

# DriveLM: Driving with Graph Visual Question Answering

Chonghao Sima<sup>4,1\*</sup> Katrin Renz<sup>2,3\*</sup> Kashyap Chitta<sup>2,3</sup> Li Chen<sup>4,1</sup>  
 Hanxue Zhang<sup>1</sup> Chengen Xie<sup>1</sup> Ping Luo<sup>4</sup> Andreas Geiger<sup>2,3†</sup> Hongyang Li<sup>1†</sup>

<sup>1</sup> OpenDriveLab <sup>2</sup> University of Tübingen <sup>3</sup> Tübingen AI Center <sup>4</sup> University of Hong Kong

<https://github.com/OpenDriveLab/DriveLM>

*In memoriam Prof. Sean Xiao'ou Tang*

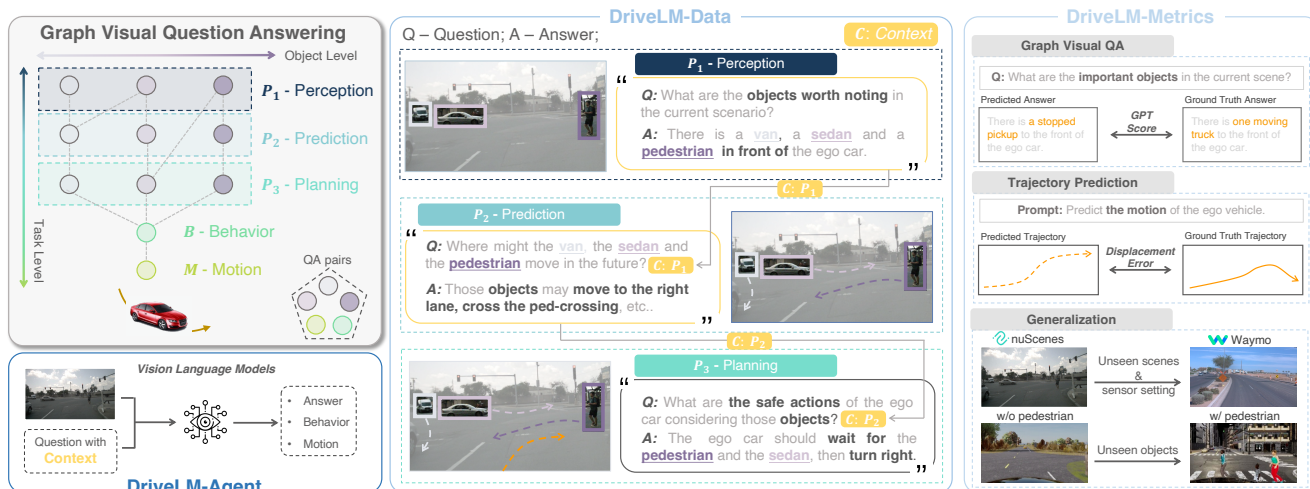


Figure 1. We present **DriveLM**: A new task, dataset, metrics, and baseline for end-to-end autonomous driving. Inspired by [8], DriveLM considers **Graph Visual Question Answering (GVQA)**, where question-answer pairs are interconnected via logical dependencies at the object-level, *i.e.*, interactions between object pairs, and the task-level, *e.g.*, perception  $\rightarrow$  prediction  $\rightarrow$  planning  $\rightarrow$  behavior (discretized action described in natural language)  $\rightarrow$  motion (continuous trajectory). We propose **DriveLM-Data** for training **DriveLM-Agent**, a baseline for GVQA. We validate its effectiveness using the **DriveLM-Metrics** on challenging settings requiring zero-shot generalization.

## Abstract

We study how vision-language models (VLMs) trained on web-scale data can be integrated into end-to-end driving systems to boost generalization and enable interactivity with human users. While recent approaches adapt VLMs to driving via single-round visual question answering (VQA), human drivers reason about decisions in multiple steps. Starting from the localization of key objects, humans estimate object interactions before taking actions. The key insight is that with our proposed task, Graph VQA, where we model graph-structured reasoning through perception, prediction and planning question-answer pairs, we obtain a suitable proxy task to mimic the human reason-

ing process. We instantiate datasets (*DriveLM-Data*) built upon *nuScenes* and *CARLA*, and propose a VLM-based baseline approach (*DriveLM-Agent*) for jointly performing Graph VQA and end-to-end driving. The experiments demonstrate that Graph VQA provides a simple, principled framework for reasoning about a driving scene, and *DriveLM-Data* provides a challenging benchmark for this task. Our *DriveLM-Agent* baseline performs end-to-end autonomous driving competitively in comparison to state-of-the-art driving-specific architectures. Notably, its benefits are pronounced when it is evaluated zero-shot on unseen objects or sensor configurations. We hope this work can be the starting point to shed new light on how to apply VLMs for autonomous driving. To facilitate future research, all code, data, and models are available to the public.

\*Equal contribution. †Equal co-advising.

# 1. Introduction

Current Autonomous Driving (AD) stacks are still lacking crucial capabilities [8, 11]. One key requirement is generalization, which involves the ability to handle unseen scenarios or unfamiliar objects. A secondary requirement pertains to the interaction of these models with humans, highlighted for example by EU regulations that mandate explainability in deployment [3]. Furthermore, unlike today’s AD models, humans do not navigate based on geometrically precise bird’s-eye view (BEV) representations [13, 26, 39]. Instead, humans implicitly perform object-centric perception, prediction, and planning (which we refer to as  $P_{1-3}$ ): a rough identification and localization of key objects, followed by reasoning about their possible movement and aggregation of this information into a driving action [49, 62].

Simultaneously, another field has been forging ahead: Vision-Language Models (VLMs) [40, 45, 73, 83]. These models have several strengths. First, they hold a base understanding of the world from internet-scale data that could potentially facilitate generalization for planning in AD. In fact, this sort of generalization has already been achieved by VLMs for simpler robotics tasks [18, 85]. Second, the use of language representations as an input and output offers a platform for human-friendly interaction with these models, unlike bounding boxes or trajectories that are more common to current methods [14, 25, 41, 58]. Finally, VLMs are able to make decisions in multiple steps linked by logical reasoning [4, 16, 75, 77, 82, 85]. Importantly, even though they reason in multiple separate steps, VLMs are end-to-end differentiable architectures, a characteristic that is highly desirable for autonomous driving [8].

Recent work towards enabling the application of VLMs to AD systems falls into two categories: scene-level or single object-level Visual Question Answering (VQA). Scene-level VQA refers to the task of describing the driving behavior by one or two supporting reasons, *e.g.*, “The car is moving into the right lane because it is safe to do so.” [34, 35]. Single object-level VQA formulates the understanding of the ego vehicle’s response to a single object by a chain of QAs in the form of “what-which-where-how-why”, *e.g.*, “The ego vehicle stops because there is a pedestrian in a white shirt crossing the intersection in front of the ego vehicle and it does not want to crash into the pedestrian.” [47, 55, 59]. Unfortunately, neither of these paradigms provides a suitable proxy task to mimic the  $P_{1-3}$  reasoning process in humans, who consider multiple objects and reason about each in multiple steps. Therefore, in this paper, we propose a new task, along with corresponding datasets and a baseline model architecture (Fig. 1).

**Task. Graph Visual Question Answering (GVQA)** involves formulating  $P_{1-3}$  reasoning as a series of question-answer pairs (QAs) in a directed graph. Its key differ-

ence to the aforementioned VQA tasks for AD is the availability of logical dependencies between QAs which can be used to guide the answering process. GVQA also encompasses questions regarding behavior and motion planning, with dedicated metrics (details in Section 2).

**Datasets. DriveLM-nuScenes** and **DriveLM-CARLA** consist of annotated QAs, arranged in a graph, linking images with driving behavior through logical reasoning. In comparison to existing benchmarks, they provide significantly more text annotations per frame (Fig. 2 and Table 1). We pair these training datasets with challenging test data for evaluating zero-shot generalization.

**Model. DriveLM-Agent** employs a trajectory tokenizer that can be applied to any general VLM [40, 45, 53, 83], coupled with a graph prompting scheme that models logical dependencies as context inputs for VLMs. The result is a simple, elegant methodology to effectively repurpose VLMs for end-to-end AD (Section 3).

Our experiments provide encouraging results. We find that GVQA on DriveLM is a challenging task, where current methods obtain moderate scores and better modeling of logical dependencies is likely necessary to achieve strong QA performance. Even so, DriveLM-Agent already performs competitively to state-of-the-art driving-specific models [26] when tested in the open-loop planning setting, despite its task-agnostic and generalist architecture. Furthermore, employing a graph structure improves zero-shot generalization, enabling DriveLM-Agent to correctly handle novel objects unseen during training or deployment on the Waymo dataset [63] after training only on nuScenes [5] data. From these results, we believe that improving GVQA holds great potential towards building autonomous driving agents with strong generalization.

## 2. DriveLM: Task, Data, Metrics

Human drivers usually decompose their decision-making process into distinct stages that follow a logical progression which encompasses the identification and localization of key objects, their possible future action and interaction, and ego planning based on all this information [22, 46]. This inspires us to propose the GVQA as the critical ingredient of DriveLM, which serves as a suitable proxy task to mimic the human reasoning process. Within this section, we illustrate the formulation of the GVQA task (Section 2.1), introduce DriveLM-Data (Section 2.2) to exemplify the instantiation of GVQA using prominent driving datasets, and overview the DriveLM-Metrics used for evaluation (Section 2.3).

### 2.1. DriveLM-Task: GVQA

We organize all the Question Answer pairs (QAs) for an image frame into a graph structure, denoted by  $G = (V, E)$ .  $V$  stands for the set of vertices, where each vertex represents a

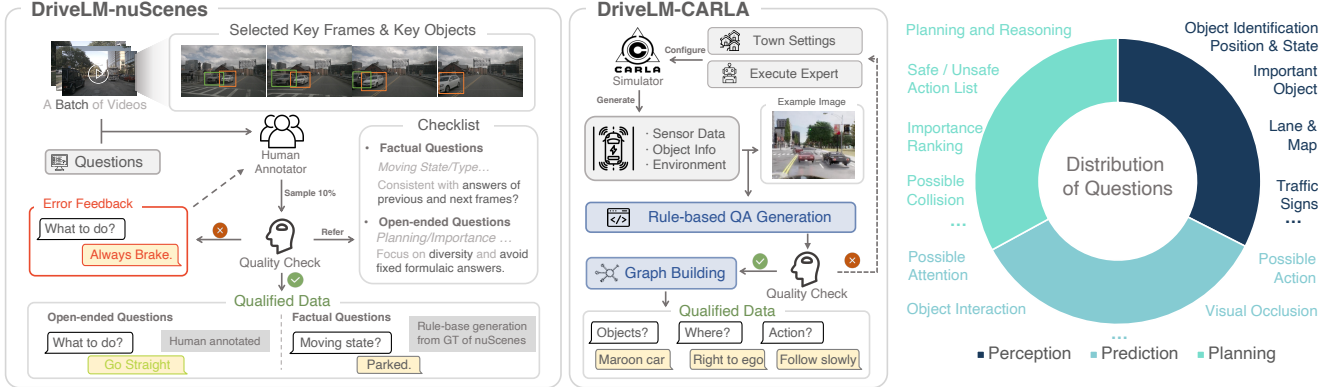


Figure 2. **(Left) Annotation Pipeline:** In DriveLM-nuScenes, we adopt a semi-rule-based QA labeling pipeline, where both the ground truth annotation in nuScenes/OpenLane-V2 and feedback from human annotators are used. A critical part of our pipeline is the multi-round quality check, which guarantees high data quality at reasonable costs. In DriveLM-CARLA, we meet the same standards while exploiting a fully rule-based QA labeling pipeline instead. **(Right) Question Distribution:** The questions in our dataset cover various specific aspects of driving tasks, most of which are annotated by human annotators, making this a suitable proxy for human-like driving reasoning.

QA pair  $v = (q, a)$  associated with one or more key objects in the scenario. The key difference between GVQA and ordinary VQA is that the QAs in GVQA have logical dependencies, which we formulate as the edges between the vertices.  $E \subseteq V \times V$ , is a set of directed edges, where each edge  $e = (v_p, v_c)$  connects the parent QA and the child QA. We formulate the edge set  $E$  by incorporating two dimensions: object-level and task-level edges. At the object level, we construct the logical edges  $e \in E$  to represent the impact of interactions between different objects. For example, the planning QA node for the sedan is influenced by the perception QA node of the pedestrian in the illustration from Fig. 1 (center). At the task-level, we establish the logical edges  $e \in E$  to capture the logical chain of different reasoning stages:

- **Perception** ( $P_1$ ): identification, description, and localization of key objects in the current scene.
- **Prediction** ( $P_2$ ): estimation of possible action/interaction of key objects based on perception results.
- **Planning** ( $P_3$ ): possible safe actions of the ego vehicle.
- **Behavior** ( $B$ ): classification of driving decision.
- **Motion** ( $M$ ): waypoints of ego vehicle future trajectory.

The concepts of perception, prediction, and planning ( $P_{1-3}$ ) are similar to those in end-to-end AD [8], while the concepts of motion and behavior are based on the ego vehicle future trajectory. Specifically, we define the motion  $M$  as the ego vehicle future trajectory, which is a set of  $N$  points with coordinates  $(x, y)$  in bird’s-eye view (BEV), denoted as  $M = \{(x_0, y_0), (x_1, y_1), \dots, (x_N, y_N)\}$ . Each point is the offset between the future position and the current position by a fixed time interval. Then, the distance for  $x, y$  at each time interval is computed as:

$$\{x, y\}_{\text{dist}} = \{(\delta_{x,1}, \delta_{y,1}), \dots, (\delta_{x,N}, \delta_{y,N})\}, \quad (1)$$

where  $\delta_{x,i} = x_i - x_{i-1}$  and  $\delta_{y,i} = y_i - y_{i-1}$ , for  $i =$

$1, 2, \dots, N$ . The goal of the behavior representation is to serve as an interface from  $P_{1-3}$  to  $M$ . To obtain a behavior representation, we map the mean of  $x_{\text{dist}}$  and  $y_{\text{dist}}$  to one of the predefined bins, where each bin corresponds to a category in either speed or steering. These are denoted as  $B_{sp}$  and  $B_{st}$  respectively. In this work, we consider 5 bins:

$$B_{sp} \in \{\text{fast}_2, \text{fast}_1, \text{moderate}, \text{slow}_1, \text{slow}_2\},$$

$$B_{st} \in \{\text{left}_2, \text{left}_1, \text{straight}, \text{right}_1, \text{right}_2\},$$

where the number in the subscript indicates the intensity. The combination of the speed and steering categories for a trajectory form its behavior category as  $B = (B_{sp}, B_{st})$ . While we use a simple definition of  $B$  as a starting point for research on driving with VLMs, we note that our formulation supports the incorporation of more abstract behaviors such as a lane changes or overtaking.

## 2.2. DriveLM-Data

In order to provide comprehensive and accurate QAs with the graph structure defined in Section 2.1, we introduce DriveLM-nuScenes and DriveLM-CARLA. Since there are significant disparities between nuScenes and CARLA, the collection methods and statistics of these datasets differ.

**DriveLM-nuScenes.** We divide the annotation process into three steps: selecting key frames from video clips, choosing key objects within these key frames, and subsequently annotating the frame-level  $P_{1-3}$  QAs for these key objects. A portion of the Perception QAs are generated from the nuScenes [5] and OpenLane-V2 [70] ground truth, while the remaining QAs are manually annotated. As we manually annotate the vast majority of data in DriveLM-nuScenes, quality is particularly crucial for this portion. When annotating, we conduct multiple rounds of rigorous quality checks. In each round, we categorize the data into

Dataset	Source Dataset	# Frames	Avg. captions / QA per annotated frame	Total captions / QA in Perception	Total captions / QA in Prediction	Total captions / QA in Planning	Logic among captions/QA pairs
nuScenes-QA [55]	nuScenes	34,149	13.5	460k**	✗	✗	None
nuPrompt [79]	nuScenes	34,149	1.0	35k*	✗	✗	None
HAD [35]	HDD	25,549	1.8	25k	✗	20k	None
BDD-X [34]	BDD	26,228	1	26k	✗	✗	None
DRAMA [47]	DRAMA	17,785	5.8	85k	✗	17k	Chain
Rank2Tell [59]	Rank2Tell	5,800	-	-	✗	-	Chain
DriveLM-nuScenes	nuScenes	4,871	<b>91.4</b>	144k*	153k	146k	<b>Graph</b>
DriveLM-CARLA	CARLA	<b>183,373</b>	20.5	<b>2.46M**</b>	<b>578k**</b>	<b>714k**</b>	<b>Graph</b>

Table 1. **Comparison of DriveLM-nuScenes & -CARLA with Existing Datasets.** \* indicates semi-rule-based labeling (w/ human annotators), \*\* indicates fully-rule-based (no human annotators), and - means publicly unavailable. DriveLM-Data significant advances annotation quantity, comprehensiveness (covering **perception, prediction and planning**), and logic (chain to **graph**).

different batches and inspect ten percent of the data in each batch. If the qualification rate of manually annotated data in this ten percent does not meet expectations, we request the annotators to re-label all data in the batch. In Fig. 2 (left), we showcase an example of the QA annotation pipeline, where all questions undergo quality checks according to our standards. As a result, DriveLM-nuScenes stands out from previously proposed datasets with its larger scale, greater comprehensiveness, and more complex structure (See Table 1). These QAs cover various aspects of the driving process, ranging from perception and prediction to planning, providing a comprehensive understanding of autonomous driving scenarios as shown in Fig. 2 (right).

