# ORION: A Holistic End-to-End Autonomous Driving Framework by Vision-Language Instructed Action Generation

Haoyu Fu[1*], Diankun Zhang[2*], Zongchuang Zhao[1*], Jianfeng Cui[2], Dingkang Liang[1†],
Chong Zhang[2], Dingyuan Zhang[1], Hongwei Xie[2†], Bing Wang[2], Xiang Bai[1]

[1] Huazhong University of Science and Technology, [2] Xiaomi EV
{hyfu, zcuangzhao, dkliang}@hust.edu.cn
https://xiaomi-mlab.github.io/Orion/

## Abstract

*End-to-end (E2E) autonomous driving methods still struggle to make correct decisions in interactive closed-loop evaluation due to limited causal reasoning capability. Current methods attempt to leverage the powerful understanding and reasoning abilities of Vision-Language Models (VLMs) to resolve this dilemma. However, the problem is still open that few VLMs for E2E methods perform well in the closed-loop evaluation due to the gap between the semantic reasoning space and the purely numerical trajectory output in the action space. To tackle this issue, we propose ORION, a hOlistic E2E autonomous dRiving framework by vIsion-language instructed actiON generation. ORION uniquely combines a QT-Former to aggregate long-term history context, a Large Language Model (LLM) for driving scenario reasoning, and a generative planner for precision trajectory prediction. ORION further aligns the reasoning space and the action space to implement a unified E2E optimization for both visual question-answering (VQA) and planning tasks. Our method achieves an impressive closed-loop performance of 77.74 Driving Score (DS) and 54.62% Success Rate (SR) on the challenge Bench2Drive datasets, which outperforms state-of-the-art (SOTA) methods by a large margin of 14.28 DS and 19.61% SR.*

## 1. Introduction

End-to-end (E2E) autonomous driving has witnessed significant advancements in recent years. Classic E2E methods [8, 18, 25, 67, 70] integrate perception [27, 43, 66], prediction [7, 15, 50], and planning [17, 44] modules through multi-task learning, as shown in Fig. 1(a). These methods optimize driving trajectories by imitating expert demon-

---

* Equal contribution. † Project leader. Work done when Haoyu Fu and Zhongchuang Zhao were interns at Xiaomi EV.
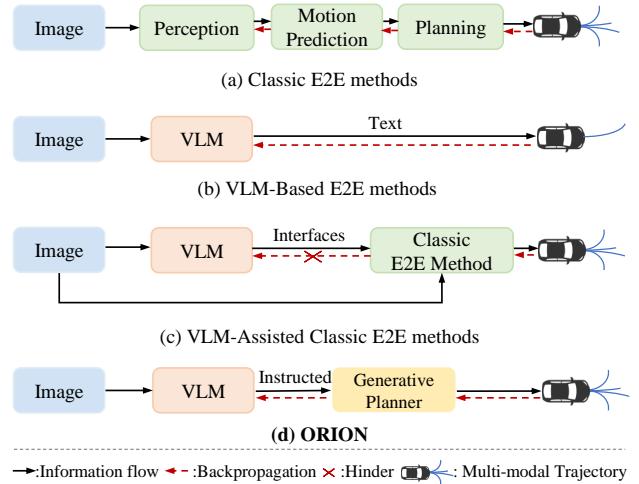


Figure 1. The comparison of different E2E paradigms. Our ORION framework establishes the differentiable connection between reasoning and action space via the generative planner.

strations, achieving promising performance in the open-loop evaluation [6, 53]. Nevertheless, these methods lack the common sense to complete complex causal reasoning. As a result, they struggle with comprehensive closed-loop benchmarks [23] that require autonomous decision-making and dynamic environmental interactions. Recently, Vision-Language Models (VLMs) [1, 10, 39, 57] have accumulated rich world knowledge and aligned vision-language space through large-scale data training, providing new insight for achieving E2E autonomous driving.

Despite these advances, leveraging VLMs for E2E autonomous driving is not trivial, as the capabilities of VLMs focus on the semantic reasoning space, while E2E methods only need the numerical planning results in the action space. Some methods [19, 49, 59, 60, 63] attempt to directly output text-based planning results using VLM, as shown in

Fig. 1(b). Although this paradigm is convenient, VLM is not well-suited for handling mathematical calculations or numerical reasoning [13, 42]. Besides, limited by the intrinsic autoregressive mechanism of VLM, this framework only infers single results, which is inconsistent with the natural uncertainty of human planning [8]. Therefore, directly using VLM for E2E autonomous driving may produce suboptimal solutions in complex scenes [62]. Other methods endeavor to bridge the gap via utilizing VLM output meta-action (*e.g.*, turn left) to assist classic E2E methods [26, 40], as shown in Fig. 1(c). They adopt a carefully crafted interface to transmit the reasoning space information into the action space. However, this paradigm decouples these two spaces, hindering collaborative optimization between the trajectory optimization and the VLM reasoning process. Thus, the capabilities of VLM for E2E planning are not fully leveraged by the above framework.

To tackle this problem, we propose a h**O**listic E2E autonomous d**R**iving framework by v**I**sion-language instructed acti**ON** generation, termed ORION. Inspired by the field of conditional generation [28, 38, 47, 48], where the semantic information controls the generation of detailed image features, we find that the generative model can construct a unified distribution of diverse types of data (*e.g.*, image, text). Therefore, considering that the reasoning space of VLM and the action space of trajectory belong to different domains, we introduce a generative planner to establish a unified latent representation for aligning the two spaces. With the help of the introduced module, we take advantage of VLMs' reasoning information to construct trajectory, facilitating the model to capture the causal relationship between scene information and driving behavior.

Furthermore, it is well-known that long-term memory is necessary for E2E autonomous driving since historical information often influences trajectory planning within the current scene. Existing VLMs for E2E methods [19, 62] typically concatenate multi-frame images for temporal modeling. They are constrained by the token length of VLM and incur significant computational overhead. Instead, motivated by OmniDrive [59], which extracts features through Q-Former-styled architecture, we introduce QT-Former, a query-based temporal module. By leveraging a memory bank and a set of history queries, QT-Former effectively stores and extracts essential historical scene information to aggregate long-term visual context, further enhancing the temporal perception ability of reasoning and action space.

We evaluate the closed-loop driving ability of ORION on the Bench2Drive dataset, which builds interactive scenarios based on the CARLA [11] simulator. ORION achieves 77.74 Driving Score (DS) and 54.62% Success Rate (SR), surpassing previous SOTA methods [24] with 14.28 driving scores and 19.61% success rates, demonstrating the powerful superiority of ORION.

**The benefits of ORION are from three aspects:** 1) Thanks to the capability of the generative model to characterize the latent distribution of data, we bridge the gap between the reasoning space of VLM and the action space of trajectories through a generative planner, enabling the VLM to understand the scene and instruct trajectory generation. 2) The QT-former in ORION effectively captures long-term temporal dependencies, enabling the model to integrate temporal vision context into reasoning and action spaces. 3) Without bells and whistles, ORION achieves excellent performance in the Bench2Drive closed-loop benchmark. Experiments also show that ORION is compatible with diverse generative models, which further demonstrate the flexibility of our proposed framework.

## 2. Related work

### 2.1. End-to-End Autonomous Driving

End-to-end autonomous driving [61, 68] aims to directly process raw sensor data to predict motion trajectories or control signals, jointly optimizing the entire system to minimize error accumulation. Recent works like UniAD [18] and VAD [25] integrate perception, prediction, and planning into a unified planning framework, making the framework ultimately planning-oriented. VADv2 [8] introduces probabilistic planning, outputting the probabilistic distribution of action and sampling one action to control the vehicle. GenAD [70] and DiffusionDrive [32] explore a new paradigm for end-to-end autonomous driving, employing the generative model to predict multi-modal trajectories. However, these methods mainly excel in open-loop evaluation, where the model could readily overfit to the ego status, as highlighted in Ego-MLP [64] and BEV-Planner [31]. Although some studies [8, 21, 22, 70] adopt closed-loop evaluation in CARLA [11] to assess robust driving ability, their performance remains suboptimal, revealing a notable gap between their open-loop and closed-loop results. Thus, we aim to construct an E2E autonomous driving system that demonstrates excellent performance in both open-loop and closed-loop evaluations.

### 2.2. Vision-Language Models (VLMs)

VLMs [1, 3, 10, 30, 35, 57] introduce visual information to large language models (LLMs) [39, 55] through various vision encoders [45, 65], demonstrating powerful visual contextual understanding and reasoning. LLaVA series [35, 36] employ visual instruction tuning to perform image-text alignment. Monkey [30] improves detail comprehension by dividing images. InternVL series [9, 10] further enhances the vision detail understanding via a dynamic resolution strategy. However, most methods map the numerous visual tokens into language space through MLP, incurring high computational costs. To alleviate this bur-
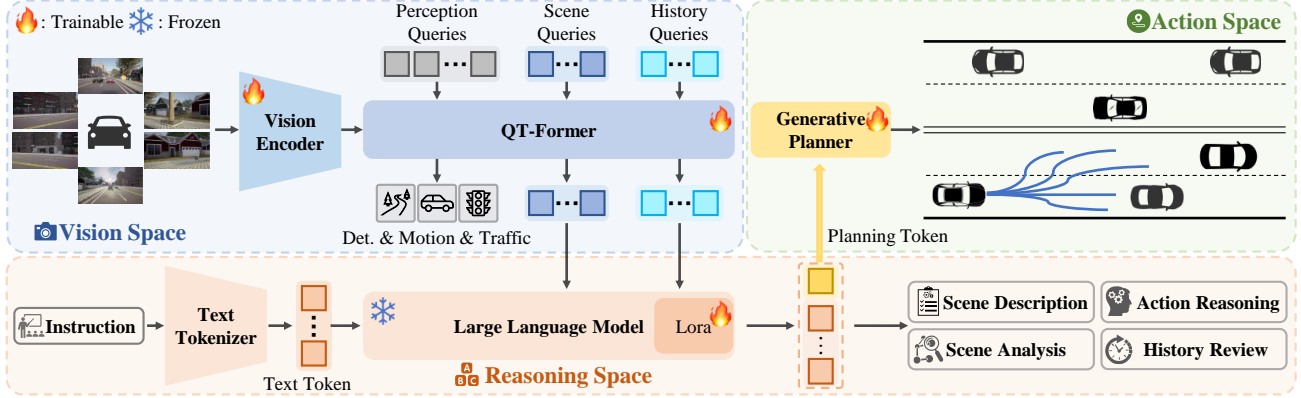
Figure 2. The pipeline of our ORION, a holistic E2E framework aligning vision-reasoning-action space. It consists of three key components: a QT-Former to extract long-term context and link the vision space of the vision encoder and the reasoning space of LLM; the LLM for performing reasoning tasks and predicting a planning token; and a generative planner that bridges reasoning and action space for generating a multi-modal trajectory conditioned by the planning token.

den, QwenVL [4] and Flamingo [2] reduce token redundancy using cross-attention, while Qwen2VL [57] enhances efficiency with dynamic resolution and multimodal rotary position embedding (M-RoPE) for simultaneously processing diverse modalities. Inspired by these, we introduce QT-Former, which leverages a set of queries and cross-attention operations to extract multi-view image features.

