote the estimated travel distance required for the
 938 background vehicle \mathcal{V}_i and the AV to reach waypoint p_k , respectively. The distance-to-goal for each
 939 vehicle is computed by determining the distance from its current position to the target waypoint p_k
 940 using the A* global path planning algorithm Hart et al. (1968). A threshold δ is introduced to define
 941 the maximum allowable difference in distance-to-goal. A background vehicle is considered critical
 942 and included in the CBV set \mathcal{C} if the absolute distance-to-goal difference between it and the AV is
 943 smaller than δ .
 944

This approach selects background vehicles whose destinations and estimated travel distances are
 sufficiently aligned with those of the AV, thereby ensuring meaningful and realistic route-level
 interactions. Once a CBV is identified, the planning path previously generated via A* during distance-
 to-goal estimation is directly adopted as its global navigation path, which is further transformed
 into the reference line for downstream CBV planning, naturally introducing route-level interactions
 between the AV and CBVs. The threshold δ serves as a tunable parameter to adjust the sensitivity of
 the CBV selection process. In this study, we set δ to 15m to achieve a balanced trade-off between
 sensitivity and selection accuracy.
 951

952 B.3 ALGORITHM FRAMEWORK 953

954 For clarity, we summarize the procedure of *RIFT* within our AV-centric closed-loop simulation frame-
 955 work in Algorithm 1. The planning model is initialized from the IL pre-trained checkpoint provided
 956 by Pluto official codebase¹, followed by RL fine-tuning within the CARLA simulator Dosovitskiy
 957 et al. (2017) to generate realistic and controllable traffic scenarios.
 958

959 B.4 TRAINING DETAILS 960

We perform RL fine-tuning on selected modules of the IL pre-trained planning model (Pluto). As
 shown in the ablation results (Section 5.4), fine-tuning only the trajectory scoring head achieves
 the best trade-off between realism and controllability. Accordingly, all fine-tuning baselines adopt
 this setting to ensure consistency and fair comparison. Our training framework is built on the open-
 source Lightning platform². Fine-tuning is conducted on 2× Bench2Drive220, while evaluation
 is performed on dev10, both from the Bench2Drive project. All experiments are conducted on
 NVIDIA GeForce RTX 4090D GPUs, with each fine-tuning run taking approximately 8 hours on a
 single GPU. Detailed training setups and hyperparameter configurations are provided in Table 4 and
 Table 5.
 961

¹<https://github.com/jchengai/pluto>

²<https://github.com/Lightning-AI/pytorch-lightning>

Algorithm 1 Procedure for *RIFT* in the AV-Centric Closed-Loop Simulation Framework.

```

1: Input: IL pre-trained planning model  $\pi_{\theta_{\text{init}}}$ , buffer  $\mathcal{D}$             $\triangleright$  IL pre-training (nuPlan Caesar et al. (2021))
2: planning model  $\pi_\theta \leftarrow \pi_{\theta_{\text{init}}}$ 
3: for interaction = 1, ...,  $I$  do                                      $\triangleright$  RL fine-tuning (CARLA Dosovitskiy et al. (2017))
4:   Update the old planning model  $\pi_{\theta_{\text{old}}} \leftarrow \pi_\theta$ 
5:   while  $\mathcal{D}$  not full do                                          $\triangleright$  Collect rollout data
6:     for step = 1, ...,  $T$  do
7:       Obtain  $G$  candidate trajectories  $\{\tau_i\}_{i=1}^G$  from  $\pi_{\theta_{\text{old}}}$  for each CBV  $\triangleright$  Policy inference
8:       Compute simulated rollouts  $\{\tilde{\tau}_i\}_{i=1}^G$  from  $\{\tau_i\}_{i=1}^G$            $\triangleright$  Forward simulation
9:       Compute reward  $\{R_i\}_{i=1}^G$ , advantage  $\{\hat{A}_i\}_{i=1}^G$  for each  $\tilde{\tau}_i$  with Equation (6)
10:      Store transition into buffer  $\mathcal{D}$ 
11:    end for
12:  end while
13:  for RIFT iteration = 1, ...,  $\mu$  do                                 $\triangleright$  Policy fine-tuning
14:    Sample mini-batches transition from the buffer  $\mathcal{D}$ 
15:    Update model  $\pi_\theta$  by maximizing the RIFT objective (Equation (9))
16:  end for
17: end for
18: Output: RL fine-tuned planning model

```

B.5 BASELINES DETAILED DESCRIPTION

To comprehensively evaluate *RIFT* in an AV-centric closed-loop simulation environment, we compare it against a range of baselines, including pure imitation learning (IL), pure reinforcement learning (RL), and various fine-tuning approaches based on IL, RL, or their combination. We initialize all fine-tuning methods from the pre-trained Pluto checkpoint and fine-tune only the trajectory scoring head to preserve trajectory-level realism. The details of each baseline are summarized below.

