

RoboTron-Sim: Improving Real-World Driving via Simulated Hard-Case

Baihui Xiao^{1†} Chengjian Feng^{1†} Zhijian Huang^{1,2}
Feng Yan¹ Yujie Zhong¹ Lin Ma^{1‡}

¹Meituan ²Shenzhen Campus of Sun Yat-sen University
<https://stars79689.github.io/RoboTron-Sim>

Abstract

Collecting real-world data for rare high-risk scenarios, long-tailed driving events, and complex interactions remains challenging, leading to poor performance of existing autonomous driving systems in these critical situations. In this paper, we propose **RoboTron-Sim** that improves real-world driving in critical situations by utilizing simulated hard cases. First, we develop a simulated dataset called *Hard-case Augmented Synthetic Scenarios (HASS)*, which covers 13 high-risk edge-case categories, as well as balanced environmental conditions such as day/night and sunny/rainy. Second, we introduce *Scenario-aware Prompt Engineering (SPE)* and an *Image-to-Ego Encoder (I2E Encoder)* to enable multimodal large language models to effectively learn real-world challenging driving skills from HASS, via adapting to environmental deviations and hardware differences between real-world and simulated scenarios. Extensive experiments on nuScenes show that RoboTron-Sim improves driving performance in challenging scenarios by $\sim 50\%$, achieving state-of-the-art results in real-world open-loop planning. Qualitative results further demonstrate the effectiveness of RoboTron-Sim in better managing rare high-risk driving scenarios.

1. Introduction

The landscape of autonomous driving (AD) has seen great breakthroughs in recent years, driven primarily by the development of end-to-end AD systems [1, 20–24, 27–29]. These systems represent a paradigm shift from traditional module-based approaches, which separate perception [30–32], prediction [33–35], and planning [36–38] into distinct components. Instead, end-to-end AD synthesizes these functionalities into a unified, integrated framework, utilizing the robustness of foundational models [19, 39, 40] and

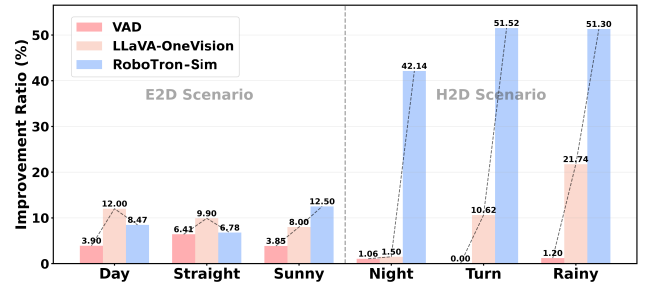


Figure 1. RoboTron-Sim achieves stronger improvements in real-world driving capabilities by leveraging simulated hard-case data. We perform representative methods to evaluate the improvements from simulated data. Results show traditional method VAD achieves minor gains while LLaVA-OneVision struggles to improve performance in challenging scenarios. RoboTron-Sim achieves over 50% improvement in hard-to-drive scenarios.

extensive datasets [15–17] to catalyze innovation. This data-driven approach not only enhances scalability [22, 42] but also promises continuous improvement as incremental data integration [39, 41]. Thus, end-to-end AD is positioned to redefine the future of autonomous mobility by offering a more adaptive solution to complex driving scenarios.

In the era of data-driven innovation, the availability of high-quality datasets [9, 44] plays a pivotal role in technological advancement. Current mainstream AD systems like UniAD [1] and VAD [49], which employ end-to-end architectures, remain constrained by supervised learning paradigms requiring fine-grained annotations such as 3D bounding boxes for agents and semantic segmentation of lane markings. The exorbitant cost of annotation has become a critical bottleneck in scaling these end-to-end methods, while the inherent long-tail distribution [45] challenge in AD tasks exacerbates the data dilemma – most collected data consists of trivial scenarios like straight-road driving [3], with safety-critical cases being extremely scarce. This distribution shift significantly limits the practical effectiveness of data-driven approaches, particularly in complex edge cases and risk-sensitive scenarios where data collection faces substantial challenges. Given the scarcity of real-

[†] Equal contribution. [‡] Corresponding author.

world data in these critical areas, Methods [24, 42] have increasingly turned to synthetic data as a potential solution. Simulation environment [10] offers the ability to generate diverse and controlled datasets, including rare and high-risk scenarios that are difficult to capture in the real world. These constraints raise a fundamental research question: *Can simulation-generated datasets be effectively leveraged to enhance model performance in real-world scenarios?*

While simulation data is cost-effective and enables training with synthetic challenging scenarios, conventional approaches (e.g. VAD) that simply mix real-world and simulated data achieve limited performance gains, even with larger datasets. This constraint primarily stems from the persistent simulation-to-reality (Sim2Real) domain gap, as evidenced by our experimental observations in Fig. 1, where models demonstrate inadequate capability in transferring knowledge acquired from synthetic scenarios to authentic driving contexts. The core challenge lies in the inherent discrepancies between simulated inputs and real-world data, hindering cross-domain knowledge transfer, as traditional methods fail to establish meaningful correspondences between virtual and real scenarios. Multimodal large language models (MLLMs), leveraging their robust reasoning [26, 48] and generalization capabilities [14, 46, 47], provide new opportunities for effectively merging different domain data, as demonstrated by the comparison of LLaVA-OneVision [14] with VAD in Fig. 1. Despite advancements, the misalignment in the geometric space between simulated and real data continues to constrain model performance. Consequently, the critical research question remains unresolved: *How can MLLMs effectively leverage synthetic data to enhance real-world autonomous driving performance?*

To address the above two questions, this paper proposes RoboTron-Sim, a multimodal large language model framework designed to bridge the Sim2Real gap by learning actionable driving knowledge from simulated hard cases. First, we implement a data stratification strategy that categorizes real-world AD data and generates targeted synthetic scenarios through automated simulation pipelines. This approach specifically augments underrepresented safety-critical cases (e.g. long-tail scenarios, hard-to-drive scenarios), effectively mitigating the imbalance distribution problem. Second, we re-engineer the multimodal input schema with driving-aware prompts that enable dual-domain awareness and geographical conditioning. By explicitly encoding data provenance (Simulation/Real-World) and location-specific traffic patterns, the model dynamically leverages LLM-embedded commonsense knowledge to adjust driving policies while maintaining domain invariance. Third, we introduce an image-to-ego encoder that explicitly injects camera parameters through a lightweight MLP-based adapter. This geometric-aware module disentangles dataset-specific sensor configurations from driving policy learning, signifi-

cantly improving cross-dataset generalization capabilities. As shown in Fig. 1, extensive experiments on nuScenes demonstrate that RoboTron-Sim achieves $\sim 50\%$ improvement in hard-case success rates compared to baseline methods, while maintaining performance in routine scenarios.

To summarize, our contributions lie in three-fold:

- To our knowledge, we present the first in-depth investigation of Sim2Real transfer limitations in MLLMs for AD, accompanied by a simple yet effective framework that strategically leverages synthetic hard-case data to enhance real-world driving competencies.
- We address the critical domain discrepancy between simulated and real-world data, including scenario-aware prompts that dynamically model data provenance and geographical context, and a geometry-aware image-to-ego encoder that disentangles sensor-specific parameters.
- Experiments on nuScenes demonstrate RoboTron-Sim achieves SOTA planning performance, particularly showing 48.1% improvement in L2 Distance and 45.8% enhancement in Collision Rate for hard scenarios.

2. Related Work

Autonomous Driving Models. The development of autonomous driving systems has been significantly advanced by end-to-end approaches that established fundamental architectural paradigms. UniAD [1] introduces unified perception-prediction-planning coordination through joint training frameworks, while VAD [49] employs vectorized scene representation methodology to optimize model computational efficiency. In recent years, MLLMs have demonstrated exceptional potential in improving end-to-end planning. RDA-Driver [12] improves decision precision via reasoning-decision alignment and optimized CoT structures for interpretability. Senna [13] and DRIVEVLM [18] combine MLLMs with end-to-end models, leveraging the spatial perception capabilities of end-to-end models and the generalization ability of MLLMs to enhance planning.