## 2.3. VLM for End-to-End Autonomous Driving

VLMs showcase excellent contextual understanding and comprehensive world knowledge, motivating their application in autonomous driving. Some methods [19, 59, 62] directly employ VLMs for environment perception and explainable trajectory prediction in text form. For example, Omnidrive [59] adopts StreamPETR [58] as Q-Formar3D to compress current scene features and connect the vision-reasoning space and then performs textual trajectory prediction. EMMA [19], trained on large-scale data, enables Gemini [3] to predict discrete textual planning with strong open-loop performance. Other studies [26, 54] integrate VLMs with representative E2E models in a fast-slow dual system. DriveVLM [54] leverages VLM to predict the low-frequency trajectory, which will be refined by an E2E model. Senna [26] further replaces the low-frequency with the meta-action, guiding the VAD [25] to predict motion. These methods only implement the open-loop evaluation. Although DriveMLM [60] and LMDrive [49] leverage the VLM to implement closed-loop evaluation, they struggle with processing complex scenarios limited by the simple CARLA Town05Long benchmark. In contrast, we propose a holistic E2E framework that employs a generative planner to bridge the reasoning space of VLM and the action space of trajectories, generating precise trajectories with interpretable action decisions in complex real-world driving scenarios of Bench2Drive.

## 3. Method

In this paper, we propose a h**O**listic end-to-end autonomous d**R**iving framework by v**I**sion-language model instructed acti**ON** generation, termed ORION. The pipeline of our ORION is shown in Fig. 2. Specifically, given the multi-view images of the current scene, the ORION first encodes the image tokens with a vision encoder. Then, QT-Former (Sec. 3.1) leverages diverse queries to aggregate long-term vision context, compress image tokens, and perceive traffic elements. The LLM (Sec. 3.2) subsequently combines the compressed scene features and historical vision information with user instructions, performing diverse understanding and reasoning tasks and generating a planning token. Finally, a generative planner (Sec. 3.3) bridges the reasoning space of LLM and the action space of trajectories, predicting multi-modal trajectory conditioned by the planning token. ORION effectively aligns the vision-reasoning-action space through these core components, achieving the collaborative optimization of scene understanding and trajectory generation in a unified space.

## 3.1. QT-Former

To achieve long-term information modeling while compressing and extracting multi-view image features $F_m$ derived from the vision encoder, we introduce QT-Former, a query-based temporal module, as shown in Fig. 3. Specifically, following Q-Former3D [59], we first set up two types of learnable queries, the scene queries $Q_s \in \mathbb{R}^{N_s \times C_q}$ and the perception queries $Q_p \in \mathbb{R}^{N_p \times C_q}$, where $N_s$ and $N_p$ are the number of scene and perception queries, respectively, and $C_q$ is the channel of queries. $Q_s, Q_p$ are processed through self-attention (SA) to exchange their information. Then they interact with image features $F_m$ with 3D positional encoding [37] $P_m$ in the cross-attention (CA) mod-
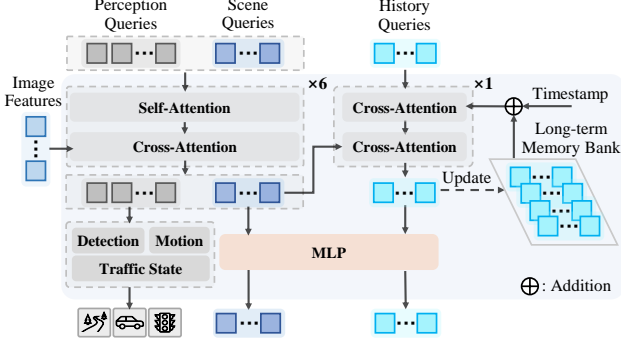
Figure 3. The detailed architecture of QT-Former. It accepts diverse queries and image features as inputs to detect traffic elements, predict motion, and aggregate long-term vision context.

ule. After that, the perception queries are fed into the multiple auxiliary heads for object detection (*e.g.*, critical objects and lanes), traffic state, and motion prediction of dynamic agents. The scene queries serve as tokens representing the key information of the current scene.

Additionally, we employ a set of history queries $Q_h \in \mathbb{R}^{N_h \times C_q}$ and a long-term memory bank $M \in \mathbb{R}^{(N_h \times n) \times C_q}$ to efficiently retrieve and store essential historical information (*e.g.*, previous road conditions and ego status), where $N_h$ is the number of history queries and $n$ is the maximum history frame length. We utilize the $Q_h$ to extract the former frame queries in $M$ with relative timestamp embedding $P_t$ through a CA block. Then $Q_h$ interacts with current scene features $Q_s$ in another CA block, enabling the extraction of relevant details about the current scenario. This process can be formulated as:

$$Q_h = \text{CA}(Q_h, M + P_t, M + P_t),$$
$$\hat{Q}_h = \text{CA}(Q_h, Q_s, Q_s),$$
(1)

where $P_t$ denotes the relative timestamp embedding.

Subsequently, the updated history queries $\hat{Q}_h$ are stored in the memory bank $M$ following the First-In-First-Out (FIFO) replacement policy, formulated as:

$$M = [\hat{Q}_h^{t-n}, \cdots, \hat{Q}_h^{t-1}, \hat{Q}_h^t],$$
(2)

where $t$ is the current frame time.

Although some methods [51, 58] also leverage the memory bank to store preceding information, they typically only store the compressed historical information without guiding for extracting the current scene information. Instead, we initialize a few numbers of the history queries to further extract the current scene features that are most closely related to historical information, enhancing the long-term memory ability of the model.

Finally, we utilize a two-layer MLP to convert the updated history queries $\hat{Q}_h$ and current scene features $Q_s$

to history tokens $x_h$ and scene tokens $x_s$ in the reasoning space of LLM.

## 3.2. Large Language Model

The LLM is pivotal in our framework because the high-quality reasoning of the current driving scenario is necessary to instruct the generative planner to generate a reasonable trajectory in action space.

As shown in Fig. 2, the user instruction is first encoded into language tokens $x_q \in \mathbb{R}^{L \times C}$ by the text tokenizer, where $L$ is the token length and $C$ is the dimension of LLM. Then, the scene tokens $x_s$ and history tokens $x_h$ are combined with the language tokens $x_q$ and fed into LLM.

Leveraging the abundant world knowledge and outstanding reasoning ability of LLM, ORION performs various text-based understanding and reasoning tasks in the driving scenario, including scene description, history information review, scene analysis, and action reasoning. Meanwhile, we design a planning QA template with a special planning token $s$ for LLM as the final QA to accumulate the understanding and reasoning context of the entire driving scenario to the $s$, formally written as:

$$s \sim p(s|x_s, x_h, x_q, x_a),$$
(3)

where $x_a$ denotes the generation answer of LLM. The embedding of the planning token $s$ will serve as a condition to control the trajectory generation.

However, there is still a lack of high-quality VQA annotations within closed-loop simulation environments to train LLMs for comprehensively understanding driving scenarios. Thus, we extend the Bench2Drive dataset via a fully automatic VQA annotation pipeline powered by Qwen2-VL [57] and propose our VQA dataset, Chat-B2D, expecting to further promote the research of VLM on closed-loop simulation. We provide detailed information on Chat-B2D and its annotation pipeline in the Appendix.

## 3.3. Generative Planner

Generative models [14, 28, 48] can effectively capture intrinsic features within data by learning the distribution of the data. Recent researches [5, 38, 47] have demonstrated semantic correlations between latent spaces of different modalities of data, where adjusting the distribution parameters of one modality space enables precise control over the generation process of another modality space.

Inspired by the generative domain, we introduce a generative planner to bridge the gap between the reasoning and action space. Specifically, we formulate the current trajectory $a$ in action space as a conditional probability distribution $p(a|s)$, where $s$ is the planning token. To construct $p(a|s)$, there are many excellent methods in the generation field (*e.g.*, variational autoencoders (VAE) [28] and diffusion model [48]).

4

As there are essential differences in the distribution between the reasoning space of VLM and the action space of trajectory, we use the VAE [28] model to align them in the Gaussian distribution. We employ two-layer MLPs to project both the state $s$ and the ground-truth trajectory $t$ into Gaussian variables $z$ in the latent space, denoted as:

$$p(z_s|s) \sim N(\mu_s, \sigma_s^2), p(z_t|t) \sim N(\mu_t, \sigma_t^2), \quad (4)$$

where $N(\mu, \sigma^2)$ denotes a Gaussian distribution with a mean of $\mu$, and standard deviation of $\sigma$. We then use Kullback-Leibler divergence loss to enforce distribution matching, represented as:

$$\mathcal{L}_{vae} = D_{KL}(p(\mathbf{z}|\mathbf{s}), p(\mathbf{z}|\mathbf{t})). \quad (5)$$

Finally, we use the GRU decoder in GenAD [70] to decode the trajectory from the latent space $z$. Significantly, the functions of VAE in this paper are not the same as VAE of GenAD. We only use a single token encoded in the reasoning space from the perspective of the ego vehicle as input, aiming to bridge the gap between reasoning space and action space. In contrast, the latter leverages features of all agents encoded in the BEV space as input, designed to learn specific patterns of the highly structured trajectories of both the ego vehicle and other agents.