- *Pluto* [Cheng et al. \(2024a\)](#) is an open-source IL-based planning framework for autonomous driving. It processes vectorized scene representations as input and outputs multimodal trajectories for downstream planning. In the AV-centric closed-loop simulation, the method directly uses a pre-trained checkpoint without additional fine-tuning.
 - *FREA* [Chen et al. \(2024b\)](#) is an RL-based approach designed to generate safety-critical yet AV-feasible scenarios. It incorporates a feasibility-aware training objective. In the AV-centric closed-loop simulation, FREA selects potential collision points along the AV’s global route as adversarial goals.
 - *PPO* [Chen et al. \(2024b\)](#) is a variant of FREA that focuses solely on generating safety-critical scenarios. Unlike FREA, it disregards the feasibility constraints of AV and treats adversariality as the only optimization objective.
 - *FPPO-RS* [Chen et al. \(2024b\)](#) is another FREA variant that integrates AV’s feasibility constraints into the reward shaping process, thereby balancing adversariality with scenario reasonability.
 - *PPO-Pluto* fine-tunes the pre-trained planning model using the PPO algorithm [Schulman et al. \(2017\)](#). The fine-tuning follows the same reward structure as detailed in Appendix B.7, aligning with *RIFT*.
 - *REINFORCE-Pluto* employs the REINFORCE algorithm [Sutton et al. \(1999\)](#) to fine-tune the pre-trained Pluto model under the same reward design as detailed in Appendix B.7.
 - *GRPO-Pluto* utilizes the basic GRPO algorithm [Shao et al. \(2024\)](#) for fine-tuning, employing the pre-trained Pluto model as the reference for KL regularization, while incorporating the standard PPO-Clip.
 - *SFT-Pluto* is a purely supervised fine-tuning approach, where PDM-Lite [Beißwenger \(2024\)](#) serves as the expert model, providing supervision at the target speed level.
 - *RTR-Pluto* [Zhang et al. \(2023a\)](#) is a hybrid framework combining imitation and reinforcement learning. While the original RTR utilizes human driving trajectories as supervision, our setting

1026 replaces this with PDM-Lite due to the lack of human-level demonstrations. The RL component
 1027 uses sparse infraction-based rewards, consistent with the original RTR, and applies PPO for
 1028 optimization.

- 1029 • *RS-Pluto* Peng et al. (2024) also adopts a hybrid IL+RL paradigm, originally trained via RE-
 1030 INFORCE using ground-truth supervision and sparse rewards to ensure safety and realism. In
 1031 our adaptation, PDM-Lite substitutes the ground-truth expert, while the rest of the methodology
 1032 remains unchanged.

1033 **B.6 FORWARD SIMULATION**

1036 Trajectory-based imitation learning often overlooks underlying system dynamics, leading to dis-
 1037crepancies between planned and executed behavior Cheng et al. (2024b). To address this issue, we
 1038 perform a forward simulation for each candidate trajectory τ_i of the CBV, yielding a rollout $\tilde{\tau}_i$. The
 1039 simulation couples a PID controller for trajectory tracking with a kinematic bicycle model for state
 1040 propagation. The PID controller is identical to that used during closed-loop execution, ensuring
 1041 behavioral consistency between training and deployment. By evaluating rollouts rather than raw
 1042 trajectories, we reduce this dynamics gap and obtain more reliable assessments.

1043 In parallel, we also forecast the motions of surrounding actors. During data collection, the cur-
 1044rent actions a^{bg} of surrounding actors are recorded. Following the rule-based forecasting scheme
 1045 in Beißwenger (2024), these actions are assumed constant over the forecast horizon and are used
 1046 to advance surrounding states. The resulting actor forecasts are combined with the CBV rollouts
 1047 to compute rewards, thereby ensuring that interaction effects with the environment are faithfully
 1048 captured in evaluation.

