



Are VLMs Ready for Autonomous Driving?

An Empirical Study from the Reliability, Data and Metric Perspectives

– Supplementary Material –

Shaoyuan Xie[†] Lingdong Kong^{‡,◇,*} Yuhao Dong^{‡,§} Chonghao Sima^{‡,▽}
Wenwei Zhang[‡] Qi Alfred Chen[†] Ziwei Liu[§] Liang Pan^{‡,✉}

[†]University of California, Irvine [‡]Shanghai AI Laboratory [◇]National University of Singapore

[§]S-Lab, Nanyang Technological University [▽]The University of Hong Kong

Code & Dataset: <https://drive-bench.github.io>

Table of Contents

A Benchmark Study	1
A.1 DriveLM-nuScenes	1
A.2 BDD-X	2
B Benchmark Setup	2
B.1. Benchmark Construction	2
B.2. Corruption Definitions	2
B.3. Overall Statistics	4
C Additional Implementation Details	4
C.1. VLM Configurations	4
C.2. VLM Prompts	5
C.3. GPT Evaluations Prompts	5
C.4. Human Evaluations	5
C.5. Denoise Model Training	6
D Detailed Experiment Results	6
D.1. GPT Scores	6
D.2. Accuracy Scores	6
D.3. ROUGE-L Scores	6
D.4. Additional RAU Results	6
E Additional Experiments	7
E.1. Visual-based Object Prompts	7
E.2. Temporal Information	8
E.3. Metric Study.	8
E.4. Case Study	8
E.5. Finding Generality	8

(*) Project lead. (✉) Corresponding author.

F. Broader Impact & Limitations	9
F.1. Broader Impact	9
F.2. Potential Limitations	10
G Public Resource Used	10

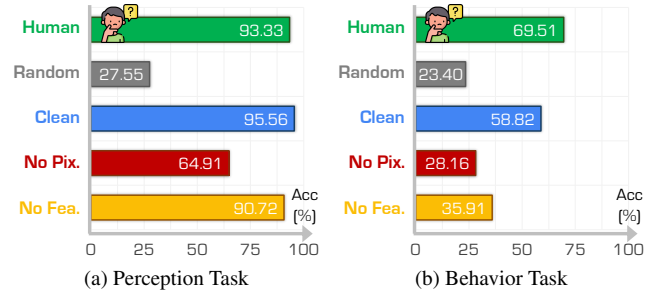


Figure .1. The accuracy scores of perception-MCQ and behavior-MCQ under different visual inputs. The results are from DriveLM-Agent [19]. **No Pix.** and **No Fea.** denote zero image pixels and zero visual features, respectively. We observe that the unbalanced dataset can severely impair the VLM fine-tuning and evaluation process, as **No Fea.** achieves over 90% accuracy without any visual information.

A. Benchmark Study

In this section, we include detailed information on the dataset distribution of the representative driving-with-language dataset. These datasets advance the development of driving with language models.

A.1. DriveLM-nuScenes

We visualize the dataset distribution in perception-MCQs and behavior-MCQs in the DriveLM-nuScenes dataset [19], as shown in Fig. A.1 and Fig. A.2, respectively. In

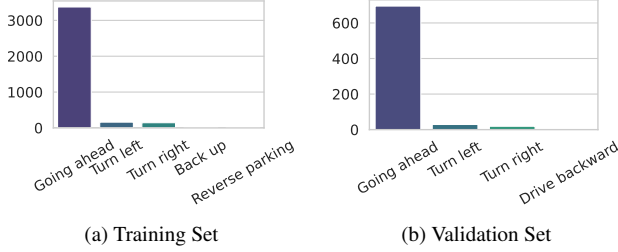


Figure A.1. **The perception-MCQ distributions** in DriveLM-nuScenes [19]. *Going Ahead* dominates both the training and testing data, leading to severe bias for both training and evaluation.

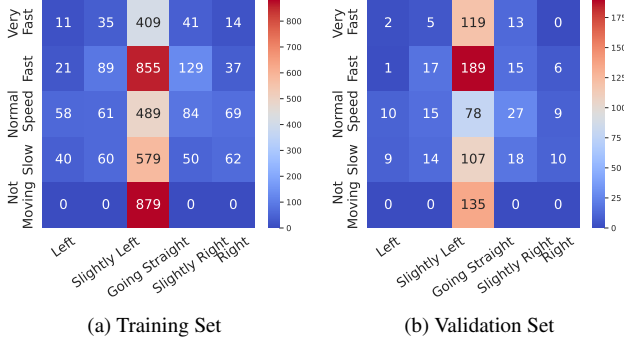


Figure A.2. **The behavior-MCQ distributions of steering and speed** in DriveLM-nuScenes [19]. The majority of actions of vehicle behaviors are “*Going Ahead*”, which has also been noted in [12].

perception-MCQs, “*Going Ahead*” dominates both the training and testing set, while in behavior-MCQs, the majority of choices for vehicle behaviors are “*Going Straight*”, which is also shown in [12].

We find that highly imbalanced data can cause several problems. When fine-tuning VLMs on this dataset, the model tends to memorize the majority of choices and thus answer with them during inference, even if the visual cues are absent, as shown in Fig. 1. The model tends to predict the majority choice even when visual information is completely absent, indicated by *No Pixel* and *No Feature*. Along with the bias of the choice distribution, DriveLM-Agent [19] without any visual input can achieve over 90% accuracy on perception-MCQs.

Furthermore, the dataset mainly adopts language metrics [13, 17] and naive GPT score: prompt with only the answer and ground truth, for evaluation. Our results in the main paper also show the limitations of these metrics in evaluating language-based driving decisions. More analysis on metrics can also be found in Sec. E.3.

A.2. BDD-X

To show the general existence of the limitations of existing benchmarks, we also study the BDD-X [7] dataset as shown

in Fig. A.3. We observe similar limitations to those in the DriveLM-nuScenes [19] dataset, where the data is highly imbalanced. Most actions of the car are “*Stop*” or “*Going Ahead*”, where the random guessing of VLMs can achieve high accuracy since we’ve shown they potentially guessed the answer based on common knowledge and general case in the main paper. The observation shows that the limitations observed in our work are not individual but general drawbacks of existing driving-with-language benchmarks.

B. Benchmark Setup

In this section, we elaborate in detail on the procedures and protocols used to establish the **DriveBench** in this work.

B.1. Benchmark Construction

We detailed the benchmark construction process in this section. Our **DriveBench** is primarily inspired by DriveLM [19] given its impact and representativeness. Given its public availability, we subsample 200 keyframes from the DriveLM-nuScenes [19] training dataset by more balanced sampling and eliminating over-challenging cases even for humans, as discussed in the main paper. These keyframes are selected to balance the ground truth distribution, which can more accurately reflect the model’s performance and prevent bias in most common cases (e.g., “*Going Ahead*”) from dominating the accuracy evaluation. The final perception-MCQs and behavior-MCQs distribution can be seen in Fig. B.1. Our dataset shows a more balanced distribution of the MCQ choices compared to the original DriveLM-nuScenes [19]. Meanwhile, we still prioritize the common choice (e.g., “*Going ahead*” in perception-MCQs and “*Going Straight*” in behavior-MCQs) given they are more common in real-world driving scenarios.