Additionally, we also attempt to replace the VAE with alternative generative models, such as the diffusion model for trajectory generation. Benefiting from the proposed method that bridges the gap between the reasoning and action space through distribution learning in latent space, our framework still demonstrates superior performance compared to other methods (detailed in Sec. 4.5).

### 3.4. Training Objectives

For the detection task of the proposed QT-Former, the detection loss is defined as $\mathcal{L}_{det} = \mathcal{L}_{cls} + \mathcal{L}_{reg}$, where $\mathcal{L}_{cls}$ is focal loss [34] and $\mathcal{L}_{reg}$ is L1 loss. For the traffic state and motion prediction, the losses are defined as $\mathcal{L}_{tra}$ and $\mathcal{L}_m = \mathcal{L}_{mcls} + \mathcal{L}_{mreg}$, respectively, where $\mathcal{L}_{tra}$ and $\mathcal{L}_{mcls}$ are focal loss, and $\mathcal{L}_{mreg}$ is L1 loss. The total loss of QT-Former is:

$$\mathcal{L}_{qt} = \mathcal{L}_{det} + \mathcal{L}_{tra} + \mathcal{L}_m. \quad (6)$$

For the LLM, we leverage the auto-regressive cross-entropy loss $\mathcal{L}_{ce}$. For the generative planner in our framework, $\mathcal{L}_{vae}$ is the Kullback-Leibler divergence loss used to align the reasoning space and action space. Following VAD [25], we adopt the collision loss $\mathcal{L}_{col}$, boundary loss $\mathcal{L}_{bd}$, and MSE loss $\mathcal{L}_{mse}$ for the planning prediction. The total loss of the generative planner is:

$$\mathcal{L}_{gp} = \mathcal{L}_{vae} + \mathcal{L}_{mse} + \mathcal{L}_{col} + \mathcal{L}_{bd}. \quad (7)$$

In summary, the total loss of the proposed ORION is:

$$\mathcal{L} = \mathcal{L}_{qt} + \mathcal{L}_{ce} + \mathcal{L}_{gp}. \quad (8)$$

## 4. Experiments

### 4.1. Dataset and Evaluation Metrics

**Dataset.** We train and evaluate ORION on the Bench2drive dataset [23], a closed-loop evaluation protocol under CARLA V2 [11] for E2E autonomous driving. It provides an official training set where we use the base set (1000 clips) for fair comparison with all the other baselines, which is divided into 950 clips for training and 50 clips for open-loop validation. Each clip captures approximately 150 meters of continuous driving within a specific traffic scene. For closed-loop evaluation, we evaluate the proposed method on the official set of 220 short routes designed by Bench2drive, spanning 44 interactive scenarios with 5 routes per scenario. Additionally, we compare our method with other SOTA baselines on nuScenes [6] open-loop evaluation, which will be provided in the Appendix.

**Evaluation Metrics.** Bench2drive includes five metrics for closed-loop evaluation: Driving Score (DS), Success Rate (SR), Efficiency, Comfortness, and Multi-Ability. The Success Rate quantifies the proportion of routes successfully completed within the allotted time. The Driving Score follows CARLA [11], incorporating both route completion status and violation penalties, where infractions reduce the score via discount factors. Efficiency and Comfortness are used to measure the speed performance and comfort of the autonomous driving system during the driving process, respectively. Multi-Ability measures 5 advanced skills independently for urban driving. For open-loop evaluation, we use the L2 distance error and the collision rate. Additionally, we use CIDEr [56], BLEU [41], and ROUGE-L [33] to evaluate the performance of ORION on VQA tasks.

### 4.2. Implementation Details

**Model Setting.** Consistent with classic E2E baselines [18, 25, 70] on Bench2Drive, ORION is a fully HD map-free method that only uses the Navigation Command (NC) as an input condition for the trajectory predictions rather than locations of lane center (*i.e.*, target point, TP). ORION is an anchor-free method that outputs 6 mode trajectory predictions corresponding to the 6 NC defined in Bench2drive.

**Training Process.** All experiments are conducted on 32 NVIDIA A800 GPUs with 80 GB of memory. Following Omnidrive [59], we adopt EVA-02-L [12] as the vision encoder. Vicuna v1.5 [69] is employed in ORION and fine-tuned using LoRA [16], with the rank dimension and alpha set to 16. The default number of scene, perception, and historical queries is 512, 600, and 16, respectively. We set the Memory Bank's stored frame number $n$ to 16. During training, data augmentations are applied to input images, which are first resized to a resolution of $640 \times 640$. More training details are provided in the Appendix.

Table 1. Closed-loop and Open-loop Results of E2E-AD Methods in Bench2Drive under base set. C/L refers to camera/LiDAR. Avg. L2 is averaged over the predictions in 2 seconds under 2Hz, similar to UniAD. * denote expert feature distillation. NC: navigation command, TP: target point, DS: Driving Score, SR: Success Rate.

| Method | Reference | Condition | Modality | Closed-loop Metric | | | | Open-loop Metric |
| | | | | DS↑ | SR(%)↑ | Efficiency↑ | Comfortness↑ | Avg. L2 ↓ |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| TCP* [61] | NeurIPS 22 | TP | C | 40.70 | 15.00 | 54.26 | 47.80 | 1.70 |
| TCP-ctrl* | NeurIPS 22 | TP | C | 30.47 | 7.27 | 55.97 | 51.51 | - |
| TCP-traj* | NeurIPS 22 | TP | C | 59.90 | 30.00 | 76.54 | 18.08 | 1.70 |
| TCP-traj w/o distillation | NeurIPS 22 | TP | C | 49.30 | 20.45 | 78.78 | 22.96 | 1.96 |
| ThinkTwice* [22] | CVPR 23 | TP | C | 62.44 | 31.23 | 69.33 | 16.22 | 0.95 |
| DriveAdapter* [21] | ICCV 23 | TP | C&L | 64.22 | 33.08 | 70.22 | 16.01 | 1.01 |
| AD-MLP [64] | arXiv 23 | NC | C | 18.05 | 0.00 | 48.45 | 22.63 | 3.64 |
| UniAD-Tiny [18] | CVPR 23 | NC | C | 40.73 | 13.18 | 123.92 | 47.04 | 0.80 |
| UniAD-Base [18] | CVPR 23 | NC | C | 45.81 | 16.36 | 129.21 | 43.58 | 0.73 |
| VAD [25] | ICCV 23 | NC | C | 42.35 | 15.00 | 157.94 | 46.01 | 0.91 |
| GenAD [70] | ECCV 24 | NC | C | 44.81 | 15.90 | - | - | - |
| MomAD[52] | CVPR25 | NC | C | 44.54 | 16.71 | 170.21 | 48.63 | 0.87 |
| DriveTransformer-Large [24] | ICLR 25 | NC | C | 63.46 | 35.01 | 100.64 | 20.78 | **0.62** |
| ORION(**Ours**) | - | NC | C | **77.74**(+14.28) | **54.62**(+19.61) | 151.48 | 17.38 | 0.68 |

Table 2. Multi-Ability Results of E2E-AD Methods under base set. * denote expert feature distillation. C/L refers to camera/LiDAR. NC: navigation command, TP: target point.

| Method | Reference | Condition | Modality | Ability (%) ↑ | | | | | |
| | | | | Merging | Overtaking | Emergency Brake | Give Way | Traffic Sign | Mean |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| TCP* [61] | NeurIPS 22 | TP | C | 16.18 | 20.00 | 20.00 | 10.00 | 6.99 | 14.63 |
| TCP-ctrl* | NeurIPS 22 | TP | C | 10.29 | 4.44 | 10.00 | 10.00 | 6.45 | 8.23 |
| TCP-traj* | NeurIPS 22 | TP | C | 8.89 | 24.29 | 51.67 | 40.00 | 46.28 | 34.22 |
| TCP-traj w/o distillation | NeurIPS 22 | TP | C | 17.14 | 6.67 | 40.00 | **50.00** | 28.72 | 28.51 |
| ThinkTwice* [22] | CVPR 23 | TP | C | 27.38 | 18.42 | 35.82 | **50.00** | 54.23 | 37.17 |
| DriveAdapter* [21] | ICCV 23 | TP | C&L | **28.82** | 26.38 | 48.76 | **50.00** | 56.43 | 42.08 |
| AD-MLP [64] | arXiv 23 | NC | C | 0.00 | 0.00 | 0.00 | 0.00 | 4.35 | 0.87 |
| UniAD-Tiny [18] | CVPR 23 | NC | C | 8.89 | 9.33 | 20.00 | 20.00 | 15.43 | 14.73 |
| UniAD-Base [18] | CVPR 23 | NC | C | 14.10 | 17.78 | 21.67 | 10.00 | 14.21 | 15.55 |
| VAD [25] | ICCV 23 | NC | C | 8.11 | 24.44 | 18.64 | 20.00 | 19.15 | 18.07 |
| DriveTransformer-Large [24] | ICLR 25 | NC | C | 17.57 | 35.00 | 48.36 | 40.00 | 52.10 | 38.60 |
| ORION (**Ours**) | - | NC | C | 25.00 | **71.11** | **78.33** | 30.00 | **69.15** | **54.72**(+16.12) |

## 4.3. Main Results

As reported in Tab. 1, the performance of ORION significantly exceeds all E2E methods on Bench2Drive, even the method with expert feature distillation. Specifically, ORION surpasses the latest SOTA method DriveTransformer [24] by +14.28 DS and +19.61% SR. It also achieves improvements of +13.52 DS and +21.54% SR over DriveAdapter [21], even if DriveAdapter distills the expert feature from Think2Drive [29] and leverages two modalities (*i.e.*, camera and LiDAR) inputs. The above promising results effectively demonstrate the superiority of our ORION.

