

RoaD: Rollouts as Demonstrations for Closed-Loop Supervised Fine-Tuning of Autonomous Driving Policies

Guillermo Garcia-Cobo^{*1} Maximilian Igl^{*1} Peter Karkus^{*1} Zhejun Zhang^{*2†}

Michael Watson¹ Yuxiao Chen¹ Boris Ivanovic¹ Marco Pavone^{1,3}

¹NVIDIA Research ²Huawei VN Research Center ³Stanford University

{guillermog, migl, pkarkus}@nvidia.com

Abstract

Autonomous driving policies are typically trained via open-loop behavior cloning of human demonstrations. However, such policies suffer from covariate shift when deployed in closed loop, leading to compounding errors. We introduce Rollouts as Demonstrations (RoaD), a simple and efficient method to mitigate covariate shift by leveraging the policy’s own closed-loop rollouts as additional training data. During rollout generation, RoaD incorporates expert guidance to bias trajectories toward high-quality behavior, producing informative yet realistic demonstrations for fine-tuning. This approach enables robust closed-loop adaptation with orders of magnitude less data than reinforcement learning, and avoids restrictive assumptions of prior closed-loop supervised fine-tuning (CL-SFT) methods, allowing broader applications domains including end-to-end driving. We demonstrate the effectiveness of RoaD on WOSAC, a large-scale traffic simulation benchmark, where it performs similar or better than the prior CL-SFT method; and in AlpaSim, a high-fidelity neural reconstruction-based simulator for end-to-end driving, where it improves driving score by 41% and reduces collisions by 54%.

1. Introduction

Autonomous vehicle (AV) policies are typically trained with behavior cloning (BC) of human demonstrations, which is scalable but inherently open loop: it assumes i.i.d. inputs and optimizes one-step accuracy under the dataset distribution. Deployed in closed loop, policies influence their own observations, creating a train-test mismatch that induces covariate shift, compounds errors, and reduces robustness to long-tail and interactive scenarios.

On the other hand, End-to-end (E2E) policies are becoming

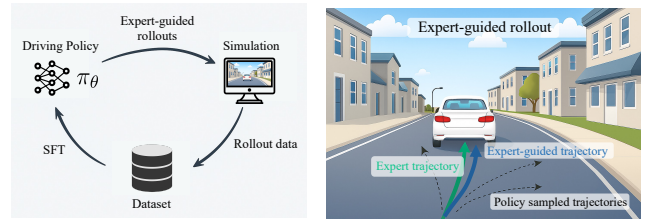


Figure 1. **RoaD closed-loop SFT with expert-guided rollout.** **Left:** Expert-guided rollouts in simulation yield additional training data that is incorporated into the dataset and used to fine-tune the policy, improving subsequent rollouts. **Right:** Expert-guided trajectories are sampled from the policy and biased toward the expert to produce higher-quality demonstrations.

ing the new norm of AV policy learning, which map sensor inputs directly to trajectories or controls. By coupling perception, prediction, and planning, they offer data efficiency, simpler deployment, and better long-horizon coordination than hand-engineered stacks [7, 12, 24, 28, 35, 36].

While RL directly optimizes closed-loop behavior, it remains impractical for end-to-end driving [17] due to brittle reward design and the cost of safe exploration and high-fidelity simulation. This leaves a gap for a scalable closed-loop training recipe for E2E driving that retains supervised simplicity and data efficiency.

Closed-loop supervised fine-tuning (CL-SFT) has recently emerged as a promising alternative to RL (see Fig. 1). The core idea of CL-SFT is to generate expert-biased on-policy rollouts in simulation and use them as additional demonstrations for fine-tuning, combining the simplicity of supervised learning with the benefits of closed-loop training. The key challenge is how to bias the rollouts towards high-quality behavior such that the fine-tuning step improves the policy. In traffic simulation, Closest Among Top-K (CAT-K) [43] instantiates this idea by selecting, at each step, the closest among a small set of policy-proposed candidate actions to the ground-truth trajectory. On these generated trajectories, CAT-K derives fine-tuning action tar-

^{*}Equal contribution, alphabetically sorted.

[†]Work performed during an internship at NVIDIA Research.

gets using an inverse dynamics model that chooses the action bringing the agent closest to ground truth.

While effective for traffic modeling, CAT-K is poorly suited to modern E2E policies because it assumes: (i) discrete actions; (ii) deterministic dynamics and known inverse dynamics; (iii) a diverse pretrained policy where at least one action sample lies close to the ground-truth trajectory at each time step; and (iv) that fresh on-policy trajectories can be generated continuously during training. In E2E driving, none of these assumptions usually hold: policies may output multi-token plans or continuous trajectories, as in the case of diffusion policies; the dynamics (including downstream controllers) are stochastic without closed-form inverse dynamics; the action distribution is typically less diverse due to safety- and, predictability-oriented training (unlike traffic agents that deliberately promote diversity); and regenerating closed-loop rollouts at every optimization step is prohibitively expensive.

In this work we introduce **Rollouts as Demonstrations (RoAD)**, a novel CL-SFT approach that addresses all limitations above (Fig. 2). First, RoAD retains the expert-biased rollout principle, but removes the need for a recovery action by treating the policy’s closed-loop rollouts directly as additional demonstrations for SFT. Empirically we find that this strategy achieves performance on par with, or better than CAT-K on large-scale traffic-simulation benchmarks. Second, we replace the Top-K selection with sampling K action candidates so that RoAD can be applied to a more general class of policies. Third, when limited action diversity prevents naive RoAD from being guided by the ground truth, we introduce a lightweight recovery-mode policy output that enables following the ground-truth trajectory even when it is not close to any of the top-k most likely actions. Finally, to reduce collection cost, we show that reusing rollout datasets across multiple optimization steps results in only marginal performance degradation, greatly improving data efficiency.

In experiments, we first validate RoAD for traffic simulation using the Waymo Open Simulation Agent Challenge (WOSAC). RoAD outperforms or matches CAT-K, even when it generates CL experiences only once with the base policy, and the performance improves further the more frequently the data is updated. We then apply RoAD to an E2E driving task to fine-tune a VLM-based policy deployed in AlpaSim [27], an E2E AV simulator that reconstructs real-world 3D scenes using SOTA 3D Gaussian splatting, 3DGUT [38]. RoAD fine-tuning improves driving scores by 41% and reduces collisions by 54% over the base model in previously unseen scenarios.

In summary, our conclusions are as follows.

- We introduce a novel CL-SFT algorithm, RoAD, that removes restrictive assumptions made by prior work.
- We apply RoAD to traffic simulation and match or outper-

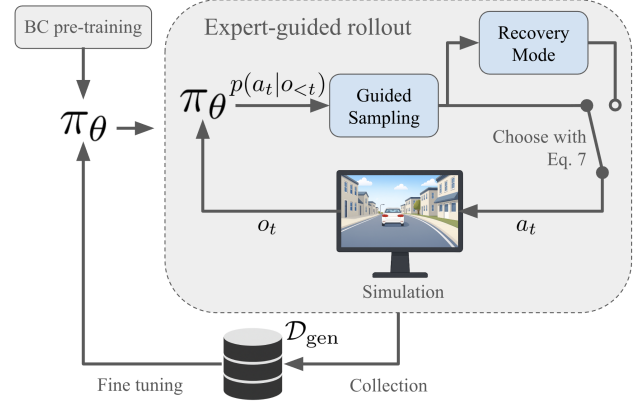


Figure 2. **RoAD method overview.** A pretrained policy generates rollouts in simulation via expert-guided sampling and an optional recovery mode. The collected trajectories form the CL-SFT dataset used to fine-tune the policy. The cycle can be repeated.

form the previous SOTA CL-SFT method.

