

Carnegie

Model-Based Policy Adaptation for Closed-Loop End-to-End Autonomous Driving







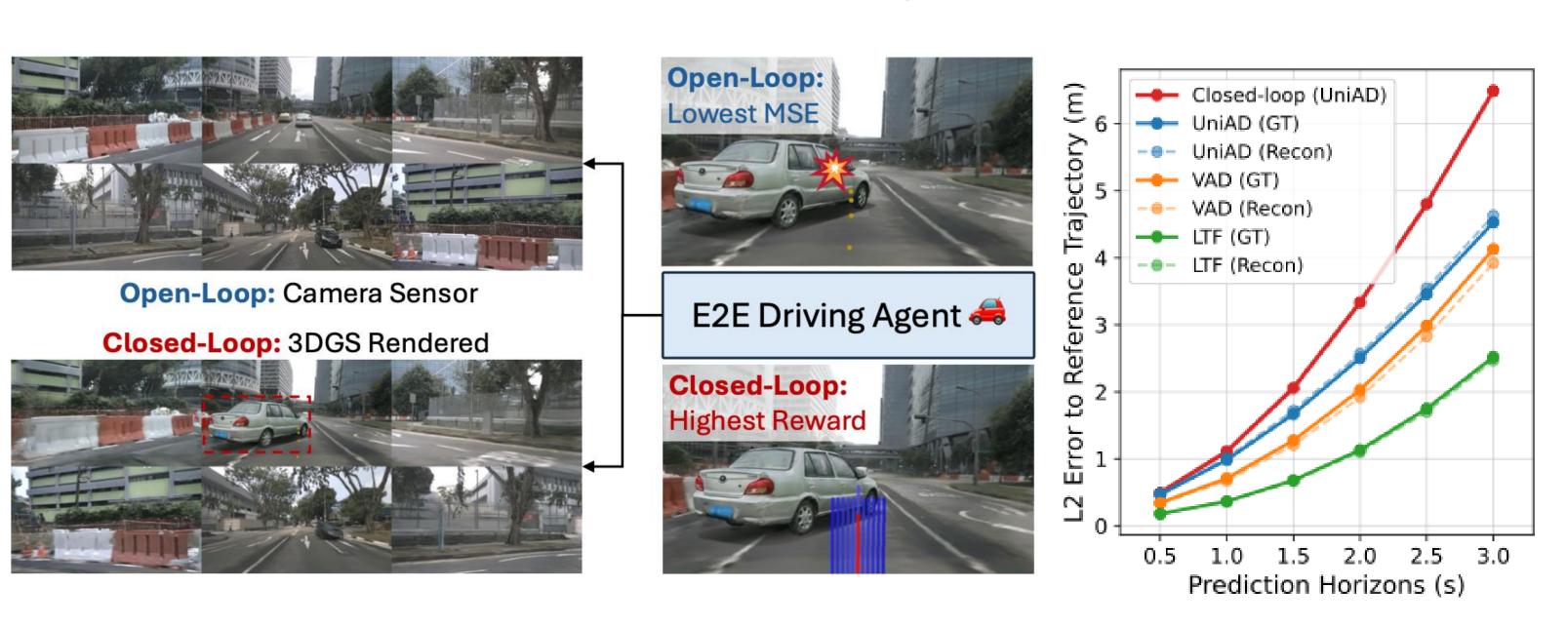


Motivation

Problem: Existing driving foundation models fall short in closed-loop evaluation, especially in the long-tailed cases.

Goal: Close the open- and closed-loop performance gap.

Approach: Use 3DGS-synthetic data to learn a policy adapter and a value model for closed-loop planning.