Additionally, the Multi-Ability results are also illustrated in Tab. 2. ORION achieves +16.12% and +12.64% performance improvements compared with DriveTransformer [24] and DriveAdapter [21] in the mean ability, respectively. Specifically, our model demonstrates outstand-ing performance in some scenarios, such as Overtaking (71.11%), Emergency Brake (78.33%), and Traffic Sign (69.15%), which shows that our model benefits from the powerful reasoning capability of VLM to understand the causal interaction between the ego vehicle, dynamic elements, and static elements (Traffic Signs) in driving scenarios. On the other hand, our model falls behind DriveAdapter in Merging and Give Way, which shows that ORION is not good at making lane-changing decisions. The phenomenon may be caused by the more diverse decision-making timing for lane-changing, making the model encounter difficulties in capturing the correct causal relationship [21].

## 4.4. Qualitative Results

The qualitative results of ORION in two canonical closed-loop evaluation scenarios of Bench2Drive are shown in
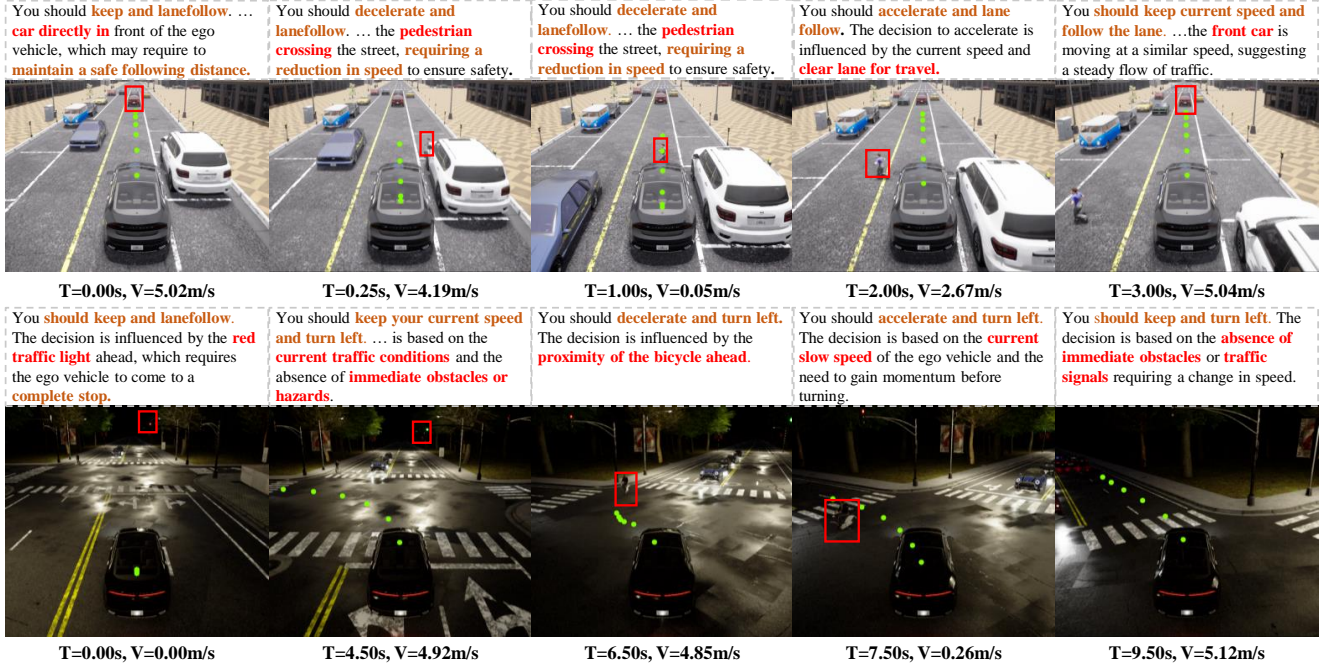
Figure 4. Qualitative results of ORION on the Bench2Drive closed-loop evaluation set. The brown, red, and green refer to the action decision, the objects that influence driving decisions, and the prediction trajectory, respectively.
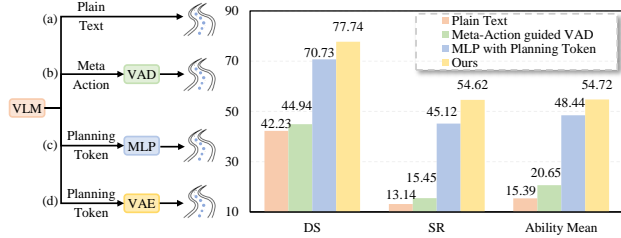


Figure 5. Advantages of the vision-language instructed action generation. DS and SR denote Driving Score and Success Rate separately. VAD [25] is a classic E2E model.

Fig. 4. It shows both the driving action reasoning and trajectory prediction outputted by our model, as well as the corresponding ego-vehicle states. We observe that ORION can capture the correct causal relationship in the scenario and make correct driving decisions, then predict the planning trajectory following the reasoning instruction, demonstrating the surprising interpretability of our method. More qualitative results can be found in the Appendix.

## 4.5. Ablation Study

**Advantages of the vision-language instructed action generation.** To validate the effectiveness of the planning generation paradigm proposed in this paper, extensive experiments are conducted to compare our paradigm with canonical trajectory prediction paradigms of VLM-based

E2E autonomous driving methods, including (a) plain text outputs [19, 59], (b) dual-system paradigm which classic E2E methods(*e.g.*, VAD [25]) output trajectory guided by elaborated design VLM interface (*e.g.*, meta-action) [26], and (c) special token decode outputs by MLP [46], as shown in the left part of Fig. 5. To ensure the fairness of the ablations, experiments of different paradigms use the same sensor inputs, vision encoder, QT-former, and VLM as our ORION and are trained by the same strategy. Only the output formats of VLMs are adjusted according to the requirements of different paradigms.

The results are illustrated in the right part of Fig. 5. The plain text paradigm performs the worst (42.23 DS, 13.14% SR, and 15.39% mean ability), indicating the limitations of plain text output in closed-loop driving scenarios, potentially due to its inadequate numerical reasoning capabilities [13, 42]. Compared with the plain text paradigm, the dual-system paradigm only obtains a slight performance improvement. Note that the reproduced results of the dual-system paradigm are very close to the official results of VAD in Tab. 1. This result may indicate that the performance of the dual-system paradigm may be bottlenecked by the insufficient capabilities of classic E2E methods. Although the effectiveness of the MLP decoder paradigm has been validated in CarLLaVA [46], our paradigm still shows a performance gain of +7.01 DS, +9.5% SR, and +6.28% mean ability. The result may be caused by the fact that the MLP is the simplest way to align features between different

7

Table 3. Ablation on diverse generative planner. DS and SR denote Driving Score and Success Rate separately.

| Generative Planner | Closed-loop | | Open-loop | | Ability |
| | DS↑ | SR(%)↑ | Avg. L2 (m) ↓ | Avg. col (%)↓ | Avg. |
|---|---|---|---|---|---|
| Diffusion | 71.97 | 46.54 | 0.73 | 0.96 | 46.68 |
| VAE (Ours) | 77.74 | 54.62 | 0.68 | 0.47 | 54.72 |

Table 4. Ablation on QT-Former designs in different frameworks. DS and SR denote Driving Score and Success Rate separately. Traffic state means using explicit traffic state supervision. T: Plain Text, G: Instructed Generator

| ID | Traffic State | Motion Pred. | Memory Bank | Output type | | Closed-loop | |
| | | | | T | G | DS ↑ | SR(%) ↑ |
|---|---|---|---|---|---|---|---|
| 1 | | | | | ✓ | 56.33 | 26.05 |
| 2 | ✓ | | | | ✓ | 74.65 | 49.31 |
| 3 | ✓ | ✓ | | | ✓ | 74.07 | 49.77 |
| 4 | ✓ | ✓ | ✓ | | ✓ | 77.74 | 54.62 |
| 5 | | | | ✓ | | 25.45 | 10.38 |
| 6 | ✓ | ✓ | ✓ | ✓ | | 42.23 | 13.14 |

Table 5. Ablation of history queries number. DS and SR denote Driving Score and Success Rate separately.

| Query Num. $N_h$ | Closed-loop | | Open-loop | |
| | DS ↑ | SR(%) ↑ | Avg. L2 (m) ↓ | Avg. col (%)↓ |
|---|---|---|---|---|
| 0 | 65.10 | 38.83 | 0.67 | 0.61 |
| 8 | 68.09 | 39.09 | 0.66 | 0.62 |
| 16 | 74.10 | 44.66 | 0.68 | 0.55 |
| 32 | 62.46 | 37.73 | 0.65 | 0.73 |

spaces, which is consistent with the viewpoint presented in this paper. Additionally, the MLP decoder struggles with handling multi-modal trajectory [8, 20], making it still significantly lag behind ORION in closed-loop evaluation.

**Analysis on different generative planners.** We then investigate the effect of employing different generative planners to bridge the reasoning-action space. Specifically, we implement the diffusion model by simply replacing the VAE, which uses K-means trajectory anchors as prior information and outputs 20 mode trajectory predictions. The results are listed in Tab. 3. Note that the VAE-based trajectory generator demonstrates a significant performance improvement over the diffusion-based. We argue the main reasons are as follows: 1) Compared with the conditional denoising process of diffusion, the latent space of VAE more directly and effectively aligns the reasoning information of VLM to the multi-modal action space. 2) The training process of VAE is inherently more stable, facilitating better alignment between the reasoning and action spaces. Surprisingly, even using diffusion, ORION still surpasses the DriveTransformer by +8.51 DS, +11.53% SR, and +8.08% mean ability. This impressive result emphasizes the effectiveness and flexibility of our framework.