Each keyframe has multiple questions related to different tasks, spanning perception, prediction, planning, and behavior. For each task, we follow the question type design in DriveLM [19], including multiple-choice questions (MCQs) and visual question answering (VQAs), as shown in Tab. A. When evaluating corruption awareness, we add information about corruption context to the question and modify the answer accordingly if necessary. We generate the corruption-related question-answering pairs by prompting GPT-4 based on original QA pairs, which we refer to as the robustness dataset, as shown in Tab. B. The corruption-related question is generated for each corruption type. For Camera Crash and Frame Lost, we design questions asking how many camera sensors fail beyond the corruption identification tasks.

B.2. Corruption Definitions

In this section, we detailed our settings for generating image corruption. **DriveBench** encompasses five distinctive

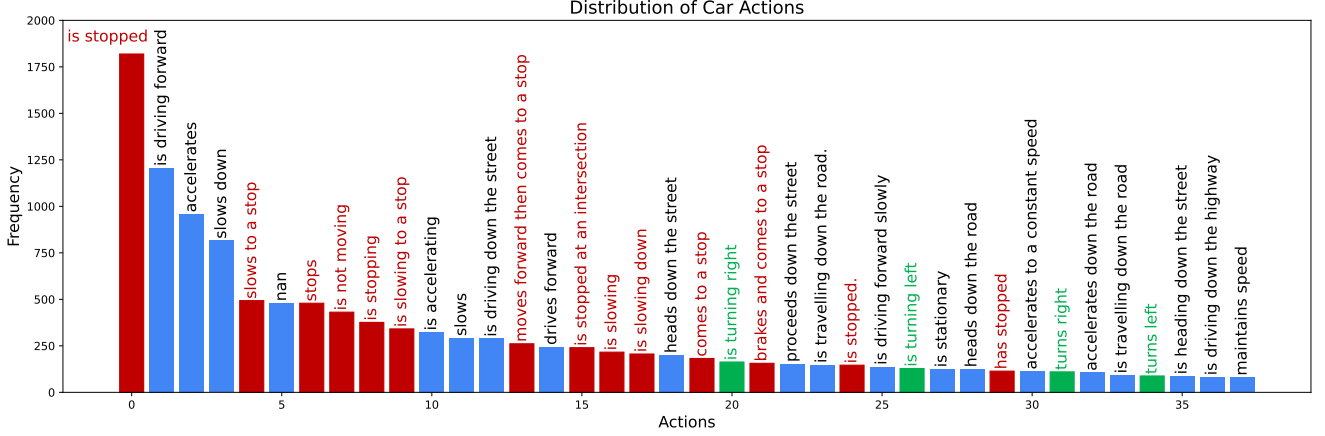


Figure A.3. **BDD-X dataset** [7]: detailed distribution of car actions. Only the actions with a frequency larger than 80 are visualized. The **Stop** actions and **Turn** actions are highlighted. We observe a similar data distribution in balance to that in DriveLM [19], where turning actions only account for a small portion of all actions.

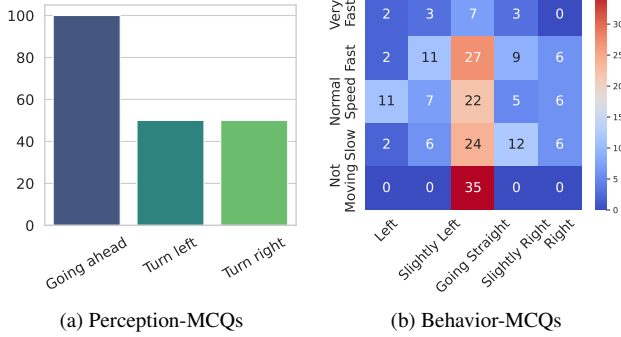


Figure B.1. **The data distributions in DriveBench**. Our dataset shows a more balanced distribution of the MCQ choices compared to the original DriveLM-nuScenes dataset [19].

Table A. **Detailed distribution** of the curated benchmark dataset. The data contains 200 keyframes in total where each keyframe has multiple driving tasks and questions.

#	Driving Task	Question Type	# Samples	Total
	◦ Perception	MCQ & VQA	400	1,261
	◦ Prediction	VQA	61	
	◦ Planning	VQA	600	
	◦ Behavior	MCQ	200	

corruption categories, each with multiple different types of corruptions reflecting the real-world scenarios.

- **Weather & Lighting Conditions (5 Types):**
The simulations of diverse environmental weather and lighting conditions are used in the driving scenarios. In this benchmark, we include the ¹Brightness, ²Dark, ³Snow, ⁴Fog, and ⁵Rain corruptions.
- **External Disturbances (2 Types):**

Table B. **Detailed distribution** of the proposed robustness benchmark dataset. The total number is summed across all the corruption types. “Corrupt. Rec.” represents corruption recognition, which asks the model to identify the current corruption types. “Corrupt. Desc.” represents questions related to the description of the current corrupted environment. The dataset is mainly used for evaluating corruption awareness of VLMs. In the main paper, we only use Corrupt.Rec-MCQs given the page limits.

#	Driving Task	Question Type	# Samples	Total
	◦ Corrupt. Rec.	MCQ	4,000	19,237
	◦ Perception	MCQ & VQA	5,475	
	◦ Prediction	VQA	799	
	◦ Planning	VQA	5,999	
	◦ Corrupt. Desc.	CAP	3,000	

The simulations of situations where camera lenses are occluded by external objects or stains. In this benchmark, we include the ⁶Water Splash and ⁷Lens Obstacle corruptions.

- **Sensor Failures (3 Types):**
The simulations of sensor failures. In this benchmark, we include the ⁸Camera Crash, ⁹Frame Lost, and ¹⁰Saturate corruptions.
- **Motion Blurs (2 Types):**
The simulations of the blurs caused by the ego vehicle’s high-speed motion. In this benchmark, we include the ¹¹Motion Blur and ¹²Zoom Blur corruptions.
- **Data Transmission Errors (3 Types):** The simulations of the errors happening during the video transmission process. In this benchmark, we include the ¹³Bit Error, ¹⁴Color Quantization, and ¹⁵H.265 Compression corruptions.
All the 15 corruption types are generated by applying

high-fidelity image processing algorithms developed in previous works [8, 9, 22, 24]. Here, we detail how each corruptions are synthesized as follows:



B.3. Overall Statistics

C. Additional Implementation Details

C.1. VLM Configurations

model that integrates vision and language understanding. With 13 billion parameters, it achieves state-of-the-art performance across multiple benchmarks, rivaling models like GPT-4.

- **LLaVA-NeXT** [15] is an evolution of the LLaVA series, enhancing multimodal capabilities by supporting multi-image, video, and 3D tasks within a unified large language model. It achieves state-of-the-art performance on a wide range of benchmarks, demonstrating strong video understanding through task transfer from images.
- **InterVL** [5] is an open-source multimodal dialogue model developed by OpenGVLab. It closely approximates the performance of proprietary models like GPT-4o, excelling in tasks that integrate visual and linguistic information, such as visual question answering and image captioning.
- **Oryx** [6] is a unified multimodal architecture created by researchers from Tsinghua University and Tencent. It is designed for spatial-temporal understanding of images, videos, and multi-view 3D scenes, offering on-demand processing of visual inputs with arbitrary spatial sizes and temporal lengths.
- **Qwen2-VL** [3, 20] is a large language model developed by Alibaba Cloud, available in both chat and pre-trained versions. It delivers high-quality language generation and understanding capabilities, optimized for tasks requiring nuanced comprehension and generation of human language.
- **DriveLM-Agent** [19] is a model from OpenDriveLab tailored for autonomous driving applications, focusing on graph-based visual question answering. It addresses challenges in driving scenarios by integrating language understanding with visual perception, enhancing decision-making processes in autonomous systems.
- **Dolphins** [16] is a multimodal language model developed by NVIDIA for driving applications. It adeptly processes inputs such as video data, text instructions, and historical control signals to generate informed outputs, facilitating a comprehensive understanding of complex driving scenarios.