1049 While subsequent rollout is open-loop, the first transition is closed-loop. This step integrates (i)
 1050 the same PID policy as in real execution, (ii) the observed current actions of surrounding actors,
 1051 and (iii) a kinematic bicycle model that approximates CARLA’s single-step dynamics. Accordingly,
 1052 the transition from (s, a, a^{bg}) to s' produces a reward consistent with the standard RL structure
 1053 $(s, a) \rightarrow s' \rightarrow r$. Subsequent rollout steps serve as open-loop estimates of longer-horizon outcomes,
 1054 enriching evaluation while preserving closed-loop fidelity at the transition boundary.

1055 **B.7 STATE-WISE REWARD MODEL SETUP**

1056 To capture diverse human driving styles, we decompose driving behaviors into distinct reward
 1057 components, following Cusumano-Towner et al. (2025). Different styles are constructed by combining
 1058 weights assigned to each reward component (detailed in Table 6), enabling a range of behaviors from
 1059 aggressive to conservative. The total driving reward is defined as:

$$R = R_{\text{collision}} + R_{\text{off-road}} + R_{\text{comfort}} + R_{\text{lane}} + R_{\text{velocity}} + R_{\text{timestep}}. \quad (17)$$

1060 The individual terms are described as follows:

- 1061 • $R_{\text{collision}} = -(\alpha_{\text{collision}} + |v|) \mathbb{1}_{\text{collision}}$: penalizes collisions, with higher penalties at higher speeds.
- 1062 • $R_{\text{off-road}} = -\alpha_{\text{boundary}} \mathbb{1}_{\text{boundary}}$: penalizes deviations from the drivable area.
- 1063 • $R_{\text{comfort}} = -\alpha_{\text{comfort}} (\mathbb{1}_{|a|>4} + \mathbb{1}_{|\omega|>4})$: penalizes excessive acceleration and angular acceleration.
- 1064 • $R_{\text{l-align}} = \alpha_{\text{l-align}} \left(\min(\cos(\theta_f), 0) + \alpha_{\text{vel-align}} \min(\cos(\theta_f) * v, 0) + 0.25 \left(1 - \frac{|\theta_f|}{\pi/2} \right) \right)$: guides
 1065 the agent to follow the correct driving direction and remain parallel to the lane markings.
- 1066 • $R_{\text{l-center}} = -\alpha_{\text{l-center}} \left(\mathbb{1}_{\cos(\theta_f)>0.5} * \left(|x_f - \alpha_{\text{center-bias}}| - \frac{0.05}{\exp(|x_f - \alpha_{\text{center-bias}}| - 0.5)} \right) \right)$: guides the
 1067 agent to prefer trajectories that remain centered within the lane.
- 1068 • $R_{\text{velocity}} = \alpha_{\text{velocity}} \max(\cos(\theta_f), 0.0) \mathbb{1}_{3<|v|<20} * |v|$: promotes forward movement and biases
 1069 the agent toward choosing routes with consistent traffic flow rather than traffic jams.
- 1070 • $R_{\text{timestep}} = -\alpha_{\text{timestep}} \mathbb{1}_{|v|>0 \vee |a|>0}$: applies a small per-step penalty, encouraging efficiency. It is
 1071 disabled when the agent is stationary to allow appropriate waiting behavior at intersections.

1072 Building on the reward definitions above, we construct a state-wise reward model StateWiseRM(\cdot),
 1073 which computes a scalar reward based on a set of interpretable features extracted from each rollout

1080

Table 4: Hyperparameters used
1081 in RIFT Training.

1082

Parameter	Value
Batch size	256
Rollout buffer capacity	4096
Fine-tune initial LR	$1 \times e^{-4}$
Minimum LR	$1 \times e^{-6}$
LR decay across iteration	0.9
LR schedule	Cosine
Num. RIFT epoch	16
Warmup Epoch of RIFT	3
AdamW weight-decay	$1 \times e^{-5}$

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

1296

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1134
 1135 **Interaction Metrics.** Following the metric design principles proposed in the WOSAC challenge [Montali et al. \(2023\)](#) and other widely adopted evaluation frameworks [Lin et al. \(2024\)](#); [Chen et al. \(2024a\)](#),
 1136 we adopt a set of well-established metrics to comprehensively evaluate agent interactions:
 1137

- 1138 • *Collision Per Kilometer (CPK)* [Chen et al. \(2024a\)](#): the average number of scenario collisions per
 1139 kilometer of driving distance.
- 1140 • *Route Progress (RP)* [Chen et al. \(2024a\)](#): the total distance traveled by all CBVs, reflecting route
 1141 completion.
- 1142 • *2D Time-to-Collision (2D-TTC)* [Guo et al. \(2023\)](#): the minimum of longitudinal and lateral
 1143 time-to-collision from the AV’s perspective, capturing the interaction risk posed by CBVs.
- 1144 • *Anticipated Collision Time (ACT)* [Venthuruthiyil & Chunchu \(2022\)](#) : a safety-critical metric
 1145 measuring the AV’s proximity to potential collisions, reflecting the interaction intensity introduced
 1146 by CBVs.