Simulation Platforms for Autonomous Driving. Several open-source simulation platforms have been developed to facilitate autonomous driving research. CARLA [10] provides a modular urban driving simulator with configurable sensors (RGB/depth/semantic segmentation) and environmental conditions, enabling comparative evaluation of modular pipelines, imitation learning, and RL-based approaches. AirSim [11] delivers high-fidelity physical simulation through Unreal Engine, featuring photorealistic rendering with dynamic shadows/reflections and cross-platform APIs for perception-controller co-simulation. To tackle autonomous corner case challenges, OASIS SIM V3.0 [43] leverages AI-powered traffic flow simulation with reinforcement learning and synthetic data to establish high-fidelity closed-loop environments, rigorously validat-

ing self-driving systems’ rare scenario handling capabilities.

3. Methodology

In this section, our goal is to utilize simulated data to enhance the driving performance of the models in challenging and crucial scenarios, which are difficult for driving and data collection in real world. To this end, we first create a simulated dataset Hard-case Augmented Synthetic Scenarios (HASS) that includes a variety of challenging and crucial scenarios in Sec. 3.1. Subsequently, we introduce two dedicated technologies, Scenario-aware Prompt Engineering (SPE) and Image-to-Ego Encoder (I2E Encoder), designed to assist the MLLM baseline in effectively leveraging simulated data to enhance its performance in Sec. 3.2.

3.1. Data Curation

The effectiveness of synthetic data hinges on its ability to mirror real-world complexity while strategically addressing data scarcity in edge cases. We present a systematic framework for generating scenario-specific synthetic data in CARLA [10], structured through three key dimensions: environmental diversity, agent behavior orchestration, and sensor-realistic multimodality.

3.1.1. Scenario-Aware Data Classification

To generate simulated data tailored to address the real-world challenging scenarios, we categorize the driving scenarios based on real-world data availability and criticality.

Common Scenarios. Common scenarios are those frequently encountered in autonomous driving environments. Our scenario taxonomy systematically organizes the common driving scenarios into two distinct classes: Easy-to-Drive (E2D) and Hard-to-Drive (H2D) based on their driving difficulty. E2D scenarios refer to the uncomplicated common cases, such as daylight lane maintenance and uncongested traffic flows, which serve as foundational references for model calibration. In contrast, H2D refers to scenarios in which driving becomes challenging due to factors like weather and lighting. Examples include night-time driving, fog-obscured, and heavy-rain conditions where optical sensors degrade significantly. Due to different occurrence frequencies of E2D and H2D scenarios and human driving preferences, the existing dataset often exhibits an imbalance between these two situations, e.g., the ratio of data volume between daytime and nighttime is approximately 7:1 in nuScenes. Nevertheless, *H2D scenarios often require more data for training*. This motivates us to synthesize the H2D data to balance and enhance the driving performance in H2D scenarios.

Long-Tail Scenarios. Long-tail scenarios encompass events that occur infrequently in natural driving environments yet exhibit high diversity and complexity. These

scenarios often involve rare combinations of environmental factors, atypical agent behaviors, or transient conditions that challenge perception and decision-making systems. Examples include near-collision events, such as a vehicle narrowly avoiding a pedestrian suddenly crossing the road, abrupt pedestrian appearances during sharp turns, or unexpected vehicle maneuvers that nearly result in accidents. Other instances might involve complex interactions at intersections, such as a car running a red light or a pedestrian stepping into the roadway despite oncoming traffic. Such scenarios defy conventional statistical modeling due to their combinatorial explosion of parameters. *While individual events may have low occurrence probabilities, they can lead to significant safety concerns.* The critical challenge lies not merely in their scarcity but in their inherent unpredictability, as these scenarios often violate implicit assumptions of spatial/temporal coherence embedded in perception models. By leveraging simulation data, which includes these specific long-tail scenarios, we can better train and evaluate the model’s ability to handle rare but critical situations, ensuring robustness and safety in real-world applications.

3.1.2. Data Collection

Generation Details. Efficient data collection in simulation scenarios relies on teacher models equipped with global environmental awareness, which leverage privileged simulator-level information (e.g., precise poses of agents, intent semantics, and full traffic signal states) to achieve robust driving performance. Unlike conventional rule-based systems constrained by limited generalization across complex scenarios, we employ Think2Drive [8], a world model-driven reinforcement learning architecture as the core data generator. To enhance reproducibility and lower the entry barrier for community-driven research in end-to-end AD, we implement a comprehensive sensor suite mirroring the nuScenes benchmark configuration, including six 900×1600 resolution cameras providing 360° coverage with overlapping fields of view.

To systematically address both routine driving patterns and challenging long-tail cases mentioned above, we procedurally develop a dataset, called Hard-case Augmented Synthetic Scenarios (HASS). These scenarios are divided into two main categories: routine maneuvers (i.e., E2D and H2D) and long-tail scenarios. Several visualization examples of the H2D and long-tail scenarios are presented in Figure 2. The long-tail scenarios are further classified into 13 high-risk edge case categories (e.g., pedestrian jaywalking, sudden vehicle cut-ins, and near-collision events), as visualized in Figure 3. These long-tail scenarios are designed to capture rare but critical driving scenarios that are often underrepresented in real-world datasets.

In addition to scenario diversity, HASS achieves balanced environmental conditions compared to traditional human-collected driving records. As shown in Table 1,

Scenario	Day	Night	Sunny	Rainy	Straight	Turn
Real Scenario	24745 (87.97%)	3385 (12.03%)	22548 (80.16%)	5582 (19.84%)	24996 (88.86%)	3134 (11.14%)
Simulated Scenario	27891 (58.65%)	19662 (41.35%)	23010 (48.38%)	24543 (51.61%)	22076 (46.42%)	25477 (53.58%)

Table 1. Comparison of the volume of data from different sources across various scenarios.



Figure 2. Visualization of long-tail and H2D scenarios in HASS. (a) Temporary Parking Ahead, (b) Roadwork Ahead, (c) Jaywalking Pedestrians, (d) Lane Invasion, (e) Opposing Lane Encroachment, (f) Parked vehicle Activation, (g)(j)(k) Rainy, (h) Turn, (i)(l) Night.

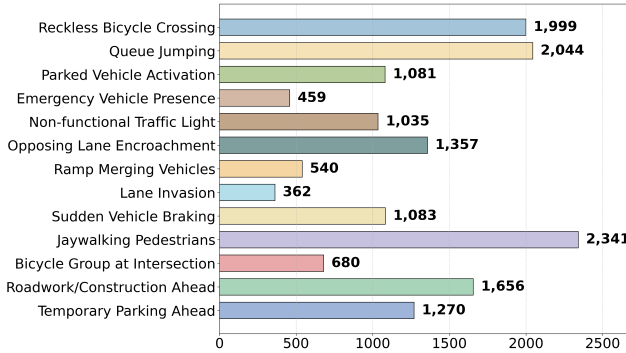


Figure 3. Types and sample counts of long-tail scenarios.

HASS shows a balanced distribution of environmental factors, including 58.65% daytime and 41.35% nighttime scenarios, as well as 48.38% sunny and 51.61% rainy conditions. Besides, HASS prioritizes interaction complexity, with 53.58% of scenarios that involve turn-dominated maneuvers, a great increase compared to 11.14% in real-world data. This deliberate design ensures that HASS not only addresses the limitations of real-world data but also provides a robust foundation for training and evaluating end-to-end AD systems in diverse and challenging conditions.

Data Alignment. Real-world scenario (e.g., nuScenes) mainly adopts a right-handed coordinate system (X: forward, Y: left, Z: upward), while CARLA defaults to a left-handed coordinate system (X: forward, Y: right, Z: upward). This discrepancy involves both axial directional conflicts and spatial offsets in sensor positioning: in nuScenes, the vehicle coordinate origin is located at the center of the roof, while in CARLA, the vehicle coordinate origin is situated in the contact plane of the wheel. These conflicting axis definitions and origin offsets would otherwise cause cross-domain data fusion failures. To bridge the gap between sim-

ulated and real-world data distributions, we transform the coordinate system of the simulation data to a right-handed coordinate system and establish the vehicle’s coordinate system origin as the center of the roof.

3.2. Model

Our early efforts in tackling the planning task involved experimenting with traditional models like VAD [49], trained on a mix of simulated and real-world data. However, as indicated by our experiments, the simulated data only yield very limited improvement (reducing L2 distance by $\sim 1\%$ for the HD scenario).

We analyzed that the primary limitation stemmed from traditional models struggling to generalize effectively across diverse and dynamic environments, largely due to the substantial domain gap between simulation and reality. Recently, MLLMs have demonstrated remarkable generalization performance in various visual tasks. This led us to explore the potential of MLLMs for addressing this task.