- We apply RoAD to E2E driving and achieve substantial improvement in closed-loop driving metrics.

2. Related work

Closed-loop training in AV. Classic closed-loop training methods such as DAGger [31] and DART [18] iteratively collect states induced by the learner and query the expert to label them, thus adapting the training distribution to the learner’s roll-out distribution. However, in driving applications expert interventions are expensive, unsafe, or infeasible during interaction. Therefore, despite efforts on reducing expert burden in variants of the algorithm [15, 30, 34, 41], repeated expert relabeling remains impractical in autonomous driving.

Reinforcement learning presents a natural closed-loop training paradigm with a large body of literature in autonomous driving [13, 21, 33, 46]. From early works such as [32] to more recent studies applying hierarchical RL [6], curriculum-based RL [8], or model-based RL [40], they face two major challenges. (1) reward design that captures the diverse and often conflicting requirements; (2) training stability, compute cost, and safety constraints for E2E driving systems. These practical issues often lead to slow convergence, brittle behavior, and poor transfer to rare scenarios when applying RL to AV.

Inspired by recent work on CAT-K [43], our method performs closed-loop training *without* relying on reinforcement rewards or on-demand expert relabeling. It derives training targets from expert demonstrations and performs supervised updates in the closed-loop setting, thereby achieving the stability and data-efficiency benefits of supervised learning. Importantly, unlike CAT-K, our method does not require a discrete action space and deterministic dynamics, which ex-

tends its applicability to domains such as E2E driving.

End-to-end driving. Modern AV systems increasingly favor end-to-end designs over modular pipelines to reduce information loss and simplify training [16]. Work in this direction was pioneered by, among others, DAVE-2 [3]. Subsequent approaches, such as ChauffeurNet [2], UniAD [11], and PARA-Drive [37] introduced intermediate BEV representations, and integrated perception, prediction, and planning within unified multi-task networks.

In recent years, the emergent class of vision-language-action (VLA) models has adapted large vision-language models (VLMs) to driving, combining visual input, language instructions, and trajectory generation [12, 14, 28, 35, 44, 45]. These models promise richer semantic understanding and human-aligned decision making. In particular, systems like EMMA [12] and Alpamayo-R1 [28] demonstrate how VLM-based architectures can ingest camera imagery (and optionally language navigation cues) and output trajectories in a unified framework. While equipped with strong semantic understanding capabilities, these models are largely trained in open-loop, and thus remain prone to covariate shift.

3. Background

Our goal is to finetune a pretrained driving policy in closed-loop to minimize the covariate shift between open-loop pre-training and closed-loop deployment. To achieve this without access to a reward function (which is hard to define for this task), we use closed-loop supervised fine-tuning (CL-SFT). In the following, we first formalize the problem and review CAT-K [43], a recent CL-SFT instantiation.

3.1. Problem formulation

We are given a policy $\pi_\theta = p(a_t | o_{<t})$ that maps a history of observation inputs $o_{<t}$ to action outputs a_t . The policy is pre-trained with behavior cloning (BC), using a dataset of expert demonstrations $\mathcal{D} = \{(o_{0:T}^{E,i}, a_{0:T}^{E,i})\}_{i=1}^{|\mathcal{D}|}$. Our goal is to perform closed-loop finetuning of π_θ to minimize the covariate shift between open-loop training and closed-loop deployment. We assume access to a stochastic simulator that generates a next observation given action, previous observation and some internal state $\mathcal{P}(o_{t+1} | o_t, a_t; \cdot)$, but no access to an on-demand expert nor to a reward function.

The contents of the observation o_t are domain specific. In traffic simulation, o_t includes the positions, velocities, and orientations of all nearby agents, together with a vectorized map (lane markings, wait lines, etc.); in E2E driving, it includes the ego vehicle’s sensor inputs (e.g., multi-view camera images) and estimated egomotion (e.g., pose, steering angle, velocity, and acceleration). We denote by $s_t \in \mathbb{R}^{D_s}$ the pose of the controlled agent. For notational simplicity, we assume control of a single agent. Since Road

operates per agent, extending to multi-agent control, as in our traffic simulation experiments, is straightforward.

Unlike prior work, we make no strong assumptions on the structure of a_t : it may represent a state delta, as is common in traffic simulation; a continuous control signal, a waypoint, or a trajectory. Single-step control inputs are executed through forward dynamics, respecting vehicle motion constraints, while waypoints and trajectories are tracked by low-level controllers. Predicting T_{pred} -step trajectories is common because such *action chunking* encourages long-horizon reasoning and often improves accuracy with open-loop training. For notational simplicity, we refer to all these outputs uniformly as a_t and use $s_{t+1} = f(s_t, a_t)$ to denote the agent state evolution over time.

Further, prior work assumed a policy with discrete modes to select top K predictions, such as next-token-prediction (NTP) traffic simulation models that encode actions in a single token. In contrast, we only assume that π_θ can generate K independent action samples, allowing for modern E2E driving policies such as Transformers with simple Gaussian outputs, NTP models with multiple tokens per action, or diffusion and flow-matching policies [28].

3.2. Closed-loop supervised fine-tuning and CatK

CL-SFT adapts a pretrained policy by behavior cloning on states encountered under its own closed-loop rollouts, aligning the training distribution with deployment and mitigating covariate shift. CAT-K [43] provides a practical instantiation with two complementary components: (i) recovery supervision, which defines action targets that move the rollout back toward the expert trajectory at the visited on-policy rollout states; and (ii) expert-proximal rollouts using top K predictions, which bias action selection during rollouts to remain close to the expert, so that the recovery supervision remains valid. Intuitively, CAT-K learns “how to get back on track” while ensuring it never drifts too far from the track in the first place.

Formally, the algorithm assumes a tokenized model or a distribution with discrete modes, which can be generally written as $\pi(a_t | o_{<t}) = \sum_{m=1}^M \pi(a | m) \pi(m | o_{<t})$, where $M \in \mathbb{N}$ denotes the vocabulary size or number of modes, $\pi(m | o_{<t})$ the token prediction or mode-selection distribution, and $\pi(a | m)$ the action decoder or action distribution within each mode. In each rollout step, the algorithm selects the top K predictions, $\Xi_t = \text{top}^K[\pi_\theta(a | o_{<t})]$, where $\text{top}^K[\pi]$ represents finding the K most likely tokens/modes under $\pi(m | o_{<t})$ and decoding/sampling the associated action. To bias rollouts toward the expert, the action “closest” to the expert is selected,

$$a_t = \arg \min_{a \in \Xi_t} d(f(s_t, a), s_{t+1}^E), \quad (1)$$

where s_t is the current agent state, f are the deterministic dynamics, s_{t+1}^E is the next expert state, and $d(\cdot, \cdot)$ is a

Algorithm 1 RoaD

```

1: Input: policy  $\pi_\theta$ , dataset  $\mathcal{D}$ , candidate action set size  $K$ , number of
   rollouts  $N_{\text{roll}}$ , number of training steps  $N_{\text{train}}$ , recovery parameters
    $(\delta_{\text{rec}}, N_{\text{rec}})$ 
2: Initialize dataset  $\mathcal{D}_{\text{gen}} = \{\}$ 
3: for  $j = 0, \dots, N_{\text{roll}}$  do
4:   Start simulation with scenario  $s_{0:T}^E \sim \mathcal{D}$ 
5:   for  $t = 1, \dots, T - 1$  do
6:     Sample candidates (Eq. (4))
7:     Choose closest to expert (Eq. (5))
8:     if trigger (Eq. (7)) then
9:       Use recovery mode output  $a_t \leftarrow a'_t$  (Eq. (8))
10:    end if
11:    Step simulator  $o_{t+1} \sim \mathcal{P}(o_{t+1}|a_t, o_t; \cdot)$ 
12:  end for
13:  Add rollout to dataset  $\mathcal{D}_{\text{gen}} \leftarrow (o_{0:T}, a_{0:T})$ 
14: end for
15: for  $j = 0, \dots, N_{\text{train}}$  do
16:   Update  $\theta$  with  $(o_t, a_t) \sim \mathcal{D}_{\text{gen}}$  and the RoaD loss (Eq. (3))
17: end for