**DriveLM-CARLA.** We collect data using CARLA 0.9.14 in the Leaderboard 2.0 framework [17] with a privileged rule-based expert [30]. We set up a series of routes in urban, residential, and rural areas and execute the expert on these routes. During this process, we collect the necessary sensor data, generate relevant QAs based on privileged information about objects and the scene, and organize the logical relationships to connect this series of QAs into a graph. We generate data and labels at 20 FPS. This process has the advantage of straightforward scalability since we only need to define route and scenario settings in CARLA and the subsequent steps can be executed automatically. The rule-based annotation pipeline is illustrated in Fig. 2 (left). Including 3.7M QAs, our DriveLM-CARLA stands out as the largest driving-language benchmark in terms of total textual content among existing benchmarks as shown in Table. 1.

### 2.3. DriveLM-Metrics

To evaluate GVQA, the DriveLM-Metrics consist of three components for evaluating motion  $M$ , behavior  $B$ , and  $P_{1-3}$ . For measuring the performance of the motion stage, we use standard metrics from the nuScenes and Waymo benchmarks: average and final displacement error, (**ADE**, **FDE**), and the **collision rate** on the predicted trajectory, following UniAD [26]. We evaluate behavior predictions by the **classification accuracy**, along with a breakdown of the overall accuracy into its steering and speed components. Finally, we measure the  $P_{1-3}$  performance using two metrics.

**SPICE** [2] is a prevailing metric used in VQA and image captioning, which calculates the structure similarity of predicted texts with ground truth while ignoring the semantic meanings. Simultaneously, we employ **GPT Score** to measure the semantic alignment of answers and complement the SPICE metric. Specifically, the question, the ground truth answer, the predicted answer, and a prompt asking for a numerical score of the answer are sent to ChatGPT-3.5 [50, 51]. We parse the text returned to get the score, where a higher score indicates better semantic accuracy.

### 3. DriveLM-Agent: A GVQA Baseline

In this section, we present DriveLM-Agent, a baseline approach for the GVQA task detailed in Section 2. DriveLM-Agent is built upon a general vision-language model and can therefore exploit underlying knowledge gained during pre-training. Our overall goal involves translating an image into the desired ego vehicle motion ( $M$ ) through the different stages of VQA ( $P_1, P_2, P_3, B$ ). For this, we choose BLIP-2 [40] as our base VLM due to its simplicity in architecture and flexibility in fine-tuning, but the proposed approach can be applied agnostically to other VLMs.

As shown in Fig. 3, DriveLM-Agent can be decomposed into several stages: (1)  $P_{1-3}$ , i.e., *perception, prediction, planning*, serve as the foundational layers to understand the scene and reason about its structure. (2) The *behavior* stage aggregates crucial information from the  $P_{1-3}$  into a description of the desired driving action in language space. (3) Finally, the *motion* stage translates the behavior into an executable driving trajectory. To implement the logical dependency between each linked QA, we propose to use context between connected nodes in the GVQA graph. We expand on this idea in the following.

#### 3.1. Prompting with Context

Directly translating images to motion as in [12, 54] is extremely challenging. Motivated by the tendency of humans to perform a multi-step reasoning process, we propose to use a similar strategy for VLM-based driving. By doing so, we facilitate the retrieval of knowledge stored in LLMs and improve explainability.

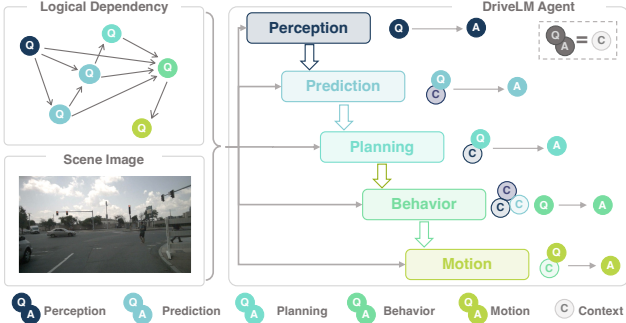


Figure 3. **DriveLM-Agent Architecture.** Given the scene image, DriveLM-Agent performs prompting with context to model the logical dependency among the five QA stages. Context is built using preceding QAs, and can have one or more sources.

More precisely, the model is designed to use answers from the previous steps in the reasoning process as the context for the following questions. For each edge  $e = (v_p, v_c) \in E$ , we append the QA from the parent node  $v_p$  to the question of the current node  $v_c$  with a prefix “Context: ”. The context can also contain QAs from multiple preceding nodes in which case we concatenate all QAs to one context sequence. It is worth noting that the context is only one possible implementation to formulate logical dependency in GVQA, which we select due to its simplicity. With this scheme, we pass forward relevant information based on the logical dependencies established by the graph.

Note that the size and structure of the graph during inference is a design choice of the algorithm, which can be adapted based on the task or available compute budget. We use this property to train on all available QAs, but perform inference on specific subgraphs, where the questions are sampled using heuristics. For more details, please refer to the supplementary material.

### 3.2. Context Aggregation through Behavior

Driving encompasses a wide array of potential situations that require an appropriate response. However, despite the diversity of these circumstances, it is interesting to note that almost all events involve decisions that can be discretized into a set of behaviors. For example, applying the brakes appropriately may address various situations such as a red light signal, a stop sign, or the presence of an object ahead of the vehicle. The focus of our behavior stage is to generate such a behavior: a statement in natural language that articulates the vehicle’s intended movement. In our methodology, the behavior is a textual description of the observed future vehicle motion, and it can also be represented as the category  $(B_{sp}, B_{st})$  which is divided into components for steering and speed (as described in Section 2.1). This description effectively serves as a reflective step wherein the model extracts and summarizes all crucial information

from the graph. To this end, we propose to use all possible sources of context for predicting behavior, *i.e.*, all the QAs in  $P_{1-3}$ . We empirically observe that having a reflective behavior step is crucial for driving with VLMs, and our design choice of using all possible context sources outperforms the naïve approach of only using  $P_3$ .

### 3.3. Trajectory Tokenization for Motion

Since it is non-trivial to output fine-grained numerical results using general VLMs, RT-2 [85] handles robotic actions based on a specialized trajectory tokenization module. We use this approach to enable DriveLM-Agent to take as input the image and behavior description and output trajectories. Specifically, we divide the coordinates of waypoints into 256 bins empirically based on the statistics of the train set trajectories. We re-define tokens in the BLIP-2 language tokenizer, establishing tokens for each bin, and fine-tune the VLM on this redefined vocabulary. For simplicity, we use the same VLM architecture (BLIP-2) to perform this task, but with independent LoRA weights and trained on a dataset consisting of only the QAs for this motion stage. Thus, it is possible to perform this functionality using a lightweight LLM [56] or driving-specific architecture that accepts a command as an input [25, 80].

## 4. Experiments

In this section, we present our experimental results that aim to address the following research questions: (1) How can VLMs be effectively repurposed for end-to-end autonomous driving? (2) Can VLMs for driving generalize when evaluated (a) with unseen sensor setups; and (b) on objects unseen during training? (3) How well do VLMs perform perception, prediction, and planning via GVQA?

**Setup.** We now briefly overview the key implementation details for the two settings used in our experiments (additional details are provided in the supplementary material). All fine-tuning is implemented with LoRA [24]. On DriveLM-nuScenes, we finetune BLIP-2 on the `train` split for 10 epochs. We use a batch size of 2 for each GPU, and the entire training process spans approximately 7 hours with 8 V100 GPUs. We train BLIP-2 on a 1/20 temporally sub-sampled `train` split of DriveLM-CARLA for 6 epochs. This takes 6 hours on 4 A100 GPUs.

### 4.1. VLMs for End-to-End Driving

In our first experiment, we aim to assess the ability of VLMs to perform open-loop planning on DriveLM-nuScenes. In particular, we investigate the impact of the context provided to the behavior and motion stages. Given sensor data (and in the case of VLM methods, a text input), the model is required to predict the ego-vehicle future trajectory in the form of waypoints.

Method	Behavior Context	Motion Context	Behavior ( $B$ )			Motion ( $M$ )	
			Acc. $\uparrow$	Speed $\uparrow$	Steer $\uparrow$	ADE $\downarrow$	Col. $\downarrow$
Command Mean	-	-	-	-	-	4.57	5.72
UniAD-Single	-	-	-	-	-	1.80	2.62
BLIP-RT-2	-	-	-	-	-	2.63	2.77
DriveLM-Agent	None	$B$	<b>61.45</b>	<b>72.20</b>	<b>84.73</b>	<b>1.39</b>	<b>1.67</b>
	Chain	$B$	50.43	60.32	75.34	2.07	2.08
	Graph	$B$	57.49	69.89	80.63	1.74	1.89
UniAD [26]	-	-	-	-	-	0.80	0.17

Table 2. **Open-loop Planning on DriveLM-nuScenes.** Using Behavior ( $B$ ) as context for Motion ( $M$ ) enables end-to-end driving with VLMs on par with UniAD-Single, a state-of-the-art driving-specific architecture.

**Baselines.** As a reference for the difficulty of the task, we provide a simple **Command Mean** baseline. Each frame in nuScenes is associated with one of 3 commands, ‘turn left’, ‘turn right’, or ‘go straight’. We output the mean of all trajectories in the training set whose command matches the current test frame command. Further, we compare our approach to the current state-of-the-art on nuScenes, UniAD [26]. Besides the author-released checkpoint, which requires video inputs, we train a single-frame version (‘UniAD-Single’) for a fair comparison to our single-frame VLMs. Finally, **BLIP-RT-2** denotes BLIP-2 [40] fine-tuned on DriveLM-Data with the trajectory tokenization scheme described in Section 3.3 for only the motion task. This acts as an indicator for the performance when using an identical network architecture as DriveLM-Agent, but no context inputs or VQA training data.

**DriveLM-Agent.** We consider 3 variants of DriveLM-Agent incorporating our proposed changes in steps: (1) a 2-stage version that predicts behavior and then motion (as described in Section 2.1), but without any  $P_{1-3}$  context for behavior prediction (‘None’); (2) a ‘Chain’ version that builds the  $P_{1-3}$  graph, but only passes the final node ( $P_3$ ) to the behavior stage; (3) the full model (‘Graph’) that uses all QAs from  $P_{1-3}$  as context for  $B$ .

**Results.** We show the results for the methods listed above in Table 2. Among the baselines, BLIP-RT-2 is unable to match UniAD-Single (though both methods perform well relative to Command Mean). This shows that the single-stage approach without any reasoning is unable to compete with the prior state-of-the-art on nuScenes. However, the proposed DriveLM-Agent, which predicts behavior as an intermediate step for motion, provides a significant boost in performance, surpassing UniAD-Single. This indicates that with the appropriate prompting, VLMs can be surprisingly competitive for end-to-end driving. Interestingly, in the experimental setting of Table 2 which does not involve generalization, the Chain and Graph versions of DriveLM-Agent do not provide any further advantage over no context. Further, single-frame VLMs fall short in comparison to the privileged video-based UniAD model, indicating that VLMs with video inputs may be necessary for this task.

Method	Behavior Context	Motion Context	Behavior ( $B$ )			Motion ( $M$ )	
			Acc. $\uparrow$	Speed $\uparrow$	Steer $\uparrow$	ADE $\downarrow$	FDE $\downarrow$
Command Mean	-	-	-	-	-	7.98	11.41
UniAD-Single	-	-	-	-	-	4.16	9.31
BLIP-RT-2	-	-	-	-	-	2.78	6.47
DriveLM-Agent	None	$B$	35.70	43.90	65.20	2.76	6.59
	Chain	$B$	34.62	41.28	64.55	2.85	6.89
	Graph	$B$	<b>39.73</b>	<b>54.29</b>	<b>70.35</b>	<b>2.63</b>	<b>6.17</b>

Table 3. **Zero-shot Generalization across Sensor Configurations.** Results on 1k randomly sampled frames from the Waymo val set after training on DriveLM-nuScenes. DriveLM-Agent outperforms UniAD-Single and benefits from graph context.

## 4.2. Generalization Across Sensor Configurations

As a more challenging setting for evaluating the models from Section 4.1, we now apply them without any further training to a new domain: the Waymo dataset [63]. Waymo’s sensor setup does not include a rear camera, so we drop this input from UniAD-Single. The VLM methods only use the front view and do not require any adaptation.

**Results.** As shown in Table 3, UniAD-Single does not cope well with the new sensor configuration, and drops below BLIP-RT-2 in performance. The multi-stage approach of DriveLM-Agent provides further improvements. In particular, the accuracy of speed predictions rises from 43.90 with no context to 54.29 with the full graph. On the other hand, the chain approach does not provide sufficient useful information, with a speed accuracy of only 41.28.

We present qualitative results for DriveLM-Agent on nuScenes and Waymo in Fig. 4. The model generally provides intuitive answers, with a few exceptions (e.g., planning on DriveLM-nuScenes, perception on Waymo). This shows the utility of GVQA towards interactive driving systems. Further, on Waymo, we see meaningful prediction and planning answers despite the imperfect perception. For more visualizations, please see the supplementary material.

## 4.3. Generalization to Unseen Objects

Next, we evaluate zero-shot generalization to novel objects. DriveLM-CARLA is collected without any pedestrians in the training or validation splits. We now generate a new test set called DriveLM-CARLA-ped, which only consists of frames where a pedestrian is present in the scene. The correct behavior is to stop for the pedestrian.

**Baselines.** For this experiment, we compare DriveLM-Agent to TransFuser++ [30], the state-of-the-art for CARLA. It uses a larger input image, an additional LiDAR sensor, and several driving-specific annotations (depth, semantics, 3D bounding boxes, HD map) in comparison to DriveLM-Agent. However, because of these task-specific inputs and outputs, TransFuser++ can only be trained on the base DriveLM-CARLA dataset and cannot incorporate general computer vision data during training, which makes generalization more challenging.

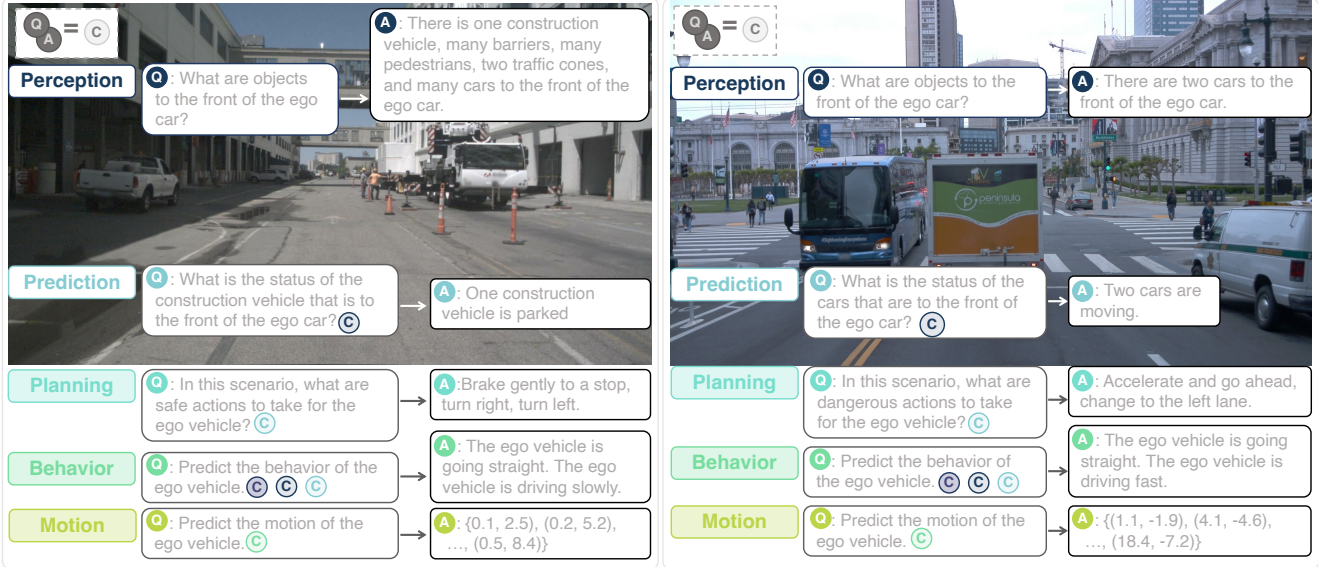


Figure 4. **Qualitative Results of DriveLM-Agent.** (Left) DriveLM-nuScenes val frame, (Right) Waymo val frame. We show the questions (Q), context (C), and predicted answers (A). DriveLM-Agent’s outputs are easy to interpret for human users.