**Effectiveness of QT-Former designs.** Tab. 4 shows the detailed ablations of each design in the introduced QT-Former. By leveraging explicit traffic state supervision (ID-2), ORION achieves 74.65 DS and 49.31% SR, which already outperforms DriveAdapter [21] and DriveTransformer [24] by a large margin and makes an improvement of +18.32 and +23.26% compared with the baseline (ID-1). This is because a better understanding of

traffic signals helps ORION directly reduce infractions in closed-loop evaluation. It is worth noting that due to the causal confusion [21], it's not trivial for previous methods to fully understand the corresponding causal relationships by simply introducing traffic state supervision, especially when encountering mixed expert behaviors before traffic signs [21, 22, 24, 61]. This result also proves that ORION can better utilize the reasoning ability of VLM to capture the causal relationship between scene information and driving behavior by aligning reasoning space and action space. This conclusion can also be verified by the results in Tab. 2, where ORION shows a significant advantage in traffic sign ability (+17.05%) compared to previous E2E methods [24].

Then, we combine the motion prediction module in the QT-Former's perception head, which gains a slight improvement of +0.4% SR and further reduces the collision rate. The slight degradation on DS may be caused by the trade-off between DS and SR in the CARLA benchmark protocol [71]. Involving a memory bank into QT-Former and supervised by QA pairs about historical information leads to an increase of +3.67 DS and +4.85% SR and boosts the final performance to 77.74 DS and 54.62% SR, which demonstrates our model can benefit from the long-temporal memory of vision tokens.

We also apply QT-former to the plain text output type (ID-6). By leveraging it, we improve the model's performance by +16.78 DS and +2.78% SR over the baseline (ID-5). Meanwhile, with the same QT-former designs, our ORION achieves further improvements of +35.51 DS and +41.48% SR compared with the plain text output mode, demonstrating the effectiveness of our framework.

**Influence of the number of history queries.** We conduct ablation experiments to further study the influence of different numbers of history queries. Here, to accelerate the training process, we only train the model using the planning trajectory and history QA pairs without other auxiliary VQA tasks. The results are detailed in Tab. 5. Increasing the history query number $N_h$ from 0 to 8 brings a significant performance boost of around 2.99 DS and 0.26% SR. Further increasing $N_h$ from 8 to 16 leads to the sweet point that achieves the best performance of 74.10 DS and 44.66% SR. However, enlarging $N_h$ from 16 to 32 shows a significant performance degradation. We argue that introducing

Table 6. Effectiveness of auxiliary VQA task training. DS and SR denote Driving Score and Success Rate separately. C/B/R refers to CIDEr/BLEU/ROUGE-L. FT: Fine Tuning

| ID | VQA FT | Planning FT | Closed-loop | | Language | | | Open-loop |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | DS ↑ | SR(%) ↑ | C↑ | B↑ | R↑ | Avg. L2 (m) ↓ |
| 1 | ✓ | | - | - | 65.65 | 50.82 | 77.65 | - |
| 2 | | ✓ | 74.10 | 44.66 | - | - | - | 0.68 |
| 3 | ✓ | ✓ | 77.74 | 54.62 | 65.77 | 52.49 | 77.58 | 0.68 |

more history queries hinders the VLM from capturing the current frame features, which are more essential than historical information in the driving scene.

**Influence between VQA task training and planning task training.** As shown in Tab. 6. The model cannot obtain both reasoning and planning capabilities with single-task training. Surprisingly, when we perform two tasks simultaneously during training, ORION achieves better performance in both planning and language metrics compared to single-task training. Specifically, the multi-task training leads to improvements of +3.64 DS and +9.66% SR in the planning task, as well as a performance gain of +0.12 CIDEr, +1.67 BLEU and competitive performance of ROUGE-L in the VQA tasks. Furthermore, the results also validate the high quality and validity of the Chat-B2D dataset produced by our auto-pipeline.

## 5. Conclusion

In this paper, we mainly focus on the challenges faced by VLM methods for end-to-end autonomous driving in aligning the reasoning space of VLM with the pure numerical action space used for planning. This dilemma makes it not trivial for existing methods to simultaneously analyze the driving scenario and output high-quality multimodal prediction trajectories. To address this problem, we propose ORION, a holistic end-to-end autonomous driving framework by vision-language instructed action generation. By leveraging a generative planner and incorporating long-term visual context, we effectively bridge the vision-reasoning-action space. Extensive experiments validate the flexibility and superiority of our proposed framework, where ORION demonstrates significant improvements in closed-loop planning evaluation, surpassing SOTA methods.

**Limitation.** Although ORION performs well in the closed-loop simulation environment on Bench2Drive [23], it is limited by the high computational complexity of the scalable VLM in real-time driving scenarios. In the future, we would like to reduce the complexity of ORION through techniques such as model compression and pruning, thereby enabling the model to achieve real-time autonomous driving.

## References

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 1, 2

[2] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. In *Proc. of Advances in Neural Information Processing Systems*, 2022. 3

[3] Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: A family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 1, 2023. 2, 3

[4] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023. 3

[5] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science.*, page 8, 2023. 4

[6] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 11621–11631, 2020. 1, 5, 2

[7] Yuning Chai, Benjamin Sapp, Mayank Bansal, and Dragomir Anguelov. Multipath: Multiple probabilistic anchor trajectory hypotheses for behavior prediction. *arXiv preprint arXiv:1910.05449*, 2019. 1

[8] Shaoyu Chen, Bo Jiang, Hao Gao, Bencheng Liao, Qing Xu, Qian Zhang, Chang Huang, Wenyu Liu, and Xinggang Wang. Vadv2: End-to-end vectorized autonomous driving via probabilistic planning. *arXiv preprint arXiv:2402.13243*, 2024. 1, 2, 8

[9] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *Science China Information Sciences*, page 220101, 2024. 2

[10] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 24185–24198, 2024. 1, 2

[11] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. Carla: An open urban driving simulator. In *Conf. on Robot Learning*, pages 1–16, 2017. 2, 5

[12] Yuxin Fang, Quan Sun, Xinggang Wang, Tiejun Huang, Xinlong Wang, and Yue Cao. Eva-02: A visual representa-

tion for neon genesis. *Image and Vision Computing*, page 105171, 2024. 5

[13] Simon Frieder, Luca Pinchetti, Ryan-Rhys Griffiths, Tommaso Salvatori, Thomas Lukasiewicz, Philipp Petersen, and Julius Berner. Mathematical capabilities of chatgpt. *Proc. of Advances in Neural Information Processing Systems*, 36, 2024. 2, 7

[14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, pages 139–144, 2020. 4

[15] Junru Gu, Chenxu Hu, Tianyuan Zhang, Xuanyao Chen, Yilun Wang, Yue Wang, and Hang Zhao. Vip3d: End-to-end visual trajectory prediction via 3d agent queries. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 5496–5506, 2023. 1

[16] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *Proc. of Intl. Conf. on Learning Representations*, 2022. 5, 2

[17] Shengchao Hu, Li Chen, Penghao Wu, Hongyang Li, Junchi Yan, and Dacheng Tao. St-p3: End-to-end vision-based autonomous driving via spatial-temporal feature learning. In *Proc. of European Conference on Computer Vision*, pages 533–549, 2022. 1

[18] Yihan Hu, Jiazhi Yang, Li Chen, Keyu Li, Chonghao Sima, Xizhou Zhu, Siqi Chai, Senyao Du, Tianwei Lin, Wenhai Wang, et al. Planning-oriented autonomous driving. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 17853–17862, 2023. 1, 2, 5, 6

[19] 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, 2, 3, 7

[20] Bernhard Jaeger, Kashyap Chitta, and Andreas Geiger. Hidden biases of end-to-end driving models. In *Porc. of IEEE Intl. Conf. on Computer Vision*, pages 8240–8249, 2023. 8

[21] Xiaosong Jia, Yulu Gao, Li Chen, Junchi Yan, Patrick Langechuan Liu, and Hongyang Li. Driveadapter: Breaking the coupling barrier of perception and planning in end-to-end autonomous driving. In *Porc. of IEEE Intl. Conf. on Computer Vision*, 2023. 2, 6, 8

[22] Xiaosong Jia, Penghao Wu, Li Chen, Jiangwei Xie, Conghui He, Junchi Yan, and Hongyang Li. Think twice before driving: Towards scalable decoders for end-to-end autonomous driving. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2023. 2, 6, 8

[23] Xiaosong Jia, Zhenjie Yang, Qifeng Li, Zhiyuan Zhang, and Junchi Yan. Bench2drive: Towards multi-ability benchmarking of closed-loop end-to-end autonomous driving. *Proc. of Advances in Neural Information Processing Systems*, 2024. 1, 5, 9, 2

[24] Xiaosong Jia, Junqi You, Zhiyuan Zhang, and Junchi Yan. Drivetransformer: Unified transformer for scalable end-to-end autonomous driving. In *Proc. of Intl. Conf. on Learning Representations*, 2025. 2, 6, 8

[25] Bo Jiang, Shaoyu Chen, Qing Xu, Bencheng Liao, Jiajie Chen, Helong Zhou, Qian Zhang, Wenyu Liu, Chang Huang, and Xinggang Wang. Vad: Vectorized scene representation for efficient autonomous driving. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 8340–8350, 2023. 1, 2, 3, 5, 6, 7

[26] Bo Jiang, Shaoyu Chen, Bencheng Liao, Xingyu Zhang, Wei Yin, Qian Zhang, Chang Huang, Wenyu Liu, and Xinggang Wang. Senna: Bridging large vision-language models and end-to-end autonomous driving. *arXiv preprint arXiv:2410.22313*, 2024. 2, 3, 7

[27] Xiaohui Jiang, Shuailin Li, Yingfei Liu, Shihao Wang, Fan Jia, Tiancai Wang, Lijin Han, and Xiangyu Zhang. Far3d: Expanding the horizon for surround-view 3d object detection. In *Proc. of the AAAI Conf. on Artificial Intelligence*, pages 2561–2569, 2024. 1

[28] Diederik P Kingma. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 2, 4, 5