C.2. VLM Prompts

We use the default system prompt for all the candidate VLMs if not specified, as shown in Fig. C.1. We prompt the VLMs to explain for MCQs to facilitate the GPT evaluation based on their explanations, as discussed in the main paper.

C.3. GPT Evaluations Prompts

We include the detailed prompt we used for GPT_{ext} evaluation. We use the prompt to evaluate perception-MCQs as shown in Fig. G.1. Given the ground truth is a single choice (e.g., answer: “A”), the DESC is used to prompt context in-

You are a smart autonomous driving assistant responsible for analyzing and responding to driving scenarios. You are provided with up to six camera images in the sequence [CAM FRONT, CAM FRONT LEFT, CAM FRONT RIGHT, CAM BACK, CAM BACK LEFT, CAM BACK RIGHT]. Each image has normalized coordinates from [0, 1], with (0, 0) at

the top left and (1, 1) at the bottom right.

Instructions:

1. Answer Requirements:

- For multiple-choice questions, provide the selected answer choice along with an explanation.

- For “is” or “is not” questions, respond with a “Yes” or “No”, along with an explanation.
- For open-ended perception and prediction questions, related objects to the camera.

2. Key Information for Driving Context:

- When answering, focus on object attributes (e.g., categories, statuses, visual descriptions) and motions (e.g., speed, action, acceleration) relevant to driving decision-making. Use the images and coordinate information to respond accurately to questions related to perception, prediction, planning, or behavior, based on the question requirements.

Figure C.1. VLM inference system prompt.

formation for a more accurate evaluation of the VLM explanation beyond the choice (e.g., the specified object is a black sedan). Limited to the current drive-with-language dataset, we extract the natural language description of critical objects in the current environment to provide more context information, e.g., if the identified object in the explanation is wrong, the score will be lower given the rubrics.

The prompt for the perception-VQA is shown in Fig. G.2. Since the ground truth for perception-VQA already included the visual description and moving status of important objects, we only prompt with PRED and GT with detailed rubrics.

C.4. Human Evaluations

In this section, we elaborate in more detail on how we conduct the human evaluation experiments in our benchmark.

Procedures. Considering the large number of questions in the dataset, we subsample 15 out of 200 keyframes from our curated dataset. To ensure no overlaps between different corruptions (i.e. different corruptions apply to the same images), which might cause information leakage, we lower the probability if the same keyframes are sampled before. Then, we design a user interface for human evaluation, focusing on perception-MCQs and behavior-MCQs. The interface is shown in Fig. G.5. The same as evaluating VLMs, we only prompt single-view images if the questions are re-

lated to only one of the cameras.

Ethic Declaration. According to the Federal Policy of Human Subjective Research¹, our research involves conducting anonymous visual recognition tasks, where participants respond to questions about visual stimuli without any interventions or demographic data collection. It qualifies for Exempt Research 2(i)² because it solely involves survey-like procedures with no physical or psychological risks to participants. Specifically, it meets the requirements of Exempt Research 2(i) as the data is recorded in a manner ensuring that participants’ identities cannot be readily ascertained, directly or through linked identifiers. It does not fall under Exempt Research 2(ii) or 2(iii) because no identifiable information is recorded, and no IRB review is required to ensure these protections. We also submit IRB review records to the corresponding institutions and receive the official confirmation of IRB review exemption.

C.5. Denoise Model Training

Experiment Setups. In this work, we choose AirNet [10] as a tool for denoising without losing the generality of the tool choice. For model training, we sample 100 keyframes in the nuScenes [18] training set, each keyframe composed of 6 camera views. Therefore, we have 600 images with resolution 1600×900 . We use the tools developed by previous works [22] to generate the corresponding corruption pairs. We train a denoise model for each of the corruptions in RoboBEV [22] benchmark, including `Bright`, `Dark`, `Fog`, `Snow`, `Color Quant`, and `Motion Blur`. Each model is trained for 100 epochs with a learning rate $1e-3$ and input patch size 224. Since the RoboBEV [22] benchmark is developed on the nuScenes [18] validation set, the training setup ensures no information leakage is encountered for the denoise model.

Qualitative Results. We show some qualitative results of the denoising model on RoboBEV [22] benchmark in Fig. C.2. The denoise model can mostly recover the original image given the visual corruptions. The only case where the denoise model is not ideal is `Color Quant` since the pixel intensity quantization process is largely irreversible. Thus, we can use VLMs as agents to apply the existing denoising tools for downstream tasks (*e.g.*, camera-based detection [11, 21]).

D. Detailed Experiment Results

In this section, we include the detailed benchmark results evaluated by multiple metrics. We also include the detailed perception-MCQs spatial distribution in Fig. G.4 for each model.

¹<https://www.federalregister.gov/d/2017-01058/p-1315>

²<https://www.federalregister.gov/d/2017-01058/p-1375>

D.1. GPT Scores

We include the detailed GPT_{cxt} scores in Tab. H, Tab. J, Tab. L, Tab. M, and Tab. O for different tasks. The observation and conclusion in the main paper are primarily derived from GPT_{cxt} scores. Therefore, we focus more on the discussion of accuracy and language scores in the following sections.

D.2. Accuracy Scores

We include accuracy scores for MCQs in addition to GPT_{cxt} scores in Tab. I and Tab. K. Compared with GPT scores, we find that the accuracy score metric is more homogeneous. For example, the LLaVA-1.5 models have 50% accuracy under all the input types, suggesting they are merely output “*Going Ahead*” for perception-MCQs (recall the distribution in Fig. B.1), which is also observed in the prediction spatial distribution in Fig. G.4e and G.4f. Moreover, we find that most models have no accuracy degradation under corruptions or even text-only inputs. This raises concerns about whether VLMs are indeed leveraging visual information to make decisions about the specified spatial location or naively guessing based on their general knowledge.

D.3. ROUGE-L Scores

We also present the detailed language scores, *i.e.*, ROUGE-L [13] here for VQA in Tab. N and Tab. P. As discussed in the main paper, we find the fine-tuning process can significantly benefit the ROUGE-L score, as indicated by the fact that DriveLM-Agent [19] outperforms other models with a large margin. On the contrary, GPT-4o [2], which generates more detailed answers, is punished by the answer length. The ROUGE-L score of GPT-4o [2] is lower than most of the models, even though the GPT_{cxt} scores are much higher.

D.4. Additional RAU Results

In this section, we provide detailed results of the proposed RAU framework. As we discussed in the main paper, RAU is orthogonal to the development of VLMs and denoising models. Therefore, the main bottlenecks of RAU are the corruption identification accuracy of VLMs and the restoration fidelity of denoising models.

RAU Accuracy. We analyze the accuracy of the agentic VLMs in identifying corruptions in Fig. D.1. We observe that the off-the-shelf VLMs (*i.e.*, InternVL2 here) are good enough to identify weather corruptions. However, the model struggles to distinguish semantic corruptions (*e.g.*, `Motion Blur`). We leave the exploration of how to improve the corruption identification accuracy with minimal cost as future work.