3.2.1. MLLM Baseline

The proposed framework, as illustrated in Figure 4, employs a multimodal architecture. It consists of three main components: a visual feature extractor, a feature adapter, and an LLM backbone. The visual feature extractor processes raw videos from multiple camera views over time, encoding them into a compact spatiotemporal representation, reducing dimensionality while preserving critical details for downstream tasks. Furthermore, the framework employs a two-layer multilayer perceptron (MLP) to project the extracted visual features into the language model’s embedding space, establishing dimensional compatibility between the visual representations and linguistic tokens while preserving semantic consistency across modalities. Consequently, the LLM decoder integrates the processed visual

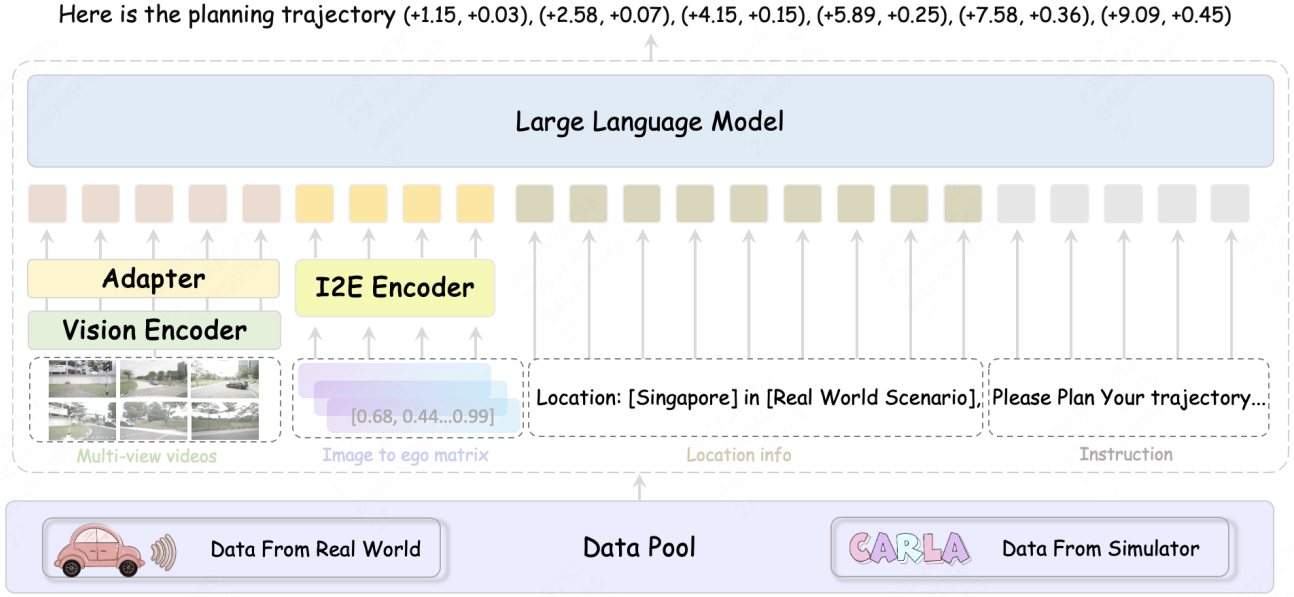


Figure 4. The overall framework of our proposed end-to-end autonomous driving system. The framework leverages simulation data to enhance performance in real-world scenarios, by integrating videos, 3D transformation, location, data source, and instruction information.


<p>Input:</p>  <p>Question:</p> <p>1: <video> 2: <video> ... 6: <video>. These 6 videos are the front view, front left view, front right view, back view, back left view, back right view of the ego vehicle. You need to make a left turn at the upcoming intersection, please provide the planning trajectory for the ego car.</p> <p>Answer:</p> <p>Here is the planning trajectory (+4.96, +0.12), (+8.93, +0.48), (+12.62, +1.03), (+16.27, +1.78), (+19.67, +2.68), (+22.94, +3.70).</p>
--

Table 2. The perspective-aware prompt for multi-view inputs.

features with tokenized text input, predicting textual outputs in an autoregressive manner. By leveraging both multimodal context and previously generated tokens, coherent and contextually relevant responses emerge.

MLLM for planning. Building on the foundational MLLM architecture, we tailor the framework for end-to-end planning by enabling it to generate actionable trajectories from multimodal inputs. As shown in Tabel 2, the model processes visual data captured from six cameras over five consecutive frames, alongside high-level command instruc-

tions (e.g., “make a left turn at the upcoming intersection,” “move forward”), ensuring a comprehensive action understanding. The visual feature extractor encodes spatiotemporal dynamics, while tokenized command instructions are fed into the LLM backbone alongside the projected visual features. To further enhance the model’s awareness of the ego vehicle’s state, we introduce velocity supervision. Consequently, the model produces two types of outputs: predicted trajectory points for the upcoming moments and anticipated vehicle speeds over the same timeframe.

3.2.2. Enhancing Driving with Simulated Data

We attempted to train the designed baseline by directly mixing the real-world data with simulated data. Although this approach has shown some improvements (reducing L2 distance from 1.03m to 0.91m), we identified that this approach alone is insufficient to fully adapt to the inherent differences between data sources, limiting the model’s ability to achieve robust generalization. To better align and leverage the complementary strengths of simulated data, we introduce two key designs aimed at enhancing the MLLM baseline to transfer diverse simulated driving skills to real-world. The resultant model is named RoboTron-Sim.

Scenario-Specific Contextual Grounding. To address the domain discrepancy between simulated and real-world environments, we implement contextual priming through Scenario-aware Prompt Engineering (SPE). Specifically, each input sequence is augmented with an environmental descriptor formatted as: “You are driving in [City Name] under [Simulation/Real-World] scenario.” This structured prompt serves dual purposes:

Model	L2(m)				Collision(%)				Boundary(%)			
	1s	2s	3s	avg	1s	2s	3s	avg	1s	2s	3s	avg
Without Ego Pose as Input												
UniAD[1]	0.59	1.01	1.48	1.03	0.16	0.51	1.64	0.77	0.35	1.46	3.99	1.93
VAD-Base[49]	0.69	1.22	1.83	1.25	0.06	0.68	2.52	1.09	1.02	3.44	7.00	3.82
BEV-Planner[3]	0.27	0.54	0.90	0.57	0.10	0.37	1.30	0.59	0.78	3.79	8.22	4.26
LLaVA[5]	1.04	1.74	2.57	1.79	0.58	1.17	1.74	1.16	-	-	-	-
Vicuna[6]	1.06	1.80	2.54	1.80	0.60	1.21	1.78	1.20	-	-	-	-
Merlin[7]	1.03	1.71	2.40	1.71	0.48	1.05	1.77	1.10	-	-	-	-
OmniDrive[25]	0.40	0.80	1.32	0.84	0.04	0.46	2.32	0.94	0.93	3.65	8.28	4.29
RoboTron-Sim	0.22	0.53	0.93	0.56	0.11	0.35	1.27	0.58	0.55	1.97	5.57	3.02
With Ego Pose as Input												
UniAD[1]	0.20	0.42	0.75	0.46	0.02	0.25	0.84	0.37	0.20	1.33	3.24	1.59
VAD-Base[49]	0.17	0.34	0.60	0.37	0.04	0.27	0.67	0.33	0.21	2.13	5.06	2.47
Ego-MLP[2]	0.15	0.32	0.59	0.35	0.00	0.27	0.85	0.37	0.27	2.52	6.60	2.93
BEV-Planner[3]	0.16	0.32	0.57	0.35	0.00	0.29	0.73	0.34	0.35	2.62	6.51	3.16
EMMA[4]	0.14	0.29	0.54	0.32	-	-	-	-	-	-	-	-
OmniDrive[25]	0.14	0.29	0.55	0.33	0.00	0.13	0.78	0.30	0.56	2.48	5.96	3.00
RoboTron-Sim	0.10	0.20	0.39	0.23	0.00	0.11	0.66	0.26	0.39	2.07	5.41	2.62

Table 3. Comparison on the open-loop planning on nuScenes dataset with and without ego pose as input. We report the L2 distance (m), collision rate (%), and boundary violation rate (%) at 1s, 2s, and 3s time horizons, along with their averages.