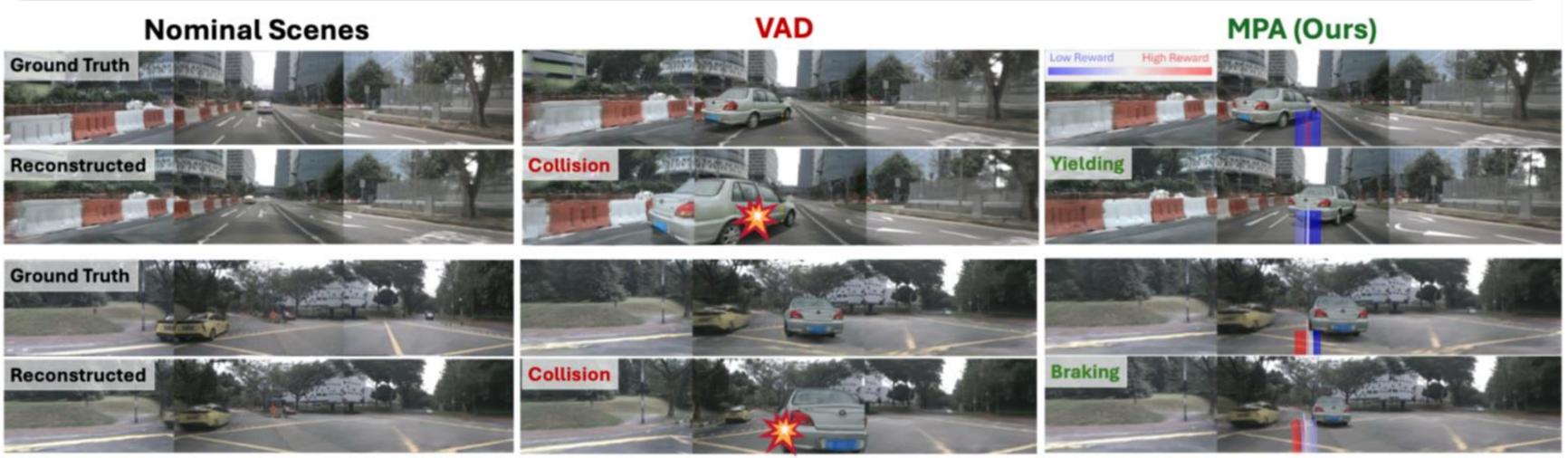
Cause of Performance Gap

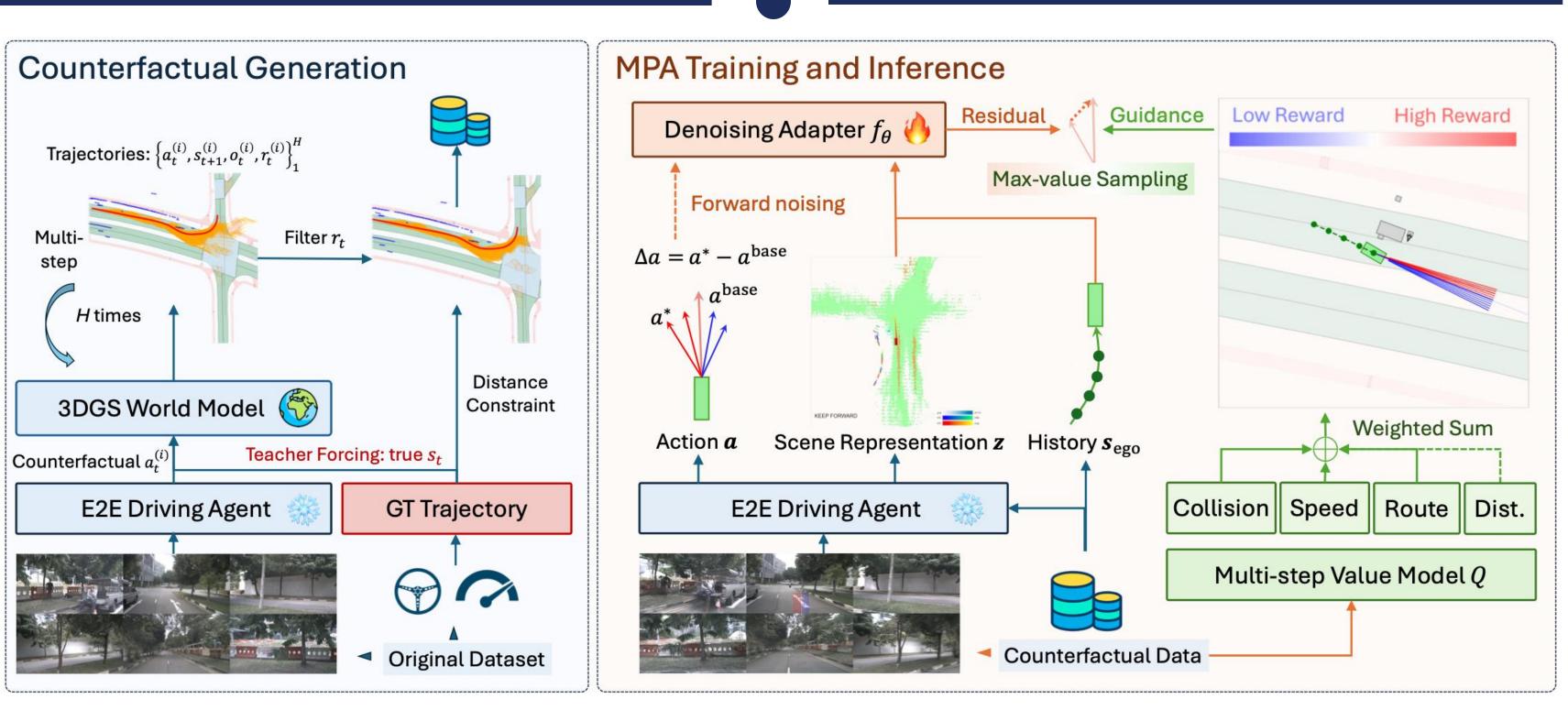
For closed-loop E2E driving, policy π^* , there would be:

- (1) Objective mismatch: covariate shift under closed-loop deployment.
- (2) Observation mismatch: between the simulated and real observations.

$$\begin{aligned} &\text{Open-loop}: \ \hat{\pi}^* = \arg\min_{\pi} \sum_{t=1}^{T} \mathbb{E}_{(s_t, a_t) \sim \pi_{\text{ref}}} \left\| a_t - \pi(s_t) \right\|_2^2, \\ &\text{Closed-loop}: \pi^* = \arg\max_{\pi} \sum_{t=1}^{T} \mathbb{E}_{s_t \sim P(s_{t-1}, a_{t-1}), \ a_{t-1} \sim \pi(o_{t-1}, s_{t-1}), \ o_{t-1} \sim P_{\text{obs}}(s_{t-1})} \left[r(s_t, a_t) \right], \\ &\text{Simulation}: \pi^* = \arg\max_{\pi} \sum_{t=1}^{T} \mathbb{E}_{s_t \sim P(s_{t-1}, a_{t-1}), \ a_{t-1} \sim \pi(o_{t-1}, s_{t-1}), \ o_{t-1} \sim \widehat{P}_{\text{obs}}(s_{t-1})} \left[r(s_t, a_t) \right]. \end{aligned}$$

Model-Based Policy Adaptation





What is MPA? Diffusion adapter + multi-principled Q head How to train MPA? With counterfactual rollouts by 3DGS.

$$\mathbb{E}_{\Delta a^{(0)},k,\epsilon} \min_{i} \left\| f_{\theta}(\Delta a^{(k)},k,z,\boldsymbol{s}_{\text{ego}},a^{\text{base}})[i] - \Delta a^{(0)} \right\|_{2}^{2}, \quad Q(o_{t},s_{t},a_{t};T) = \sum_{t=1}^{T} \gamma^{t} r(s,a_{t})$$
 where $\Delta a^{(k)} = \sqrt{\bar{\alpha}_{k}} \Delta a^{(0)} + \sqrt{1-\bar{\alpha}_{k}} \epsilon$, with $\epsilon \sim \mathcal{N}(0,\mathbf{I})$. $Q = \sum_{i \in \{\text{collision, dist, ...}\}} w_{i} \times Q_{i}$