BEV Detector Performance. We provide the complete BEV detector performance in Tab. C and Tab. D. We find that when equipped with RAU, the model can be much

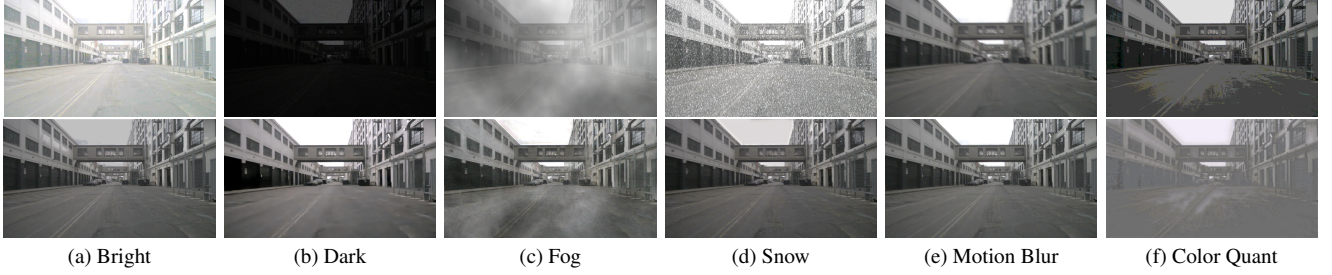


Figure C.2. **Denoise qualitative results.** The first row shows the image from the RoboBEV [22] dataset, and the second row shows the image after applying the denoise model trained for each corruption.

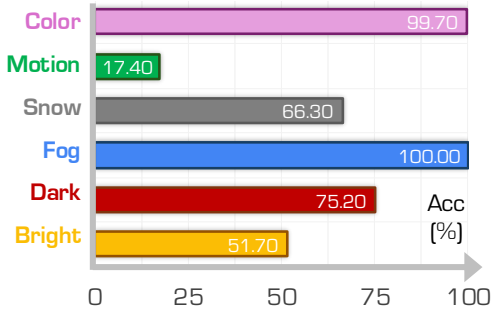


Figure D.1. **Corruption identification accuracy** of RAU on nuScenes validation set. We observe that the off-the-shelf VLMs (*i.e.*, InternVL2 here) are good enough to identify weather corruptions.

Table C. **Detail Performance** of BEVFormer [11] under corruption. We highlight the improved performance when equipped with RAU. For Color Quant, we find the performance degradation might be due to imperfect denoising (shown in Fig. C.2).

Method	Corruption	NDS	mAP	mATE	mASE	mAOE	mAVE	mAAE
BEVFormer	Clean	51.7	41.6	67.3	27.4	37.2	39.4	19.8
BEVFormer	Bright	41.7	33.1	74.3	28.6	46.8	78.1	20.4
BEVFormer	Dark	25.3	14.1	87.3	32.7	67.4	101.1	29.6
BEVFormer	Fog	40.6	31.6	76.1	28.4	47.1	80.0	20.7
BEVFormer	Snow	16.9	6.2	93.3	42.5	82.5	109.0	43.6
BEVFormer	Motion	22.9	10.3	90.0	32.8	75.4	97.5	27.2
BEVFormer	Color	36.4	25.5	80.1	28.8	50.1	81.1	23.3
BEVFormer _{RAU}	Bright	42.8	34.5	74.4	28.3	43.9	77.2	21.2
BEVFormer _{RAU}	Dark	29.2	18.2	83.7	30.5	65.4	93.6	26.5
BEVFormer _{RAU}	Fog	41.3	32.5	76.4	28.4	45.9	77.4	21.5
BEVFormer _{RAU}	Snow	41.3	33.0	75.2	28.3	47.8	78.9	21.8
BEVFormer _{RAU}	Motion	28.4	14.8	83.8	30.7	62.4	89.3	23.5
BEVFormer _{RAU}	Color	29.7	17.4	87.1	30.0	54.8	91.1	27.1

more robust towards image corruption, especially under severe visual degradation (*e.g.*, Snow). For DETR3D [21], the performance under Bright and Fog stays similar with and without RAU. The reason might be DETR3D [21] is more robust compared to BEVFormer [11], the mRR is 0.71 vs. 0.59 without RAU (shown in Tab. 6 in the main paper). Thus, the improvement when equipped with RAU is not as significant as BEVFormer [11]. Additionally, we find the performance under Color Quant is even lower. The reason is due to the imperfection of the denoise model and the

Table D. **Detail Performance** of DETR3D [21] under corruption. We highlight the improved performance when equipped with RAU. For Color Quant, we find the performance degradation might be due to imperfect denoising (shown in Fig. C.2).

Method	Corruption	NDS	mAP	mATE	mASE	mAOE	mAVE	mAAE
DETR3D	Clean	43.4	34.9	71.6	26.8	38.0	84.2	20.0
DETR3D	Bright	41.8	32.3	73.3	27.2	40.6	82.6	19.8
DETR3D	Dark	28.9	15.7	82.6	29.5	51.5	108.7	25.8
DETR3D	Fog	40.5	30.5	75.1	27.0	41.6	82.8	20.4
DETR3D	Snow	17.2	5.0	92.3	40.9	80.5	128.9	39.6
DETR3D	Motion	22.2	9.0	92.8	32.0	71.9	110.4	26.3
DETR3D	Color	34.0	23.1	85.0	28.0	47.3	91.9	23.6
DETR3D _{RAU}	Bright	41.1	31.7	77.7	27.1	39.5	83.1	20.2
DETR3D _{RAU}	Dark	29.7	17.7	87.7	28.7	50.7	101.4	24.1
DETR3D _{RAU}	Fog	39.0	28.8	82.1	27.2	41.7	82.3	20.8
DETR3D _{RAU}	Snow	39.4	29.5	81.7	27.1	40.6	82.4	21.3
DETR3D _{RAU}	Motion	26.7	12.8	87.6	30.0	55.9	100.0	23.9
DETR3D _{RAU}	Color	28.8	15.9	89.6	29.2	52.9	97.7	22.1

corruption, as Color Quant is a highly irreversible process, and the denoise model [10] struggles to restore the images, as shown in Fig. C.2.

E. Additional Experiments

E.1. Visual-based Object Prompts

In this section, we consider another evaluation setup: we specify the target object by visual-based prompts, where we visualize the bounding box around the target object instead of using the numerical coordinates in the image.

Experimental Setups. We conduct the experiments using perception-MCQs. We make a slight change to the text prompt (in Fig. C.1) for this evaluation. Specifically, we replace the “What is the moving status of the object at (480,520)?” with “What is the moving status of the object inside red bounding box?” and visualize the corresponding object with a red bounding box in the image. We choose Qwen2VL_{7B} for the experiments.

Results. The results can be seen in Tab. F. We find that the performance remains largely unchanged across corruptions. Given that *Going Ahead* accounts for half the answer distribution (recall the distribution in Fig. B.1), the model constantly predicts *Going Ahead* when both the visual and text information is completely absent (*i.e.*, text-only input

Table E. **Ablation study of temporal input.** We evaluate multiple frame inputs on perception-MCQs. The results suggest that temporal information can benefit the model prediction under clean input. However, the model’s bias towards general knowledge is still observed under text-only inputs, the same as those observed in single-frame input.