```

distance metric on states, e.g., a weighted ℓ_2 over position, heading, and speed. For each state, recovery actions are defined by projecting the expert continuation onto the action vocabulary,

$$\hat{a}_t = \arg \min_{a \in \{1, \dots, |M|\}} d(f(s_t, a), s_{t+1}^E), \quad (2)$$

and θ is updated with behavior cloning on the rollout states: $\mathcal{L}_{\text{BC}}(\theta) = -\frac{1}{NT} \sum_{t=0}^{T-1} \sum_{i=1}^N \log \pi_\theta(\hat{a}_t^i | o_{<t})$, where N represents the number of controlled agents.

CAT-K is highly effective in achieving this goal for traffic simulation, but is limited when applied to E2E driving, due to the assumptions of deterministic dynamics and known inverse dynamics to construct recovery action targets, and it’s reliance on single-step, discrete policies to efficiently compute the top-K operator.

4. Method

Our goal is a CL-SFT recipe that works with modern E2E driving policies and reduces the covariate shift between open-loop training and closed-loop deployment without requiring a reward function. Our proposed method, rollouts as demonstrations (RoaD), keeps CAT-K’s bias-toward-expert idea but removes its main constraints while remaining simple and data-efficient.

4.1. Rollouts as demonstrations (RoaD)

The key idea of RoaD is to treat the policy’s own expert-guided, closed-loop rollouts as additional supervision for fine-tuning. Formally, let $\mathcal{R}_{s_{0:T}^E}^{\mathcal{P}}[\pi_\theta]$ denote the expert-guided rollout operator for π_θ given the simulator \mathcal{P} and expert (GT) trajectory $s_{0:T}^E$. We accumulate generated rollouts in a dataset,

$$\mathcal{D}_{\text{gen}} = \left\{ (o_{0:T}, a_{0:T}) \mid o_{0:T}, a_{0:T} \sim \mathcal{R}_{s_{0:T}^E}^{\mathcal{P}}[\pi_\theta], s_{0:T}^E \sim \mathcal{D} \right\}$$

and fine-tune the policy by behavior cloning:

$$\mathcal{L}_{\text{RoaD}}(\theta) = - \sum_{(o^i, a^i) \in \mathcal{D}_{\text{gen}}} \sum_t \sum_{i=1}^N \log \pi_\theta(a_t^i | o_{<t}). \quad (3)$$

The expert guidance is designed to produce trajectories that are simultaneously near on-policy (i.e. sampled from π_θ), but also higher-quality than unassisted rollouts. Because this data is collected on-policy, it covers states the policy is likely to encounter, reducing covariate shift between training and deployment. In practice, $\mathcal{R}_{s_{0:T}^E}^{\mathcal{P}}[\pi_\theta]$ can be implemented by biasing the policy’s output toward the expert continuation, for example using Top- K selection (Eq. (1)) or Sample- K (see Sec. 4.2).

Crucially, compared to CAT-K, RoaD does not require the construction of target recovery actions which are challenging to construct under stochastic or non-invertible dynamics, and are often low-quality for policies that output future trajectories rather than single-step actions. Instead, it uses the future trajectory itself as the target. In the following, we discuss three further modifications which makes RoaD applicable to E2E driving.

4.2. Sample- K expert-guided rollouts

To preserve expert guidance without discrete top- K enumeration, we draw K action candidates from the current policy distribution (e.g., trajectory samples or diffusion/flow-matching draws),

$$\Xi_t = \{a_t^{(k)}\}_{k=1}^K \sim \text{i.i.d. } \pi_\theta(\cdot | o_{<t}) \quad (4)$$

and select the candidate closest to the expert continuation under a generalized distance metric (Eq. (6)):

$$a_t = \arg \min_{a \in \Xi_t} d^g(a, s_{t:T}^E). \quad (5)$$

This Sample- K relaxation maintains the “closest-to-expert” bias while accommodating continuous policy outputs and large vocabularies.

When a_t represent trajectories, the distance function $d(\cdot, \cdot)$ can be implemented as a trajectory-level (generalized) distance between predicted trajectories and the future expert trajectory. A concrete choice is a weighted step-wise distance over a comparison horizon H_t ,

$$d^g(a_t, s_{t:T}^E) = \sum_{k=1}^{H_t} w_k d(\tilde{s}_{t+k}(a_t), s_{t+k}^E), \quad (6)$$

where $\tilde{s}_{t+k}(a_t)$ denotes the predicted state at step $t+k$ implied by action a_t , and $w_k \geq 0$ are arbitrary weights.

4.3. Recovery-mode policy output

E2E driving policies often exhibit limited action diversity as they are trained to drive safely and predictably, preventing

WOSAC leaderboard, test split Method	# model params	RMM ↑	Kinematic metrics ↑	Interactive metrics ↑	Map-based metrics ↑	min ADE ↓
SMART-tiny RoaD (ours)	7 M	0.7847	0.4932	0.8106	0.9178	1.3042
SMART-tiny CAT-K [43]	7 M	0.7846	0.4931	0.8106	0.9177	1.3065
SMART-large [39]	102 M	0.7614	0.4786	0.8066	0.8648	1.3728
SMART-tiny [39]	7 M	0.7591	0.4759	0.8039	0.8632	1.4062

Table 1. **WOSAC leaderboard [20] for traffic simulation comparing CL-SFT approaches.** RMM stands for Realism Meta Metric, the key metric used for ranking. RoaD fine-tuning significantly improves over the base model (SMART-tiny), it outperforms a much larger model from the same model family (SMART-large), and it is on par with the SOTA CL-SFT method, CAT-K.

WOSAC local val. split Method	Data update frequency	RMM ↑	Kinematic metrics ↑	Interactive metrics ↑	Map-based metrics ↑	min ADE ↓
RoaD fine-tuning	always	0.7673	0.4871	0.8107	0.8715	1.3004
RoaD fine-tuning	every 2 epochs	0.7669	0.4865	0.8098	0.8720	1.2893
RoaD fine-tuning	one-off	0.7664	0.4865	0.8093	0.8712	1.2983
SMART-tiny base model	–	0.7653	0.4831	0.8081	0.8716	1.3240
CAT-K fine-tuning [43]	always	0.7616	0.4583	0.8105	0.8720	1.3105
SMART-tiny base model (from [43])	–	0.7581	0.4512	0.8076	0.8697	1.3152

Table 2. **Ablation of data collection frequency for traffic simulation, WOSAC 2% validation split.** RoaD fine-tuning leads to significant improvement even when closed-loop data is only collected once, achieving similar levels of improvements over the base model as through CAT-K fine-tuning. The more frequently the data is updated the larger the performance gain. Note that results for CAT-K were taken from [43], where likely a different SMART-tiny checkpoint was used as a base model.

naive sampling from reliably producing a candidate near the expert. To address this, we introduce an optional recovery-mode policy output that nudges the policy toward the expert when all sampled actions are too far from the expert.