Method	DriveLM-CARLA ( <i>B</i> )			DriveLM-CARLA-ped ( <i>B</i> )		
	Acc. $\uparrow$	Spd. $\uparrow$	Str. $\uparrow$	Acc. $\uparrow$	Spd. $\uparrow$	Str. $\uparrow$
TransFuser++ [30]	<b>70.19</b>	<b>73.29</b>	<b>90.68</b>	8.72	8.72	100.00
DriveLM-Agent	59.63	61.50	78.26	4.59	4.59	100.00
+ Pedestrian QA	52.17	55.28	77.64	<b>27.04</b>	<b>27.04</b>	100.00
<i>DriveLM-Agent (GT)</i>	60.25	65.22	80.12	20.92	20.92	100.00
+ Pedestrian QA	60.25	65.22	80.12	92.35	92.35	100.00

Table 4. **Generalization in DriveLM-CARLA.** All methods underperform on DriveLM-CARLA-ped with the novel pedestrian object, but DriveLM-Agent can be significantly improved by including a pedestrian-specific question in its GVQA graph.

**DriveLM-Agent.** Taking advantage of the more general architecture of a VLM, we include samples from COCO [44] and GQA [28] along with DriveLM-CARLA during training for DriveLM-Agent. We compare several versions: (1) we investigate the addition of a new  $P_1$  question during inference, ‘*Is there a person crossing the road?*’ (+ Pedestrian QA). (2) As an upper bound, we directly input the ground truth  $P_{1-3}$  graph to the model during inference, instead of the model’s predictions. For more details, please refer to the supplementary material.

**Results.** We present our findings in Table 4. We observe that TransFuser++ struggles on DriveLM-CARLA-ped relative to DriveLM-CARLA, with a drop in accuracy from 70.19 to just 8.72. DriveLM-Agent experiences a similar drop from 59.63 to 4.59. However, adding the pedestrian QA significantly boosts performance on the generalization setting to 27.04, albeit with slightly reduced accuracy on regular scenes. This is mainly attributed to the fact that the VLM is not able to detect all pedestrians correctly. This indicates that the large performance gains of recently published VLMs [21, 53] can support even better

generalization ability in the domain of driving. Additionally, when the pedestrian QA is provided in the privileged setting that assumes access to perfect context for each question in the graph, DriveLM-Agent achieves a near-perfect score (20.92  $\rightarrow$  92.35) on the frames with pedestrians. Note that DriveLM-CARLA-ped only contains pedestrians crossing a straight road, so all models obtain a 100% accuracy on the steering class (which is always straight).

#### 4.4. Performance for $P_{1-3}$ via GVQA

In our final experiment, we establish baseline results for the  $P_{1-3}$  stages of GVQA, studying the impact of context. We use two VLMs, the off-the-shelf BLIP-2 [40] model (which is not fine-tuned on DriveLM), as well as the proposed DriveLM-Agent.

**Baselines.** We consider the lower bound of no context (‘None’), which corresponds to training and evaluation with the same setting as standard VQA (image and question in, answer out). As an upper bound for each architecture, we perform GVQA but input the ground truth (‘GT’) context to the model at test time instead of its own prior predictions.

**Results.** Our results are summarized in Table 5. Firstly, we observe that DriveLM-nuScenes is significantly more challenging for both models, as indicated by the lower scores on this dataset relative to DriveLM-CARLA in all context settings. This is likely due to the higher diversity in human answers obtained for DriveLM-nuScenes, as opposed to the rule-based generation in CARLA. On both datasets, we observe that DriveLM-Agent, which is fine-tuned on DriveLM, significantly outperforms BLIP-2 which is applied in a zero-shot manner. We also observe the SPICE and GPT

Context	DriveLM-nuScenes ( $P_{1-3}$ )				DriveLM-CARLA ( $P_{1-3}$ )			
	BLIP-2 [40]		DriveLM-Agent		BLIP-2 [40]		DriveLM-Agent	
	SPICE $\uparrow$	GPT $\uparrow$	SPICE $\uparrow$	GPT $\uparrow$	SPICE $\uparrow$	GPT $\uparrow$	SPICE $\uparrow$	GPT $\uparrow$
None	4.34	42.97	42.56	71.39	13.89	66.77	87.52	86.84
Graph	7.71	45.21	49.54	72.51	12.79	59.98	86.86	84.51
<i>GT</i>	<i>8.19</i>	<i>41.10</i>	<i>50.29</i>	<i>72.94</i>	<i>39.46</i>	<i>61.40</i>	<i>87.70</i>	<i>87.15</i>

Table 5. **Baseline  $P_{1-3}$  Results.** DriveLM-Agent outperforms BLIP-2, but is unable to benefit from the availability of context.

score metrics to be misaligned, in particular for BLIP-2. Interestingly, on DriveLM-CARLA, BLIP-2 with GT context gains over 25 points in terms of SPICE accompanied by a small drop in GPT score compared to the no context version. This shows that current VLMs are able to mimic the sentence structure and style of the provided context, but face challenges in performing logical reasoning. Overall, we conclude that DriveLM-Agent can obtain a reasonable baseline performance on  $P_{1-3}$  question answering without context. However, specialized architectures or prompting schemes beyond naive concatenation may be necessary to make better use of the logical dependencies in GVQA.

## 5. Related Work

**Generalization in Autonomous Driving.** The inadequacy of generalization to the “long tail” of corner cases poses significant safety concerns to AD systems [8, 65, 66]. To tackle this issue, prior research primarily makes efforts in data-driven methods [1, 6, 23, 64, 71]. For example, TrafficSim [64] collects more data for safety-critical cases by simulation. An emerging direction involves leveraging semantic information to supervise the detection of unseen or anomalous objects [19, 72]. These efforts alleviate the problem of insufficient generalization. Even so, the zero-shot performance of AD systems is currently not satisfactory. In this paper, we bring a new approach towards better generalization: learning logical reasoning using Graph VQA.

**Embodied Planning with LLMs.** Recent work endeavors to leverage the formidable reasoning and generalization capacity of LLMs [20, 36, 68] for embodied AI systems [18, 27, 29, 33, 42, 57, 85]. PaLM-E [18] trains an LLM for various embodied tasks including sequential robotic manipulation planning. CaP [42] provides a robot-centric formulation of language model generated programs executed on real systems. RT-2 [85] represents robot actions as language tokens, training vision-language models to output robot policies. These methods showcase the capabilities of LLMs in embodied planning tasks, inspiring us to apply them to address the current shortcomings in generalization in AD, which is far less explored.

**Language-grounded Driving.** Several concurrent methods attempt to incorporate multi-modal inputs into LLMs for AD tasks [7, 32, 48, 60, 72, 76, 81]. Specifically, GPT-Driver [48] and LLM-Driver [7] encode the perceived scene state into prompts, relying on LLMs to formulate reason-

able plans. DriveGPT4 [81] projects raw sensor data into tokens and utilizes LLMs for end-to-end prediction of control signals and explanations. Despite these preliminary attempts, there is untapped potential in addressing generalization in AD through LLMs. Our work combines VLMs with training over graph-structured QAs from DriveLM. This enables us to show benefits on zero-shot end-to-end planning, which was not demonstrated by these concurrent studies.

## 6. Discussion

Even though DriveLM exhibits promising generalization, there are concerning limitations of this work.

- **Driving-specific Inputs.** DriveLM-Agent directly applies the VLM’s vision module, taking a low-resolution front-view frame as its input. Currently, driving-specific sensors such as LiDAR cannot be processed. This results in our model lacking temporal information and 360-degree scene understanding. Extending DriveLM-Agent to observe images from multiple views is straightforward as the graph formulation allows different input frames for different nodes. We leave it for future work to explore options for multi-view and multi-frame inputs.
- **Closed-loop Planning.** Our approach is currently evaluated under an open-loop scheme. In this setting, incorporating the ego vehicle’s status as input can significantly enhance the metrics, but its effectiveness may not translate well to the real world, and hence we only consider methods that do not do so. Extending our work to a closed-loop setting with an affordable budget in training time and computational cost is a promising direction to explore. With the usage of CARLA we provide a promising foundation for more research in the direction of closed-loop planning with VLMs.
- **Efficiency Constraints.** Inheriting the drawbacks of LLMs, our model suffers from long inference times, especially as we require multiple rounds of predictions based on the graph structure (roughly  $4\times$  slower than UniAD). This may impact practical implementation. Exploring how to conduct model quantization, distillation and pruning is a worthy direction for future research.

**Conclusion.** We show how VLMs can be leveraged as end-to-end autonomous driving agents with improved generalization over task-specific driving stacks. For this we propose the task of Graph VQA together with new datasets and metrics. Equipped with these tools, we build a baseline approach that has a simple architecture and obtains promising results. We believe that this approach can accelerate progress in the field of autonomous driving by enabling it to directly benefit from better VLMs.

**Broader Impact.** Our goal is to make progress towards autonomous driving, which will have profound impact if successful. We recognize that by bringing VLMs into this area,



we accept their ethical implications, such as hallucinations and high resource use. Yet, by improving the interactivity between humans and autonomous driving systems, we can build confidence in the technology. This could hasten its acceptance and lead to safer transportation in the long term.

**Acknowledgements.** The OpenDriveLab team is part of the Shanghai AI Lab and kindly supported by National Key R&D Program of China (2022ZD0160104) and NSFC (62206172). This work was also supported by the BMBF (Tübingen AI Center, FKZ: 01IS18039A), the DFG (SFB 1233, TP 17, project number: 276693517), and by EXC (number 2064/1 – project number 390727645). We thank the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for supporting K. Renz and K. Chitta. Our gratitude goes to Tai Wang for the valuable feedback, Jens Beißwenger for assisting with the CARLA setup, Qingwen Bu for refining the figures, Jiajie Xu for refining the DriveLM-nuScenes and cleaning the DriveLM-Agent codebase, and Jiazhi Yang, Shenyuan Gao, Yihang Qiu, Tianyu Li, Yunsong Zhou, Zetong Yang for the fruitful discussion.

## References

- [1] Shivam Akhauri, Laura Y Zheng, and Ming C Lin. Enhanced transfer learning for autonomous driving with systematic accident simulation. In *IROS*, 2020. 8
- [2] Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. Spice: Semantic propositional image caption evaluation. In *ECCV*, 2016. 4, 23
- [3] Shahin Atakishiyev, Mohammad Salameh, Housam Babiker, and Randy Goebel. Explaining autonomous driving actions with visual question answering. *arXiv preprint arXiv:2307.10408*, 2023. 2
- [4] Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Michal Podstawski, Hubert Niewiadomski, et al. Graph of Thoughts: Solving Elaborate Problems with Large Language Models. *arXiv preprint arXiv:2308.09687*, 2023. 2, 32
- [5] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yuxin Pan, Giancarlo Baldan, and Oscar Beijbom. nuScenes: A multi-modal dataset for autonomous driving. In *CVPR*, 2020. 2, 3, 15, 16
- [6] Dian Chen and Philipp Krähenbühl. Learning from all vehicles. In *CVPR*, 2022. 8
- [7] Long Chen, Oleg Sinavski, Jan Hünermann, Alice Karnsund, Andrew James Willmott, Danny Birch, Daniel Maund, and Jamie Shotton. Driving with llms: Fusing object-level vector modality for explainable autonomous driving. *arXiv preprint arXiv:2310.01957*, 2023. 8
- [8] Li Chen, Penghao Wu, Kashyap Chitta, Bernhard Jaeger, Andreas Geiger, and Hongyang Li. End-to-end autonomous driving: Challenges and frontiers. *arXiv preprint arXiv:2306.16927*, 2023. 1, 2, 3, 8
- [9] Xiaojun Chen, Shengbin Jia, and Yang Xiang. A review: Knowledge reasoning over knowledge graph. *Expert Syst. Appl.*, 2020. 32
- [10] Yunpeng Chen, Marcus Rohrbach, Zhicheng Yan, Yan Shuicheng, Jiashi Feng, and Yannis Kalantidis. Graph-based global reasoning networks. In *CVPR*, 2019. 32
- [11] Pranav Singh Chib and Pravendra Singh. Recent advancements in end-to-end autonomous driving using deep learning: A survey. *IEEE T-IV*, 2023. 2
- [12] Kashyap Chitta, Aditya Prakash, and Andreas Geiger. Neat: Neural attention fields for end-to-end autonomous driving. In *ICCV*, 2021. 4
- [13] Kashyap Chitta, Aditya Prakash, Bernhard Jaeger, Zehao Yu, Katrin Renz, , and Andreas Geiger. Transfuser: Imitation with transformer-based sensor fusion for autonomous driving. *IEEE T-PAMI*, 2023. 2
- [14] Daniel Dauner, Marcel Hallgarten, Andreas Geiger, and Kashyap Chitta. Parting with misconceptions about learning-based vehicle motion planning. In *CoRL*, 2023. 2
- [15] Vikrant Dewangan, Tushar Choudhary, Shivam Chandhok, Shubham Priyadarshan, Anushka Jain, Arun Singh, Siddharth Srivastava, Krishna Murthy Jatavallabhula, and Madhava Krishna. Talk2BEV: Language-enhanced bird’s-eye view maps for autonomous driving. *arXiv preprint arXiv:2310.02251*, 2023. 32
- [16] Ruomeng Ding, Chaoyun Zhang, Lu Wang, Yong Xu, Minghua Ma, Wei Zhang, Si Qin, Saravan Rajmohan, Qingwei Lin, and Dongmei Zhang. Everything of thoughts: Defying the law of penrose triangle for thought generation. *arXiv preprint arXiv:2311.04254*, 2023. 2
- [17] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. CARLA: An open urban driving simulator. In *CoRL*, 2017. 4, 20
- [18] Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, et al. PaLM-E: An embodied multimodal language model. In *ICML*, 2023. 2, 8
- [19] Amine Elhafsi, Rohan Sinha, Christopher Agia, Edward Schmerling, Issa AD Nesnas, and Marco Pavone. Semantic anomaly detection with large language models. *Auton. Robot.*, 2023. 8
- [20] Luciano Floridi and Massimo Chiriatti. GPT-3: Its nature, scope, limits, and consequences. *MIND MACH*, 2020. 8
- [21] Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan Lu, Conghui He, Xianguyue Yue, Hongsheng Li, and Yu Qiao. LLaMA-Adapter v2: Parameter-efficient visual instruction model. *arXiv preprint arXiv:2304.15010*, 2023. 7
- [22] John A Groeger. *Understanding driving: Applying cognitive psychology to a complex everyday task*. Routledge, 2013. 2
- [23] Niklas Hanselmann, Katrin Renz, Kashyap Chitta, Apratim Bhattacharyya, and Andreas Geiger. King: Generating safety-critical driving scenarios for robust imitation via kinematics gradients. In *ECCV*, 2022. 8
- [24] 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 *CoRL*, 2021. 5, 24