[29] Qifeng Li, Xiaosong Jia, Shaobo Wang, and Junchi Yan. Think2drive: Efficient reinforcement learning by thinking with latent world model for autonomous driving (in carla-v2). In *Proc. of European Conference on Computer Vision*, pages 142–158, 2024. 6

[30] Zhang Li, Biao Yang, Qiang Liu, Zhiyin Ma, Shuo Zhang, Jingxu Yang, Yabo Sun, Yuliang Liu, and Xiang Bai. Monkey: Image resolution and text label are important things for large multi-modal models. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 26763–26773, 2024. 2

[31] Zhiqi Li, Zhiding Yu, Shiyi Lan, Jiahan Li, Jan Kautz, Tong Lu, and Jose M Alvarez. Is ego status all you need for open-loop end-to-end autonomous driving? In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 14864–14873, 2024. 2

[32] Bencheng Liao, Shaoyu Chen, Haoran Yin, Bo Jiang, Cheng Wang, Sixu Yan, Xinbang Zhang, Xiangyu Li, Ying Zhang, Qian Zhang, et al. Diffusiondrive: Truncated diffusion model for end-to-end autonomous driving. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2025. 2

[33] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Proc. Annual Meeting of the Association for Computational Linguistics Workshop*, pages 74–81, 2004. 5

[34] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Porc. of IEEE Intl. Conf. on Computer Vision*, pages 2980–2988, 2017. 5

[35] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *Proc. of Advances in Neural Information Processing Systems*, 2023. 2

[36] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, 2024. 2

[37] Yingfei Liu, Tiancai Wang, Xiangyu Zhang, and Jian Sun. Petr: Position embedding transformation for multi-view 3d object detection. In *Proc. of European Conference on Computer Vision*, pages 531–548, 2022. 3

[38] Yixin Liu, Kai Zhang, Yuan Li, Zhiling Yan, Chujie Gao, Ruoxi Chen, Zhengqing Yuan, Yue Huang, Hanchi Sun, Jianfeng Gao, et al. Sora: A review on background, technology, limitations, and opportunities of large vision models. *arXiv preprint arXiv:2402.17177*, 2024. 2, 4

[39] Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, et al. Deepseek-vl: towards real-world vision-language understanding. *arXiv preprint arXiv:2403.05525*, 2024. 1, 2

[40] Jianbiao Mei, Yukai Ma, Xuemeng Yang, Licheng Wen, Xinyu Cai, Xin Li, Daocheng Fu, Bo Zhang, Pinlong Cai, Min Dou, et al. Continuously learning, adapting, and improving: A dual-process approach to autonomous driving. In *Proc. of Advances in Neural Information Processing Systems*, 2024. 2

[41] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proc. Annual Meeting of the Association for Computational Linguistics*, pages 311–318, 2002. 5

[42] Shuai Peng, Ke Yuan, Liangcai Gao, and Zhi Tang. Mathbert: A pre-trained model for mathematical formula understanding. *arXiv preprint arXiv:2105.00377*, 2021. 2, 7

[43] Jonah Philion and Sanja Fidler. Lift, splat, shoot: Encoding images from arbitrary camera rigs by implicitly unprojecting to 3d. In *Proc. of European Conference on Computer Vision*, pages 194–210, 2020. 1

[44] Aditya Prakash, Kashyap Chitta, and Andreas Geiger. Multimodal fusion transformer for end-to-end autonomous driving. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 7077–7087, 2021. 1

[45] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *Proc. of Intl. Conf. on Machine Learning*, pages 8748–8763, 2021. 2

[46] Katrin Renz, Long Chen, Ana-Maria Marcu, Jan Hünermann, Benoit Hanotte, Alice Karnsund, Jamie Shotton, Elahe Arani, and Oleg Sinavski. Carllava: Vision language models for camera-only closed-loop driving. *arXiv preprint arXiv:2406.10165*, 2024. 7

[47] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 2, 4

[48] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Proc. of Intl. Conf. on Medical Image Computing and Computer Assisted Intervention*, pages 234–241, 2015. 2, 4

[49] Hao Shao, Yuxuan Hu, Letian Wang, Guanglu Song, Steven L Waslander, Yu Liu, and Hongsheng Li. Lmdrive: Closed-loop end-to-end driving with large language models. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 15120–15130, 2024. 1, 3

[50] Shaoshuai Shi, Li Jiang, Dengxin Dai, and Bernt Schiele. Motion transformer with global intention localization and local movement refinement. *Proc. of Advances in Neural Information Processing Systems*, 35:6531–6543, 2022. 1

[51] Enxin Song, Wenhao Chai, Guanhong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Haozhe Chi, Xun Guo, Tian Ye, Yanting Zhang, et al. Moviechat: From dense token to sparse memory for long video understanding. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 18221–18232, 2024. 4

[52] Ziying Song, Caiyan Jia, Lin Liu, Hongyu Pan, Yongchang Zhang, Junming Wang, Xingyu Zhang, Shaoqing Xu, Lei Yang, and Yadan Luo. Don't shake the wheel: Momentum-aware planning in end-to-end autonomous driving. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2025. 6

[53] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 2446–2454, 2020. 1

[54] 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. 3, 2

[55] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 2

[56] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 4566–4575, 2015. 5

[57] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 1, 2, 3, 4

[58] Shihao Wang, Yingfei Liu, Tiancai Wang, Ying Li, and Xiangyu Zhang. Exploring object-centric temporal modeling for efficient multi-view 3d object detection. In *Porc. of IEEE Intl. Conf. on Computer Vision*, pages 3621–3631, 2023. 3, 4

[59] Shihao Wang, Zhiding Yu, Xiaohui Jiang, Shiyi Lan, Min Shi, Nadine Chang, Jan Kautz, Ying Li, and Jose M Alvarez. Omnidrive: A holistic llm-agent framework for autonomous driving with 3d perception, reasoning and planning. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2024. 1, 2, 3, 5, 7

[60] Wenhai Wang, Jiangwei Xie, ChuanYang Hu, Haoming Zou, Jianan Fan, Wenwen Tong, Yang Wen, Silei Wu, Hanming Deng, Zhiqi Li, et al. Drivemlm: Aligning multi-modal large language models with behavioral planning states for autonomous driving. *arXiv preprint arXiv:2312.09245*, 2023. 1, 3

[61] Penghao Wu, Xiaosong Jia, Li Chen, Junchi Yan, Hongyang Li, and Yu Qiao. Trajectory-guided control prediction for end-to-end autonomous driving: A simple yet strong baseline. In *Proc. of Advances in Neural Information Processing Systems*, 2022. 2, 6, 8

[62] Shuo Xing, Chengyuan Qian, Yuping Wang, Hongyuan Hua, Kexin Tian, Yang Zhou, and Zhengzhong Tu. Openemma: Open-source multimodal model for end-to-end autonomous driving. In *Proc. of IEEE Winter Conf. on Applications of Computer Vision*, pages 1001–1009, 2025. 2, 3

[63] Zhenhua Xu, Yujia Zhang, Enze Xie, Zhen Zhao, Yong Guo, Kwan-Yee K Wong, Zhenguo Li, and Hengshuang Zhao. Drivegpt4: Interpretable end-to-end autonomous driving via large language model. *IEEE Robotics and Automation Letters*, 2024. 1

[64] Jiang-Tian Zhai, Ze Feng, Jinhao Du, Yongqiang Mao, Jiang-Jiang Liu, Zichang Tan, Yifu Zhang, Xiaoqing Ye, and Jingdong Wang. Rethinking the open-loop evaluation of end-to-end autonomous driving in nuscenes. *arXiv preprint arXiv:2305.10430*, 2023. 2, 6

[65] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Porc. of IEEE Intl. Conf. on Computer Vision*, pages 11975–11986, 2023. 2

[66] Diankun Zhang, Zhijie Zheng, Haoyu Niu, Xueqing Wang, and Xiaojun Liu. Fully sparse transformer 3-d detector for lidar point cloud. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–12, 2023. 1

[67] Diankun Zhang, Guoan Wang, Runwen Zhu, Jianbo Zhao, Xiwu Chen, Siyu Zhang, Jiahao Gong, Qibin Zhou, Wenyuan Zhang, Ningzi Wang, et al. Sparsead: Sparse query-centric paradigm for efficient end-to-end autonomous driving. *arXiv preprint arXiv:2404.06892*, 2024. 1

[68] Zhejun Zhang, Alexander Liniger, Dengxin Dai, Fisher Yu, and Luc Van Gool. End-to-end urban driving by imitating a reinforcement learning coach. In *Porc. of IEEE Intl. Conf. on Computer Vision*, 2021. 2

[69] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Proc. of Advances in Neural Information Processing Systems*, 36:46595–46623, 2023. 5

[70] Wenzhao Zheng, Ruiqi Song, Xianda Guo, Chenming Zhang, and Long Chen. Genad: Generative end-to-end autonomous driving. In *Proc. of European Conference on Computer Vision*, pages 87–104, 2024. 1, 2, 5, 6

[71] Julian Zimmerlin, Jens Beißwenger, Bernhard Jaeger, Andreas Geiger, and Kashyap Chitta. Hidden biases of end-to-end driving datasets. *arXiv preprint arXiv:2412.09602*, 2024. 8

# ORION: A Holistic End-to-End Autonomous Driving Framework by Vision-Language Instructed Action Generation

## Supplementary Material

We provide supplementary material to complement the main paper, arranged as follows:

- Appendix A: Details on the Chat-B2D dataset.
- Appendix B: Traning Details.
- Appendix C: More results.

## A. Details on the Chat-B2D dataset

To compensate for the absence of a high-quality scene text annotation dataset and promote the application of VLM in the closed-loop simulated driving scenario, we carefully design an automated annotation pipeline to extend the Bench2Drive dataset [23] to support VQA pairs, named Chat-B2D, covering diverse tasks.