Concretely, when the chosen action a_t is a trajectory, we linearly interpolate between a_t and the expert continuation, acting as guidance rather than a discrete override. Let the prediction horizon be F . We reuse the notation $\tilde{s}_{t+k}(a_t)$ from Eq. (6) to denote the predicted state at step $t+k$ implied by a_t . Recovery is triggered when the generalized distance to the expert exceeds a threshold:

$$d^g(a_t, s_{t:T}^E) > \delta_{\text{rec}}. \quad (7)$$

Upon triggering, we define a weight vector $\lambda \in [0, 1]^F$ (e.g., a linear schedule $\lambda_k = \min(1, k/N_{\text{rec}})$) and blend the trajectories as

$$\tilde{s}_{t+k}(a'_t) = (1 - \lambda_k) \tilde{s}_{t+k}(a_t) + \lambda_k s_{t+k}^E, \quad k = 1, \dots, F, \quad (8)$$

which defines the recovery trajectory a'_t . The weight vector λ subsumes all parameters of the schedule; in practice we use a simple linear ramp over the first N_{rec} steps.

4.4. CL-SFT with off-policy data

CAT-K regenerates rollouts at each gradient step, which is feasible in BEV traffic simulation, but prohibitive for E2E driving due to the high cost of rendering sensor inputs. To reduce this collection cost, we evaluate reusing the same rollout dataset across multiple optimization steps, similar to a replay buffer in off-policy RL, including the extreme case of generating only a single dataset at the start of

fine-tuning. Empirically, we find that rollout data reuse incurs only small degradation, making RoaD practical when high-fidelity rollouts are expensive to obtain.

5. Experiments

We validate our method for traffic simulation on the WOSAC benchmark [22], and for E2E driving using the AlpSim simulator [27] and the NVIDIA Physical AI - AV NuRec Dataset [26].

5.1. Traffic simulation

We first validate RoaD for traffic simulation using the WOSAC benchmark. Note that our primary goal here is not to outperform CAT-K, but to show that the simplified RoaD approach can achieve comparable performance to CAT-K while also being applicable to E2E driving due to fewer restrictive assumptions.

5.1.1. Experimental setup

We follow the experimental setup of [43] and use RoaD to fine-tune the SMART-tiny model on the WOMD dataset. The model receives 1 second of trajectory history for all agents, it outputs delta x-y actions, and at test time it is rolled out for 8 seconds with 0.1s time steps. Note that for this experiment we do not use the recovery mode as traffic models are naturally diverse enough.

Metrics. Evaluation follows the WOSAC protocol. For each scenario, we generate 32 rollouts for all agents in the scene and compare the resulting joint behavior distribution to human driven trajectories. We report the following metrics which, excluding minADE, measure the distri-

<i>AV NuRec dataset</i> Method	Driving score \uparrow	Collision rate \downarrow	Offroad rate \downarrow	Distance traveled (m) \uparrow
Fine-tuning with RoaD (ours)	0.6300 \pm 0.0090	0.0239 \pm 0.0000	0.2098 \pm 0.0029	147.2 \pm 0.54
Fine-tuning with re-rendered expert trajectories	0.4985 \pm 0.0046	0.0464 \pm 0.0051	0.2583 \pm 0.0044	151.9 \pm 0.36
Continued large-scale training with BC	0.4215 \pm 0.0092	0.0627 \pm 0.0033	0.2783 \pm 0.0039	143.7 \pm 0.76
Base model pre-trained with BC	0.4443 \pm 0.0210	0.0525 \pm 0.0051	0.2833 \pm 0.0084	149.0 \pm 1.08

Table 3. **End-to-end simulation results over the AV NuRec dataset.** RoaD fine-tuning significantly increases the driving score, and it outperforms both continued open-loop training with real data, as well as fine-tuning with re-rendered expert trajectories.

<i>AV NuRec dataset</i> Method	Driving score \uparrow	Collision rate \downarrow	Offroad rate \downarrow	Distance traveled (m) \uparrow
RoaD	0.6300 \pm 0.0090	0.0239 \pm 0.0000	0.2098 \pm 0.0029	147.2 \pm 0.54
RoaD (no expert guidance)	0.4847 \pm 0.0027	0.0576 \pm 0.0022	0.2543 \pm 0.0011	151.2 \pm 0.24
RoaD (no recovery)	0.5030 \pm 0.0107	0.0518 \pm 0.0044	0.2493 \pm 0.0013	151.4 \pm 0.42
RoaD (1 rollout)	0.5898 \pm 0.0108	0.0341 \pm 0.0006	0.2091 \pm 0.0045	143.4 \pm 0.29
RoaD (3 rollouts, default)	0.6300 \pm 0.0090	0.0239 \pm 0.0000	0.2098 \pm 0.0029	147.2 \pm 0.54
RoaD (9 rollouts)	0.6317 \pm 0.0070	0.0264 \pm 0.0027	0.2083 \pm 0.0054	148.3 \pm 0.41
RoaD (fine-tune once, default)	0.6300 \pm 0.0090	0.0239 \pm 0.0000	0.2098 \pm 0.0029	147.2 \pm 0.54
RoaD (fine-tune twice)	0.6613 \pm 0.0130	0.0420 \pm 0.0035	0.1967 \pm 0.0043	157.8 \pm 0.37
RoaD (1k steps)	0.6171 \pm 0.0124	0.0246 \pm 0.0035	0.2152 \pm 0.0011	148.0 \pm 0.66
RoaD (4.2k steps, default)	0.6300 \pm 0.0090	0.0239 \pm 0.0000	0.2098 \pm 0.0029	147.2 \pm 0.54
RoaD (10k steps)	0.6095 \pm 0.0245	0.0424 \pm 0.0029	0.2043 \pm 0.0066	150.2 \pm 0.15
RoaD (K=16)	0.5789 \pm 0.0039	0.0322 \pm 0.0023	0.2207 \pm 0.0029	146.4 \pm 0.18
RoaD (K=32)	0.5898 \pm 0.0060	0.0290 \pm 0.0006	0.2196 \pm 0.0033	146.6 \pm 0.16
RoaD (K=64, default)	0.6300 \pm 0.0090	0.0239 \pm 0.0000	0.2098 \pm 0.0029	147.2 \pm 0.54
RoaD (K=128)	0.6396 \pm 0.0119	0.0304 \pm 0.0039	0.2047 \pm 0.0045	150.4 \pm 0.10
Base model	0.4443 \pm 0.0210	0.0525 \pm 0.0051	0.2833 \pm 0.0084	149.0 \pm 1.08

Table 4. **Ablation study for E2E driving.** The setting is identical to the main experiments apart from the ablated property. *steps* refer to the number of optimization steps for fine-tuning. *K* denotes the number of trajectory samples for expert-guided rollouts. Results indicate that both expert guidance and recovery mode are important in the algorithm; and performance gains are observed over a wide range of hyperparameters.

butional similarity between the policy and the data. The principal metric on the leaderboard is **Realism Meta Metric (RMM)**, which combines three distributional metrics: **kinematic metrics**, e.g. velocities and accelerations; **interactive metrics**, e.g., collisions; and **map-based metrics**, e.g. off-road driving. For the exact definition of the metrics we refer to [23]. Additionally we also report **minADE**, i.e., minimum Average Displacement Error, a widely used metric for trajectory prediction.

5.1.2. Results

Main results on WOSAC. Tab. 1 provides results on the public WOSAC leaderboard. RoaD fine-tuning significantly improves over the base model (SMART-tiny), it outperforms a much larger model from the same model family (SMART-large), and it is on-par with the SOTA CL-SFT method, CAT-K [43]. We note that multiple works on the leaderboard, concurrently developed with ours, such as SMART-R1 [29] achieve higher RMM using a combination of CL-SFT with CAT-K, and RL fine-tuning. However, these approaches require highly specialized rewards derived from the WOSAC evaluation metrics, and a large number of

environment interactions, making them unsuitable for E2E driving. A snapshot of the complete leaderboard at the time of submission is included in the Appendix.