1147
 1148 **Map Metrics.** Map metrics evaluate adherence to road geometry, reflecting how well agents remain
 1149 within drivable areas and comply with map constraints.

- 1150 • *Off-Road Rate (ORR)* [Chen et al. \(2024a\)](#): the percentage of time CBVs spend off-road on average.
 1151

1152 C AV EVALUATION DETAILS

1153 C.1 AV METHODS IMPLEMENTATION

1154 To assess the effectiveness of *RIFT* in generating reliable and interactive scenarios for AV evaluation
 1155 in the AV-centric closed-loop simulation environment, we evaluate the following representative and
 1156 stable AV methods:

- 1157 • *PDM-Lite* [Beißwenger \(2024\)](#): A rule-based privileged expert method that achieves state-of-the-art
 1158 performance on the CARLA Leaderboard 2.0 by leveraging components such as the Intelligent
 1159 Driver Model and the kinematic bicycle model. This open-source method serves as a strong baseline
 1160 for comparison.
- 1161 • *PlanT* [Renz et al. \(2022\)](#): An explainable, learning-based planning method that operates on an
 1162 object-level input representation and is trained through imitation learning.
- 1163 • *UniAD* [Hu et al. \(2023\)](#): A planning-oriented unified framework integrating perception, prediction,
 1164 mapping, and planning into one end-to-end model using query-based interfaces.
- 1165 • *VAD* [Jiang et al. \(2023a\)](#): A fast, end-to-end vectorized driving paradigm representing scenes with
 1166 vectorized motion and map elements for efficient, safe planning.

1167 C.2 AV EVALUATION METRICS

1168 As detailed in Appendices B.1 and B.4, we develop an AV-centric closed-loop simulation environment,
 1169 including a training and evaluation pipeline based on Bench2Drive. The AV closed-loop
 1170 evaluation metrics proposed in Bench2Drive extend the original metrics of the CARLA Leaderboard
 1171 by emphasizing the specific strengths and weaknesses of different methods across various aspects,
 1172 such as merging and overtaking, thereby making them suitable for evaluating performance under
 1173 predefined scenarios. However, as noted in Appendix B.1, replacing predefined scenarios with CBV-
 1174 generated traffic scenarios precludes the evaluation of specific AV capabilities. To systematically
 1175 assess the quality of traffic scenarios generated by different CBV methods, we follow the practice
 1176 of KING [Hanselmann et al. \(2022\)](#) and introduce PDM-Lite [Beißwenger \(2024\)](#)—a rule-based
 1177 privileged planner—as a reference AV. By measuring its performance under various CBV methods,
 1178 we evaluate:

- 1179 • *Feasibility*, via PDM-Lite’s Driving Score (DS)—a high DS indicates the PDM-Lite can complete
 1180 its route without severe collisions or rule violations, implying the generated traffic scenario is
 1181 feasible.
- 1182 • *Naturalness*, via a newly proposed metric, Blocked Rate (BR)—a low BR suggests that CBVs do
 1183 not unrealistically obstruct the AV, reflecting naturalistic behavior.



Figure 5: Representative closed-loop interactions between *RIFT*-generated traffic flows and end-to-end autonomous driving algorithms. UniAD (top) and VAD (bottom) are shown interacting with surrounding vehicles orchestrated by *RIFT*, which preserves realistic driving styles while enabling dynamic CBV–AV interactions. The controlled background vehicle (CBV) is highlighted in purple, the autonomous vehicle (AV, end-to-end) in red, and other background vehicles (BVs) in blue.

These metrics enable a principled comparison of traffic quality generated by different CBV methods. Furthermore, to assess the capacity of generated traffic scenarios to expose AV limitations, we test multiple learning-based AV methods under an identical CBV method and quantify their relative performance drop compared to PDM-Lite [Beißwenger \(2024\)](#). The relative driving score degradation (ΔDS) reflects how effectively the traffic scenario stresses the AV policy, with larger drops indicating stronger capability in revealing planning weaknesses.