### A.1. Data Annotation Pipeline

As shown in Fig. A1, the automated annotation pipeline consists of three steps:

**Critical object selection.** Unlike mainline self-driving perception modules that process all detected objects equally, we emphasize identifying the crucial object that potentially affects the ego vehicle's driving behavior, grounded in human driving strategies. Our selection criteria include: 1) Objects have potential collisions within three seconds. 2) Leading vehicles in current and adjacent lanes. 3) Active traffic signals. 4) The vulnerable road users (VRUs), such as pedestrians/cyclists.

**Description generation.** We extract video clips comprising the current and five preceding frames. Subsequently, these clips, along with the ego vehicle's status and the ground truth information (*e.g.*, 2D/3D coordinates and velocity, *etc.*) of selected crucial objects, serve as input to Qwen2VL-72B [57] for multi-task generation: 1) the scene description; 2) attributes of key objects and their impact on the ego vehicle; 3) operational meta-commands and action reasoning for autonomous navigation.

**History Information.** During the generation process, we incorporate a queue mechanism to preserve essential historical information. The stored information comprises two principal components: 1) Environmental dynamics that capture spatiotemporal variations of critical scene elements, and 2) Ego-motion characteristics derived from comparative analysis between current speed/action and their historical counterparts across previous frames.

The generated description and collected historical information are combined with predefined question templates to create VQA pairs. Tab. A3 displays the detailed crafted prompt, and Tab. A4 shows the question templates.
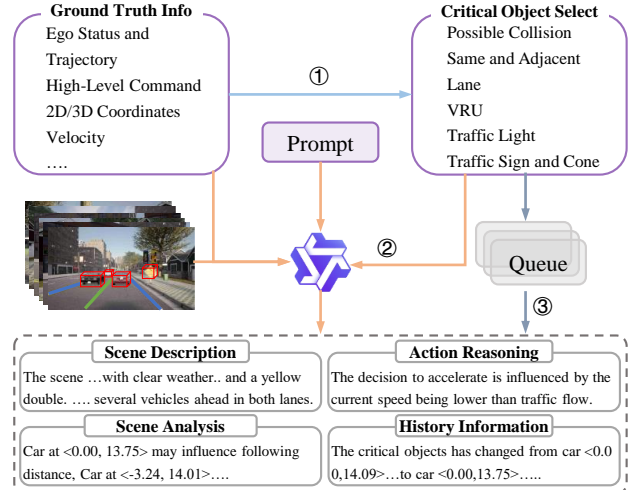


Figure A1. The automated annotation pipeline for the Chat-B2D dataset.

### A.2. Chat-B2D Attribute

Through the carefully crafted prompts and the above generation pipeline, we have automatically conducted a large-scale, high-quality VQA dataset for the Bench2Drive [23], creating Chat-B2D. This dataset, including a total of 2.11M VQA pairs for training and 0.12M for validation, supports four primary categories: 1) Scene description, which provides a comprehensive overview of the driving scenarios, including weather, time of day, traffic situations, and road characteristics. 2) Behavior description of critical objects detailing their current state and intentions. 3) Meta-driving decisions and action reasoning of the ego car, such as turning left and lane following. 4) Recall of essential historical information.

## B. Training Details

To accelerate the alignment of the vision-reasoning-action space and gradually enhance the reasoning and planning capabilities of our ORION, we adopt a three-stage training strategy. In each stage, the model inherits the weights from the previous stage and continues training. Additionally, we train the model for six epochs per stage with a total batch size of 32. The three-stage training strategy is as follows:

**3D Vision-Language Alignment:** In this first stage, we primarily train the QT-Former and the VLM while freezing the generative planner. By training on VQA pairs from

Table A1. Comparison of the Open-loop planning in nuScene. †: The ego status and planning trajectory are both processed by LLM in textual modality. ‡: The high-level command is not used during the training and testing phases.

| Method | VLM-Based | Ego Status | | L2 (m) ↓ | | | | Collision (%) ↓ | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | BEV | Planner | 1s | 2s | 3s | Avg. | 1s | 2s | 3s | Avg. |
| ST-P3 | - | - | - | 1.33 | 2.11 | 2.90 | 2.11 | 0.23 | 0.62 | 1.27 | 0.71 |
| UniAD [18] | - | - | - | 0.48 | 0.96 | 1.65 | 1.03 | 0.05 | 0.17 | 0.71 | 0.31 |
| UniAD | - | ✓ | ✓ | 0.20 | 0.42 | 0.75 | 0.46 | 0.02 | 0.25 | 0.84 | 0.37 |
| VAD-Base [25] | - | - | - | 0.69 | 1.22 | 1.83 | 1.25 | 0.06 | 0.68 | 2.52 | 1.09 |
| VAD-Base | - | ✓ | - | 0.41 | 0.70 | 1.06 | 0.72 | 0.04 | 0.43 | 1.15 | 0.54 |
| VAD-Base | - | ✓ | ✓ | 0.17 | 0.34 | 0.60 | 0.37 | 0.04 | 0.27 | 0.67 | 0.33 |
| Ego-MLP [64] | - | - | ✓ | 0.15 | 0.32 | 0.59 | 0.35 | 0.00 | 0.27 | 0.85 | 0.37 |
| BEV-Planner [31] | - | - | - | 0.30 | 0.52 | 0.83 | 0.55 | 0.10 | 0.37 | 1.30 | 0.59 |
| BEV-Planner++ | - | ✓ | ✓ | 0.16 | 0.32 | 0.57 | 0.35 | 0.00 | 0.29 | 0.73 | 0.34 |
| DriveVLM† [54] | ✓ | - | - | 0.18 | 0.34 | 0.68 | 0.40 | 0.10 | 0.22 | 0.45 | 0.27 |
| DriveVLM-Dual [54] | ✓ | ✓ | - | 0.15 | 0.29 | 0.48 | 0.31 | 0.05 | 0.08 | 0.17 | 0.10 |
| OmniDrive‡ [59] | ✓ | - | - | 1.15 | 1.96 | 2.84 | 1.98 | 0.80 | 3.12 | 7.46 | 3.79 |
| OmniDrive | ✓ | - | - | 0.40 | 0.80 | 1.32 | 0.84 | 0.04 | 0.46 | 2.32 | 0.94 |
| OmniDrive++ | ✓ | ✓ | ✓ | 0.14 | 0.29 | 0.55 | 0.33 | 0.00 | 0.13 | 0.78 | 0.30 |
| Senna [26] | ✓ | - | - | 0.37 | 0.54 | 0.86 | 0.59 | 0.09 | 0.12 | 0.33 | 0.18 |
| Senna | ✓ | ✓ | ✓ | **0.11** | **0.21** | **0.35** | **0.22** | 0.04 | **0.08** | **0.13** | **0.08** |
| EMMA† [19] | ✓ | - | - | 0.14 | 0.29 | 0.54 | 0.32 | - | - | - | - |
| ORION (**Ours**) | ✓ | ✓ | - | 0.17 | 0.31 | 0.55 | 0.34 | 0.05 | 0.25 | 0.80 | 0.37 |

Chat-B2D, we focus on aligning the vision space with the reasoning space.

**Language-Action Alignment:** In this stage, we unfreeze the generative planner and train the entire model except for the LLM, which is trained by LoRA [16], to predict planning trajectories without auxiliary VQA pairs. This stage primarily focuses on transmitting world knowledge from the reasoning space to the action space.

**End-to-End Fine-tuning:** We follow the training settings from the previous stage, with the only difference being the incorporation of joint training on VQA and planning tasks. This step further facilitates the alignment of the vision-reasoning-action space.

## C. More Results

### C.1. Experiments on nuScenes dataset

**nuScenes Dataset.** nuScenes [6] is a popular autonomous driving benchmark typically used for detection and open-loop planning evaluation. The dataset contains 1000 scenes from Singapore and Boston, with 700 scenes for training, 150 scenes for validation, and 150 scenes for testing. Each scene spans 20 seconds and is annotated at 2 Hz. nuScenes utilizes the L2 error and collision rate as planning metrics.

**Results on nuScenes.** We compare the ORION with previous SOTA end-to-end autonomous driving methods on the nuScenes dataset. Here, for a fair comparison with other VLM-Based methods, we modify ORION by replacing QT-Former with the Q-Former from OmniDrive [59], and without the explicit ego status in the generative planner. As shown in Tab. A1, our ORION achieves comparable perfor-

Table A2. Ablation study of training strategy. V/L/A indicates vision/language/action space. DS and SR denote Driving Score and Success Rate separately. C/B/R refers to CIDEr/BLEU/ROUGE-L.

| ID | V→L | L→A | V→L→A | Closed-loop | |
| --- | --- | --- | --- | --- | --- |
| | | | | DS ↑ | SR(%) ↑ |
| 1 | | ✓ | | 57.96 | 26.32 |
| 2 | ✓ | ✓ | | 65.10 | 38.83 |
| 3 | ✓ | ✓ | ✓ | 74.65 | 49.31 |

mance to classic SoTA methods [18, 25, 31] without VLM. However, compared with other VLM-Based methods, our ORION is suboptimal. We argue that this is due to the latent space of VAE being more suitable for multimodal trajectory distributions of Bench2Drive [23]. In contrast, the nuScene dataset follows a uni-modal Gaussian distribution (with straight trajectories accounting for about 70%).

Additionally, as highlighted in BEV-Planner [31] and Ego-MLP [64], even a simple MLP decoder with ego status can achieve strong open-loop planning performance on nuScenes. Thus, in the main paper, we primarily focus on evaluating ORION's closed-loop performance on the Bench2Drive dataset.