Re-using CL experience. To assess the effect of reusing previously generated CL-SFT data, we ablate the data refresh frequency in Tab. 2, evaluating locally on 2% of the WOMD validation set following [43]. As expected, more frequent refreshes yield higher performance, though the incremental gains are modest. Importantly, even when closed-loop data is generated only once at the start of fine-tuning, RoaD already delivers a substantial improvement. Given the high cost of data rendering, this motivates our default E2E setup (Sec. 5.2) of generating CL-SFT data only once. For completeness, we also ablate repeated data generation and observe additional, albeit smaller, improvements in the E2E experiment.

5.2. End-to-end driving

Our main result is that CL-SFT with RoaD can significantly improve closed-loop performance in E2E driving.

Local scene set	3D-GS	NeRF
Method	Driving score \uparrow	Driving score \uparrow
Road (ours)	0.75 ± 0.23	0.58 ± 0.09
Re-rendered expert trajectories	0.42 ± 0.07	0.35 ± 0.05
Base model	0.28 ± 0.05	0.33 ± 0.04

Table 5. **Sim2sim transfer results.** Policies are fine-tuned with 3DGS generated data, and evaluated in previously unseen 75 scenarios reconstructed either as 3DGS (default setting) or as a NeRF (sim2sim transfer). As expected, performance reduces when transferring fine-tuned policies to a new simulation environment, but fine-tuning with RoaD improves over the base model even in the transfer setting.

5.2.1. Experimental setup.

End-to-end VLA policy We employ a VLA-based policy structured similar to [28]. The policy takes in 1.6s ego motion history, and a sequence of timestamped images from two onboard cameras (front facing wide-angle and tele camera), and generates 6.4s trajectory sample output, which is then tracked by a downstream controller when executed in closed-loop. The policy is trained with a large-scale dataset comprising of 20,000 hours of human driving data from 25 countries, covering a variety of scenarios including highway and urban driving, weather conditions, day and night times. A 1700+ hour subset of this dataset is publicly available [25].

Simulation environment. For E2E driving experiments, we employ AlpaSim [27], a closed-loop simulator built on SOTA neural scene reconstruction [38]. The system reconstructs real-world driving logs as temporal 3D Gaussian Splatting (3D-GS) scenes and renders novel camera views when the ego vehicle diverges from the recorded path. We employ custom controllers that track predicted trajectories with separate lateral and longitudinal control, using a 200 ms control delay and ego-motion noise. The vehicle dynamics is governed by a dynamically extended bicycle model. The controller, control delay and dynamics model are designed to imitate real-world driving as closely as possible. All other traffic participants, including vehicles and pedestrians, replay their logged trajectories. Qualitative examples from AlpaSim [27] are shown in Fig. 4.

Fine-tuning. To fine-tune the VLA policy we generate 20s long simulated CL data using 8251 3D-GS scenes reconstructed from real-world driving logs, 3x rollouts per scene by default. We fine-tune for 4.2k steps (approximately one epoch of non-overlapping trajectory data) with frozen encoders to mitigate overfitting to visual artifacts.

Metrics To evaluate models, we use 920 openly accessible challenging 3D-GS scenes from the NVIDIA Physical AI - AV NuRec Dataset [26], and generate 3 rollouts per scenes. In Tabs. 3 to 5, mean values are computed over all scenes and rollouts, standard deviations are computed by taking the mean across scenes and computing the standard

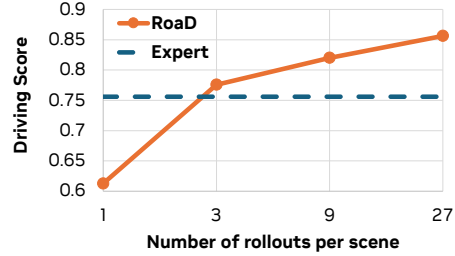


Figure 3. **The impact of generating multiple rollouts per scene.** In this experiment the same set of AV NuRec scenes are used for training and evaluation. RoaD performance improves monotonically as more rollouts are added to its SFT dataset (orange), while fine-tuning with resimulated expert demonstrations cannot make use of multiple rollouts.

deviation across rollouts, estimating the evaluation uncertainty. The standard deviation over the full RoaD training, including data generation, fine-tuning, and evaluation with three rollouts per scene, is too expensive to perform for every model. For the driving score of our main result in Tab. 3 we found it to be ± 0.0057 .

We report the following metrics. The **driving score** measures the average distance traveled (in kilometers) between incident events, where an incident corresponds to either off-road driving or a collision. The **collision rate** denotes the proportion of scenarios in which the ego vehicle is involved in a close encounter or collision for which it is deemed responsible, i.e., excluding rear-end and side contacts. The **off-road rate** captures the fraction of scenarios where the ego vehicle leaves the drivable area; this value appears relatively high because in the AV NuRec Dataset only the region bounded by lane markings is considered drivable. Finally, the **distance traveled** denotes the distance traveled by the ego vehicle in meters.

Each simulation terminates after the first incident. Evaluation in reconstructed scenes is inherently sensitive to rendering artifacts, particularly when the ego vehicle diverges from the logged path. To reduce the impact of such artifacts, we exclude any events where the ego deviates by more than 4 m from the original trajectory. Nonetheless, a portion of recorded incidents remain attributable to visual artifacts or imperfect scene reconstructions.

5.2.2. Results

Main results. The main results for E2E driving are reported in Tab. 3. RoaD fine-tuning increases the driving score in previously unseen scenarios by 41% and reduces collisions by 54%. RoaD outperform fine-tuning with expert demonstrations re-rendered in the same simulation environment, indicating that the performance gains of RoaD are not only from adjusting to simulation artifacts. RoaD also outperforms continued large-scale open-loop training using real-world driving data, indicating that the performance gains are not simply due to further training steps.

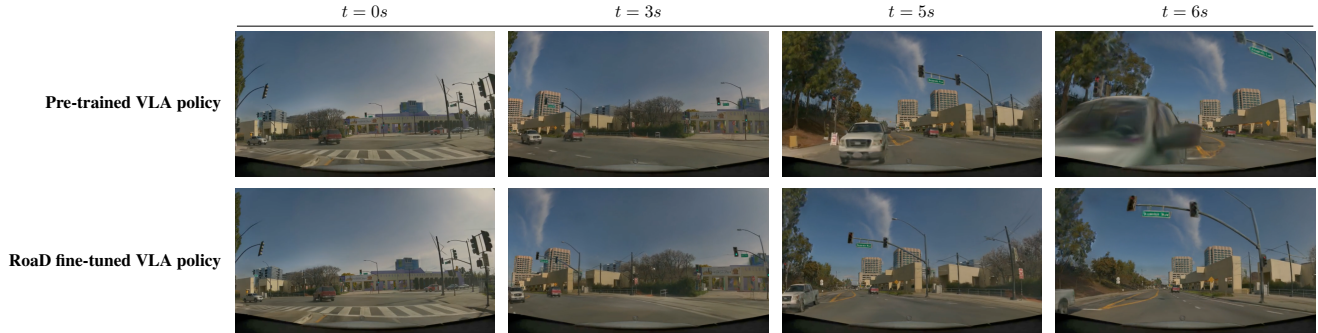


Figure 4. **Qualitative comparison of policy rollouts in our E2E simulator.** **Top row:** before fine-tuning, the policy navigates this intersection poorly, ends up in a wrong lane and fails to avoid a collision with a stationary vehicle. **Bottom row:** after fine-tuning with RoAD, the policy handles the intersection correctly and avoids any collision.