- **Domain Awareness:** Explicitly informs the model about the data characteristics (e.g., sensor noise levels) through categorical labeling of simulation/real-world contexts.
- **Geographical Conditioning:** Embeds location-specific priors (e.g., traffic rules) via city name specification, enabling adaptive processing of regional driving patterns.

For instance, the prompt “*You are driving in Town13 under simulation scenario*” activates the model’s knowledge of both the city’s unique driving habits (left-hand or right-hand driving) and traffic rules. This added context helps the MLLM baseline adjust its decision-making, ignoring unrealistic details from simulations when working with real-world data, while making use of useful patterns learned from synthetic scenarios to improve planning stability.

Geometry-Aware Visual Alignment. Due to variations in vehicles and cameras, together with different sensor installation positions in simulated and real-world scenarios, the intrinsic and extrinsic parameters of the camera typically differ between the two systems. This discrepancy creates a critical cross-domain gap, leading to degraded performance when transferring between simulation and reality. To enhance adaptability to different camera parameters, we introduce an Image-to-Ego Encoder (I2E Encoder) to explicitly incorporate geometric transformations into the input processing, enhancing the visual features with 3D spatial information, as shown in Figure 4. Specifically, we first compute the image-to-ego transformation matrix using the

camera’s intrinsic and extrinsic parameters for each view, ensuring a consistent spatial representation across different viewpoints. Subsequently, we employ the I2E Encoder, implemented as a two-layer MLP, to integrate this geometric information into the MLLM baseline. This encoder maps the transformation matrix into an embedding space that captures the spatial context of each camera view. The resulting embeddings are then concatenated with the tokenized text input, allowing the model to incorporate spatial reasoning directly into its decision-making process.

4. Experiment

4.1. Experimental Setting

4.1.1. Dataset

RoboTron-Sim is trained using a hybrid data strategy combining:

- **Real-world Data:** 28,130 samples from nuScenes[9].
- **Simulated Data:** 47,553 purpose-built samples generated in CARLA simulator [10].

For evaluation, we use nuScenes validation set. Please refer to the supplementary material for more details.

4.1.2. Evaluation Metrics

Following BEV-Planner [3], we evaluate via L2 Distance, Collision Rate, and Boundary Violation Rate. Please refer to the supplementary material for more details.

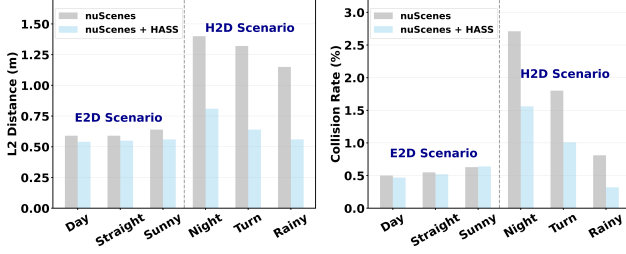


Figure 5. Comparison of RoboTron-Sim on different datasets.

4.1.3. Training Details

Our model architecture builds upon LLaVA-OneVision [14], with continuous five-frame sequences from six camera views as input. The training was performed on 16 NVIDIA A100 GPUs for 3 epochs, taking approximately 20 hours when using nuScenes and HASS.

4.2. Main Results

Open-loop Planning. To ensure a fair comparison of planning capabilities across different methods, we evaluate our method under two settings: (1) without ego pose as input and (2) with ego pose as input. As shown in Tab. 3, when ego pose is excluded, the framework achieves state-of-the-art trajectory prediction accuracy with the lowest average L2 distance (0.56m) and collision rate (0.58%) while maintaining a competitive boundary violation rate (3.02% vs. 4.29% in OmniDrive). Incorporating ego pose further enhances these metrics, yielding substantial improvements to 0.23m L2 distance and 0.26% collision rate, achieving the SOTA. The boundary violation rate also decreases to 2.62%, demonstrating consistent performance gains across all safety-critical metrics. This highlights RoboTron-Sim’s adaptability to both sensor-limited and sensor-rich conditions while maintaining superior safety margins.

Scenario-Specific Improvement. Here we evaluate the performance improvement of RoboTron-Sim in each scenario. Specifically, we train the proposed model on two kinds of data: (1) nuScenes, and (2) nuScenes + HASS. The results presented in Fig. 5 expose striking performance disparities across environments. In E2D scenarios (Day + Straight + Sunny) where existing models already demonstrate mature capabilities, the performance ceiling leaves limited optimization space. Conversely, in H2D scenarios (Night + Turn + Rainy), the integration of HASS drives remarkable progress: the nighttime collision rate is reduced by **42.4%** (from 2.71% to 1.56%), while turning maneuver precision shows **51.5%** improvement in L2 metrics. It is worth noting that RoboTron-Sim achieves **51.3%** lower collision rates than the baseline approach. This flipped contrast underscores the value of HASS: while maintaining stability in routine scenarios, it primarily empowers breakthrough advancements in hard cases via synthetic scenario.

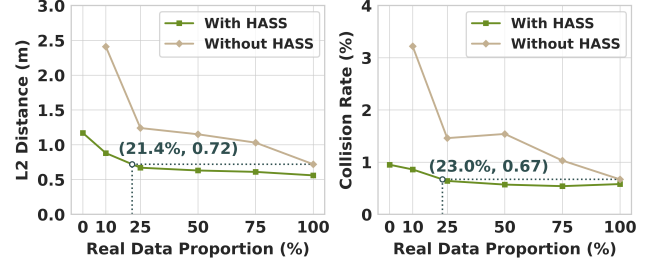


Figure 6. Performance comparison between methods **with** and **without** HASS under different ratios of nuScenes data.

Data Robustness. We conduct comparative experiments in multiple data regimes to verify the robustness of the data: (1) different proportions of nuScenes data (from 100% to 0%), and (2) different proportions of nuScenes data + full HASS data. As shown in Fig. 6, RoboTron-Sim maintains stable performance even as the proportion of real-world data decreases from 100% to 10%. In contrast, purely real-data-based models experience severe performance degradation in limited-data conditions, with both L2 distance and collision rate increasing significantly as real samples diminish. By incorporating HASS, performance degradation is alleviated, ensuring minimal performance loss even when real data is scarce. Notably, we find that by combining simulated data with only around 20% real samples, we can achieve the same level of performance as using 100% real data. This highlights the crucial role of simulated data in compensating for real-world data limitations while preserving robust trajectory planning across diverse data distributions.

4.3. Ablation Study

Model Compatibility. To investigate model compatibility with simulated data augmentation, we conduct cross-architecture evaluations on L2 distance across driving scenarios. As evidenced in Table 4, the effectiveness of HASS varies significantly across architectures: VAD exhibits inherent incompatibility with HASS, showing marginal L2 improvements (3-4% in daytime/sunny conditions) but showing negligible gains (below 1.5%) in H2D scenarios. While MLLMs demonstrate preliminary adaptability with moderate L2 reductions (8-12% in E2D), their gains decline in H2D scenarios, exposing limitations in modeling vehicle dynamics. In contrast, RoboTron-Sim achieves paradigm-shifting enhancements through Sim2Real alignment, it reduces L2 distance by **51.5%** in complex turns and **42-51%** in night/rain conditions. This empowers knowledge transfer from synthetic domains while preserving real-world physical constraints, unlocking the model’s untapped potential.

Ablation on Model Designs. Tab. 6 (rows 1-2) shows SPE’s enhancements: L2 is reduced by 5.5%, collisions drops by 16.0%, and most notably, boundary violations see a substantial decrease of 38.2% (3.22%→2.68%), which

Method	Data	E2D			H2D		
		Day	Straight	Sunny	Night	Turn	Rainy
VAD	nuScenes	0.77	0.78	0.78	0.94	0.87	0.83
	nuScenes + HASS	0.74(↓ 3.9%)	0.78(↓ 0.0%)	0.75(↓ 3.8%)	0.93(↓ 1.1%)	0.87(↓ 0.0%)	0.82(↓ 1.2%)
MLLM	nuScenes	1.00	1.01	1.00	1.33	1.13	1.15
	nuScenes + HASS	0.88(↓ 12.0%)	0.91(↓ 9.9%)	0.92(↓ 8.0%)	1.31(↓ 1.5%)	1.02(↓ 9.7%)	0.90(↓ 21.7%)
RoboTron-Sim	nuScenes	0.59	0.59	0.64	1.40	1.32	1.15
	nuScenes + HASS	0.54(↓ 8.5%)	0.55(↓ 6.8%)	0.56(↓ 12.5%)	0.81(↓ 42.1%)	0.64(↓ 51.5%)	0.56(↓ 51.3%)

Table 4. L2 Performance gains of HASS on different models in E2D and H2D scenarios. Red arrows indicate improvement (**lower values are better**), gray denotes no significant change. “MLLM” means the baseline obtained by finetuning LLaVA-OneVision on driving data.