### C.2. More Ablation Studies on Bench2Drive

**Ablation of training pipeline.** To facilitate the vision-language-action space alignment of our model, we implement a progressive space alignment training strategy. We validate the effectiveness of the training pipeline, and the
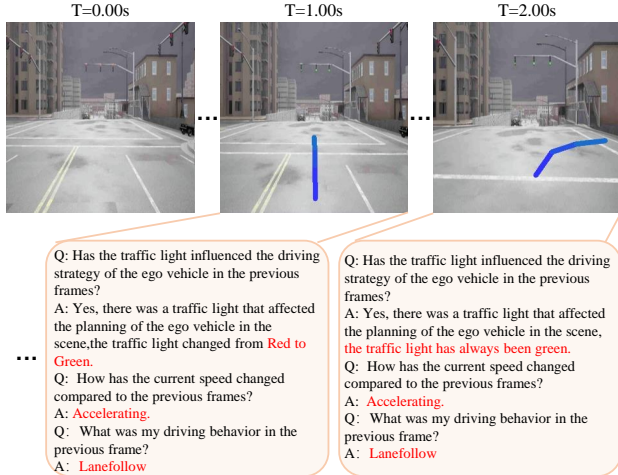
Figure A2. Qualitative results of historical information memory and retrieval on Bench2Drive open-loop validation set.

results are presented in Tab. A2. Here, the QT-Former of our model does not incorporate collision loss or long-term memory bank with history queries. Specifically, through our second-stage training(ID-2), ORION achieves a significant improvement by +7.14 DS and +12.51% SR compared to direct training planning without the first stage (ID-1). After completing the third-stage training (ID-3), our model further improved the performance and achieved optimal (74.65 DS and 49.32 SR), demonstrating the effectiveness of our training strategy.
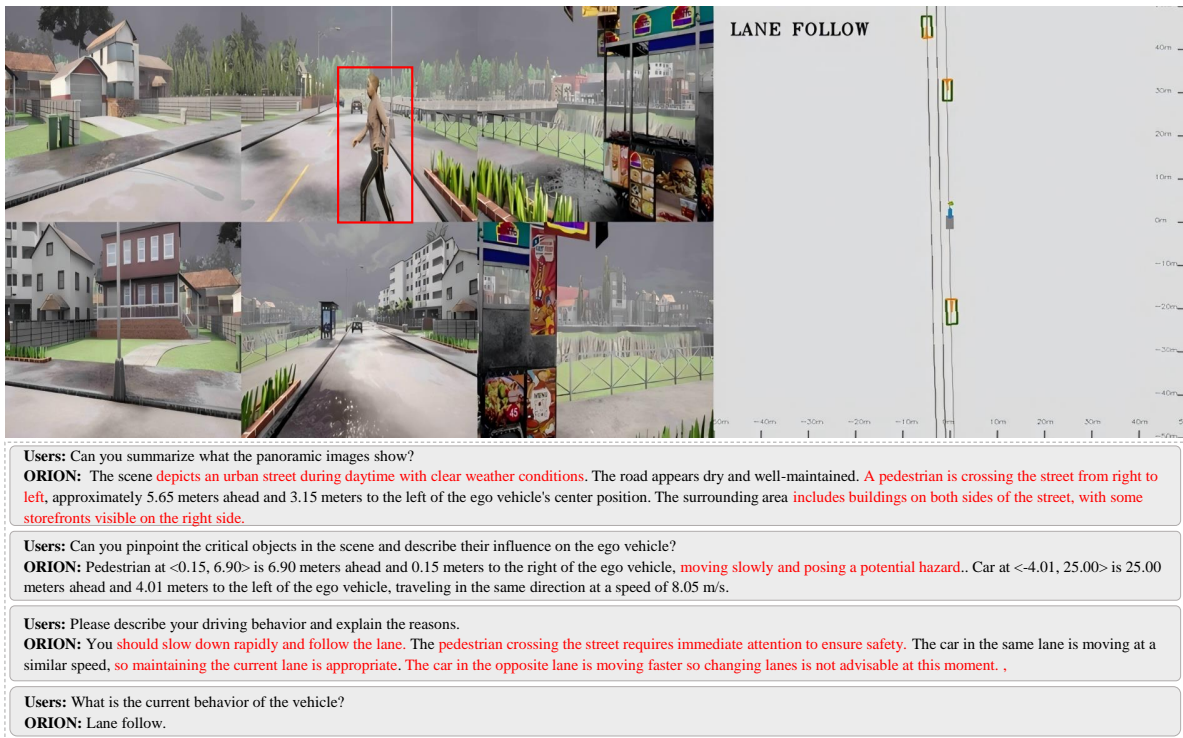
## C.3. More Qualitative Results

**Historical information memory and retrieval.** Benefiting from the introduced long-term memory bank and history queries in QT-Former, our ORION could store and retrieve historical information, as illustrated in Fig. A2. Our model could perceive critical elements (*e.g.*, traffic light) changes in previous and current times.

**Scene understanding and action reasoning.** Fig. A3 shows scene understanding and action reasoning results of ORION. It could be observed that ORION could not only accurately perceive detailed scene information but also identify key objects influencing the ego vehicle's behavior and infer appropriate motion decisions. Even in extreme situations (e.g., a pedestrian suddenly crosses the road in Fig. A3(b)), our model maintains robust performance, highlighting its superior reasoning and decision-making ability.

(a)



(b)

Figure A3. Qualitative results for scene understanding and action reasoning on Bench2Drive open-loop validation. From top to bottom, each sub-figure displays the multi-view input and traffic conditions in Bird's Eye View (BEV) of the current scene, the scene understanding, and the reasoning result. The red rectangles indicate the critical objects influencing the action of the ego vehicle, while the red text highlights our method's correct scene comprehension.

Table A3. Prompts fed into Qwen2VL to generate corresponding response.

**Prompt 1: Scene Description**
Suppose you are driving, generate a description of the driving scene which includes the key factors for driving planning, including the traffic conditions, weather, time of day and road conditions, traffic signs, and traffic lights that affect the driving of the ego vehicle if it exists, indicating smooth surfaces or the presence of obstacles; The description should be concise, and accurate to facilitate informed decision-making. Please make sure the traffic light colors you provide are accurate; otherwise, give 'unknown.'

---

**Prompt 2: Critical Objects Analysis**
I will provide you with several critical objects that are most important to my short-term command in the last image of the video. I provide you with 2d coordinates, which are two points of the top-left and bottom-right coordinates, and the 3d position and velocity information of these critical objects: {objects_desc}. Please describe their action and explain why they are most important, including their speed, position, heading, and influence on ego vehicle. Please associate these objects with the objects in the image and please remember the ego vehicle is located at the **center of the bottom edge of the picture**.

---

**Prompt 3: Expert Meta-Decision**
Besides, I will provide you speed, historical trajectory and future driving behaviors of ego vehicle, which can be divided into SPEED decisions and COMMAND decisions, SPEED includes keep, accelerate, decelerate, while COMMAND includes left, right, straight, lane follow, change lane left, change lane right. Your current speed is {ego_vel} m/s, historical trajectory is {ego_his_trajs}. The next SPEED decision is {speed_decision}, the next COMMAND decision is {path_decision}. Please analyze the reasons for the future driving behaviors or the reason why ego vehicle can {path_decision} based on the driving environment, including the behavior of other traffic participants, especially the critical objects, road conditions, and traffic light status.

---

**Example:**
You should refer to the following example and format the results like {"description": "xxx","critical_objects": "xxx", "action": "{speed_decision} and {path_decision}"}}:
{{ "description": "The scene captures a moment of urban life framed by a red traffic light in mid-transition. To the right, a pedestrian crossing, ..., waiting for the signal to change. Directly ahead, ... On the left, the sidewalk bustles with people of all ages, ... Behind this foreground of orderly traffic and pedestrian movement, ..."
"critical_objects: "["Car at <-0.24, 7.56> directly in front of ego vehicle, ...", "Car at <-2.64, 10.00> ..., moving at a slower speed, may influence left change."]"
"action": "Slow down and right lane change. - The decision to change lanes is influenced by the need to overtake Car at <-0.24, 7.56> in front of the ego vehicle. - There are no traffic lights for the vehicle,... - Pedestrians are visible on the sidewalk to the right, ..." }}
If it has no critical_objects, you should refer to the following example and format the results like {{"description": "xxx", "critical_objects": [], "action": "xxx"}}.

Table A4. A list of question templates for diverse VQA tasks.

**Type 1: Scene Description**
1. What can you tell about the current driving conditions from the images?
2. What can be observed in the panoramic images provided?
3. Can you provide a summary of the current driving scenario based on the input images?
4. What can you observe from the provided images regarding the driving conditions?
5. Please describe the current driving conditions based on the images provided.
6. Can you describe the current weather conditions and the general environment depicted in the images?
7. Please describe the current driving conditions based on the input images.
8. Could you summarize the current driving conditions based on the input images?
9. Please provide an overview of the current driving conditions based on the images.
10. Can you summarize what the panoramic images show?
11. Can you describe the overall conditions and environment based on the images?
12. Could you describe the overall environment and objects captured in the images provided?

**Type 2: Critical Objects Analysis**
1. Where are the critical objects in the scene and what impact do they have on the ego vehicle?
2. Identify the significant objects in the scene and their specific impacts on the ego vehicle.
3. Can you pinpoint the critical objects in the scene and describe their influence on the ego vehicle?
4. Which objects in the scene are critical, and what effects do they have on the ego vehicle's movement?
5. Please describe the critical objects in the scene, their positions, and the influence they have on the ego vehicle.

**Type 3: Interpretable Action of Ego Vehicle**
1. Please describe your driving behavior and explain the reasons.
2. What is the current behavior of the vehicle?

**Type 4: Historical Information**
1. What are the differences between the current scene and the past scene in terms of critical objects?
2. How do the critical objects in the current scene differ from those in the past scene?
3. What changes have occurred in the critical objects between the current and past scenes?
4. What are the differences in critical objects between the present scene and the previous scene?
5. What distinctions exist between the critical objects of the current scene and those of the past scene?
6. In the past few frames, has a traffic light affected the driving strategy of the ego vehicle?
7. Within the recent frames, has the driving strategy of the ego vehicle been influenced by a traffic light?
8. In the last few frames, has the driving strategy of the ego vehicle been impacted by a traffic light?
9. Has the driving strategy of the ego vehicle been affected by a traffic light in the past few frames?
10. Has the traffic light influenced the driving strategy of the ego vehicle in the previous frames?
11. How has the current speed changed compared to the previous frames?
12. What was my driving behavior in the previous frame?