Fig. 4 shows a qualitative example: while the base policy encounters a collision, after fine-tuning with RoAD, the policy handles the intersection correctly without collision.

Ablation studies. Ablation results in Tab. 4 indicate that both expert guidance during rollouts and recovery-mode policy outputs are important for best performance in E2E driving. Furthermore, RoAD is not strongly sensitive to its hyperparameters, including the number of rollouts generated per scene, re-collecting CL data and further fine-tuning the policy, changing the number of optimization steps used for fine-tuning, or the number of trajectory samples (K) during CL data generation. In all alternative settings, RoAD improves upon the base model. While some hyper-parameter choices can further increase performance, in particular, increasing K and re-collecting CL data for additional fine-tuning, these also increase the computational costs of RoAD. On the other hand, we found that fine-tuning for too many optimization steps can slightly reduce performance, likely due to a lack of co-training with the original large-scale training data in our experiments.

Data scaling. Given the high cost of scene reconstruction for E2E CL-SFT (i.e. generating 3D-GS artifacts), scalability of RoAD with the number of rollouts per scene is an important question. To this end, in Fig. 3 we vary the number of rollouts generated per scene for CL-SFT. RoAD performance improves monotonically as more rollouts are added to its SFT dataset, while fine-tuning with re-simulated expert demonstrations cannot make use of multiple rollouts.

Sim2sim transfer. Finally, in our experiments so far we have generated CL fine-tuning data in the same simulation environment where the policy is evaluated. Given a gap between simulation and the real-world, the policy may overfit to artifact of the simulation and in turn it may degrade in real-world deployment. While addressing sim2real gap is not in the scope of this work, to shed some light on this issue, we perform a sim2sim transfer experiment, where the policies are fine-tuned with 3DGS generated data, and

evaluated in either 3DGS (default setting) or NeRF reconstructions (sim2sim transfer). For this experiment, we use an in-house scenarios set consisting of 75 scenarios curated for dense ego-agent interactions. Results are reported in Tab. 5. As expected, performance reduces when transferring fine-tuned policies to a new simulation environment, but fine-tuning with RoAD improves over the base model even in the transfer setting, indicating that RoAD has potential to improve real-world driving performance, despite possible sim2real gaps.

6. Conclusions

We presented RoAD, a simple closed-loop supervised fine-tuning (CL-SFT) method that treats the policy’s own expert-guided rollouts as additional demonstrations. By avoiding discrete recovery targets and introducing a lightweight recovery mode, RoAD removes key assumptions that limit prior CL-SFT approaches and makes the recipe applicable to modern E2E driving policies, allowing closed-loop training without the need for reward functions.

Across vectorized traffic simulation and high-fidelity E2E driving, RoAD consistently improves closed-loop performance over strong baselines. Because RoAD can achieve substantial improvements even when closed-loop data is only collected once, it is a promising, data-efficient, approach for training E2E driving policies in closed-loop.

Limitation of all CL-SFT approaches include the reliance of a pre-trained policy with sufficiently high performance, the assumption that the expert trajectory remains good behavior despite small deviations by the actor, and a distance metric for expert-guided rollouts. Further, our method relies on a high-fidelity simulator such as AlpaSim. Results on sim2sim transfer suggest that RoAD has potential to improve real-world driving performance, despite possible sim2real gaps. Future work may more explicitly address sim-to-real transfer and reduce overfitting to simulation, e.g., by co-training on simulated and real images, or introducing feature similarity bottlenecks [9].

References

- [1] Ehsan Ahmadi and Hunter Schofield. Rlftsim: Multi-agent traffic simulation via reinforcement learning fine-tuning. Technical report, 2025. Technical Report. 11
- [2] Mayank Bansal, Alex Krizhevsky, and Abhijit Ogale. Chauffeurnet: Learning to drive by imitating the best and synthesizing the worst. In *Robotics: Science and Systems (RSS) Conference*, 2019. 3
- [3] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon Goyal, Lawrence D. Jackel, Mathew Monfort, Urs Müller, Jiakai Zhang, Xin Zhang, Jake Zhao, and Karol Zieba. End to end learning for self-driving cars. *arXiv preprint arXiv:1604.07316*, 2016. Preprint. 3
- [4] Holger Caesar, Juraj Kabzan, Kok Seang Tan, Whye Kit Fong, Eric M. Wolff, Alex H. Lang, Luke Fletcher, Oscar Beijbom, and Sammy Omari. NuPlan: A closed-loop ml-based planning benchmark for autonomous vehicles. In *Conference on Computer Vision and Pattern Recognition (CVPR) ADP3 Workshop*, 2021. 11
- [5] Wei Cao, Marcel Hallgarten, Tianyu Li, Daniel Dauner, Xunjiang Gu, Caojun Wang, Yakov Miron, Marco Aiello, Hongyang Li, Igor Gilitschenski, Boris Ivanovic, Marco Pavone, Andreas Geiger, and Kashyap Chitta. Pseudo-simulation for autonomous driving. In *Conference on Robot Learning (CoRL)*, 2025. 11
- [6] Jianyu Chen, Shengbo Eben Li, and Masayoshi Tomizuka. Interpretable end-to-end urban autonomous driving with latent deep reinforcement learning. *IEEE Transactions on Intelligent Transportation Systems*, 23(6):5068–5078, 2021. 2
- [7] Li Chen, Penghao Wu, Kashyap Chitta, Bernhard Jaeger, Andreas Geiger, and Hongyang Li. End-to-end autonomous driving: Challenges and frontiers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. 1
- [8] Felipe Codevilla, Eder Santana, Antonio M López, and Adrien Gaidon. Exploring the limitations of behavior cloning for autonomous driving. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9329–9338, 2019. 2
- [9] Lan Feng, Yang Gao, Eloi Zablocki, Quanyi Li, Wuyang Li, Sichao Liu, Matthieu Cord, and Alexandre Alahi. Rap: 3d rasterization augmented end-to-end planning. *arXiv preprint arXiv:2510.04333*, 2025. 8
- [10] Ke Guo, Haochen Liu, Xiaojun Wu, and Chen Lv. Decomp-gail: Learning realistic traffic behaviors with decomposed multi-agent generative adversarial imitation learning, 2025. 11
- [11] Yihan Hu, Jiazhi Yang, Li Chen, Keyu Li, Chonghao Sima, Xizhou Zhu, Siqi Chai, Senyao Du, Tianwei Lin, Wenhai Wang, Lewei Lu, Xiaosong Jia, Qiang Liu, Yu Qiao, and Hongyang Li. Planning-oriented autonomous driving (uniad). In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, page –, 2023. 3
- [12] Jyh-Jing Hwang, Runsheng Xu, Hubert Lin, Wei-Chih Hung, Jingwei Ji, Kristy Choi, Di Huang, Tong He, Paul Covington, Benjamin Sapp, et al. EMMA: End-to-end multimodal model for autonomous driving. *arXiv preprint arXiv:2410.23262*, 2024. 1, 3
- [13] David Isele, Reza Rahimi, Akansel Cosgun, Kaushik Subramanian, and Kikuo Fujimura. Navigating occluded intersections with autonomous vehicles using deep reinforcement learning. In *2018 IEEE international conference on robotics and automation (ICRA)*, pages 2034–2039. IEEE, 2018. 2
- [14] Sicong Jiang, Zilin Huang, Kangan Qian, Ziang Luo, Tianze Zhu, Yang Zhong, Yihong Tang, Menglin Kong, Yunlong Wang, Siwen Jiao, Hao Ye, Zihao Sheng, Xin Zhao, Tuopu Wen, Zheng Fu, Sikai Chen, Kun Jiang, Diange Yang, Seongjin Choi, and Lijun Sun. A survey on vision-language-action models for autonomous driving. *arXiv preprint arXiv:2506.24044*, 2025. Preprint. 3
- [15] Gregory Kahn, Adam Villafior, Bosen Ding, Pieter Abbeel, and Sergey Levine. Self-supervised deep reinforcement learning with generalized computation graphs for robot navigation. In *International Conference on Robotics and Automation*, pages 1–8, 2018. 2
- [16] Peter Karkus, Boris Ivanovic, Shie Mannor, and Marco Pavone. Diffstack: A differentiable and modular control stack for autonomous vehicles. In *6th Annual Conference on Robot Learning*, 2022. 3
- [17] B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. *IEEE transactions on intelligent transportation systems*, 23(6):4909–4926, 2021. 1
- [18] Michael Laskey, Jonathan Lee, Roy Fox, Anca Dragan, and Ken Goldberg. Dart: Noise injection for robust imitation learning. In *Conference on robot learning*, pages 143–156. PMLR, 2017. 2
- [19] Longzhong Lin, Xuewu Lin, Kechun Xu, Haojian Lu, Lichao Huang, Rong Xiong, and Yue Wang. Revisit mixture models for multi-agent simulation: Experimental study within a unified framework, 2025. 11
- [20] Waymo LLC. Waymo open sim agent challenge (wosac) 2024 leaderboard. <https://waymo.com/open/challenges/2024/sim-agents/>, 2024. Accessed: 2025-03-14. 5
- [21] Xiaobai Ma, Jiachen Li, Mykel J Kochenderfer, David Isele, and Kikuo Fujimura. Reinforcement learning for autonomous driving with latent state inference and spatial-temporal relationships. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6064–6071. IEEE, 2021. 2
- [22] Nico Montali, John Lambert, Paul Mougin, Alex Kuefler, Nicholas Rhinehart, Michelle Li, Cole Gulino, Tristan Emrich, Zoey Yang, Shimon Whiteson, et al. The waymo open sim agents challenge. *Advances in Neural Information Processing Systems*, 36:59151–59171, 2023. 5
- [23] Nico Montali, John Lambert, Paul Mougin, Alex Kuefler, Nicholas Rhinehart, Michelle Li, Cole Gulino, Tristan Emrich, Zoey Yang, Shimon Whiteson, et al. The waymo open sim agents challenge. *Advances in Neural Information Processing Systems (NeurIPS)*, 36, 2024. 6
- [24] NVIDIA. Autonomous vehicles safety report, 2025. Available at <https://images.nvidia.com/aem->