The evaluation metrics are summarized as follows:

- *Driving Score (DS)*: $R_i P_i$ — The main metric of the leaderboard, calculated as the product of route completion and the infraction penalty. Here, R_i represents the percentage of completion of the i -th route, and P_i denotes the infraction penalty. The maximum value is 100.
- *Block Rate (BR)*: The average number of occurrences where a CBV fails to navigate its route normally and obstructs the AV’s progress.
- *Relative Driving Score Degradation (ΔDS)*: The reduction in Driving Score of a learning-based AV compared to PDM-Lite under the same CBV method, indicating how effectively the scenario reveals weaknesses in AV planning.

C.3 END-TO-END AV VISUALIZATION

To further validate the feasibility of *RIFT* as a closed-loop evaluation framework, we extend its application beyond traffic scenario generation to testing end-to-end autonomous driving algorithms. In contrast to conventional adversarial approaches that often introduce unrealistic or overly aggressive behaviors, *RIFT* generates traffic flows that preserve the realism of human driving styles, engage the autonomous vehicle in genuine interactive behaviors, and maintain feasibility by ensuring that the resulting scenes, though diverse and stress-inducing, remain solvable for the end-to-end method. This balance enables *RIFT* to evaluate robustness under credible conditions while avoiding degenerate or unsolvable scenarios.

Figure 5 presents representative interactions between *RIFT*-generated traffic flows and two representative end-to-end driving models, UniAD and VAD. As shown, *RIFT* adapts seamlessly to different driving policies, producing realistic and interactive scenes where the autonomous vehicle must negotiate with surrounding traffic. These results underscore the capability of *RIFT* to provide realistic yet interactive closed-loop evaluations, highlighting its potential as a versatile tool for testing the robustness of end-to-end AV systems.



Figure 6: Qualitative illustration of *RIFT*'s style-level controllability under different reward configurations. CBV is marked in purple, AV (PDM-Lite) is in red, and BVs are in blue.

D ADDITIONAL RESULTS

D.1 DETAILED QUALITATIVE RESULTS OF STYLE-LEVEL CONTROLLABILITY

As discussed in Section 5.4, we investigate the style-level controllability of *RIFT* under different reward configurations. The aggressive variant applies a reduced collision penalty and places greater emphasis on driving efficiency (Table 6), encouraging assertive behaviors such as overtaking. In contrast, the normal configuration imposes a higher collision penalty to promote safer and more conservative driving behaviors.

Quantitative results in Table 3 show that the aggressive variant achieves greater driving efficiency at the expense of more frequent collisions and off-road events. To complement these findings, Figure 6 presents a qualitative comparison in a single-lane intersection scenario where a leading BV halts at a stop sign. The aggressive CBV variant attempts an overtaking maneuver, resulting in a collision, whereas the normal CBV variant yields and waits, demonstrating distinct behavioral patterns induced by different reward preferences. These results highlight the controllability of *RIFT* in modulating driving style according to user-specified reward configuration.

D.2 DETAILED ANALYSIS IN DRIVING COMFORT

Metrics. To further evaluate the driving comfort of different CBV methods, we define several comfort metrics based on Bench2Drive, which assesses agent comfort through acceleration and jerk profiles. Specifically, we measure comfort using the following metrics:

- *Uncomfortable Rate (UCR)*: the percentage of simulation time during which CBVs experience discomfort.
- *Driving Jerk (Jerk)*: the time derivative of acceleration, quantifying the abruptness of acceleration changes and the smoothness of CBV rollouts.

1296 **Table 7: Comparison of CBV Comfort Metrics across**
 1297 **Various AV Methods.** Each metric is evaluated across three
 1298 random seeds.

Method	PDM-Lite		PlanT	
	UCR ↓	Jerk ↓	UCR ↓	Jerk ↓
Pluto	56.45 ± 4.14	-0.16 ± 3.72	50.26 ± 2.17	-0.42 ± 3.38
PPO	74.76 ± 2.71	-0.51 ± 4.61	74.90 ± 1.21	0.40 ± 4.83
FREA	72.40 ± 1.72	0.29 ± 4.61	73.48 ± 3.83	-0.15 ± 4.91
FPPO-RS	68.33 ± 1.90	-0.07 ± 3.96	66.67 ± 0.82	-0.15 ± 3.95
SFT-Pluto	68.14 ± 4.91	-0.06 ± 4.06	59.78 ± 4.72	-0.11 ± 4.00
RS-Pluto	70.31 ± 4.07	0.32 ± 4.12	65.18 ± 2.11	-0.16 ± 4.07
RTR-Pluto	55.58 ± 4.76	-0.19 ± 3.37	45.12 ± 2.66	-0.14 ± 3.34
PPO-Pluto	58.29 ± 2.70	-0.32 ± 3.70	54.85 ± 5.82	-0.07 ± 3.40
REINFORCE-Pluto	68.10 ± 1.22	0.23 ± 3.96	64.94 ± 5.36	-0.11 ± 3.96
GRPO-Pluto	78.58 ± 0.59	0.22 ± 4.62	77.13 ± 0.65	-0.23 ± 4.58
<i>RIFT</i> -Pluto (ours)	76.90 ± 2.82	0.59 ± 4.12	72.41 ± 4.02	0.21 ± 4.44