- [25] 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 *ECCV*, 2022. 2, 5, 23
- [26] Yihan Hu, Jiazhi Yang, Li Chen, Keyu Li, Chonghao Sima, Xizhou Zhu, Siqi Chai, Senyao Du, Tianwei Lin, et al. Planning-oriented autonomous driving. In *CVPR*, 2023. 2, 4, 6, 23
- [27] Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, et al. Inner monologue: Embodied reasoning through planning with language models. In *CoRL*, 2022. 8
- [28] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *CVPR*, 2019. 7
- [29] Brian Ichter, Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel Ho, Julian Ibarz, Alex Irpan, et al. Do as i can, not as i say: Grounding language in robotic affordances. In *CoRL*, 2023. 8
- [30] Bernhard Jaeger, Kashyap Chitta, and Andreas Geiger. Hidden biases of end-to-end driving models. In *ICCV*, 2023. 4, 6, 7, 21
- [31] Frederick Jelinek, John D Lafferty, and Robert L Mercer. *Basic methods of probabilistic context free grammars*. Springer, Berlin, Heidelberg, 1992. 23
- [32] Bu Jin, Xinyu Liu, Yupeng Zheng, Pengfei Li, Hao Zhao, Tong Zhang, Yuhang Zheng, Guyue Zhou, and Jingjing Liu. Adapt: Action-aware driving caption transformer. In *ICRA*, 2023. 8
- [33] Siddharth Karamcheti, Suraj Nair, Annie S. Chen, Thomas Kollar, Chelsea Finn, Dorsa Sadigh, and Percy Liang. Language-Driven representation learning for robotics. In *RSS*, 2023. 8
- [34] Jinkyu Kim, Anna Rohrbach, Trevor Darrell, John Canny, and Zeynep Akata. Textual explanations for self-driving vehicles. In *ECCV*, 2018. 2, 4, 32
- [35] Jinkyu Kim, Teruhisa Misu, Yi-Ting Chen, Ashish Tawari, and John Canny. Grounding human-to-vehicle advice for self-driving vehicles. In *CVPR*, 2019. 2, 4, 32
- [36] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. In *NeurIPS*, 2022. 8
- [37] Alex M Lamb, Anirudh Goyal ALIAS PARTH GOYAL, Ying Zhang, Saizheng Zhang, Aaron C Courville, and Yoshua Bengio. Professor forcing: A new algorithm for training recurrent networks. In *NeurIPS*, 2016. 24
- [38] Alon Lavie and Abhaya Agarwal. METEOR: An automatic metric for MT evaluation with high levels of correlation with human judgments. In *ACL Workshop*, 2007. 22, 25
- [39] Hongyang Li, Chonghao Sima, Jifeng Dai, Wenhai Wang, Lewei Lu, Huijie Wang, Jia Zeng, Zhiqi Li, Jiazhi Yang, Hanming Deng, et al. Delving into the devils of bird’s-eye-view perception: A review, evaluation and recipe. *IEEE T-PAMI*, 2023. 2
- [40] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, 2023. 2, 4, 6, 7, 8
- [41] Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Yu Qiao, and Jifeng Dai. Bevformer: Learning bird’s-eye-view representation from multi-camera images via spatiotemporal transformers. In *ECCV*, 2022. 2
- [42] Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. Code as policies: Language model programs for embodied control. In *ICRA*, 2023. 8
- [43] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *ACL Workshop*, 2004. 22, 25
- [44] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 7
- [45] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023. 2
- [46] Charles C Macadam. Understanding and modeling the human driver. *Veh. Syst. Dyn*, 2003. 2
- [47] Srikanth Malla, Chiho Choi, Isht Dwivedi, Joon Hee Choi, and Jiachen Li. DRAMA: Joint risk localization and captioning in driving. In *WACV*, 2023. 2, 4, 32
- [48] Jiageng Mao, Yuxi Qian, Hang Zhao, and Yue Wang. GPT-Driver: Learning to drive with gpt. *arXiv preprint arXiv:2310.01415*, 2023. 8
- [49] David Marr. *Vision: A computational investigation into the human representation and processing of visual information*. The MIT Press, 2010. 2
- [50] OpenAI. OpenAI: Introducing ChatGPT. <https://openai.com/blog/chatgpt>, 2022. 4
- [51] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, et al. Training language models to follow instructions with human feedback. In *NeurIPS*, 2022. 4
- [52] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *ACL*, 2002. 22, 25
- [53] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*, 2023. 2, 7
- [54] Aditya Prakash, Kashyap Chitta, and Andreas Geiger. Multi-modal fusion transformer for end-to-end autonomous driving. In *CVPR*, 2021. 4
- [55] Tianwen Qian, Jingjing Chen, Linhai Zhuo, Yang Jiao, and Yu-Gang Jiang. NuScenes-QA: A multi-modal visual question answering benchmark for autonomous driving scenario. *arXiv preprint arXiv:2305.14836*, 2023. 2, 4, 32
- [56] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI blog*, 2019. 5
- [57] Krishan Rana, Jesse Haviland, Sourav Garg, Jad Abou-Chakra, Ian Reid, and Niko Suenderhauf. SayPlan: Grounding large language models using 3d scene graphs for scalable robot task planning. In *CoRL*, 2023. 8
- [58] Katrin Renz, Kashyap Chitta, Otniel-Bogdan Mercea, Almut Sophia Koepke, Zeynep Akata, and Andreas Geiger.

- Plant: Explainable planning transformers via object-level representations. In *CoRL*, 2022. 2
- [59] Enna Sachdeva, Nakul Agarwal, Suhas Chundi, Sean Roelofs, Jiachen Li, Behzad Dariush, Chiho Choi, and Mykel Kochenderfer. Rank2Tell: A multimodal driving dataset for joint importance ranking and reasoning. *arXiv preprint arXiv:2309.06597*, 2023. 2, 4, 32
- [60] Ari Seff, Brian Cera, Dian Chen, Mason Ng, Aurick Zhou, Nigamaa Nayakanti, Khaled S Refaat, Rami Al-Rfou, and Benjamin Sapp. MotionLM: Multi-agent motion forecasting as language modeling. In *ICCV*, 2023. 8
- [61] Jiaxin Shi, Hanwang Zhang, and Juanzi Li. Explainable and explicit visual reasoning over scene graphs. In *CVPR*, 2019. 32
- [62] Elizabeth S Spelke and Katherine D Kinzler. Core knowledge. *Dev Sci*, 2007. 2
- [63] Pei Sun, Henrik Kretschmar, 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 *CVPR*, 2020. 2, 6
- [64] Simon Suo, Sebastian Regalado, Sergio Casas, and Raquel Urtasun. Trafficsim: Learning to simulate realistic multi-agent behaviors. In *CVPR*, 2021. 8
- [65] Ardi Tampuu, Tambet Matiisen, Maksym Semikin, Dmytro Fishman, and Naveed Muhammad. A survey of end-to-end driving: Architectures and training methods. *IEEE T-NNLS*, 2020. 8
- [66] Siyu Teng, Xuemin Hu, Peng Deng, Bai Li, Yuchen Li, Yunfeng Ai, Dongsheng Yang, Lingxi Li, Zhe Xuanyuan, Fenghua Zhu, et al. Motion planning for autonomous driving: The state of the art and future perspectives. *IEEE T-IV*, 2023. 8
- [67] Nikzad Benny Toomarian and Jacob Barhen. Learning a trajectory using adjoint functions and teacher forcing. *Neural networks*, 1992. 24
- [68] Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 8
- [69] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *CVPR*, 2015. 22, 25
- [70] Huijie Wang, Tianyu Li, Yang Li, Li Chen, Chonghao Sima, Zhenbo Liu, Bangjun Wang, Peijin Jia, Yuting Wang, Shengyin Jiang, et al. OpenLane-V2: A topology reasoning benchmark for unified 3d HD mapping. In *NeurIPS Datasets and Benchmarks*, 2023. 3, 15
- [71] Jingkan Wang, Ava Pun, James Tu, Sivabalan Manivasagam, Abbas Sadat, Sergio Casas, Mengye Ren, and Raquel Urtasun. Advsim: Generating safety-critical scenarios for self-driving vehicles. In *CVPR*, 2021. 8
- [72] Tsun-Hsuan Wang, Alaa Maalouf, Wei Xiao, Yutong Ban, Alexander Amini, Guy Rosman, Sertac Karaman, and Daniela Rus. Drive anywhere: Generalizable end-to-end autonomous driving with multi-modal foundation models. *arXiv preprint arXiv:2310.17642*, 2023. 8
- [73] Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, and Jifeng Dai. VisionLLM: Large language model is also an open-ended decoder for vision-centric tasks. *arXiv preprint arXiv:2305.11175*, 2023. 2
- [74] Xiang Wang, Dingxian Wang, Canran Xu, Xiangnan He, Yixin Cao, and Tat-Seng Chua. Explainable reasoning over knowledge graphs for recommendation. In *AAAI*, 2019. 32
- [75] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-Consistency improves chain of thought reasoning in language models. In *ICLR*, 2023. 2
- [76] Wayve. Lingo-1. <https://wayve.ai/thinking/lingo-natural-language-autonomous-driving/>, 2023. 8
- [77] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In *NeurIPS*, 2022. 2
- [78] Dongming Wu, Wencheng Han, Tiancai Wang, Xingping Dong, Xiangyu Zhang, and Jianbing Shen. Referring Multi-Object tracking. In *CVPR*, 2023. 32
- [79] Dongming Wu, Wencheng Han, Tiancai Wang, Yingfei Liu, Xiangyu Zhang, and Jianbing Shen. Language prompt for autonomous driving. *arXiv preprint arXiv:2309.04379*, 2023. 4, 32
- [80] 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 *NeurIPS*, 2022. 5
- [81] Zhenhua Xu, Yujia Zhang, Enze Xie, Zhen Zhao, Yong Guo, Kenneth KY Wong, Zhenguo Li, and Hengshuang Zhao. DriveGPT4: Interpretable end-to-end autonomous driving via large language model. *arXiv preprint arXiv:2310.01412*, 2023. 8, 23
- [82] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of Thoughts: Deliberate problem solving with large language models. *arXiv preprint arXiv:2305.10601*, 2023. 2
- [83] Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, Peng Gao, and Yu Qiao. LLaMA-Adapter: Efficient fine-tuning of language models with zero-init attention. *arXiv preprint arXiv:2303.16199*, 2023. 2
- [84] Xingcheng Zhou, Mingyu Liu, Bare Luka Zagar, Ekim Yurtsever, and Alois C. Knoll. Vision language models in autonomous driving and intelligent transportation systems. *arXiv preprint arXiv:2310.14414*, 2023. 32
- [85] Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker, et al. RT-2: Vision-language-action models transfer web knowledge to robotic control. In *CoRL*, 2023. 2, 5, 8, 24

## Overview

In the appendices below, we first delve deeper into various discussions, along with additional details around the annotation process of DriveLM-nuScenes & -CARLA, GVQA metrics, context setting & trajectory tokenization in DriveLM-Agent, and more ablation results of DriveLM-Agent on DriveLM-nuScenes and Waymo. Finally, we provide additional results and visualizations that further complement the findings from the main text.

For readers who want to focus on specific topics, we provide a summary below:

---

### Appendix A – *Motivating Questions*

We index a list of “motivating” questions that may arise from reading the main text and that we expand on further here (e.g., “why adapt general VLMs to driving”). These questions are open to be explored and thus our answers here are *intuitive* and *empirical*.

### Appendix B – *DriveLM-nuScenes*

We provide the DriveLM-nuScenes dataset composition, introduce the detailed annotation pipeline and conduct statistics of QA categories.

### Appendix C – *DriveLM-CARLA*

We provide a detailed description of the composition of the dataset, how the graph looks like and explain the data generation and annotation process.

### Appendix D – *DriveLM-Metrics*

We explain the details of the metrics for each task in the GVQA, illustrate their difference, and provide the reasons of proposing GPT-score as the main metric used in  $P_{1-3}$  VQA tasks.

### Appendix E – *DriveLM-Agent*

We introduce the detailed design of the prompting with context and the trajectory tokenization, including the difference of context in training and evaluation, the pattern for the trajectory as sentence, and the hyperparameters in the tokenizer.

### Appendix F – *Experiments*

We provide more experiments on DriveLM-nuScenes and Waymo, including the effects of more context design on the zero-shot ability, the performance with more conventional VQA metrics and the model efficiency comparison.

### Appendix G – *Qualitative Results*

We show qualitative examples of the context, questions, and answers on nuScenes, Waymo, and CARLA. Additionally, we contrast predicted and ground truth answers together with their SPICE and GPT Score on nuScenes to provide some intuition for those metrics.

### Appendix H – *Other Related Work*

We provide more related works from two new perspectives. One is reasoning over graph structure which is similar to our idea of graph-structure reasoning, the other is more vision-language benchmarks for autonomous driving.

# Appendix

<b>A Motivating Questions</b>	<b>14</b>
<b>B DriveLM-nuScenes</b>	<b>15</b>
B.1. Dataset Composition . . . . .	15
B.2. Collection Methodology . . . . .	15
B.3. Statistics and Facts . . . . .	17
<b>C DriveLM-CARLA</b>	<b>20</b>
C.1. Dataset Composition . . . . .	20
C.2. Collection Methodology . . . . .	20
<b>D DriveLM-Metrics</b>	<b>22</b>
D.1. $P_{1-3}$ VQA Metrics . . . . .	22
D.2. Behavior Task Metrics . . . . .	23
D.3. Motion Task Metrics . . . . .	23
<b>E DriveLM-Agent</b>	<b>24</b>
E.1. Prompting with Context . . . . .	24
E.2. Trajectory Tokenization Details . . . . .	24
<b>F. Experiments</b>	<b>24</b>
F.1. Implementation Details . . . . .	24
F.2. Results with More Metrics in VQA . . . . .	25
F.3. Ablation on Zero-shot Generalization across Sensor Configurations . . . . .	25
F.4. Computational Complexity . . . . .	26
<b>G Qualitative Results</b>	<b>26</b>
G.1. DriveLM-nuScenes . . . . .	26
G.2. Waymo . . . . .	26
G.3. DriveLM-CARLA . . . . .	26
<b>H More Related Work</b>	<b>32</b>

## A. Motivating Questions

**Q1.** *In what situations could we expect planning with VLMs to outperform conventional end-to-end autonomous driving?*

One of the key challenges of autonomous driving is to generalize to the long-tail of scenarios, that are rarely encountered but have critical importance. Considering the large-scale pre-training of VLMs, their acquired knowledge of the world, and the reasoning ability of the LLM, it is anticipated that planning with VLMs work better, particularly in situations that are novel or unseen in the context of driving scenarios but encountered during pre-training in unrelated contexts.

**Q2.** *Why adapt general VLMs to driving rather than adding language inputs to driving-specific models?*

General VLMs benefit from billion-scale pre-training data for vision-language tasks extracted from the internet, which can be adapted to the driving domain through fine-tuning on small autonomous driving datasets like DriveLM. Conversely, driving-specific models are only pre-trained on small autonomous driving datasets, and adding language inputs to these with data from outside the self-driving domain is non-trivial. Combining the advantages of VLMs and driving-specific models is however an interesting direction to explore.

**Q3.** *Can open-loop evaluation of planning provide meaningful results?*

When performing open-loop evaluation, providing the ego history as an input to the planning module prevents fair comparisons, as this signal alone is sufficient for achieving low errors on existing benchmarks. DriveLM alleviates this issue by evaluating key frames, where the intention of the ego vehicle changes, and the ego history is not strongly indicative of the future behavior or motion. Additionally, we consider baselines in our analysis that do not input the ego history to the planning module. Finally, we introduce DriveLM-CARLA as a means to show closed-loop planning results in the future.

**Q4.** *Why are there currently no closed-loop planning results on CARLA?*

Running 4B parameter models at 20 FPS as required by CARLA needs more engineering effort. This could be solved by using distillation, quantization, and caching techniques in LLM inference. Another approach would be to execute only the final motion stage of DriveLM-Agent at 20 FPS, while the other GVQA stages are executed at a lower frame rate.

**Q5.** *Is DriveLM-Agent efficient enough to be applicable to real-world autonomous driving?*

We comment on the run-time of DriveLM-Agent in Table 6. Without any optimization, the approach is around 1 order of magnitude slower than UniAD. However, with the optimizations proposed for closed-loop results on CARLA (see Q4), practical applications of VLMs in driving should be possible.

**Q6.** *Why is VQA more suited than alternative techniques to train internet-scale models (such as generative modeling) for the downstream application of autonomous driving?*

Both perception and planning in driving require reasoning and involve zero-shot generalization. VLMs potentially have the reasoning ability inherited from LLMs, making VQA a promising direction for bringing the benefits of web-scale training to autonomous driving.

**Q7.** *Do today's VLMs understand and reason about the visual world as well as LLMs understand text-based worlds?*

This is not known but deserves to be explored because VLMs approach the problem of generalization in a data-driven way, which has been proved successful repeatedly.

**Q8.** *Why does the proposed graph reasoning scheme not provide strong improvements in VQA?*

It is possible that the simple prompting scheme, relatively small base VLMs, or insufficiently strong logical dependencies in the dataset (or a combination of these factors) contribute to the lack of clear improvements. DriveLM-CARLA provides a platform to carefully study these factors and inform the annotation of future datasets for GVQA.

Method	#Params	#Trainable	FLOPs	FPS
UniAD-Single	131.9M	58.8M	1.7T	1.8
DriveLM-Agent	3.955B	12.9M	24.2T	0.16

Table 6. **Computational complexity and runtime.** Compared to UniAD-Single, DriveLM-Agent has fewer trainable parameters. For now, the efficiency is not good, and requires engineering efforts.