dam/en-zz/Solutions/auto-self-driving-safety-report.pdf. 1

- [25] NVIDIA. Physical AI autonomous vehicles dataset. <https://huggingface.co/datasets/nvidia/PhysicalAI-Autonomous-Vehicles>, 2025. 7
- [26] NVIDIA. Physical AI autonomous vehicles NuRec dataset. <https://huggingface.co/datasets/nvidia/PhysicalAI-Autonomous-Vehicles-NuRec>, 2025. 5, 7
- [27] NVIDIA, :, Yulong Cao, Riccardo de Lutio, Sanja Fidler, Guillermo Garcia-Cobo, Zan Gojcic, Maximilian Igl, Boris Ivanovic, Peter Karkus, Janick Martinez, Marco Pavone, Aaron Smith, Michal Tyszkiewicz, Michael Watson, Qi Wu, and Le Zhang. Alpasim: A modular, lightweight, and data-driven research simulator for end-to-end autonomous driving, 2025. 2, 5, 7, 11
- [28] NVIDIA, Yan Wang, Wenjie Luo, Junjie Bai, Yulong Cao, Tong Che, Ke Chen, Yuxiao Chen, Jenna Diamond, Yifan Ding, Wenhao Ding, Liang Feng, Greg Heinrich, Jack Huang, Peter Karkus, Boyi Li, Pinyi Li, Tsung-Yi Lin, Dongran Liu, Ming-Yu Liu, Langechuan Liu, Zhijian Liu, Jason Lu, Yunxiang Mao, Pavlo Molchanov, Lindsey Pavao, Zhenghao Peng, Mike Ranzinger, Ed Schmerling, Shida Shen, Yunfei Shi, Sarah Tariq, Ran Tian, Tilman Wekel, Xinshuo Weng, Tianjun Xiao, Eric Yang, Xiaodong Yang, Yurong You, Xiaohui Zeng, Wenyuan Zhang, Boris Ivanovic, and Marco Pavone. Alpamayo-R1: Bridging reasoning and action prediction for generalizable autonomous driving in the long tail. *arXiv preprint arXiv:2511.00088*, 2025. 1, 3, 7, 11
- [29] Muleilan Pei, Shaoshuai Shi, and Shaojie Shen. Advancing multi-agent traffic simulation via rl-style reinforcement fine-tuning. *arXiv preprint arXiv:2509.23993*, 2025. 6, 11
- [30] Stephane Ross and J. Andrew Bagnell. Reinforcement and imitation learning via interactive no-regret learning. In *Proceedings of the 17th International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2014. 2
- [31] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 627–635. JMLR Workshop and Conference Proceedings, 2011. 2
- [32] Ahmad EL Sallab, Mohammed Abdou, Etienne Perot, and Senthil Yogamani. Deep reinforcement learning framework for autonomous driving. *Electronic Imaging*, 2017(19):70–76, 2017. 2
- [33] Dhruv Mauria Saxena, Sangjae Bae, Alireza Nakhaei, Kikuo Fujimura, and Maxim Likhachev. Driving in dense traffic with model-free reinforcement learning. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5385–5392. IEEE, 2020. 2
- [34] Wen Sun, Arun Venkatraman, Geoffrey J. Gordon, Byron Boots, and J. Andrew Bagnell. Deeply AggreVaTeD: Differentiable imitation learning for sequential prediction. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, pages 3309–3318. PMLR, 2017. 2
- [35] Xiaoyu Tian, Junru Gu, Bailin Li, Yicheng Liu, Yang Wang, Zhiyong Zhao, Kun Zhan, Peng Jia, Xianpeng Lang, and Hang Zhao. Drivevlm: The convergence of autonomous driving and large vision-language models. *arXiv preprint arXiv:2402.12289*, 2024. 1, 3
- [36] Wayve. Pioneering the end-to-end ai driving model, 2025. Available at <https://wayve.ai/technology/#AV2.0>. 1
- [37] Xinshuo Weng, Boris Ivanovic, Yan Wang, Yue Wang, and Marco Pavone. Para-drive: Parallelized architecture for real-time autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15449–15458, 2024. 3
- [38] Qi Wu, Janick Martinez Esturo, Ashkan Mirzaei, Nicolas Moenne-Loccoz, and Zan Gojcic. 3dgt: Enabling distorted cameras and secondary rays in gaussian splatting. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 26036–26046, 2025. 2, 7
- [39] Wei Wu, Xiaoxin Feng, Ziyang Gao, and Yuheng Kan. Smart: Scalable multi-agent real-time simulation via next-token prediction. *Advances in Neural Information Processing Systems (NeurIPS)*, 2025. 5
- [40] Zhuo Xu, Jianyu Chen, and Masayoshi Tomizuka. Guided policy search model-based reinforcement learning for urban autonomous driving. *arXiv preprint arXiv:2005.03076*, 2020. 2
- [41] Jiakai Zhang and Kyunghyun Cho. Query-efficient imitation learning for end-to-end autonomous driving. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2017. 2
- [42] Zhiyuan Zhang, Xiaosong Jia, Guanyu Chen, Qifeng Li, and Junchi Yan. Trajtok: Technical report for 2025 waymo open sim agents challenge, 2025. 11
- [43] 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. 1, 2, 3, 5, 6
- [44] Xingcheng Zhou, Mingyu Liu, Ekim Yurtsever, Bare Luka Zagar, Walter Zimmer, Hu Cao, and Alois C. Knoll. Vision language models in autonomous driving: A survey and outlook. *IEEE Transactions on Intelligent Vehicles*, PP(99): 1–20, 2024. 3
- [45] Zewei Zhou, Tianhui Cai, Seth Z. Zhao, Yun Zhang, Zhiyu Huang, Bolei Zhou, and Jiaqi Ma. Autovla: A vision-language-action model for end-to-end autonomous driving with adaptive reasoning and reinforcement fine-tuning. In *arXiv preprint arXiv:2506.13757*, 2025. 3
- [46] Meixin Zhu, Yinhai Wang, Ziyuan Pu, Jingyun Hu, Xuesong Wang, and Ruimin Ke. Safe, efficient, and comfortable velocity control based on reinforcement learning for autonomous driving. *Transportation Research Part C: Emerging Technologies*, 117:102662, 2020. 2