Method	Temporal	• Clean	• T.O.
Qwen2VL _{7B}	$t = 2$	61.0	51.5
Qwen2VL _{7B}	-	59.0	56.5

of visual-prompt evaluation). The results prove our findings in the main paper that VLMs tend to leverage general knowledge to predict the answer when visual information is degraded. We also observe the accuracy difference between visual-based object prompts and text-based object prompts, the accuracy is 50.0% vs. 56.5%. The results further prove our finding in the main paper that VLMs tend to leverage text cues to guess the answer since neither is provided with visual information.

E.2. Temporal Information

In this section, we study how temporal frames influence the performance of VLMs and their reliability.

Experimental Setups. We augment the input frames by providing the previous two key frames defined in the nuScenes [4] dataset. We only evaluate the clean input and text-only input, given that it is costly to generate 15 types of corruption on all the history frames. For text-only, we use the same temporal length as clean input, with all the frames being black. We use text-prompt as it is the default evaluation protocol in this paper. We conduct the experiments using perception-MCQs and choose Qwen2VL_{7B} for the experiments.

Results. The results can be seen in Tab. E. We found that temporal information can potentially benefit the perception, leading to higher accuracy in predicting the moving status of other vehicles. However, we find that providing temporal information doesn’t mitigate the hallucination issues of VLMs, as the text-only accuracy still achieves 51.5%. Therefore, our conclusions are still held given the temporal information. We leave the exploration of temporal-aware evaluation under OoD corruptions as future work.

E.3. Metric Study.

In this section, we provide a detailed analysis of different metrics and study how metrics are related given the same question-answer pairs. We also investigate how the prompts influence the GPT evaluation.

Language Metrics. We visualize ROUGE-L vs. GPT_{cxt} score given the same question-answer pairs in Fig. E.2a. We observe that ROUGE-L remains around 0.2 while GPT scores vary across a large range. The results reveal how

the language score fails to reflect underlying key information. Interestingly, DriveLM-Agent, which is fine-tuned on the DriveLM-nuScenes dataset [19], achieves the highest ROUGE-L score (*i.e.*, around 0.55). However, the GPT evaluation towards the same set of answers doesn’t indicate such a large advantage. The observation indicates that the main improvement of in-distribution fine-tuning on the current small-scale driving dataset largely comes from the answering template.

Accuracy. In terms of accuracy, we also study how accuracy and GPT_{cxt} scores are related. The results are presented in Fig E.2b. The GPT evaluation highly aligns with accuracy since we prompt the GPT to assign certain scores if the answer is correct. The deviation is because GPT assigns another portion of scores to the coherence of explanation dimensions, capturing nuanced differences between answers.

GPT Evaluation. A critical question remains: is GPT evaluation currently the optimal approach? The answer is nuanced. GPT-based scoring can effectively capture human preferences and emphasize critical elements in driving scenarios, yet this capability is highly contingent on the provision of comprehensive driving contextual information. We empirically compare how the same response is scored given different information, shown in Fig. E.1. When GPT evaluation is prompted solely with GT and model response, the resulting scores are highly homogeneous, while the inclusion of specific rubrics, questions, and specific driving context yields greater score diversity. We provide an example in Figs. G.16 and G.18. The results suggest GPT scores are sensitive to prompts, and it is critical to provide enough environmental information for the LLM Judge to evaluate the answer safety tied to specific scenarios.

E.4. Case Study

We provide a detailed case study at the end of the Appendix as listed below:

- **Failure Case:** please refer to Fig. G.3.
- **Driving Tasks:** please refer to Fig. G.7 and Fig. G.8.
- **Comparison of VLMs:** please refer to Fig. G.9, Figs. G.6 and G.10.
- **Comparison w/ and w/o Corruption:** please refer to Fig. G.11, Fig. G.12, Fig. G.13, Figs. G.14 and G.15.
- **GPT Evaluate Prompt:** please refer to Fig. G.17.

E.5. Finding Generality

We also conduct experiments on other datasets beyond DriveLM [19] to study the generality of our findings. Specifically, we use the nuScenes-OIA dataset [23]. We crafted one corruption from each corruption category, including Snow, Color Quantization, Motion Blur, Saturate, and Lens Obstacle. Following their setups, we adopt the F-1 score to compute the VLMs’

Table F. **Ablation study of different prompt modality** (i.e., visual-prompt vs. text-prompt) of Qwen2VL_{7B} on perception-MCQs. *Visual-prompt* means visualizing a bounding box around the target object to ask questions. *Text-prompt* means to provide the numerical coordinate of the target object to ask a question. We find VLMs tend to predict general knowledge (e.g., *Going Ahead*) when all visual and text-cues are excluded. The comparison between prompts further proves that VLMs tend to leverage text cues to guess the answer.

Method	Clean	T.O.	Brightness	Dark	Snow	Fog	Rain	Lens Obstacle	Water Splash	Camera Crash	Frame Lost	Saturate	Motion Blur	Zoom Blur	Bit Error	Color Quant	H.265 Compression
Visual-prompt	56.0	50.0	57.0	52.0	49.5	52.0	55.5	54.5	55.5	54.5	48.5	55.5	57.0	48.5	49.0	48.5	55.5
Text-prompt	59.0	56.5	60.0	59.0	60.0	59.5	59.0	59.0	59.0	58.0	56.0	59.0	59.5	55.0	54.5	57.0	59.0

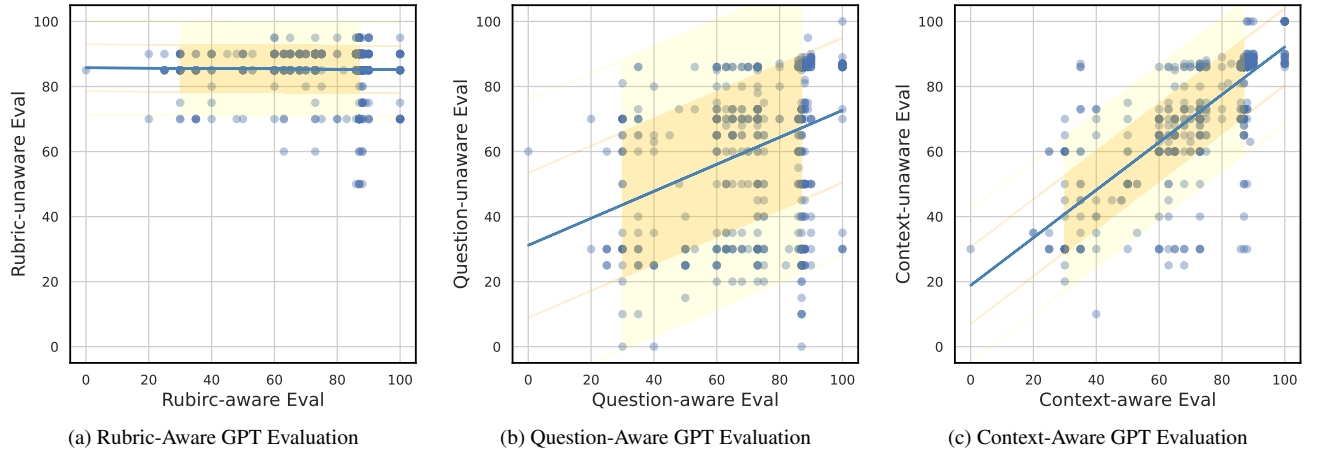


Figure E.1. **Comparisons among different evaluation types (rubric, question-aware, and context-aware).** The GPT scores vary depending on the rubric, question, and physical driving context. With more information added, the results become more distinguishable.