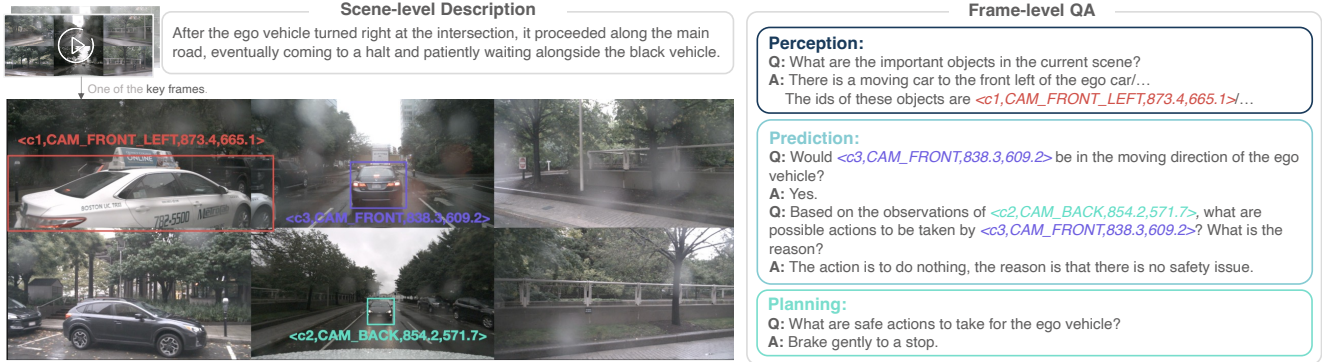


Figure 5. **Overall Composition of DriveLM-nuScenes.** The dataset comprises scene-level descriptions and frame-level QA, which can be divided into three parts: *Perception*, *Prediction*, and *Planning*. Objects are encoded using *c tags*, which contain identifiers, camera affiliations, and center coordinates of its 2D bounding box in the corresponding camera frame.

## B. DriveLM-nuScenes

In this section, we introduce the details of DriveLM-nuScenes, including the dataset composition, collection methodology, and statistics.

### B.1. Dataset Composition

DriveLM-nuScenes comprises a **training set of 4072 frames** and a **validation set of 799 frames**, consisting of scene-level descriptions and frame-level QA accompanied by 2D bounding boxes within multi-view images from the nuScenes dataset. The scene-level description delineates the behavior of the ego vehicle throughout the entire video clip. The frame-level QA encompasses three distinct categories: perception, prediction, and planning.

- **Perception** involves queries related to the thorough examination of the entire frame. Apart from several questions in this question set that are manually annotated, we design prompts to generate questions about the observational facets of objects within the scene, leveraging ground truth from nuScenes [5] and OpenLane-V2 [70].
- **Prediction** encompasses a series of inquiries regarding the projection of the forthcoming state of key objects and the ego vehicle in the current frame, and the underlying reasoning process behind the prediction. Because the predictions are intricate and challenging, we *manually annotate* the answers.
- **Planning** contains questions related to planning subsequent actions of the ego vehicle in the current frame. As “Planning” is the same challenging as “prediction”, we design the prompt for the reasoning process and *manually annotate* the answers to the questions.

For the key objects referred to in the QA, we encode them as **c tags** in the format  $\langle c, CAM, x, y \rangle$ , where  $c$  is the identifier,  $CAM$  indicates the camera where the object’s center point is situated, and  $x, y$  represent the horizontal and vertical coordinates of the 2D bounding box in the respective camera’s coordinate system. We also provide a dictionary in each key frame, recording more basic information about the key objects such as the size of the bounding box, the category, the moving state, and the visual description. The overview of data organization forms is shown in Fig. 5.

### B.2. Collection Methodology

During the annotation process, we employ individuals with driving experience for the labeling task. We provide annotators with the stitched results from the six cameras of nuScenes as source data. As shown in Fig. 6 (left), we divide the annotation process into three steps: selecting key frames from video clips, choosing key objects within these key frames, and subsequently annotating the frame-level QAs in the key frames. Following this, we conduct multiple rounds of quality checks

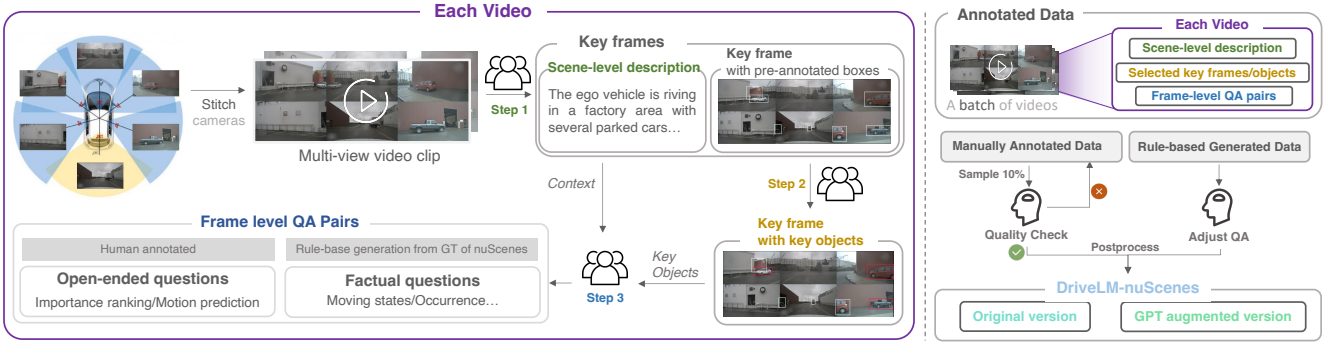


Figure 6. **(Left) Pipeline of the three-steps annotation process.** For each video, we ask the annotators to annotate the key frames, key objects, and QA attributes step by step. **(Right) The quality check and post-processing progress.** We divided the annotated data into batches, where each batch contains 8 video clips and their related annotations. We conduct rigorous quality checks, and after the post-processing, we finally get two versions of our DriveLM-nuScenes dataset.

to ensure the data reliability and perform post-processing procedures on the qualified data as shown in Fig. 6 (right). The specific details of this pipeline will be introduced below.

**Key Frame Selection.** In this process, we ask annotators to review the entire video clip to pinpoint key frames rich in scene information and potentially indicative of future state changes. Simultaneously, annotators are instructed to label the ego vehicle’s behavior throughout the video clip. This segment serves as the foundation for our scene-level description.

**Key Object Selection.** In this annotation step, we instruct annotators to identify objects in key frames that are relevant to the ego vehicle’s driving, denoted as key objects. To ensure accuracy, we provide pre-annotated bounding boxes based on ground truth categories from nuScenes [5]. Annotators also have the flexibility to designate objects not present in the ground truth as key objects if they are deemed significant.

**QA Labeling.** In the QA labeling process, we have two sets of questions, factual questions and open-ended questions. For the factual questions, we generate the answers with a rule-based method. For the open-ended questions, we instruct annotators to manually annotate the meticulously designed questions. Options are provided for most manually annotated questions, and we include an "Other - Fill in the Blank" option for answer choices in such cases to ensure flexibility. We have also incorporated free-form questions, allowing annotators to generate their own inquiries about the current frame.

**Quality Check.** We prioritize the quality of our data. In addition to establishing clear criteria and implementing autonomic checking strategies at each annotation step, we conduct rigorous manual quality checks. We organize the final data into batches, with each batch comprising 8 video clips, along with their scene-level descriptions, key frames with key objects selected from the 8 video clips, and corresponding QA pairs for each key frame. We provide explicit standards to quality check inspectors, instructing them to assess data eligibility based on these criteria. For manually annotated data, if the accuracy of the manual annotations falls below expectations for a particular batch, we compile feedback on the issues encountered and request annotators to re-annotate the entire batch. For data generated from ground truth, we instruct quality inspectors to manually adjust the inconsistent or unreasonable QA pairs.

**Post Processing.** Since our annotators are Chinese speakers, we need to translate the labeled data into English after obtaining it. Initially, we establish mappings between Chinese and English using a vocabulary. For texts that are not successfully mapped, we utilize GPT-3.5 for translation, and then perform manual checks and corrections on the GPT outputs. We also provide a version augmented by GPT-3.5, utilizing the prompt as shown in Table 7.

```

Messages = [
  { "role": "system", "content": f"" You are an English improver. "" },
  { "role": "user", "content": f"" I have a question and answer that I need you to help me modify and embellish, please make a few simple changes to the content in written language and keep the meaning same, you only need to answer the changes to: {QA}"" } ]

```

Table 7. **Prompt for GPT-refined version of DriveLM-nuScenes.** We try 50 different prompt and select this pattern as the final one to do the refinement.



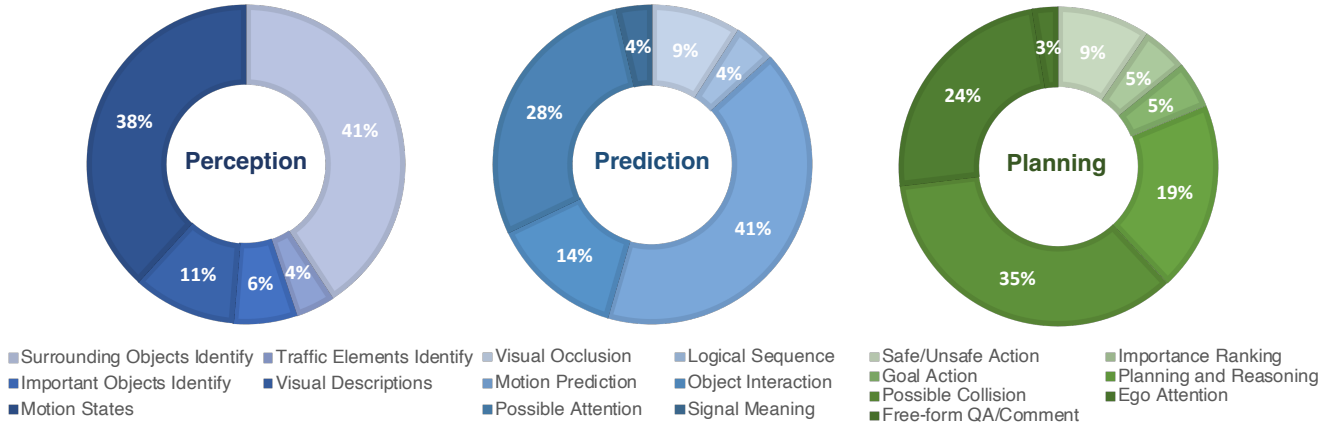


Figure 7. The distribution of question types according to different tasks in DriveLM-nuScenes. We categorize questions into perception, prediction, and planning tasks, each further subdivided into more specific question types.

### B.3. Statistics and Facts

In this section, we conduct a distribution analysis of our DriveLM-nuScenes QA categories at both the task level and object level. Additionally, for the task level, we provide the templates for all our QA under this classification criterion. The results indicate the richness of our QA categories, covering various aspects of autonomous driving. Moreover, the abundance of logical relationships is sufficient to construct a graph-structured QA.

**Task level.** Our DriveLM-nuScenes orchestrates a benchmark that encompasses various aspects of autonomous driving, connecting the whole stages of the human driving logic. To delve deeper into this aspect, we present the detailed QA types distribution at task level in Fig. 7. For a better understanding, we also provide examples of QA templates in all of the P3 stages in Table 8.

#### Perception

##### Surrounding Objects Identify

Q: Please describe the current scene.

A: There are two moving cars behind the ego car and two barriers in front of it.

Q: What are objects to the front left/back right/... of the ego car?

A: There are two barriers to the front left of the ego car.

Q: Are there traffic cones/moving cars/... to the front right/back left/... of the ego car?

A: No.

##### Traffic Elements Identify

Q: Is there any traffic element in the front view?

A: Yes, there are some traffic elements in the front view.

Q: Identify all the traffic elements in the front view, categorize them, determine their status, and predict the bounding box around each one. The output should be a list formatted as (c, s, x1, y1, x2, y2), where c represents the category, s denotes the status, and x1, y1, x2, y2 are the offsets of the top-left and bottom-right corners of the box relative to the center point.

A: There are three traffic elements in the front view. The information of these traffic elements are [(road sign, go straight, 907.58, 590.67, 992.54, 630.95)...].

##### Important Objects Identify

Q: What are the important objects in the current scene? Those objects will be considered for the future reasoning and driving decision.

A: There is a parked truck to the back of the ego car... The ids of these objects are <c1,CAM\_BACK,827.5,484.2>...

Q: What is the relative positioning of the important objects in the current scene?

A: <c3,CAM\_FRONT,689.2,527.5> is to the front of <c1,CAM\_BACK,827.5,484.2>...

Q: Which lanes are each important object on in the scene?

A: <c2,CAM\_FRONT,820.8,473.3> is on the ego lane...

#### Visual Description

Q: What is the visual description of <c2,CAM\_FRONT\_LEFT,415.8,580.8>/...?

A: Pedestrian riding a bicycle.

#### Motion State

Q: What is the status of the cars/pedestrians/... that are to the front/front right/... of the ego car?

A: Many cars are parked.

Q: What is the observed status of object <c1,CAM\_FRONT,920.0,509.2>/...?

A: Moving.

Q: What is the moving status of object <c1,CAM\_FRONT,920.0,509.2>/...?

A: Going ahead.

#### Prediction

##### Visual Occlusion

Q: Which object is most likely to be occluded by <c1,CAM\_FRONT,707.5,472.5>/...? Would this object affect the ego vehicle? Based on this object, what action of the ego vehicle is dangerous?

A: The object in front of <c1,CAM\_FRONT,840.8,507.5>, yes, accelerating forward.

##### Logical Sequence

Q: What object should the ego vehicle notice first when the ego vehicle is getting to the next possible location? What is the state of the object that is first noticed by the ego vehicle and what action should the ego vehicle take? What object should the ego vehicle notice second when the ego vehicle is getting to the next possible location? What is the state of the object perceived by the ego vehicle as second and what action should the ego vehicle take? What object should the ego vehicle notice third? What is the state of the object perceived by the ego vehicle as third and what action should the ego vehicle take?

A: Firstly notice that <c2,CAM\_FRONT,514.7,462.2>, the state of it is traffic sign, so the ego vehicle should slow down and go ahead. Secondly notice that <c3,CAM\_FRONT,950.3,613.1>, the state of it is traffic sign, so the ego vehicle should slow down and go ahead. Thirdly notice that <c1,CAM\_FRONT,707.5,472.5>, the state of it is going ahead, so the ego vehicle should slow down and go ahead.

##### Motion Prediction

Q: Would <c1,CAM\_FRONT,920.0,509.2>/... be in the moving direction of the ego vehicle?

A: Yes.

Q: What is the future state of <c1,CAM\_FRONT,920.0,509.2>/...?

A: Keep going straight.

Q: Will <c2,CAM\_FRONT,1223.3,598.3>/... be in the moving direction of <c1,CAM\_BACK,514.2,503.3>/...?

A: No.

##### Object Interaction

Q: Will <c2,CAM\_FRONT,1223.3,598.3>/... change its motion state based on <c1,CAM\_BACK,514.2,503.3>/...?

A: No.

Q: Based on the observations of <c1,CAM\_BACK,514.2,503.3>/..., what are possible actions to be taken by <c2,CAM\_FRONT,1223.3,598.3>/...? What is the reason?

A: The action is to keep going at the same speed, the reason is there is no safety issue.

Q: Based on the observation of <c4,CAM\_FRONT,1071.2,346.2>/..., what actions may <c1,CAM\_FRONT,1126.7,515.0>/... take?

A: The action is to keep going at the same speed, the reason is there is no safety issue.

##### Possible Attention

Q: In this scenario, what object is most likely to consider <c3,CAM\_FRONT,400.1,717.2>/...?

A: The ego vehicle.

Q: Would <c1,CAM\_BACK,514.2,503.3>/... take <c3,CAM\_FRONT,400.1,717.2>/... into account?

A: No.

Q: What object would consider <c1,CAM\_FRONT,985.8,516.7>/... to be most relevant to its decision?

A: The ego vehicle.

Q: Except for the ego vehicle, what object would consider <c1,CAM\_FRONT,985.8,516.7>/... to be most relevant to its decision?

A: <c2,CAM\_FRONT,1217.5,511.7>.

#### Signal Meaning

Q: What does <c2,CAM\_BACK\_LEFT,400.8,654.2>/... mean?

A: No entry.

Q: What kind of traffic sign is <c2,CAM\_BACK\_LEFT,400.8,654.2>/...?

A: Traffic cone.

#### Planning

##### Safe/Unsafe Action

Q: In this scenario, what are safe actions to take for the ego vehicle?

A: Decelerate gradually without braking, keep going at the same speed.

Q: In this scenario, what are dangerous actions to take for the ego vehicle?

A: Accelerate and go ahead, brake suddenly, drive backward, turn right.

##### Importance Ranking

Q: What is the priority of the objects that the ego vehicle should consider? (in descending order)

A: <c2,CAM\_FRONT,514.7,462.2>, <c3,CAM\_FRONT,950.3,613.1>, <c1,CAM\_FRONT,707.5,472.5>.

##### Goal Action

Q: What is the target action of the ego vehicle?

A: Go straight.

##### Planning and Reasoning

Q: What actions could the ego vehicle take based on <c1,CAM\_FRONT,920.0,509.2>/...? Why take this action and what's the probability?

A: The action is to decelerate gradually without braking, the reason is to keep a safe distance, high.

Q: Based on <c3,CAM\_FRONT,1591.1,441.8>/... in this scene, what is the most possible action of the ego vehicle?