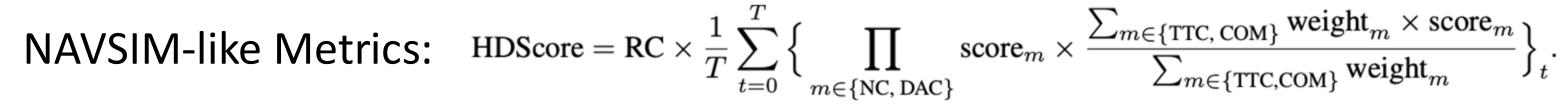
Closed-loop: $\pi^* = \underset{\pi}{\operatorname{arg\,max}} \sum_{t=1}^{\infty} \mathbb{E}_{s_t \sim P(s_{t-1}, a_{t-1}), \ a_{t-1} \sim \pi(o_{t-1}, s_{t-1}), \ o_{t-1} \sim P_{obs}(s_{t-1})} [r(s_t, a_t)]$, How to use MPA? Inference-time Scaling with Q-values

$$\Delta \hat{a}^* = rg \max_{\Delta a \in a^{ ext{adapt}}} Q(o_t, oldsymbol{s}_{ ext{ego}}, a^{ ext{base}} + \Delta a; T), \quad \hat{a}^* = a^{ ext{base}} + \Delta \hat{a}^*.$$