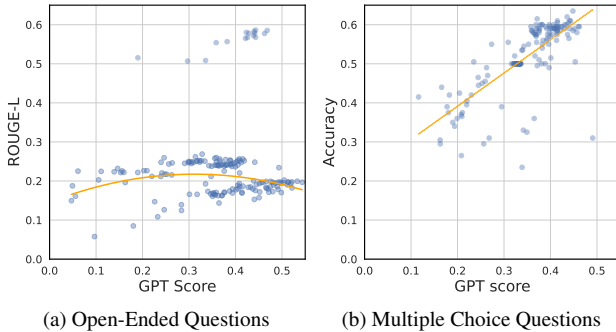


Figure E.2. **Correlations when using different metrics.** We study how well accuracy or ROUGE-L [13] matches the GPT scores for VQA and MCQs, respectively. We find that ROUGE-L [13] fails to reflect semantic information (e.g., key object) that is critical in driving. Contrarily, accuracy aligns well with the GPT score for MCQ, while the GPT score can further capture nuanced differences in explanation when the answer is correct.

performance on action prediction. The results are shown in

Table G. **BDD-OIA [23] Evaluation.** We report the F-1 score of action prediction following the dataset setups. We find similar results on the BDD-OIA dataset, where the model performance under corruption remains close to that under clean inputs.

Model	Clean	T.O.	Snow	Color	Motion	Saturate	Lens
LLaVA1.5 _{7B}	27.00	17.64	27.92	28.45	26.31	25.84	27.38
LLaVA1.6 _{7B}	15.57	15.44	15.43	15.43	15.79	15.64	15.49
Qwen2-VL _{7B}	20.67	15.44	16.16	20.81	28.18	20.99	19.45
Phi3.5	29.09	15.44	15.66	17.36	27.71	20.69	25.32

Tab. G. We find similar results on the BDD-OIA dataset, where the model performance under corruption remains close to that under clean inputs.

F. Broader Impact & Limitations

In this section, we discuss the broader implications of our study and acknowledge its potential limitations.

F.1. Broader Impact

Our research focuses on evaluating the reliability of VLMs in autonomous driving, emphasizing three critical perspec-

tives: the model’s robustness, data quality, and evaluation metrics. The findings reveal a concerning tendency of VLMs to fabricate explanations, particularly under conditions of visual degradation. This issue is not limited to autonomous driving but is likely relevant to other VLM-embodied systems, such as robotics and other safety-critical cyber-physical systems. For example, VLM-based robots could generate misleading task explanations or actions based on hallucinated understanding, potentially compromising safety and operational reliability.

The implications of our work extend beyond autonomous driving, calling for a reassessment of benchmark and metric designs to better evaluate the trustworthiness of VLMs in real-world applications. Current benchmarks often fail to account for the complexity and variability of real-world scenarios, particularly in environments where system malfunctions could result in life-threatening consequences. Our study highlights the urgency of addressing these gaps to develop robust, reliable, and interpretable VLMs that can be safely integrated into such systems.

Finally, the design of benchmarks, testbeds, and evaluation metrics that accurately capture the reliability and safety implications of applying VLMs to real-world physical systems is of paramount importance. These tools must go beyond traditional performance metrics to consider the nuanced requirements of autonomous systems, such as contextual understanding, interpretability, and robustness against adversarial conditions.

F.2. Potential Limitations

While this study provides valuable insights, it is essential to recognize its limitations to contextualize the findings:

- The experimental results are derived exclusively from the DriveLM [19] dataset due to the prohibitive computational cost of large-scale VLM inference and GPT-based evaluations. While DriveLM is a comprehensive dataset, its scope may limit the generalizability of our findings to other driving benchmarks or real-world settings. Future work should expand the analysis to additional datasets and environments to validate the observed trends.
- The lack of detailed contextual data in the DriveLM dataset poses a constraint on our evaluations. For instance, the GPT-based assessments rely on limited visual descriptions of key objects, which may not comprehensively capture the broader situational context required for accurate and nuanced evaluations. Expanding datasets to include richer temporal and spatial contexts could improve evaluation fidelity.
- This study primarily investigates the language-based explanations generated by VLMs. While these insights are crucial for understanding VLM reliability, it remains unclear whether the observations generalize to action models that generate vehicle trajectories or other non-

language outputs. Exploring how VLMs’ visual grounding affects action-oriented tasks, such as trajectory prediction or manipulation control, represents an important direction for future research.

- The study evaluates a finite set of 12 VLMs across specific tasks, metrics, and settings. Although the insights are significant, the scalability of these findings to emerging VLM architectures or more diverse driving scenarios warrants further investigation.

By addressing these limitations in future studies, we aim to build a more comprehensive understanding of the challenges and opportunities in applying VLMs to autonomous driving and other safety-critical domains.

G. Public Resource Used

In this section, we acknowledge the use of the following public resources during the course of this work:

- nuScenes³ CC BY-NC-SA 4.0
- nuScenes-devkit⁴ Apache License 2.0
- RoboBEV⁵ Apache License 2.0
- AirNet⁶ Apache License 2.0
- DriveLM⁷ Apache License 2.0
- Phi-3.5-vision-instruct⁸ MIT License
- Phi-3-mini-4k-instruct⁹ MIT License
- LLaVA-1.5-7B-hf¹⁰ Apache License 2.0
- LLaVA-1.5-13B-hf¹¹ Apache License 2.0
- LLaVA-v1.6-mistral-7B¹² Apache License 2.0
- InternVL2-8B¹³ Apache License 2.0
- Qwen2-VL-7B¹⁴ Apache License 2.0
- Qwen2-VL-72B¹⁵ Apache License 2.0

³<https://www.nuscenes.org/nuscenes>

⁴<https://github.com/nutonomy/nuscenes-devkit>

⁵<https://github.com/Daniel-xsy/RoboBEV>

⁶<https://github.com/XLearning-SCU/2022-CVPR-AirNet>

⁷<https://github.com/OpenDriveLab/DriveLM>

⁸<https://huggingface.co/microsoft/Phi-3.5-vision-instruct>

⁹<https://huggingface.co/microsoft/Phi-3-mini-4k-instruct>

¹⁰<https://huggingface.co/llava-hf/llava-1.5-7b-hf>

¹¹<https://huggingface.co/llava-hf/llava-1.5-13b-hf>

¹²<https://huggingface.co/liuhaotian/llava-v1.6-mistral-7b>

¹³<https://huggingface.co/OpenGVLab/InternVL2-8B>

¹⁴<https://huggingface.co/Qwen/Qwen2-VL-7B-Instruct>

¹⁵<https://huggingface.co/Qwen/Qwen2-VL-72B-Instruct>

Please evaluate the multiple-choice answer on a scale from 0 to 100, where a higher score reflects precise alignment with the correct answer and well-supported reasoning. Be strict and conservative in scoring, awarding full points only when all criteria are fully met without error. Deduct points for minor inaccuracies, omissions, or lack of clarity. Distribute the **Total Score** across the following criteria:

1. Answer Correctness (50 points):

- Exact Match (50 points): Assign 50 points if the predicted answer exactly matches the correct answer.
- No Match (0 points): Assign 0 points if the predicted answer does not match the correct answer, regardless of explanation quality.

2. Object Recognition (10 points):

- Award up to 5 points for accurately identifying all relevant object(s) in the scene.
- Award up to 5 points for correct descriptions of the identified object(s), including attributes like colors, materials, sizes, or shapes.
- Guideline: Deduct points for any missing, misidentified, or irrelevant objects, particularly if they are crucial to the driving context. Deduct points if any important visual details are missing, incorrect, or overly generalized, especially if they affect comprehension or recognition.