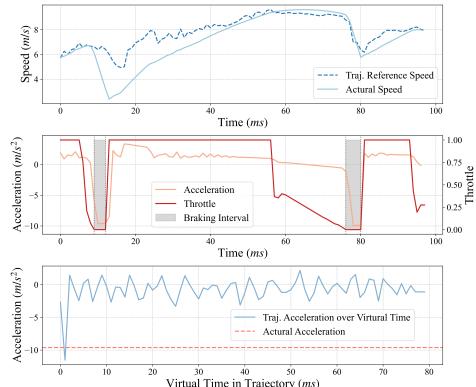


Figure 7: Controller Performance.

To determine whether a CBV’s current state is considered comfortable, we adopt the Frame Variable Smoothness (FVS) criterion from Bench2Drive:

$$\text{Frame Variable Smoothness (FVS)} = \begin{cases} \text{True} & \text{if lower bound} \leq p_i \leq \text{upper bound}. \\ \text{False} & \text{otherwise} \end{cases} \quad (22)$$

$p \in \text{smoothness vars}, 0 \leq i \leq \text{total frames}$

The smoothness variables include longitudinal acceleration (expert bounds: [-4.05, 2.40]), maximum absolute lateral acceleration (expert bounds: [-4.89, 4.89]), and maximum jerk magnitude (expert bounds: [-8.37, 8.37]).

Main Results. The quantitative results of the comfort metrics are presented in Table 7. All CBV methods exhibit notable levels of driving discomfort. Although the more conservative methods identified in Section 5.2 achieve relatively lower levels of discomfort, a high baseline of discomfort persists between methods.

To investigate the underlying causes of discomfort, we further decouple the planned trajectories from the executed control actions. In CARLA, most CBV methods rely on PID controllers to transform high-level trajectory waypoints into executable driving commands, including throttle, steering, and brake. As shown in Figure 7, the upper panel illustrates the speed tracking curve, while the middle panel presents the raw throttle signal and corresponding acceleration profile.

Because trajectory generation is performed state-wise, predicting only the immediate next action, the reference states may vary discontinuously over time. These discontinuities are amplified by the PID controller, whose binary throttle/brake responses induce abrupt changes in acceleration, ultimately leading to discomfort during vehicle operation. Such execution-level instabilities are a major contributor to the discomfort observed across CBV methods.

While many CBV methods attempt to mitigate discomfort through fine-tuning strategies that incorporate post-action feedback via reward shaping or expert action alignment, *RIFT* adopts a different approach. It employs a state-wise reward model (see Appendix B.7) that quantifies comfort within the trajectory’s virtual forward simulation.

To further analyze this, we visualize both the actual acceleration after executing a selected trajectory and the corresponding virtual-time acceleration (shown in Figure 7). The results reveal that while virtual-time acceleration aligns with actual motion at the beginning of the trajectory, it underestimates acceleration variations in later segments of the trajectory. This leads to an overly conservative estimation of trajectory-level discomfort, resulting in insufficient supervision during training and reflected in *RIFT*’s comfort performance in Table 7.

In summary, the discomfort exhibited by CBV methods can be attributed to two primary sources:

- **Tracking instability**, caused by discontinuities in planned trajectories and the limited control fidelity of PID controllers. Discrete, state-wise planning combined with low-resolution, often binary control outputs amplifies acceleration fluctuations and leads to uncomfortable motion.

- 1350
 1351 • **Inadequate comfort modeling**, particularly in state-wise reward formulations such as that adopted
 1352 by *RIFT* and GPRO. These formulations fail to capture long-term trajectory-level discomfort,
 1353 leading to insufficient supervision during training and suboptimal comfort performance.