Data	E2D			H2D		
	Day	Straight	Sunny	Night	Turn	Rainy
nuScenes	0.59	0.59	0.64	1.40	1.32	1.15
nuScenes + GASS	0.55(↓ 6.8%)	0.50(↓ 15.3%)	0.52(↓ 18.8%)	1.00(↓ 28.6%)	1.21(↓ 8.3%)	0.99(↓ 13.9%)
nuScenes + HASS	0.54(↓ 8.5%)	0.55(↓ 6.8%)	0.56(↓ 12.5%)	0.81(↓ 42.1%)	0.64(↓ 51.5%)	0.56(↓ 51.3%)

Table 5. L2 Performance comparison of different datasets in E2D and H2D scenarios (**lower values are better**).

SPE	I2E Encoder	L2(m)	Collision(%)	Boundary(%)
×	×	0.91	0.94	3.22
✓	×	0.86	0.79	2.68
✓	✓	0.56	0.58	3.02

Table 6. Ablation study on model designs.

proves SPE’s ability to enhance the model’s sensitivity to lane discipline and road geometry constraints. The last two rows of Tab. 6 show that integrating I2E Encoder reduces L2 distance by 34.9% and collision rate by 26.6%. This demonstrates that explicit geometric grounding effectively bridges the Sim2Real domain gap by establishing camera-independent spatial representations, enabling consistent perception across synthetic and physical sensors.

Effectiveness of HASS. To evaluate the impact of different data synthesis strategies, we conduct experiments using two kinds of synthetic data: (1) Hard-case Augmented Synthetic Scenarios (HASS), where we adjust the synthesis proportion to focus on hard cases, and (2) General Augmented Synthetic Scenarios (GASS), where the data is synthesized following the scene distribution in nuScenes. As shown in Tab. 5, the method trained on GASS performs well across E2D scenarios but struggles in H2D conditions, such as night and turn scenarios, where higher L2 errors are observed (1.00m and 1.21m). In contrast, the model trained on HASS, which increases the representation of hard cases, achieves great improvements in H2D scenarios. Please refer to the supplementary material for more details.

4.4. Case Study

We visualize predicted trajectories from both the baseline and RoboTron-Sim under ego-pose-free conditions, together with ground truth. In Fig. 7a, when facing an oncoming vehicle crossing lanes and a stationary ego-lane ob-

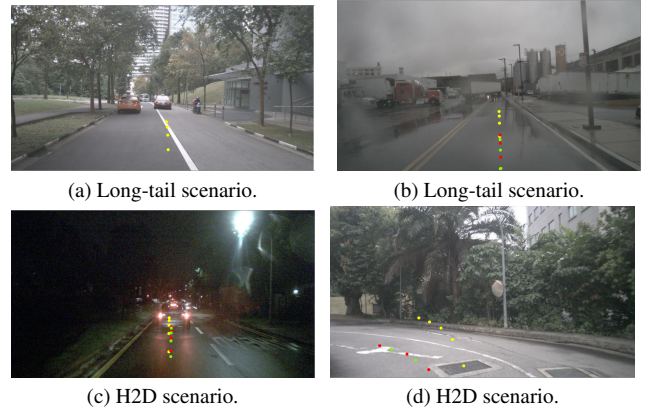


Figure 7. Visual comparison of trajectories in challenging scenarios. **Ground-truth** trajectories are marked in red, **baseline** predictions in yellow, and **RoboTron-Sim**’s predictions in green.

stacle, **RoboTron-Sim** accurately executes a safe stop like the **ground truth**, whereas the **baseline** dangerously continues forward. This demonstrates our model’s superior risk awareness in complex situations. RoboTron-Sim also yields safer sharp turns (Fig. 7d). These visualizations (and Fig. 5) confirm RoboTron-Sim’s effectiveness in safety-critical scenarios.

5. Conclusion

In this paper, we present RoboTron-Sim, an approach that improves real-world driving robustness by leveraging long-tail and hard-to-drive cases. To our knowledge, this is the first solution addressing sim-to-real gaps in MLLMs for autonomous driving. We introduce the HASS dataset and the RoboTron-Sim model with SPE and I2E, boosting real-world performance in rare high-risk driving scenarios.

References

- [1] Y. Hu, J. Yang, L. Chen, K. Li, C. Sima, X. Zhu, S. Chai, S. Du, T. Lin, W. Wang *et al.*, “Planning-oriented autonomous driving,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 17 853–17 862. 1, 2, 6
- [2] J.-T. Zhai, Z. Feng, J. Du, Y. Mao, J.-J. Liu, Z. Tan, Y. Zhang, X. Ye, and J. Wang, “Rethinking the open-loop evaluation of end-to-end autonomous driving in nuscenese,” *arXiv preprint arXiv:2305.10430*, 2023. 6
- [3] Z. Li, Z. Yu, S. Lan, J. Li, J. Kautz, T. Lu, and J. M. Alvarez, “Is ego status all you need for open-loop end-to-end autonomous driving?” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 14 864–14 873. 1, 6, 11
- [4] J.-J. Hwang, R. Xu, H. Lin, W.-C. Hung, J. Ji, K. Choi, D. Huang, T. He, P. Covington, B. Sapp *et al.*, “Emma: End-to-end multimodal model for autonomous driving,” *arXiv preprint arXiv:2410.23262*, 2024. 6
- [5] H. Liu, C. Li, Q. Wu, and Y. J. Lee, “Visual instruction tuning,” *Advances in neural information processing systems*, vol. 36, pp. 34 892–34 916, 2023. 6
- [6] W.-L. Chiang, Z. Li, Z. Lin, Y. Sheng, Z. Wu, H. Zhang, L. Zheng, S. Zhuang, Y. Zhuang, J. E. Gonzalez *et al.*, “Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, march 2023,” URL <https://lmsys.org/blog/2023-03-30-vicuna>, vol. 3, no. 5, 2023. 6
- [7] E. Yu, L. Zhao, Y. Wei, J. Yang, D. Wu, L. Kong, H. Wei, T. Wang, Z. Ge, X. Zhang *et al.*, “Merlin: Empowering multimodal llms with foresight minds,” in *European Conference on Computer Vision*. Springer, 2024, pp. 425–443. 6
- [8] Q. Li *et al.*, “Think2drive: Efficient reinforcement learning by thinking with latent world model for autonomous driving (in carla-v2),” in *European Conference on Computer Vision*. Cham: Springer Nature Switzerland, 2024. 3
- [9] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, “nuscenese: A multimodal dataset for autonomous driving,” in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020, pp. 11 621–11 631. 1, 6, 11
- [10] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, “Carla: An open urban driving simulator,” in *Conference on robot learning*. PMLR, 2017, pp. 1–16. 2, 3, 6, 11
- [11] S. Shah, D. Dey, C. Lovett, and A. Kapoor, “Airsim: High-fidelity visual and physical simulation for autonomous vehicles,” in *Field and Service Robotics: Results of the 11th International Conference*. Springer, 2018, pp. 621–635. 2
- [12] Z. Huang, T. Tang, S. Chen, S. Lin, Z. Jie, L. Ma, G. Wang, and X. Liang, “Making large language models better planners with reasoning-decision alignment,” in *European Conference on Computer Vision*. Springer, 2024, pp. 73–90. 2
- [13] B. Jiang, S. Chen, B. Liao, X. Zhang, W. Yin, Q. Zhang, C. Huang, W. Liu, and X. Wang, “Senna: Bridging large vision-language models and end-to-end autonomous driving,” *arXiv preprint arXiv:2410.22313*, 2024. 2
- [14] B. Li, Y. Zhang, D. Guo, R. Zhang, F. Li, H. Zhang, K. Zhang, P. Zhang, Y. Li, Z. Liu *et al.*, “Llava-onevision: Easy visual task transfer,” *arXiv preprint arXiv:2408.03326*, 2024. 2, 7, 12
- [15] Y. Li, W. Zhang, K. Chen, Y. Liu, P. Li, R. Gao, L. Hong, M. Tian, X. Zhao, Z. Li *et al.*, “Automated evaluation of large vision-language models on self-driving corner cases,” *arXiv preprint arXiv:2404.10595*, 2024. 1
- [16] A.-M. Marcu, L. Chen, J. Hünermann, A. Karnsund, B. Hanotte, P. Chidananda, S. Nair, V. Badrinarayanan, A. Kendall, J. Shotton *et al.*, “Lingoqa: Video question answering for autonomous driving,” in *Eur. Conf. Comput. Vis.*, 2024.
- [17] X. Ding, J. Han, H. Xu, X. Liang, W. Zhang, and X. Li, “Holistic autonomous driving understanding by bird’s-eye-view injected multi-modal large models,” in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2024, pp. 13 668–13 677. 1
- [18] X. Tian, J. Gu, B. Li, Y. Liu, Y. Wang, Z. Zhao, K. Zhan, P. Jia, X. Lang, and H. Zhao, “Drivevlm: The convergence of autonomous driving and large vision-language models,” *arXiv preprint arXiv:2402.12289*, 2024. 2
- [19] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altschmidt, S. Altman, S. Anadkat *et al.*, “Gpt-4 technical report,” *arXiv preprint arXiv:2303.08774*, 2023. 1
- [20] L. Wen, D. Fu, X. Li, X. Cai, T. Ma, P. Cai, M. Dou, B. Shi, L. He, and Y. Qiao, “Dilu: A knowledge-driven approach to autonomous driving with large language models,” *arXiv preprint arXiv:2309.16292*, 2023. 1
- [21] Z. Xu, Y. Zhang, E. Xie, Z. Zhao, Y. Guo, K.-Y. K. Wong, Z. Li, and H. Zhao, “Drivegpt4: Interpretable end-to-end autonomous driving via large language model,” *IEEE Robotics and Automation Letters*, 2024.
- [22] Z. Huang, T. Tang, S. Chen, S. Lin, Z. Jie, L. Ma, G. Wang, and X. Liang, “Making large language models better planners with reasoning-decision alignment,” in *European Conference on Computer Vision*. Springer, 2025, pp. 73–90. 1
- [23] C. Cui, Y. Ma, X. Cao, W. Ye, and Z. Wang, “Drive as you speak: Enabling human-like interaction with large language models in autonomous vehicles,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024, pp. 902–909.
- [24] H. Shao, Y. Hu, L. Wang, G. Song, S. L. Waslander, Y. Liu, and H. Li, “Lmdrive: Closed-loop end-to-end driving with large language models,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 15 120–15 130. 1, 2
- [25] S. Wang, Z. Yu, X. Jiang, S. Lan, M. Shi, N. Chang, J. Kautz, Y. Li, and J. M. Alvarez, “OmniDrive: A holistic llm-agent framework for autonomous driving with 3d perception, reasoning and planning,” *arXiv preprint arXiv:2405.01533*, 2024. 6
- [26] H. Liu, C. Li, Q. Wu, and Y. J. Lee, “Visual instruction tuning,” *Advances in neural information processing systems*, vol. 36, 2024. 2
- [27] Z. Yang, X. Jia, H. Li, and J. Yan, “A survey of large language models for autonomous driving,” *arXiv preprint arXiv:2311.01043*, 2023. 1