A: Decelerate gradually without braking.

##### Possible Collision

Q: What is the probability of colliding with <c1,CAM\_FRONT,920.0,509.2>/... after the ego vehicle goes straight and keeps the same speed/accelerates and goes straight/...?

A: Low.

Q: What actions taken by the ego vehicle can lead to a collision with <c1,CAM\_FRONT,920.0,509.2>/...?

A: Accelerate and go straight.

##### Ego Attention

Q: What is the traffic signal that the ego vehicle should pay attention to?

A: None.

Q: Is <c1,CAM\_FRONT,920.0,509.2>/... an object that the ego vehicle should consider in the current scene?

A: Yes.

Q: Is it necessary for the ego vehicle to take <c3,CAM\_FRONT,400.1,717.2>/... into account?

A: Yes.

##### Free-form QA/Comment

Q: What impact does this situation have on driving vehicles?

A: The road scene is complex, please slow down.  
 Q: What’s your comment on this scene?  
 A: Pedestrians at the intersection, please be careful and give way.  
 ...

Table 8. **Question templates of DriveLM-nuScenes at task level.** The categories in the table correspond to those in Fig. 7.

**Object level.** We also conduct some statistics at the object level since QAs in our DriveLM-nuScenes revolve around key objects. Fig. 8 (left) shows the distribution of our key object types. Given the substantial differences in questions associated with traffic elements compared to other categories, we separately conduct statistics for QA types related to traffic elements and the remaining categories. The results are depicted in Fig. 8 (right).

### C. DriveLM-CARLA

In this section, we introduce the details of DriveLM-CARLA, including the dataset composition and collection methodology.

#### C.1. Dataset Composition

DriveLM-CARLA consists of automatically generated frame-level question-answer pairs that are structured with an interconnected graph. The graph structure can be seen in Fig. 9. In the current version, the dataset consists of questions about the road layout, stop signs, traffic lights, and vehicles. In future versions, the dataset can be extended to more categories like static objects, weather, other signs, and others.

Utilizing the driving simulator CARLA for the data generation process allows for scalable annotations and data without any manual effort involved. Additionally, the dataset supports a variety of sensor outputs from CARLA, including semantic segmentation, depth map, LiDAR, and others, which can be employed to train different network architectures. Each question within the graph is designed in a way to facilitate situational reasoning, which could be instrumental in answering subsequent questions. As with DriveLM-nuScenes, each question can be categorized into perception, prediction, or planning. For each QA-pair, besides the corresponding question and answer, we also save the object ID in case the QA-pair is about an object. This ID is consistent over time, enabling object tracking and temporal reasoning in future studies. In addition, relationships to parent and child questions within the graph are documented to allow efficient traversal of the graph.

#### C.2. Collection Methodology

In this section, we provide details about the data collection and the annotation process.

**Simulator settings.** We utilize the CARLA Simulator (version 0.9.14) with Leaderboard 2.0 [17] to generate our dataset. Leaderboard 2.0 introduces two new large maps along with a suite of new scenarios, enhancing the diversity of the training and evaluation environments. Town 12 serves as the training town, while Town 13 is reserved for evaluation. Each town

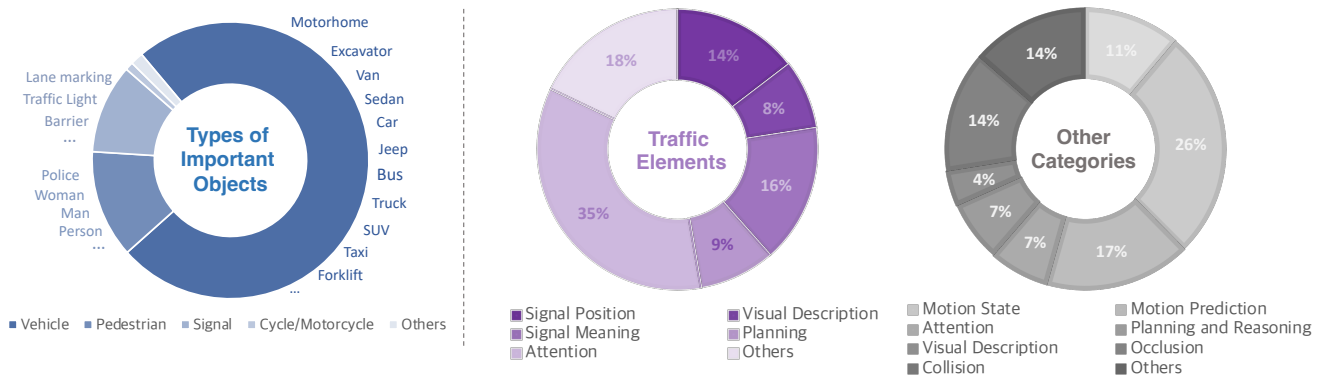


Figure 8. **(Left) The distribution of key objects in DriveLM-nuScenes.** The sub-categories are extracted from the visual description. **(Right) The distribution of question types related to different key objects in DriveLM-nuScenes.** Since the questions associated with traffic elements differ significantly from other categories, we separately conduct statistics for QA types related to traffic elements and the remaining categories.

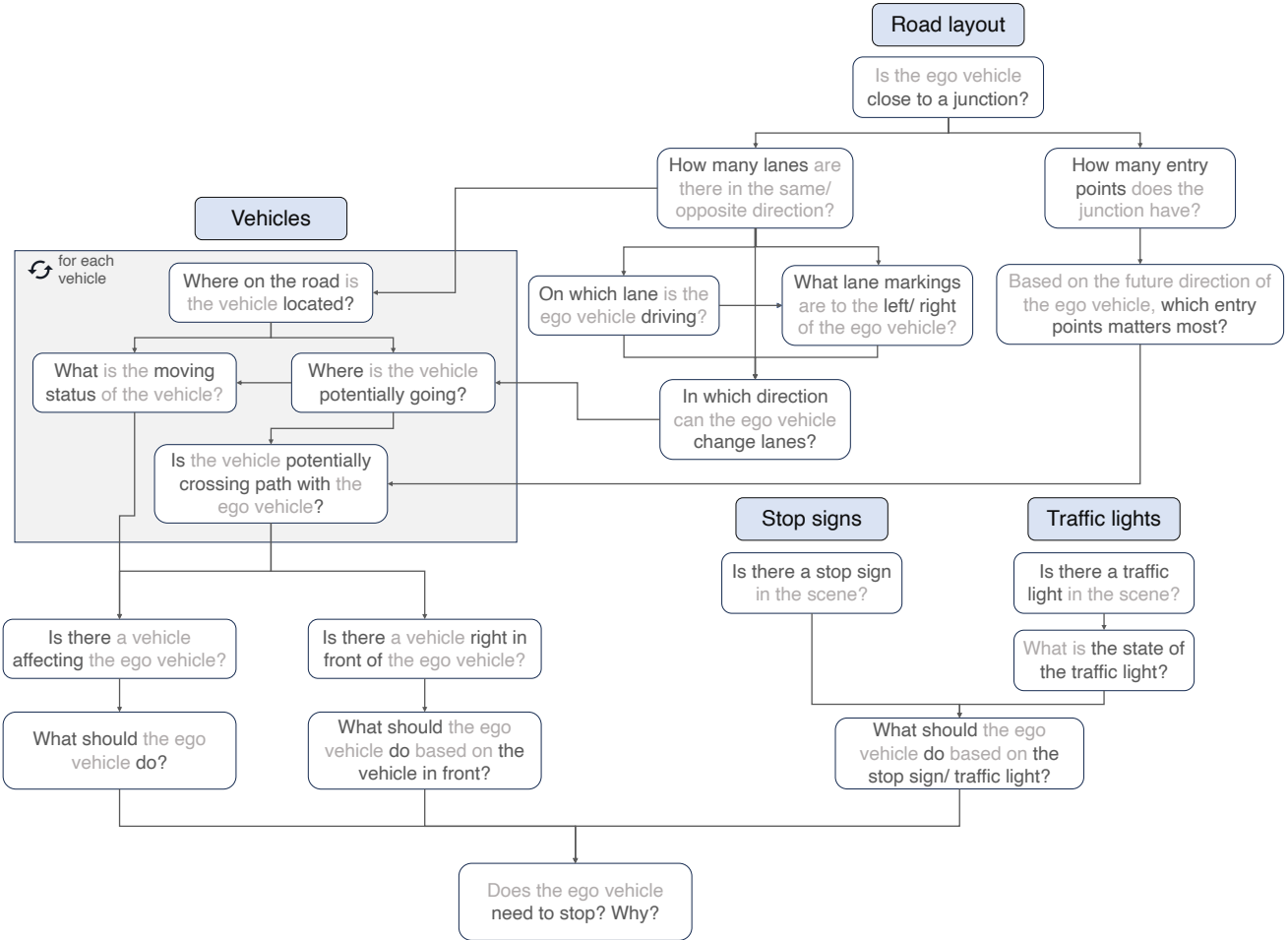


Figure 9. **Detailed flow on CARLA graph.** We show the full graph of DriveLM-CARLA. The graph consists of questions and answers about the road layout, traffic lights, stop signs, and vehicles.

covers an area of 10 x 10 square meters, encompassing varied environments such as rural, residential, and urban landscapes to replicate real-world driving conditions.

The CARLA team provides a total of 90 training routes, spanning 780.6 kilometers, and 20 evaluation routes, measuring 247.6 kilometers. Every route incorporates multiple driving scenarios. We segmented these routes into shorter segments, approximately 200 meters in length and filter routes that start and end at the same position. During the training phase, all predefined scenarios are omitted to test generalization ability and because the expert was designed for Leaderboard 1.0 and cannot solve most of the new scenarios. We use the pedestrian scenarios to assess the generalizability of DriveLM-Agent.

The traffic manager within Leaderboard 2.0 initializes random background traffic around the ego vehicle comprised exclusively of 'car' entities. To enrich the dataset with greater diversity, we introduce additional vehicle classes including 'trucks', 'vans', 'bicycles', and 'motorcycles'. Moreover, we implemented randomized weather configurations for each training and evaluation route to mimic realistic driving conditions. However, night-time settings were excluded from our study due to the inadequate illumination in certain map regions. Low-light conditions significantly impede the correctness of the automatic labeling process since it is hard to obtain information about the visibility of certain objects in the image.

**Expert.** We integrate the rule-based expert from Transfuser++ [30] and slightly modify it to function within the new version of the CARLA simulator (version 0.9.14). The expert follows the provided waypoints, typically the center lane of the road, driving with a pre-defined target speed. In each timestep, the expert evaluates the environment for potential hazards that could impede safe driving. Those hazards include traffic signals such as red lights and stop signs as well as dynamic obstacles like pedestrians crossing the road, and vehicles that may pose a risk of collision. Upon identification of any such hazard, the

expert starts braking, to ensure safety. With the transition to the newer CARLA version, the vehicle dynamics changed. These changes necessitated fine-tuning the parameters of the PID controller, which is responsible for the vehicle’s speed and steering control. The tuning aimed to mitigate any oscillatory behaviors that were particularly evident during braking and when navigating sharp turns.

**Data collection.** We execute the expert on each of the routes and gather a comprehensive set of sensor data. The sensor data includes:

- *RGB image*
- *LiDAR point cloud*
- *semantic segmentation images*
- *depth maps*
- *Bird’s Eye View (BEV) semantic segmentation*

While DriveLM-Agent leverages only RGB images, retraining TransFuser++ needs the additional data for the auxiliary tasks.

In addition, we extract privileged information from the simulator about the status of the static and dynamic objects in the scene, as follows:

- *Ego vehicle:* 3D bounding box, speed, brake, id
- *Other vehicles:* 3D bounding box, number of lidar points inside BB, distance to ego, speed, steer, throttle, brake, id, color, vehicle type, number of wheels, traffic light state, lane information (i.e., on which road and lane is the vehicle driving), vehicle in junction or not, distance to next junction, next high-level command
- *Pedestrians:* 3D bounding box, number of lidar points inside BB, gender, age, distance to ego, speed, id, lane information
- *Traffic lights:* 3D bounding box, distance to ego, state, affects ego vehicle
- *Stop signs:* 3D bounding box, distance to ego, affects ego vehicle
- *Static cars (parked cars):* 3D bounding box, lane information
- *Landmarks (e.g., speed signs):* 3D bounding box, distance to ego, id, text, value
- *Weather:* weather parameters

**Language labels.** Based on the information we extract from the simulator we create questions and answers with hand-crafted sentence templates. For more linguistic diversity and to prevent overfitting to those sentence structures those sentences could be further augmented with current state-of-the-art language models like GPT-4. However, in this work, we use a version of the dataset that is not augmented.

## D. DriveLM-Metrics

In this section, we offer a detailed introduction to DriveLM-Metrics. DriveLM-Metrics can be broadly categorized into three parts:  $P_{1-3}$  VQA Metrics, Behavior Task Metrics, and Motion Task Metrics.

### D.1. $P_{1-3}$ VQA Metrics

We assess the performance of  $P_{1-3}$  using common VQA metrics, and we introduce the *GPT score* for a more semantically comprehensive evaluation of our QA results. Additionally, given the graph structure of our QA, we propose the *Completeness* score to provide a thorough assessment.

**BLEU** [52] measures the similarity between a generated text and one or more reference texts. It operates by comparing n-grams in the generated text to those in the reference texts, with higher precision indicating a better match. However, the BLEU score exhibits insensitivity to semantic nuances and variations in word order.

**ROUGE\_L** [43] calculates scores with the longest common sub-sequence of the model outputs and the reference answers. Similar to the BLEU metric, ROUGE is used to assess the level of matching between generated results and standard references, with the key difference being that ROUGE is based on recall.

**METEOR** [38] takes into account precision, recall, stemming, synonymy, stemming, and word order. It establishes alignment between model outputs and references, computes the 1-gram matching between them, and then applies penalties based on chunk blocks, providing a more nuanced evaluation.

**CIDEr** [69] combines elements from BLEU and vector space models. The underlying concept involves treating each sentence as a document, calculating its n-gram TF-IDF vector, and using cosine similarity to measure the semantic consistency between candidate and reference sentences. CIDEr captures matches between n-grams of different lengths and differentiates the importance of various n-grams through TF-IDF weighting.

**SPICE** [2] first parses the text into a syntactic dependency tree using Probabilistic Context-Free Grammar [31], then maps the dependency tree into a scene graph in a rule-based manner. The scene graph describes the objects, attributes, and their relationship in the original text, and the SPICE score is computed as the F-score of the generated scene graphs from prediction and ground truth.

**GPT Score** is a metric provided by ChatGPT. Traditional metrics mainly assess word-level performance and may not capture semantic nuances, potentially yielding unexpected evaluation outcomes. Leveraging ChatGPT’s robust reasoning capabilities, we employ it to gauge prediction quality and derive a more rational score. ChatGPT is prompted to assign a numerical score between 0 and 100, with higher scores indicative of enhanced prediction accuracy. The detailed prompt for GPT score evaluation is shown in Table 9.

```

Messages = [
  { "role": "system", "content": f"" An evaluator who rates my answer based on the correct answer. ""
},
  { "role": "user", "content": f"" Rate my answer based on the correct answer out of 100, with higher
scores indicating that the answer is closer to the correct answer, and you should be accurate to single digits like 62,
78, 41, etc. This is the correct answer: {GT}. This is my answer: {Pred}. ""}]

```

Table 9. **Prompt for GPT score.** This differs from the prompt used in DriveGPT4 [81], but the resulting score is similar.

**Completeness** provides a score that accounts for how many ground truth questions are correctly answered associated with a frame. For each QA, if the score of the predicted answer is above a threshold, then this QA is considered “correctly answered” and is a correct prediction, otherwise an incorrect prediction. We then calculate the accuracy, which is the ratio of the number of correct predictions to the total number of predictions. In our setting, we utilize the SPICE score and set the threshold at 0.5.

### D.2. Behavior Task Metrics

We evaluate behavior predictions by classification accuracy, along with a breakdown of the overall accuracy into its steering and speed components.

**Classification Accuracy** is the metric we use to evaluate Behavior Prediction Task, comprising *accuracy of behavior*, *behavior speed*, and *behavior steer*. Specifically, the ground truth of the ego vehicle future trajectory is a set of  $N$  points with coordinates  $(x, y)$  under the bird’s-eye-view, noted as  $\{(x_0, y_0), (x_1, y_1), \dots, (x_N, y_N)\}$ . Each point denotes the offset of the future position to the current position by a number of fixed interval times. Then, the distance for  $x, y$  at each time interval is independently computed as:

$$\{x, y\}_{dist} = ((\{x, y\}_1 - \{x, y\}_0), \dots, (\{x, y\}_N - \{x, y\}_{N-1})) \tag{2}$$

The mean of  $x_{dist}$  and  $y_{dist}$  are mapped to one of the predefined bins, where each bin corresponds to a category in either speed or steering, noted as  $B_{speed}$  and  $B_{steer}$  respectively. Finally, the speed and steering categories for this trajectory form the behavior category as  $(B_{speed}, B_{steer})$ . We compare them with the behaviors of our DriveLM-Agent outputs and calculate the related accuracy.