1354 **D.3 VISUALIZATION OF THE AV-CENTRIC CLOSED-LOOP SIMULATION**

1355 To qualitatively evaluate the robustness of *RIFT* across diverse AV-centric scenarios, we provide
 1356 additional temporal visualizations of closed-loop simulations. As shown in Figure 8, the traffic scene
 1357 consists of the autonomous vehicle (AV, controlled by PDM-Lite), background vehicles (BVs), and
 1358 critical background vehicles (CBVs), which interact dynamically over time.

1359 The visualizations demonstrate the ability of *RIFT* to generate temporally coherent, realistic, and
 1360 controllable trajectories across a variety of traffic situations. Even under complex and evolving
 1361 closed-loop conditions, *RIFT* maintains stable multimodal behavior, highlighting its effectiveness in
 1362 simulating realistic and controllable traffic scenarios around the AV.

1363 **E DISCUSSION AND BROADER IMPLICATIONS**

1364 **E.1 USE OF LARGE LANGUAGE MODELS (LLMs)**

1365 The large language model (LLM) was employed as a general-purpose writing assistant during the
 1366 preparation of this manuscript. Its use was limited to:

- 1367 • Language refinement: improving grammar, syntax, and overall readability to ensure clarity and
 1368 professionalism.
 1369 • Style adjustments: suggesting more concise and precise phrasing while preserving the original
 1370 meaning and technical content.

1371 The LLM was not involved in research ideation, experimental design, data collection, analysis, or
 1372 interpretation of results. All intellectual contributions and scientific conclusions are solely those of
 1373 the authors. This disclosure is provided in accordance with the conference guidelines on LLM usage.

1374 **E.2 LIMITATIONS AND FUTURE WORK.**

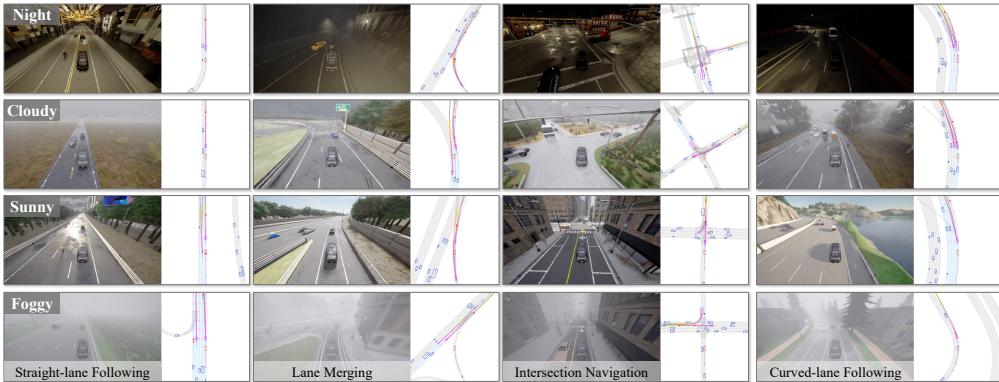
1375 In the current framework, the generation head is frozen during RL fine-tuning, and its reliability
 1376 stems from the robust trajectory generation capability learned during IL pre-training. However,
 1377 without reliable expert demonstrations, the realism and robustness of generated trajectories cannot be
 1378 further improved during RL fine-tuning. This limitation highlights a key avenue for future research:
 1379 developing methods—potentially leveraging RL or other self-improvement paradigms—that can
 1380 enhance the trajectory generation quality without relying on expert demonstrations.

1381 **E.3 SOCIAL IMPACT**

1382 **Positive Societal Impacts.** This work presents a practical framework that bridges the gap between
 1383 realism and controllability in traffic simulation. By decoupling pre-training and fine-tuning, our
 1384 method enables models pre-trained on real-world datasets to adapt effectively to physics-based
 1385 simulators, preserving trajectory-level realism and route-level controllability while improving long-
 1386 horizon closed-loop performance. This paradigm establishes a viable pathway for transitioning
 1387 data-driven approaches to physics-based simulators, enabling more reliable closed-loop testing and
 1388 training. Consequently, it advances safer and more robust autonomous systems.

1389 **Negative Societal Impacts.** While fine-tuning in physics-based simulators improves closed-loop
 1390 performance, it may also lead to overfitting to the specific characteristics of the simulator. As a
 1391 result, the learned policy could struggle to generalize beyond the simulated environment, giving rise
 1392 to a sim-to-real gap. This gap poses challenges for real-world deployment, as models that perform
 1393 well in simulation may not retain the same level of reliability when applied to actual autonomous
 1394 driving systems. Such discrepancies can affect the testing and training stages, highlighting the need
 1395 for further work to ensure real-world transferability.