3. Object Location and Orientation (15 points):

- Score up to 5 points for a precise description of the object's location, orientation, or position relative to the ego vehicle.
- Award up to 5 points for acknowledging environmental factors, such as lighting, visibility, and other conditions that influence perception.
- Score up to 5 points based on how well the answer reflects an understanding of situational context, such as obstacles, traffic flow, or potential hazards.
- Guideline: Deduct points for inaccuracies or omissions in spatial information that could affect scene understanding. Deduct points if the answer fails to consider factors impacting object visibility or situational awareness. Deduct points for overlooked or misinterpreted contextual factors that may impact driving decisions.

4. Environmental Condition Awareness (15 points):

- Award up to 15 points if the explanation considers environmental conditions (e.g., weather or sensor limitations) that could impact perception.
- Guideline: Deduct points if relevant environmental conditions are ignored or inadequately addressed.

5. Clarity of Reasoning (10 points):

- Award up to 5 points for clear, logically structured reasoning that is easy to understand.
- Assign up to 5 points for grammatical accuracy and coherent structure.
- Guideline: Deduct points for vague or confusing explanations that hinder comprehension. Deduct points for grammar or syntax issues that impact clarity or logical flow.

Assign 0 points from criteria 2 to 5 if no explanation is provided.

Here is the multiple-choice question: **QUESTION**

Here is the ground truth object visual description: **DESC**

Here is the correct answer: **GT**

Here is the predicted answer and explanation (if any): **PRED**

Please fill in the following scoring sheet, and then provide a brief summary supporting the score:

1. Answer Correctness (50 points):
2. Object Recognition (10 points):
3. Object Location and Orientation (15 points):
4. Environmental Condition Awareness (15 points):
5. Clarity of Reasoning (10 points):

Total Score:

Brief Summary:

Figure G.1. GPT evaluation prompts for **Multiple-Choice Questions** (MCQs) in our benchmark.

Please evaluate the predicted answer on a scale from 0 to 100, where a higher score reflects precise alignment with the correct answer and well-supported reasoning. Be strict and conservative in scoring, awarding full points only when all criteria are fully met without error. Deduct points for minor inaccuracies, omissions, or lack of clarity. Distribute the **Total Score** across the following criteria:

1. Action Alignment (20 points):

- Assign up to 20 points based on how accurately the predicted action (e.g., forward, turn left, turn right) matches the correct answer.
- Guideline: Award full points only for exact matches or highly similar actions. Deduct points for any inaccuracies or missing elements. Assign 0 points if no action prediction is provided.

2. Motion Precision (20 points):

- Award up to 20 points based on how closely the predicted motion (e.g., speed up, decelerate) aligns with the correct motion in the answer.
- Guideline: Deduct points if the predicted motion fails to match the type or intensity of the correct answer. Ensure that the intended speed or deceleration aligns accurately with the driving context. Assign 0 points if no motion prediction is provided.

3. Driving Context Appropriateness (15 points):

- Score up to 15 points for the relevance of the predicted answer to the driving context implied by the correct answer, emphasizing logical alignment with the situation.
- Guideline: Award higher scores only if the answer fully reflects an accurate understanding of the driving context. Deduct points if the action or motion is illogical or does not align with the scenario's requirements.

4. Situational Awareness (15 points):

- Award up to 15 points for demonstrated awareness of environmental factors (e.g., traffic participants, obstacles) relevant to the action or motion.
- Guideline: Deduct points if the answer misses key situational details that may lead to unsafe or incorrect predictions.

5. Conciseness and Clarity (20 points):

- Assess the clarity and brevity of the predicted answer. Answers should be concise, clear, and easy to understand, effectively communicating the intended actions and motions.
- Guideline: Deduct points for verbosity, ambiguity, or lack of focus that could hinder quick comprehension.

6. Grammar (10 points):

- Evaluate the grammatical accuracy and structure of the answer. Assign up to 5 points for clarity and logical flow, and up to 5 points for grammatical accuracy.
- Guideline: Deduct points for grammar or syntax issues that reduce readability or coherence.

Here is the predicted answer: **PRED**

Here is the correct answer: **GT**

Please fill in the following scoring sheet, and then provide a brief summary supporting the score:

1. Action Alignment (20 points):
2. Motion Precision (20 points):
3. Driving Context Appropriateness (15 points):
4. Situational Awareness (15 points):
5. Conciseness and Clarity (20 points):
6. Grammar (10 points):

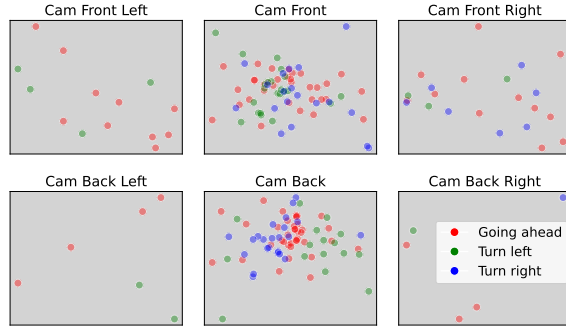
Total Score:

Brief Summary:

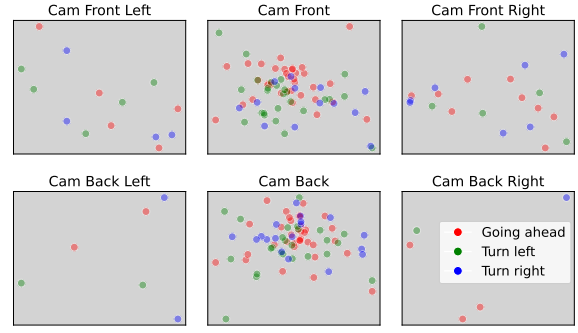
Figure G.2. GPT evaluation prompts for **Visual Question Answering (VQA)** in our benchmark.



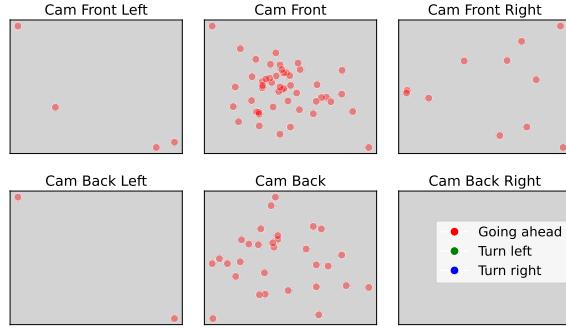
Figure G.3. **GPT-4o failure cases.** (a): GPT-4o reasons the moving status of the pedestrian by the moving position related to the frame, instead of the coordinate of the moving object itself, thus leading to wrong perception results. (b): The model struggles to distinguish the correct direction based on the coordinates of the target object. (c): GPT-4o reasons that the moving status of the SUV by the relative location of the object to the current frame causes the wrong perception results. (d): GPT-4o fails to perceive the orientation of the car. (e): The dataset contains examples that need multiple frames to reason successfully, GPT-4o fails to address these examples with a single image input. (f): GPT-4o reasons that the moving status of the SUV by the relative location of the object to the current frame causes the wrong perception results.