- [28] L. Chen, P. Wu, K. Chitta, B. Jaeger, A. Geiger, and H. Li, “End-to-end autonomous driving: Challenges and frontiers,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- [29] Z. Zhang, X. Bo, C. Ma, R. Li, X. Chen, Q. Dai, J. Zhu, Z. Dong, and J.-R. Wen, “A survey on the memory mechanism of large language model based agents,” *arXiv preprint arXiv:2404.13501*, 2024. 1
- [30] J. Huang, G. Huang, Z. Zhu, Y. Ye, and D. Du, “Bevdet: High-performance multi-camera 3d object detection in bird-eye-view,” *arXiv preprint arXiv:2112.11790*, 2021. 1
- [31] T. Liang, H. Xie, K. Yu, Z. Xia, Z. Lin, Y. Wang, T. Tang, B. Wang, and Z. Tang, “Befusion: A simple and robust lidar-camera fusion framework,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 10 421–10 434, 2022.
- [32] Z. Liu, H. Tang, A. Amini, X. Yang, H. Mao, D. L. Rus, and S. Han, “Befusion: Multi-task multi-sensor fusion with unified bird’s-eye view representation,” in *2023 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2023, pp. 2774–2781. 1
- [33] J. Gu, C. Hu, T. Zhang, X. Chen, Y. Wang, Y. Wang, and H. Zhao, “Vip3d: End-to-end visual trajectory prediction via 3d agent queries,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 5496–5506. 1
- [34] J. Gao, C. Sun, H. Zhao, Y. Shen, D. Anguelov, C. Li, and C. Schmid, “Vectornet: Encoding hd maps and agent dynamics from vectorized representation,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 11 525–11 533.
- [35] F. Da and Y. Zhang, “Path-aware graph attention for hd maps in motion prediction,” in *2022 International conference on robotics and automation (ICRA)*. IEEE, 2022, pp. 6430–6436. 1
- [36] O. Scheel, L. Bergamini, M. Wolczyk, B. Osiński, and P. Ondruska, “Urban driver: Learning to drive from real-world demonstrations using policy gradients,” in *Conference on Robot Learning*. PMLR, 2022, pp. 718–728. 1
- [37] A. Sadat, S. Casas, M. Ren, X. Wu, P. Dhawan, and R. Urtasun, “Perceive, predict, and plan: Safe motion planning through interpretable semantic representations,” in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIII 16*. Springer, 2020, pp. 414–430.
- [38] L. Gao, Z. Gu, C. Qiu, L. Lei, S. E. Li, S. Zheng, W. Jing, and J. Chen, “Cola-hrl: Continuous-lattice hierarchical reinforcement learning for autonomous driving,” in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2022, pp. 13 143–13 150. 1
- [39] Z. Yang, X. Jia, H. Li, and J. Yan, “Llm4drive: A survey of large language models for autonomous driving,” *arXiv preprint arXiv:2311.01043*, 2023. 1
- [40] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners,” *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020. 1
- [41] P. Wu, L. Chen, H. Li, X. Jia, J. Yan, and Y. Qiao, “Policy pre-training for autonomous driving via self-supervised geometric modeling,” *arXiv preprint arXiv:2301.01006*, 2023. 1
- [42] X. Jia, Z. Yang, Q. Li, Z. Zhang, and J. Yan, “Bench2drive: Towards multi-ability benchmarking of closed-loop end-to-end autonomous driving,” *arXiv preprint arXiv:2406.03877*, 2024. 1, 2
- [43] Z. Yang, Z. Zhang, Z. Zheng, Y. Jiang, Z. Gan, Z. Wang, Z. Ling, J. Chen, M. Ma, B. Dong, *et al.*, “Oasis: Open agent social interaction simulations with one million agents,” *arXiv preprint arXiv:2411.11581*, 2024. 2
- [44] P. Sun, H. Kretzschmar, X. Dotiwalla, A. Chouard, V. Patnaik, P. Tsui, J. Guo, Y. Zhou, Y. Chai, B. Caine *et al.*, “Scalability in perception for autonomous driving: Waymo open dataset,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 2446–2454. 1
- [45] A. Jain, L. Del Pero, H. Grimmer, and P. Ondruska, “Autonomy 2.0: Why is self-driving always 5 years away?” *arXiv preprint arXiv:2107.08142*, 2021. 1
- [46] Z. Huang, C. Feng, F. Yan, B. Xiao, Z. Jie, Y. Zhong, X. Liang, and L. Ma, “Drivemm: All-in-one large multimodal model for autonomous driving,” *arXiv preprint arXiv:2412.07689*, 2024. 2
- [47] F. Yan, F. Liu, L. Zheng, Y. Zhong, Y. Huang, Z. Guan, C. Feng, and L. Ma, “Robomm: All-in-one multimodal large model for robotic manipulation,” *arXiv preprint arXiv:2412.07215*, 2024. 2
- [48] Y. Zhong, C. Feng, F. Yan, F. Liu, L. Zheng, and L. Ma, “P3nav: A unified framework for embodied navigation integrating perception, planning, and prediction,” *arXiv preprint arXiv:2503.18525*, 2025. 2
- [49] B. Jiang, S. Chen, Q. Xu, B. Liao, J. Chen, H. Zhou, Q. Zhang, W. Liu, C. Huang, and X. Wang, “Vad: Vectorized scene representation for efficient autonomous driving,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 8340–8350. 1, 2, 4, 6, 12

A. Experimental Details

A.1. Datasets

RoboTron-Sim is trained using a hybrid data strategy combining:

- **Real-world Data:** 28,130 samples from nuScenes[9].
- **Simulated Data:** 47,553 purpose-built samples from our Hard-case Augmented Synthetic Scenarios(HASS) dataset, generated in CARLA simulator[10], designed to address the inherent imbalance in real-world data distribution. While the dataset covers a broad range of driving situations, it places particular emphasis on addressing challenging cases, including H2D scenarios and Long-Tail scenarios. Partial results are illustrated in Figure 8.