RoaD: Rollouts as Demonstrations for Closed-Loop Supervised Fine-Tuning of Autonomous Driving Policies

Supplementary Material

A. Leaderboard Snapshot

In Fig. 5 we show a snapshot of the nuPlan WOSAC leaderboard, with our *SMART-tiny-CLSFT-RoaD* entry highlighted (red box). We briefly discuss the other high-performing methods to clarify how they are either orthogonal to our approach or specialized to the WOSAC task, and thus do not directly transfer to E2E driving.

SMART-tiny-DecompGAIL [10], *SMART-tiny-RLFTSim* [1], and *SMART-RI* [29] use reinforcement learning (RL) to optimize policies for the WOSAC task of matching the data distribution. *SMART-tiny-DecompGAIL* does so by using GAIL with PPO, while *SMART-tiny-RLFTSim* and *SMART-RI* directly optimize the WOSAC metric using RL (with small differences in implementation). However, these methods do not translate well to E2E driving, where we typically lack a well-defined reward function and RL tends to be too data-inefficient given the high cost of high-fidelity simulation.

By contrast, *TrajTok* [42] proposes an improved tokenizer for traffic models. This contribution is orthogonal to our approach. However, it is only applicable to tokenizing short actions (e.g., one-step actions), and hence does not translate to E2E driving, where policies typically predict multiple seconds into the future.

Finally, *unimotion* [19] proposes an alternative to the CatK rollout approach, whereby it finds the closest action to the ground truth not only among the top- K actions, but among all actions. As a result, it does not require additional recovery actions (since it already tracks the ground-truth trajectory as closely as possible), which yields a setup more similar to our RoaD approach, where rollouts are taken directly as demonstrations. However, this exhaustive search makes the method unsuitable for E2E driving: it requires predicting and evaluating *all* possible actions of the policy, which is only feasible for small action spaces. This excludes multi-token trajectory predictions, whose action space grows exponentially with the horizon length, and flow-matching policies, whose action space is continuous. In their work, focussing on traffic models, they use either a fixed set of up to 2024 actions or a set of up to 16 actions predicted from action queries.

B. Experimental details: end-to-end driving

B.1. Simulation

For data generation, we run AlpaSim [27] at 30 Hz to match the frequency of the ground-truth logs and reuse the same

dataloader as for pre-training. For evaluation, we run the simulator at 10 Hz, which matches the model’s training frequency.

At each step, we render two camera views (front-facing wide-angle and telephoto). The policy predicts 6.4s trajectories, which are tracked by a downstream controller. As in the main text, the controller models a 200 ms control delay and ego-motion noise, and uses a dynamically extended bicycle model for the ego-vehicle dynamics.

Scenes are reconstructed using 3D Gaussian Splatting (3D-GS). Reconstruction quality degrades with distance from the recorded trajectory. This is negligible for rollout generation, where expert-guided rollouts remain close to the log, but can reduce visual fidelity during evaluation if the ego vehicle deviates too far. To avoid such artifacts, we discard rollouts whose ego trajectory deviates by more than 4 m from the recorded trajectory.

All other traffic participants, including vehicles and pedestrians, replay their logged trajectories and do not react to the ego vehicle. Consequently, they cannot avoid rear-end collisions if the ego drives more slowly than in the recording. Following prior work [4, 5, 28], we therefore count only “at-fault” collisions for the ego: rear-end collisions caused by following vehicles are ignored, while lateral collisions are still included.

B.2. RoaD fine-tuning

For RoaD rollout generation, we sample $K=64$ candidate trajectories from the policy at a temperature of 0.8. To select the executed trajectory, we compute the ground-truth distance d^g as the average distance between the four corners of the ego vehicle along the predicted and ground-truth trajectories over the first 20 time steps (2s). The same distance metric is used to decide whether to trigger the recovery mode, with a threshold of $\delta_{\text{rec}}=3$ m. When recovery is triggered, we linearly interpolate between the predicted and ground-truth trajectories over $N_{\text{rec}}=30$ time steps and follow the ground-truth trajectory thereafter. Recovery is disabled in the last 4s of the episode because the controller requires at least 4s of input trajectory.

Method Name	Realism Meta metric ↓	Kinematic metrics	Interactive metrics	Map-based metrics	minADE	Date (Pacific Daylight Time)
SMART-tiny-DecompGAIL	0.7864	0.4919	0.8152	0.9176	1.4209	9/22/2025, 6:33:53 AM
SMART-R1	0.7858	0.4944	0.8110	0.9201	1.2885	5/22/2025, 7:03:10 PM
SMART-tiny-RLFTSim	0.7857	0.4927	0.8129	0.9183	1.3252	7/3/2025, 8:56:51 PM
SMART-R1	0.7855	0.4940	0.8109	0.9194	1.2990	5/22/2025, 11:20:11 AM
TrajTok	0.7852	0.4887	0.8116	0.9207	1.3179	5/21/2025, 1:45:05 PM
unimotion	0.7851	0.4943	0.8105	0.9187	1.3036	5/8/2025, 1:52:07 PM
SMART-tiny-CLSFT-Road	0.7847	0.4932	0.8106	0.9178	1.3042	6/17/2025, 3:59:51 PM
SMART-tiny-CLSFT	0.7846	0.4931	0.8106	0.9177	1.3065	4/11/2025, 7:55:44 PM
MDG	0.7844	0.4928	0.8099	0.9183	1.3123	11/12/2025, 10:18:36 PM
SMART-tiny-RLFTSim	0.7844	0.4893	0.8128	0.9164	1.3470	5/15/2025, 7:38:00 AM
MDG	0.7842	0.4913	0.8102	0.9182	1.3074	8/12/2025, 9:11:43 PM
R1Sim	0.7839	0.4913	0.8107	0.9168	1.3421	6/13/2025, 5:08:01 AM
SMART-clsttlocal	0.7839	0.4909	0.8105	0.9170	1.3347	7/8/2025, 4:23:04 AM
ntu	0.7839	0.4901	0.8128	0.9145	1.4490	8/19/2025, 2:32:51 PM
R1Sim	0.7839	0.4916	0.8106	0.9166	1.3430	6/3/2025, 3:59:12 AM

Figure 5. Snapshot of the WOSAC leaderboard with our SMART-tiny-CLSFT-Road entry highlighted (red box).