### D.3. Motion Task Metrics

For measuring the performance of the motion stage, we use standard metrics from the nuScenes and Waymo benchmarks: average and final displacement error, (ADE, FDE), and the collision rate of the predicted trajectory.

**ADE** stands for Average Displacement Error, indicating the average L2 distance between the predicted trajectory and the ground truth trajectory over all predicted horizons. It is the average of the errors at 1<sup>st</sup>, 2<sup>nd</sup> and 3<sup>rd</sup> second.

**FDE** stands for Final Displacement Error, which measures the Euclidean distance between the predicted endpoint and the true endpoint at the last predicted step (the 3<sup>rd</sup> second).

**Collision Rate** accounts for the ratio of how many test frames the predicted trajectory collides with objects in over all test frames. The number reported in Table 2 of the main paper is the average of the collision rate at 1<sup>st</sup>, 2<sup>nd</sup> and 3<sup>rd</sup> second.

Note that the calculation of **ADE**, **FDE** and **Collision Rate** follows the setting used in UniAD [26] but not ST-P3 [25]. For example, in terms of the FDE and the collision rate at 3<sup>rd</sup> second, the UniAD setting will consider the error/collision rate at only this timestep, while the ST-P3 setting will consider the error/collision rate as an average over 0.5, 1, 1.5, 2, 2.5,

3 seconds. For more details, please refer to the UniAD repo [discussion](#). Additionally, please note that errors reported on the full nuScenes validation dataset (in prior work) is not directly comparable to results reported on the DriveLM-nuScenes val split, a challenging subset of this consisting of only keyframes with intention changes.

## E. DriveLM-Agent

In this section, we introduce the details of DriveLM-Agent, including the graph prompting scheme and the trajectory tokenization process.

### E.1. Prompting with Context

In terms of the implementation, the content of context differs during the training and inference of DriveLM-Agent, following the teacher-forcing setting [37, 67] generally adopted in recurrent networks. During training, for each edge  $e \in E$  in the frame, we pick the child QA. The child questions in the edges are appended with the ground truth parent QA as the context. All QA pairs are used during training, including those without context. The objective used is next token prediction, the standard approach for language modeling. During inference, the model is applied interactively in multiple rounds to get the required context predictions as inputs for each child question. Specifically, the model is prompted with the five stages of questions in the sequential order of  $P_1, P_2, P_3, B, M$ . In this order, the model can only infer the questions in the succeeding stages after getting the predicted answer from the preceding stages.

### E.2. Trajectory Tokenization Details

To generate action sequences (*i.e.*, ego future trajectories) directly with the language model we use for building the graph, we adopt the method of RT-2 [85]. This process entails the discretization and tokenization of the continuous trajectory.

Initially, we analyze the distribution of the future trajectories within the nuScenes dataset. To effectively convert the continuous  $(x, y)$  coordinate space into a discrete set of actions, we partition each coordinate axis into 256 discrete intervals. This granularity ensures a sufficient level of detail while maintaining a manageable number of tokens for the language model.

Each discretized bin corresponds to a unique token within the vocabulary of the language model. We extract the token identifiers (IDs) for numeric tokens within the vocabulary. To ensure coherence and preserve the ability to express numerical values, we omit single-digit tokens from this mapping process. Out of the remaining numeric tokens, a subset of 256 token IDs is selected to represent the trajectory data. In addition to these, we introduce two special tokens designated for marking the start and end of a trajectory sequence – the start-of-trajectory (SOT) token and the end-of-trajectory (EOT) token, respectively. This tokenization scheme enables us to encode complex trajectory information as a sequence of tokens that a language model can process. Using this mapped vocabulary, the language model can generate predicted future trajectory sequences by outputting a series of tokens, which are then translated back into the coordinate space.

## F. Experiments

In this section, we introduce the details of experiments, including implementation details of each subsection in **Section 4** in the main paper, more metrics on the VQA part, more ablation and comparison on computational complexity.

### F.1. Implementation Details

Here we explain the implementation details for the training and validation settings used in our experiments of each of each subsections in **Section 4** in the main paper.

**Fine-tuning Details.** We configure the learning rate as 0.0001, no learning-rate scheduler, random seed as 1234, and other settings following the default LoRA [24] configuration. For the BLIP-2 model, we use a maximal sequence length of 400, and other hyperparameters remain the same as the official BLIP-2 implementation.

**Implementation Details for Experiment in Section 4.1 & 4.4.** During training, we utilize all QAs as input per frame, with a subset of them having contexts (questions from  $P_{2,3}, B$ , and  $M$ ). The contexts are extracted from ground truth, following the teacher-forcing setting [37, 67] generally adopted in recurrent networks. As for inference, due to the variant complexity of the scenarios, the count of  $P_{1-3}$  QA per frame is highly imbalanced across the dataset, with a variance of over 260 on DriveLM-nuScenes. To balance the impact of this, we compute the GVQA Scores on only a subset of QA associated with each frame. To extract the QA subset, we design a set of QA patterns for each stage based on the questions generally associated with that stage. We ensure that for all our validation frames, each stage has at least one question matched with the designed pattern. In this process, except for the questions in stage  $P_1$ , all questions in other stages have context from the



Context	DriveLM-nuScenes													
	Off-the-Shelf BLIP-2							DriveLM-Agent						
	BLEU-4 ↑	METEOR ↑	CIDEr ↑	ROUGE_L ↑	SPICE ↑	GPT ↑	Comp. ↑	BLEU-4 ↑	METEOR ↑	CIDEr ↑	ROUGE_L ↑	SPICE ↑	GPT ↑	Comp. ↑
None	0.022	3.317	0.1185	7.205	4.336	42.97	1.064	51.89	35.81	2.463	66.79	42.56	71.39	30.04
Graph	0.022	3.882	0.0771	7.397	7.710	45.21	0.859	53.09	36.19	2.786	65.58	49.54	72.51	31.66
<i>GT</i>	<i>0.022</i>	<i>4.397</i>	<i>0.0758</i>	<i>8.033</i>	<i>8.192</i>	<i>41.10</i>	<i>1.315</i>	<i>53.06</i>	<i>36.64</i>	<i>3.069</i>	<i>66.69</i>	<i>50.29</i>	<i>72.94</i>	<i>32.41</i>

Table 10. **Graph-structured reasoning facilitates improved VQA with VLMs.** Completeness measures how many questions are correctly answered in one frame of data. The improvement trends are not consistent across different conventional metrics, thus we need the GPT score as the main metrics as it evaluates the performance more comprehensively.

Method	Behavior Context	Motion Context	Behavior ( <i>B</i> )			Motion ( <i>M</i> )	
			Acc. ↑	Speed ↑	Steer ↑	ADE ↓	FDE ↓
Command Mean	-	-	-	-	-	7.98	11.41
UniAD-Single	-	-	-	-	-	4.16	9.31
BLIP-RT-2	-	-	-	-	-	2.78	6.47
DriveLM-Agent	None	<i>B</i>	35.70	43.90	65.20	2.76	6.59
	$P_1$	<i>B</i>	38.20	43.67	70.74	2.67	6.41
	$P_2$	<i>B</i>	39.52	44.20	<b>78.67</b>	<b>2.62</b>	6.19
	$P_3$	<i>B</i>	34.62	41.28	64.55	2.85	6.89
	$P_{1-3}$	<i>B</i>	<b>39.73</b>	<b>54.29</b>	70.35	2.63	<b>6.17</b>
<i>BLIP-RT-2</i>	-	<i>B*</i>	<i>100.0</i>	<i>100.0</i>	<i>100.0</i>	<i>2.41</i>	<i>5.79</i>

Table 11. **Zero-shot Generalization across Sensor Configurations.** *B\** denotes using ground truth behavior QA as context for motion task. Results on 1k randomly sampled frames from the Waymo val set after training on DriveLM-nuScenes. Our key observation is that the higher the accuracy of the behavior task, the better the performance of the motion task.

previous stage’s QA, where the answers are derived from the prediction in the preceding steps. Two specific graph structure examples can be found in Fig. 10 and Fig. 11.

**Implementation Details for Experiment in Section 4.2.** The model is trained in the same scheme as in Section 4.1 & 4.4, and the training set is the DriveLM-nuScenes train split. During the inference, as there are no annotation on the Waymo (not even questions), we devise the question in a rule-based manner. Specifically, we re-use the general-purpose questions in the perception stage from DriveLM-nuScenes for Waymo as the starting questions. Then we try to find if there is any objects in the answer that is matched in the DriveLM-nuScenes annotation, such as “pedestrians”, “cars”, “trucks” and so on. Then, we generate the questions based on those matched objects automatically, which serve as the following questions in the prediction and the planning stages. A specific graph structure example can be found in Fig. 12.

**Implementation Details for Experiment in Section 4.3.** For the generalization experiment in Section 4.3 we add two new questions to the graph: (1) *Is there a person in the scene?* and (2) if the answer to the first question is yes we ask *What should the vehicle do based on the pedestrian that is crossing the road?* if the answer is no we ask *What should the ego vehicle do?*. The answer to the second question gets concatenated to the context of the final behavior question. Three examples can be found in Fig. 13.

## F.2. Results with More Metrics in VQA

In Table 10, we provide the performance under BLEU-4 [52], METEOR [38], CIDEr [69] and ROUGE-L [43] of the Table 5 in the main paper. One key observation is that different metrics reflect different trends in the performance, and the improvement is not consistent across all these metrics. This brings us the motivation to use the GPT Score as the main metric in the VQA evaluation part.

## F.3. Ablation on Zero-shot Generalization across Sensor Configurations

In Table 11, we provide more settings of context in the zero-shot generalization across sensor configurations as in Table 3 in the main paper. One key observation is that the higher the accuracy of the behavior task, the better the performance of the motion task. With more context in the behavior task, the accuracy improvement mainly originates from the improvement of the speeding accuracy, which finally affects the FDE score.

## F.4. Computational Complexity

In Table 6, we provide a comparison of the computational complexity of DriveLM-Agent to UniAD-Single. A future direction would be caching the vision tokens and batching the different question patterns, which can speed up the inference time fundamentally.

## G. Qualitative Results

In this section, we show the qualitative results of the experiments, including VQA on our DriveLM-nuScenes, generalization results across sensor data on Waymo, and generalization results to unseen objects on DriveLM-CARLA.

### G.1. DriveLM-nuScenes

This section shows the qualitative examples for the DriveLM-nuScenes. In Fig. 10, we showcase a detailed example of GVQA reasoning process on DriveLM-nuScenes, encompassing  $P_{1-3}$  QA and the behavior task. We compare the predicted answers with ground truth and provide SPICE scores and GPT scores. In this figure, the second question in the prediction stage represents a typical error. Due to the input of single-frame images, our model often struggles to accurately determine the correct movement status of objects. This judgment is indeed challenging even for humans. Furthermore, in Fig. 11, we present additional qualitative results to showcase our model’s performance.

### G.2. Waymo

This section demonstrates the generalizability of our model across sensor configurations. Fig. 12 illustrates the results of our model, trained on DriveLM-nuScenes, when applied to inference on Waymo. As we do not annotate data on Waymo, the questions are manually defined, and ground truth is not provided. These results showcase the robust generalization capability of our model.

### G.3. DriveLM-CARLA

In this section, we provide qualitative examples for the CARLA dataset.

**Generalization to the unseen pedestrian scenario.** Fig. 13 shows the generated behaviors for the generalization test set on the unseen pedestrian scenario. The first example, illustrated on the top of Fig. 13, demonstrates a successful situation where DriveLM-Agent accurately recognizes a pedestrian. It subsequently infers the appropriate action to undertake, which in this case, is to stop the vehicle. The behavior generation is able to interpret this context, resulting in the correct behavior pattern, as evidenced by the ego vehicle coming to a complete stop. The other two examples represent the predominant failure modes of DriveLM-Agent in scenarios involving pedestrians. The middle example of Fig. 13 shows the case where the model still detects the pedestrian. However, it fails to translate this detection into the correct behavior. The final example, shown at the bottom, highlights a more severe limitation where the model completely overlooks the pedestrian. In such instances, DriveLM-Agent acts as if the pedestrian is non-existent, which consequently results in it not executing any evasive or stopping maneuvers, posing a significant risk in a real-world scenario.

**Graph Visual Question Answering.** This section presents two examples of the graph visual question answering tasks using the CARLA dataset to evaluate the performance of DriveLM-Agent (Fig. 14). We only show a subset of the evaluated questions. In the first example the ego vehicle drives behind another vehicle. The primary task is to follow the road and adjust the speed in accordance with the leading vehicle. Our results indicate that DriveLM-Agent demonstrates proficient scene understanding by accurately identifying all important objects in the scene. Despite the ground truth data indicating occasional inaccuracies in object color labeling by the CARLA simulator, DriveLM-Agent maintains reliable performance in object recognition. Additionally, the model can identify the vehicle in front and reason about what to do based on the leading vehicle.

The second example takes place at an intersection regulated by a stop sign. DriveLM-Agent identifies all objects and can reason about the situation. It correctly identified that it needs to stop not simply due to the stop sign, but primarily because of a motorcycle positioned ahead. This implies that DriveLM-Agent is capable of prioritizing dynamic obstacles over traffic control devices under certain circumstances.

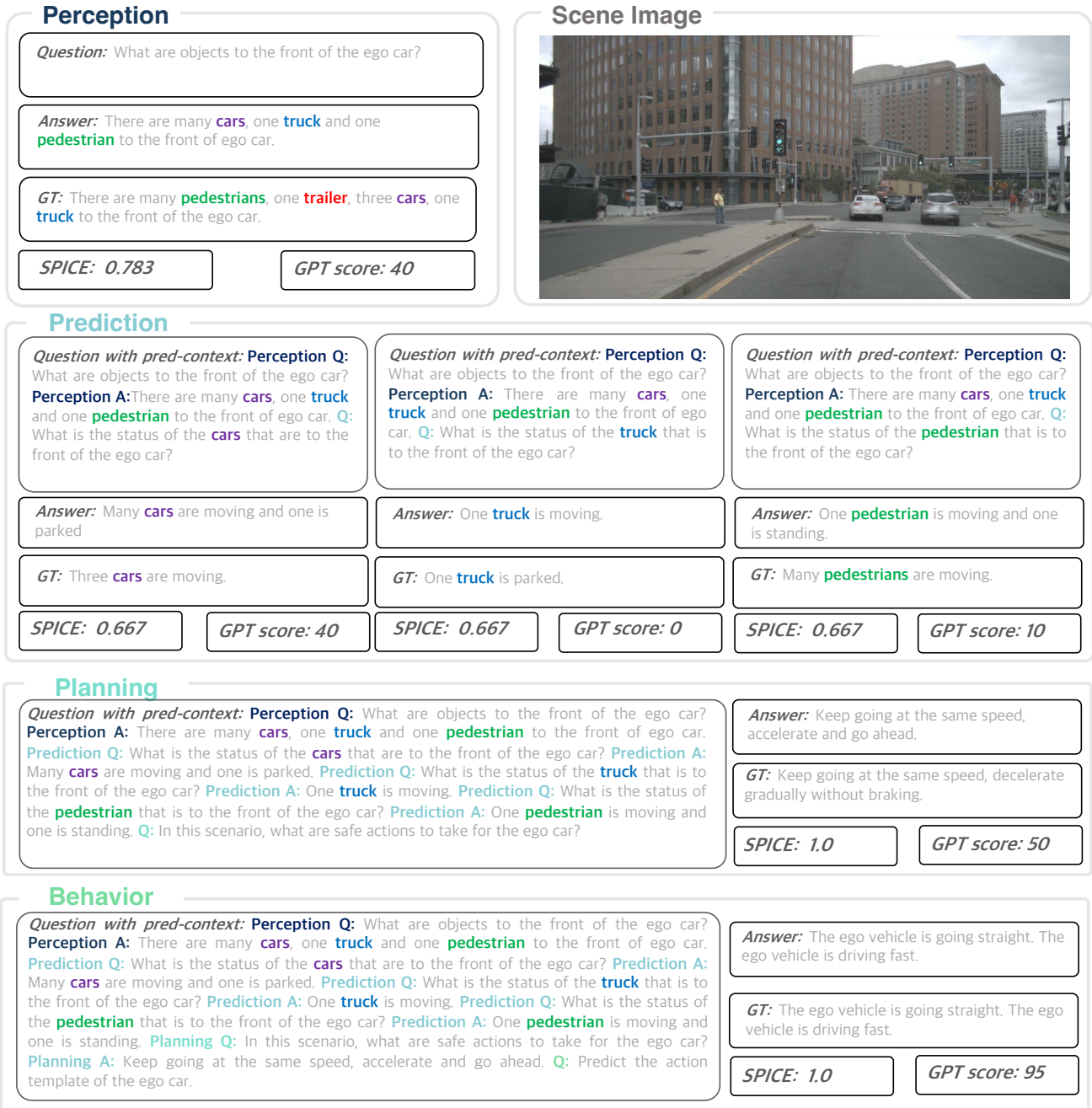


Figure 10. Detailed qualitative results on DriveLM-nuScenes. The graph prompting process can be divided into different tasks, and different QAs in each task revolve around different objects.