A.2. Evaluation Metrics

Following BEV-Planner [3], we evaluate via L2 Distance, Collision Rate, and Boundary Violation Rate.

- **Trajectory Accuracy (L2 Distance):**

$$L2 = \frac{1}{T} \sum_{t=1}^T \|\hat{p}_t - p_t^{gt}\|_2 \quad (1)$$

where \hat{p}_t and p_t^{gt} denote the predicted and ground-truth positions at timestep t over a $T = 3s$ horizon.

- **Safety Metrics (Collision Rate):**

Computes the percentage of predicted trajectories that result in collisions with other agents or obstacles.

$$Collision = \frac{1}{T} \sum_{t=1}^T \frac{N_{collision,t}}{N_{total,t}} \times 100\% \quad (2)$$

where $N_{collision}$ is the number of predicted trajectories leading to collisions, and N_{total} is the total number of evaluated trajectories at timestep t over a $T = 3s$ horizon.

- **Boundary Violation Rate:**

$$Boundary = \frac{1}{T} \sum_{t=1}^T \frac{N_{violation,t}}{N_{total,t}} \times 100\% \quad (3)$$

where $N_{violation}$ counts trajectories exceeding road boundaries, and N_{total} is the total evaluated trajectories at timestep t over $T = 3s$. Calculated by comparing ego segmentation masks with drivable area labels.

B. More Results

B.1. Robustness of HASS

We investigate the performance trend divergence between simulated data augmentation and real data-only scenarios across multiple orders of magnitude in real data volume in RoboTron-Sim, with quantitative comparisons presented in Table 7 and Table 8. Table 7 presents quantitative results

nuScenes	HASS	L2(m)	Collision(%)
0%	100%	1.24	0.99
10%	100%	0.87	0.89
25%	100%	0.67	0.64
50%	100%	0.63	0.57
75%	100%	0.61	0.54
100%	100%	0.56	0.58

Table 7. Performance variation with nuScenes blending ratio under full HASS integration.

nuScenes	HASS	L2(m)	Collision(%)
10%	0%	2.41	3.22
25%	0%	1.24	1.46
50%	0%	1.15	1.54
75%	0%	1.03	1.03
100%	0%	0.72	0.67

Table 8. Performance scaling with nuScenes blending ratio (HASS Excluded).

with full simulated data integration, while Table 8 provides detailed metrics when trained without any simulated data, using real-world data exclusively. The experimental results demonstrate enhanced stability of overall performance through simulated data augmentation.

B.2. Effectiveness of HASS

We generate two distinct datasets based on nuScenes scenarios: General Augmented Synthetic Scenarios (GASS) for common driving conditions and Hard-case Augmented Synthetic Scenarios (HASS) for challenging situations, aiming to investigate which synthetic data generation mechanisms yield more meaningful performance improvements. The evaluation results categorized by individual scenarios are presented in Table 9, while the aggregated metrics for H2D scenarios (Night+Turn+Rain) are summarized in Table 10.

B.3. Model Generalization

To verify the model generalization in the **planning task**, we further evaluate model performance on the NAVSIM (NV) benchmark using the predictive driver model score (PDMS), which is based on five factors: no at-fault collisions (NC), drivable area compliance (DAC), time-to-collision (TTC), comfort (Comf.), and ego progress (EP). Table 11 demonstrates that RoboTron-Sim delivers comparable or superior performance compared to existing methods, with the integration of HASS achieving a PDMS of 85.6 and setting a new SOTA result on NV benchmark.

We also explore the robustness of the model on the **VQA**

Data	Day		Night		Straight		Turn		Sunny		Rainy	
	L2	Col	L2	Col	L2	Col	L2	Col	L2	Col	L2	Col
nuScenes	0.59	0.50	1.40	2.71	0.59	0.55	1.32	1.80	0.64	0.63	1.15	0.81
nuScenes + GASS	0.55	0.42	1.00	2.53	0.50	0.49	1.21	1.89	0.52	0.58	0.99	0.79
nuScenes + HASS	0.54	0.47	0.81	1.56	0.55	0.52	0.64	1.01	0.56	0.64	0.56	0.32

Table 9. Performance comparison across various training data in each scenario.

Training Data	E2D		H2D	
	L2 (m)	Collision (%)	L2 (m)	Collision (%)
nuScenes	0.61	0.56	1.29	1.77
nuScenes + GASS	0.52	0.50	1.07	1.74
nuScenes + HASS	0.55	0.54	0.67	0.96

Table 10. Performance comparison in H2D and E2D scenarios.

Method	Data	NC↑	DAC↑	TTC↑	Comf.↑	EP↑	PDMS↑
Human	-	100.0	100.0	100.0	99.9	87.5	94.8
Ego-MLP	NV	93.0	77.3	83.6	100.0	62.8	65.6
UniAD	NV	97.8	91.9	92.9	100.0	78.8	83.4
ParaDrive	NV	97.9	92.4	93.0	99.8	79.3	84.0
RoboTron-Sim	NV	98.0	93.0	93.3	99.8	79.9	84.6
RoboTron-Sim NV+HASS		98.2	93.6	93.8	99.9	81.1	85.6

Table 11. Performance on NAVSIM benchmark, † indicates that RoboTron-Sim is trained without HASS.

task. The VQA data curated for HASS encompasses three categories of questions: (1) Descriptive questions, such as “What is the color of the traffic light ahead?”, are answered directly using data generated by the simulator; (2) Hypothetical questions, such as “If you turn right at this intersection, what would you encounter?”, are annotated using GPT-4o based on environment visuals and predefined rules; (3) Reasoning questions, such as “Why are you slowing down here?”, are generated by GPT-4o based on driving videos and trajectories to enhance the understanding of vehicle behavior. We conduct separate validations on the BDD-X and LingoQA datasets. As shown in Table 12, with HASS integration, RoboTron-Sim achieves SOTA performance on both benchmarks (e.g., improving METEOR from 52.23 to 56.30 on BDD-X, and increasing CIDEr from 61.3 to 62.2 on LingoQA).

B.4. Model Compatibility

We conduct a comparative analysis of three models: VAD [49] (representing classical end-to-end models), LLaVA-OneVision [14] (as a representative multimodal large language model), and our RoboTron-Sim, evaluating their performance gains when augmenting real-world data with simulated data. To systematically investigate model

Method	Data	BLEU	METEOR	CIDEr
QwenVL	BDD-X	25.89	46.54	19.91
LLaVA-1.5	BDD-X	25.97	45.08	21.62
Senna	BDD-X	31.04	50.44	34.31
RoboTron-Sim	BDD-X	32.54	52.23	37.19
RoboTron-Sim	BDD-X+HASS	33.25	56.30	38.17
LLaVA	LingoQA	12.5	18.5	57.0
Vicuna-7B	LingoQA	10.1	15.2	51.0
BLIP-2	LingoQA	13.0	17.4	60.1
LingoQA	LingoQA	15.0	18.6	59.5
RoboTron-Sim	LingoQA	15.5	18.5	61.3
RoboTron-Sim	LingoQA+HASS	16.6	19.0	62.2

Table 12. Performance on NAVSIM benchmark, † indicates that RoboTron-Sim is trained without HASS.

compatibility with simulated data augmentation, we conduct cross-architecture evaluations on L2 distance. As evidenced in Table 13, VAD exhibits fundamental compatibility limitations, with marginal L2 reductions ($\downarrow 2.5\%$ E2D, $\downarrow 1.1\%$ H2D). Although MLLM demonstrates preliminary compatibility, showing gradual improvements, the gains remain constrained in the hard cases ($\downarrow 9.0\%$ E2D, $\downarrow 7.3\%$ H2D). In stark contrast, our RoboTron-Sim achieves breakthrough enhancements ($\downarrow 48.1\%$) in H2D case while maintaining stable performance in E2D case. This empowers knowledge transfer from synthetic domains while preserving real-world physical constraints, unlocking the model’s untapped potential.