Experiments and Analysis

Closed-Loop Evaluation on the nuScenes dataset and HUGSIM simulator

290 Training Scenarios, 70 In-domain, 70 unseen, and 10 safety-critical scenarios

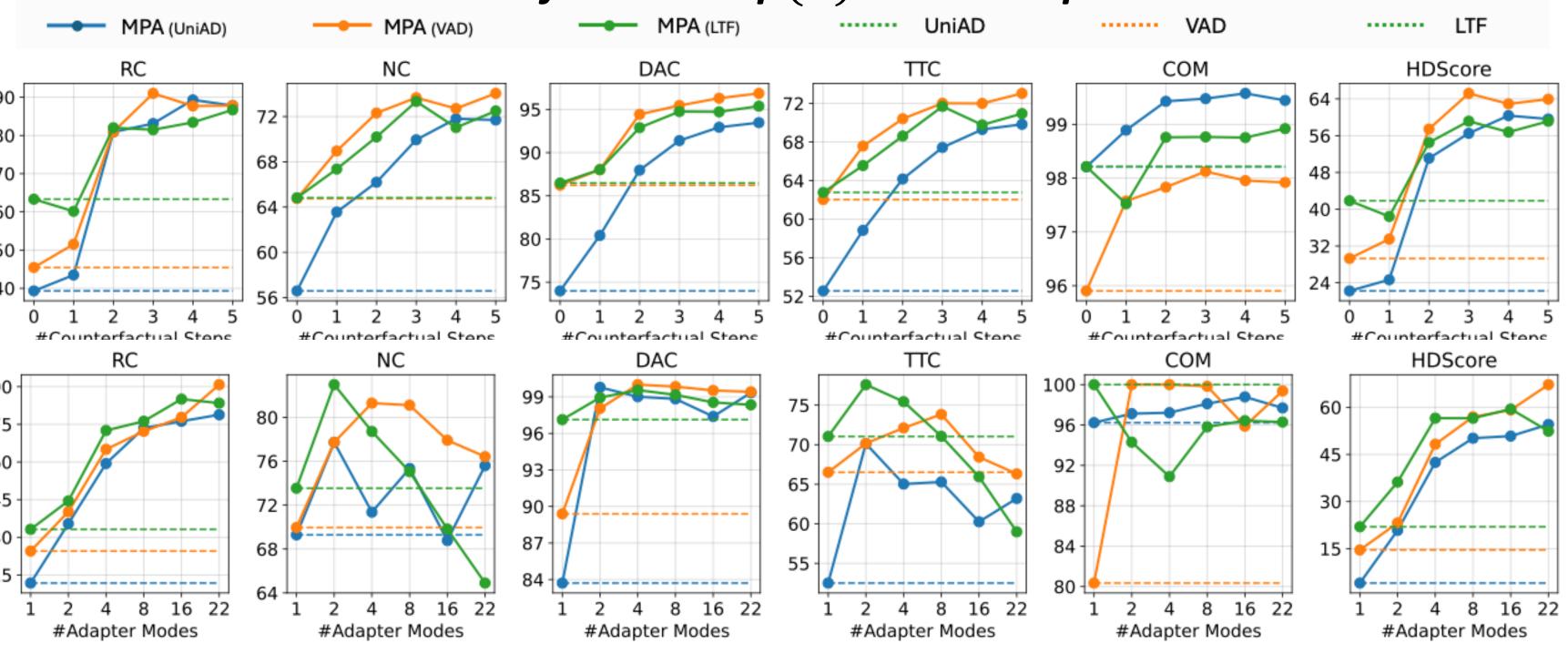


Model	Eg	o Statu	is C	amera	Curati	on RC	NC	DAG	\mathbf{C} \mathbf{T}'	TC	COM	HDScore	Observation: MPA
UniAD		\checkmark		√	X	39.4	56.9	75.1	1 52	2.1	98.7	19.4	
VAD		\checkmark		\checkmark	X	50.1	68.4	87.2	2 60	6.1	90.2	31.9	manages to
LTF		✓		✓	X	65.2	71.3	92.1	l 6'	7.6	<u>98.4</u>	46.7	outperform both the
AD-MLP		\checkmark		X	\checkmark	13.4	80.2	86.2	2 79	9.4	90.1	6.5	•
BC-Safe		\checkmark		\checkmark	\checkmark	57.0	59.8	87.9	5:	5.2	89.4	33.6	pretrained baselines
Diffusion		\checkmark		\checkmark	\checkmark	71.8	67.4	88.	l 64	4.5	91.5	45.1	and baselines solely
MPA (UniAD)		\checkmark		\checkmark	\checkmark	93.6	<u>76.4</u>	92.8	3 72	2.8	91.8	<u>66.4</u>	and baselines solely
MPA (VAD)		\checkmark		\checkmark	\checkmark	94.9	75.4	93.0	5 <u>7</u> 2	2.5	92.8	67.0	trained on the
MPA (LTF)		\checkmark		\checkmark	\checkmark	93.1	70.8	90.9	9 6'	7.9	94.9	60.0	
	Unseen Nominal Scenes Safety-Critical Scenes									counterfactual			
Model	RC	NC	DAC	TTC	COM	HDScore	RC	NC	DAC	TTC	COM	HDScore	dataset.
UniAD 3	39.3	56.6	74.0	52.6	98.2	22.2	11.4	76.2	82.1	57.8	95.9	4.5	The closed-loop
VAD 4	45.4	64.8	86.2	62.0	95.9	29.3	25.4	77.0	88.3	73.2	88.4	16.0	The closed-loop

 NC
 DAC
 TTC
 COM
 HDScore
 RC
 NC
 DAC
 TTC
 COM
 HDScore
 dataset.

 56.6
 74.0
 52.6
 98.2
 22.2
 11.4
 76.2
 82.1
 57.8
 95.9
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.5
 4.2
 4.9
 4.3
 4.9
 4.3
 4.9
 4.3
 4.9
 4.3
 4.9
 4.3
 4.9
 4.3
 4.9
 4.3
 4.9
 4.3
 4.9
 4.3
 4.3
 4.9
 4.3
 4.5
 4.3
 4.5
 4.3

Performance w.r.t. #Counterfactual Step(T) and #Sample Modes



Short Takeaways: Using <u>counterfactual</u> data to learn policy <u>adapter</u> and <u>Q-value</u> critic leads to better closed-loop performance in E2E autonomous driving!