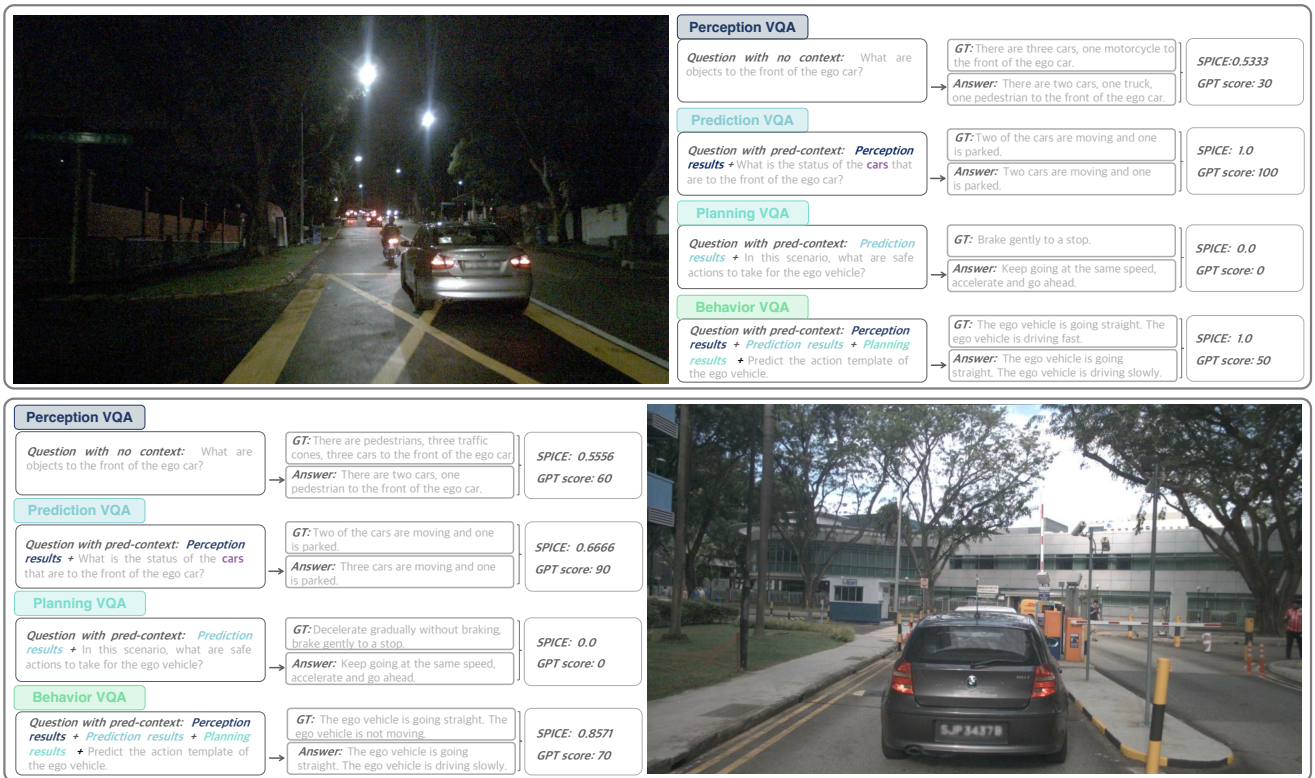


Figure 11. More qualitative results on DriveLM-nuScenes. The examples in the figure illustrate the robust ability of our DriveLM-Agent to perform VQA tasks in driving scenarios.

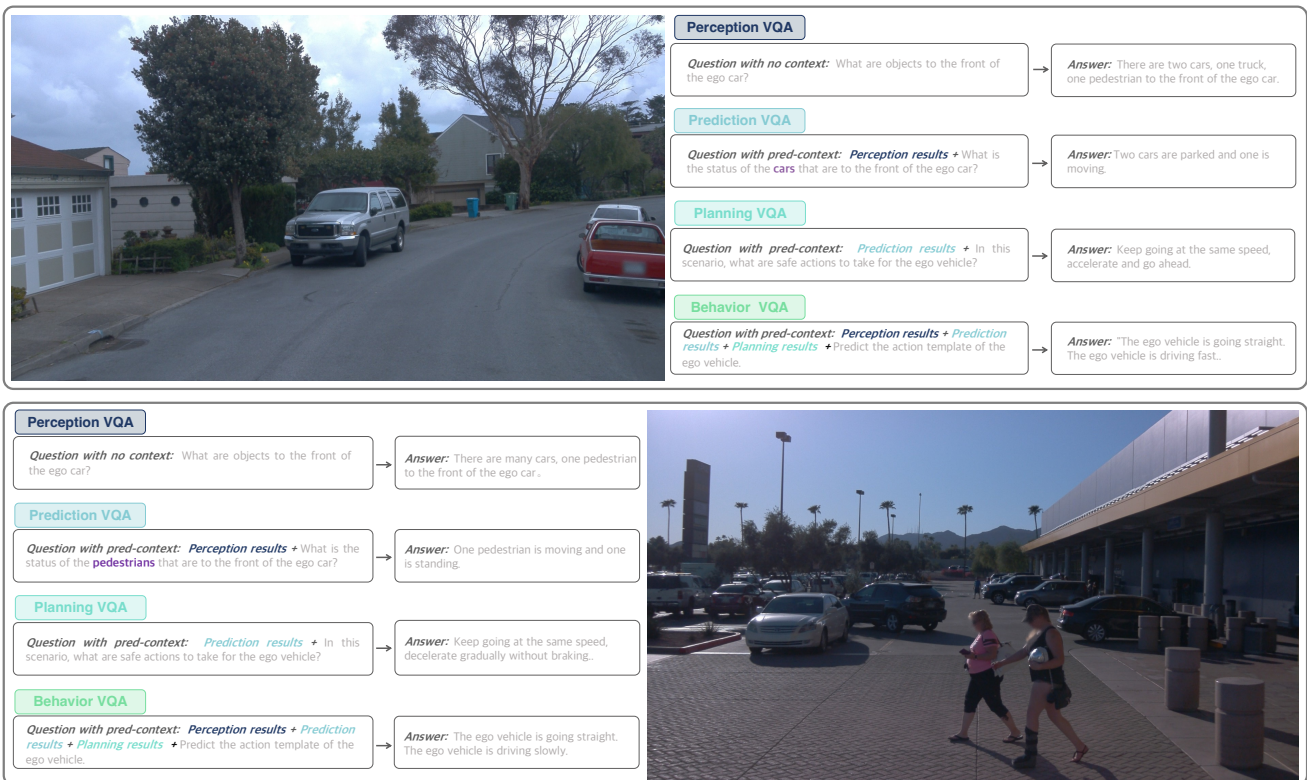


Figure 12. **Qualitative results on Waymo.** We present two examples showcasing the generability of DriveLM-Agent to new sensor configurations, demonstrating the strong generalization capability of DriveLM-Agent.

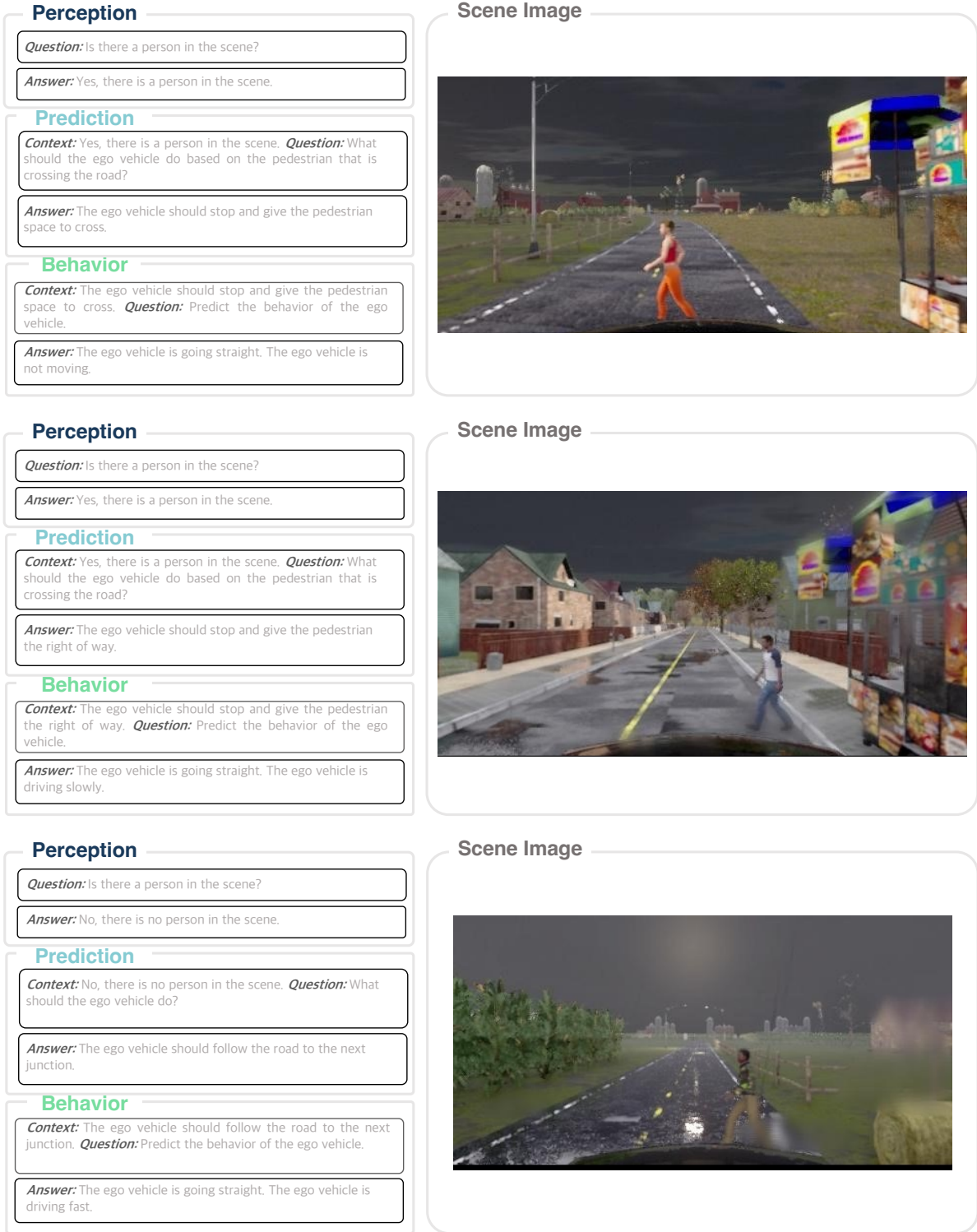


Figure 13. **Qualitative results on the generalization test set in DriveLM-CARLA.** We show three examples of DriveLM-Agent handling the pedestrian scenario. The first example shows a success case and the second and third show two common failure cases of the model.

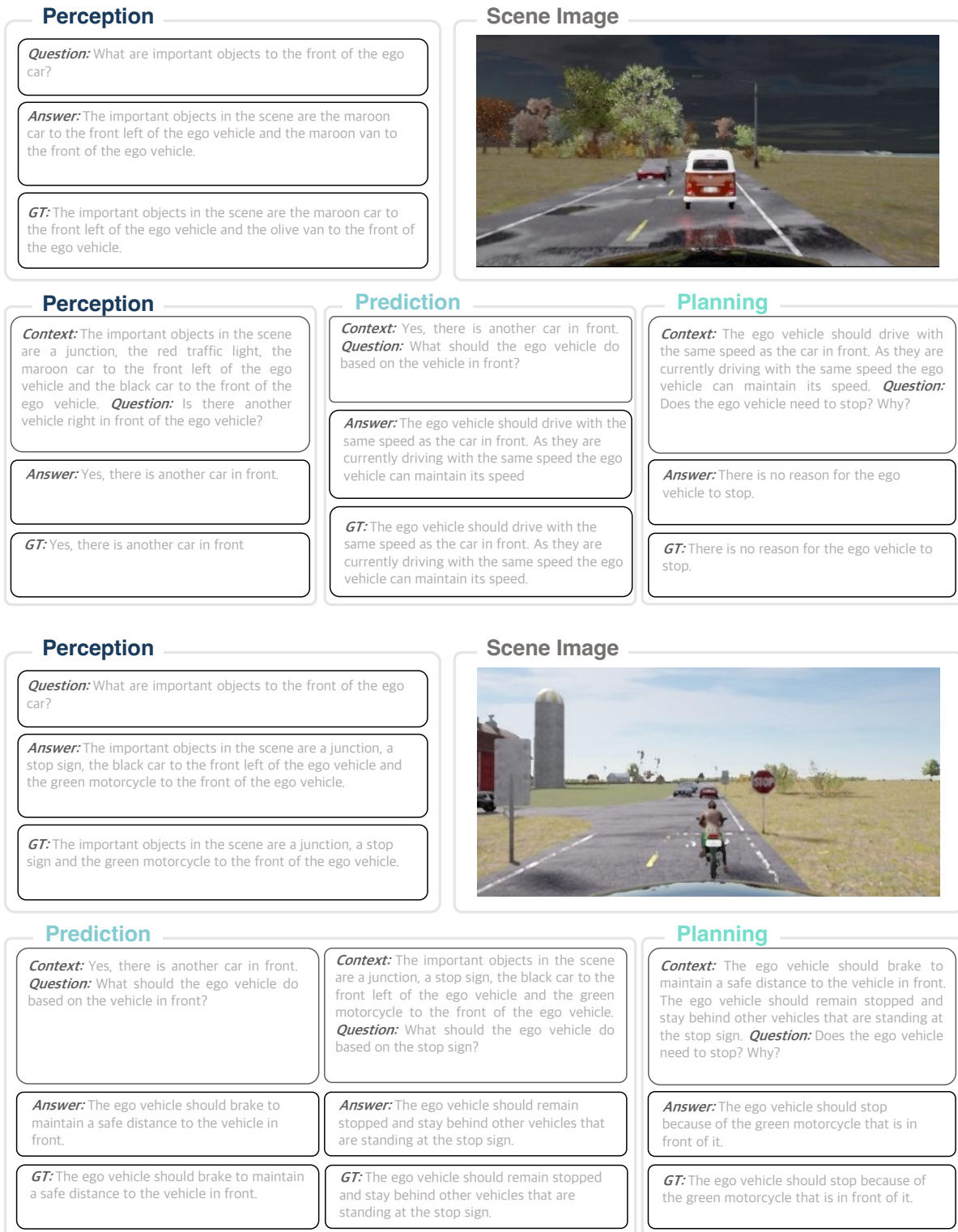


Figure 14. **Qualitative VQA results in DriveLM-CARLA.** The first example shows how DriveLM-Agent deals with a situation with a vehicle directly in front of the ego vehicle. The second example shows an intersection with a stop sign and other traffic participants waiting in front of the ego vehicle at the stop sign.

## H. More Related Work

In this section, we supplement the related work mentioned in the main paper.

**Reasoning Over Graph Structure.** Reasoning is one of the basic forms of simulated human thinking, enabling the derivation of new judgments from one or several existing judgements [9]. Many prior reasoning works have been grounded in graph-based approaches [4, 10, 61, 74]. XNMs [61] employs scene graphs for explainable and explicit reasoning with structured knowledge. KPRN [74] utilizes knowledge graph for reasoning and applies it to recommender systems. GoT [4] models LLM-generated information as an arbitrary graph and brings the LLM reasoning closer to brain mechanisms. Inspired by these successful attempts, we try to link the stages of perception, prediction, and planning in autonomous driving through a graph, enabling the model to grasp the reasoning process and deduce unseen scenarios based on learned graph structure.

**Vision-language Benchmarks for Driving.** An increasing number of vision-language datasets have been proposed for AD systems [15, 34, 35, 47, 55, 59, 78, 79, 84]. NuScenes-QA [55] and NuPrompt [79] provide perceptual information as text by describing the positions and states of surrounding objects. BDD-X [34] provides reasons for the ego vehicle’s actions in natural language descriptions. DRAMA [47] and Rank2Tell [59] identify crucial objects and provide corresponding driving suggestions. However, these datasets focus on scene-level context or individual objects. DriveLM fills this gap in the literature by organizing language annotations from object-level and task-level with a graph structure.