B.5. Deployment Costs

We compare the key deployment metrics for the models on RTX-4090, as shown in Table 14. It shows that RoboTron-Sim is applicable to smaller models like RoboTron-Sim-0.5B (replacing the LLM from Qwen2-7B to Qwen1.5-0.5B), achieving comparable performance to RoboTron-Sim-7B and exhibiting deployment efficiency akin to traditional end-to-end model. This alignment of low deployment costs and performance improvement makes RoboTron-Sim practical for a wide range of real-world applications.

Method	Data	L2 Distance (m)	
		E2D	H2D
VAD	nuScenes	0.78	0.88
	nuScenes + HASS	0.76(↓ 2.5%)	0.87(↓ 1.1%)
MLLM	nuScenes	1.00	1.23
	nuScenes + HASS	0.91(↓ 9.0%)	1.14(↓ 7.3%)
RoboTron-Sim	nuScenes	0.61	1.29
	nuScenes + HASS	0.57(↓ 6.6%)	0.67(↓ 48.1%)

Table 13. L2 Distance performance gains of HASS across different models in E2D and H2D scenarios. To rigorously evaluate the model’s inherent capability to comprehend dynamic environments without relying on ego-pose dependencies, we conducted ablation studies by removing ego-pose inputs from both MLLM and RoboTron-Sim architectures.

Model	Latency	E2D			H2D		
		Day	Straight	Sunny	Night	Turn	Rainy
VAD	115.3ms	0.77	0.78	0.78	0.94	0.87	0.83
RoboTron-Sim-7B	612.8ms	0.54	0.55	0.56	0.81	0.64	0.56
RoboTron-Sim-0.5B	141.4ms	0.57	0.62	0.60	0.81	0.69	0.64

Table 14. Comparison of deployment costs.

C. Visualization

C.1. In Hard-to-Drive(H2D) Scenarios

We conduct trajectory visualization comparisons among the baseline method, RoboTron-Sim, and ground truth (GT) using representative long-tail cases from the nuScenes test set (including turn, night, and similar challenging scenarios), as shown in Figure 9.

C.2. In Long-Tail Scenarios

We conduct trajectory visualization comparisons among the baseline method, RoboTron-Sim, and ground truth (GT) using representative long-tail cases from the nuScenes test set (including lane invasion, temporary parking ahead, and similar challenging scenarios), as shown in Figure 10.



(a) Scenario 1



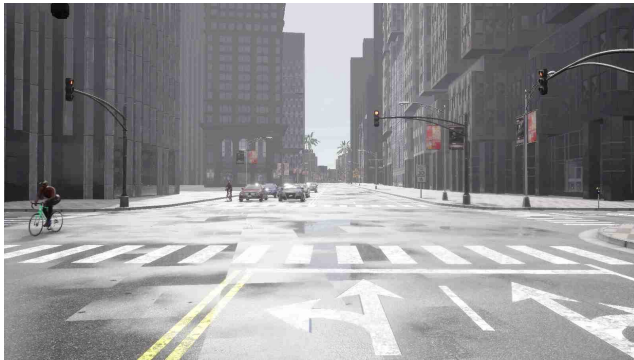
(b) Scenario 2



(c) Scenario 3



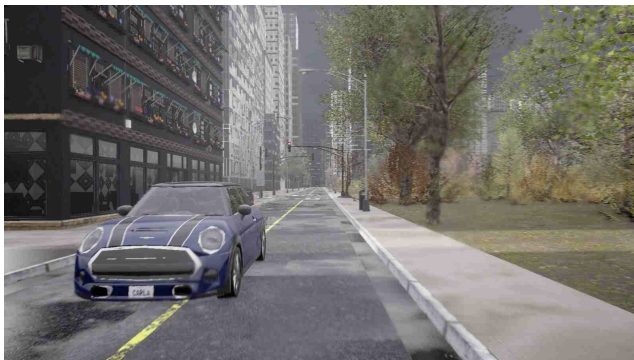
(d) Scenario 4



(e) Scenario 5



(f) Scenario 6

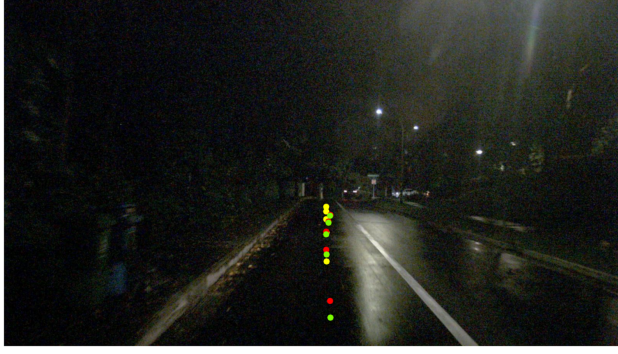


(g) Scenario 7



(h) Scenario 8

Figure 8. Visualization of HASS.



(a) Scenario 1



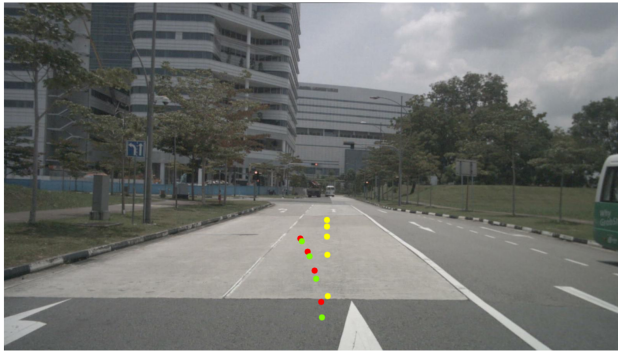
(b) Scenario 2



(c) Scenario 3



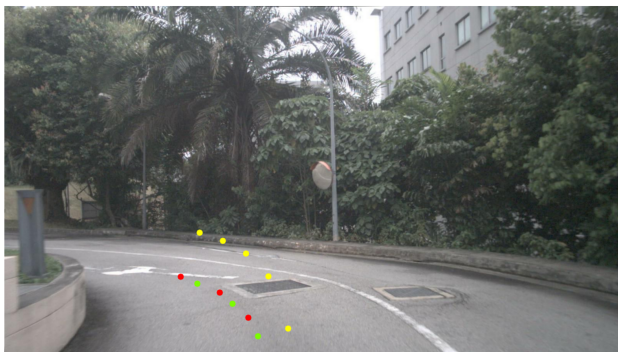
(d) Scenario 4



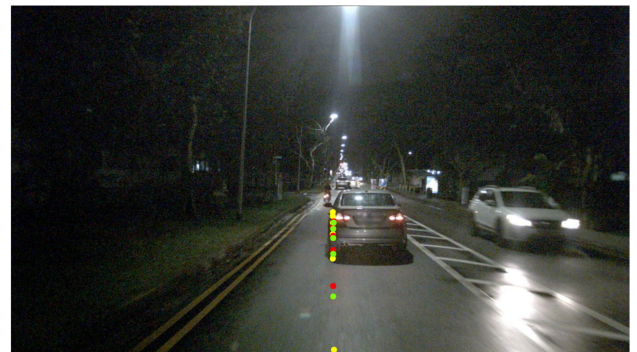
(e) Scenario 5



(f) Scenario 6



(g) Scenario 7



(h) Scenario 8

Figure 9. Visual comparison of planning trajectories in H2D scenarios. Ground-truth trajectories are marked in red, baseline predictions in yellow, and RoboTron-Sim's predictions in green.



(a) Scenario 1



(b) Scenario 2



(c) Scenario 3



(d) Scenario 4



(e) Scenario 5



(f) Scenario 6

Figure 10. Visual comparison of planning trajectories in Long-Tail scenarios. Ground-truth trajectories are marked in red, baseline predictions in yellow, and RoboTron-Sim's predictions in green.