1404



1405

1406

1407

1408

1409

1410

1411

1412

1413

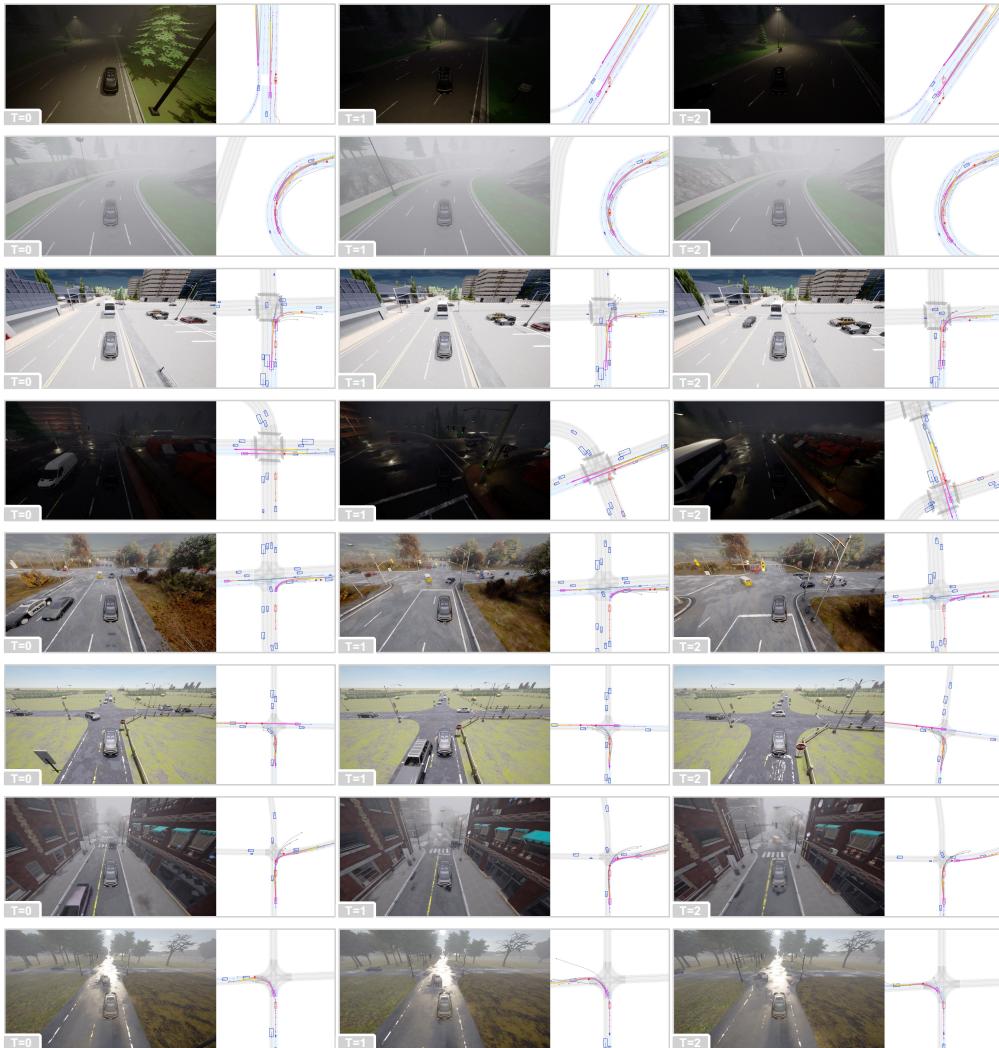
1414

1415

1416

1417

1418

(a) Robustness of *RIFT* across diverse AV-centric traffic scenarios.

1419

1420

1421

1422

1423

1424

1425

1426

1427

1428

1429

1430

1431

1432

1433

1434

1435

1436

1437

1438

1439

1440

1441

1442

1443

1444

1445

1446

1447

1448

1449

1450

1451

1452

1453

1454

1455

1456

1457

(b) Temporal stability of *RIFT* in closed-loop simulation.Figure 8: Visualizations of *RIFT* in diverse AV-centric scenarios. (a) Robustness of *RIFT* across diverse AV-centric traffic scenarios. (b) Temporal stability of *RIFT* in closed-loop simulation. CBV is marked in purple, AV (PDM-Lite) is in red, and BVs are in blue.