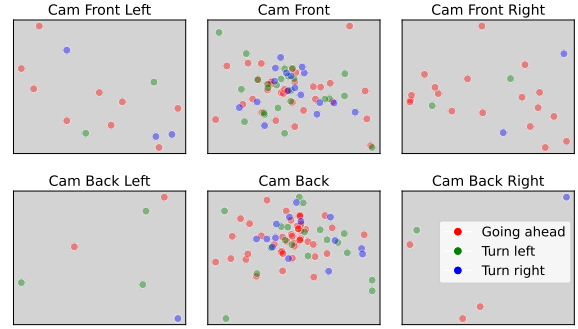
(a) Ground Truth



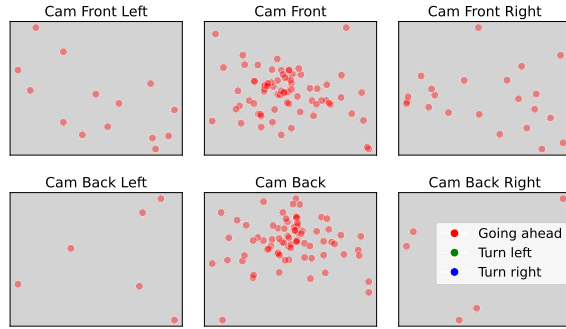
(b) GPT-4o [2]



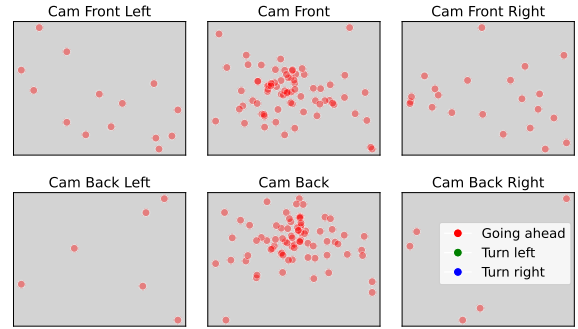
(c) Phi-3 [1]



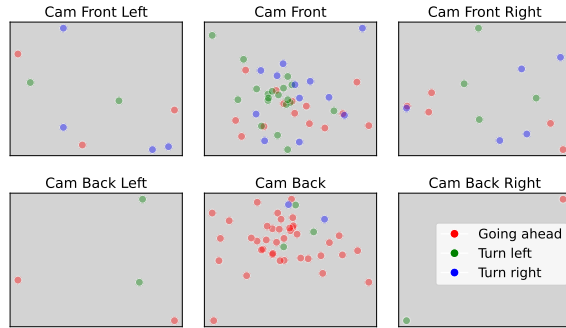
(d) Phi-3.5 [1]



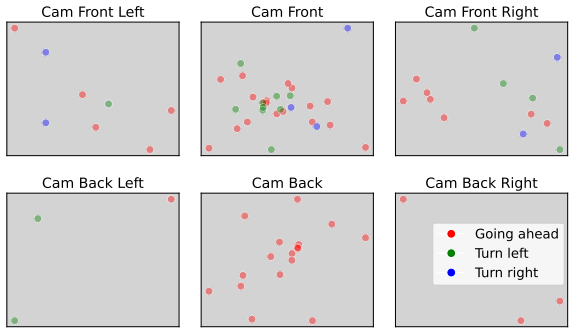
(e) LLaVA-1.5_{7B} [14]



(f) LLaVA-1.5_{13B} [14]



(g) Qwen2-VL_{72B} [20]



(h) InterVL_{8B} [5]

Figure G.4. **Prediction spatial distributions from VLMs.** The locations represent the object positions in the image within each camera, which is input to the model as a text description. We only visualize the data point where the model response aligns with the provided multiple choices (e.g., *Going ahead*, *Turn Left*, and *Turn Right*).

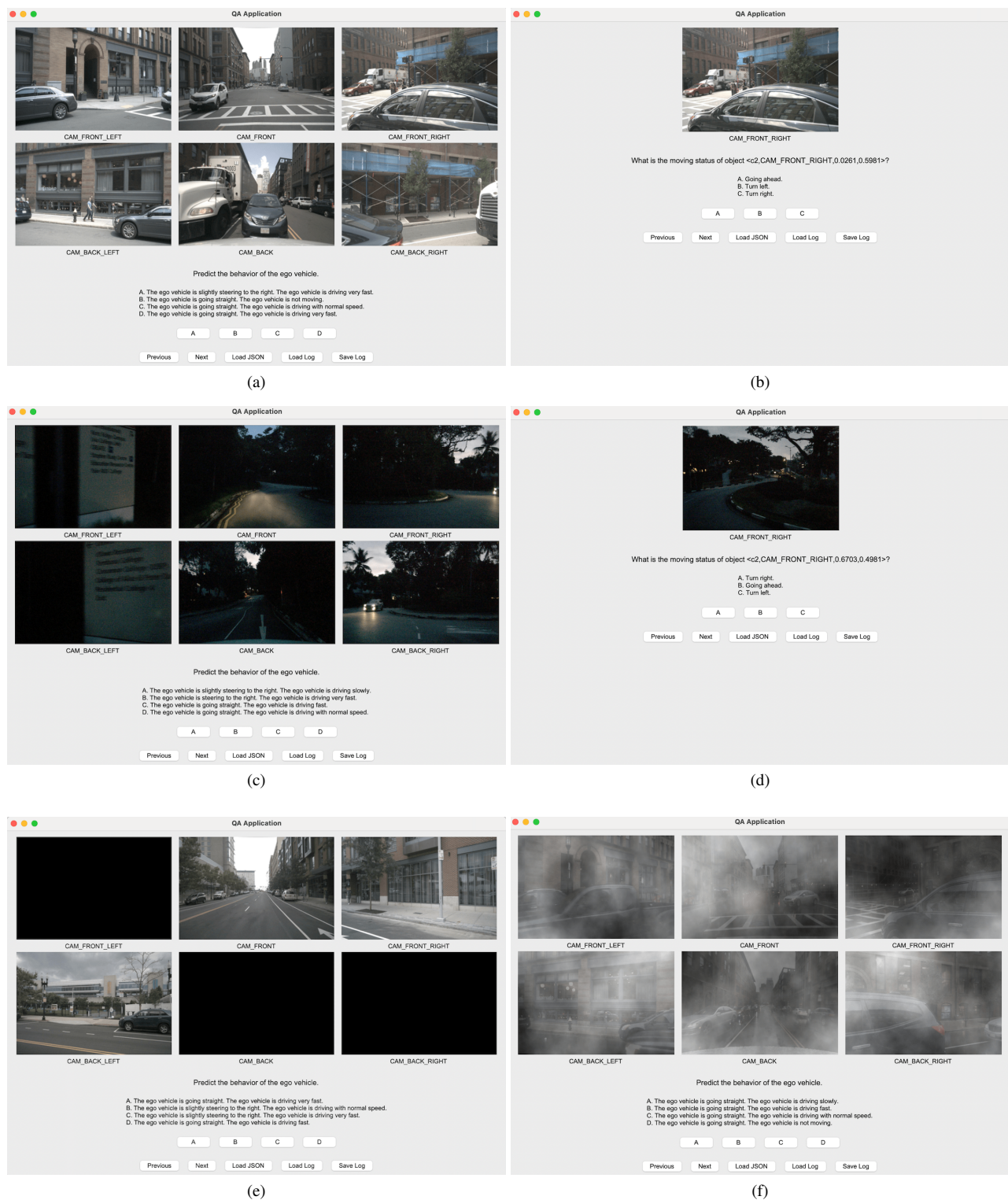


Figure G.5. Illustrative examples from our **human evaluation interfaces**.

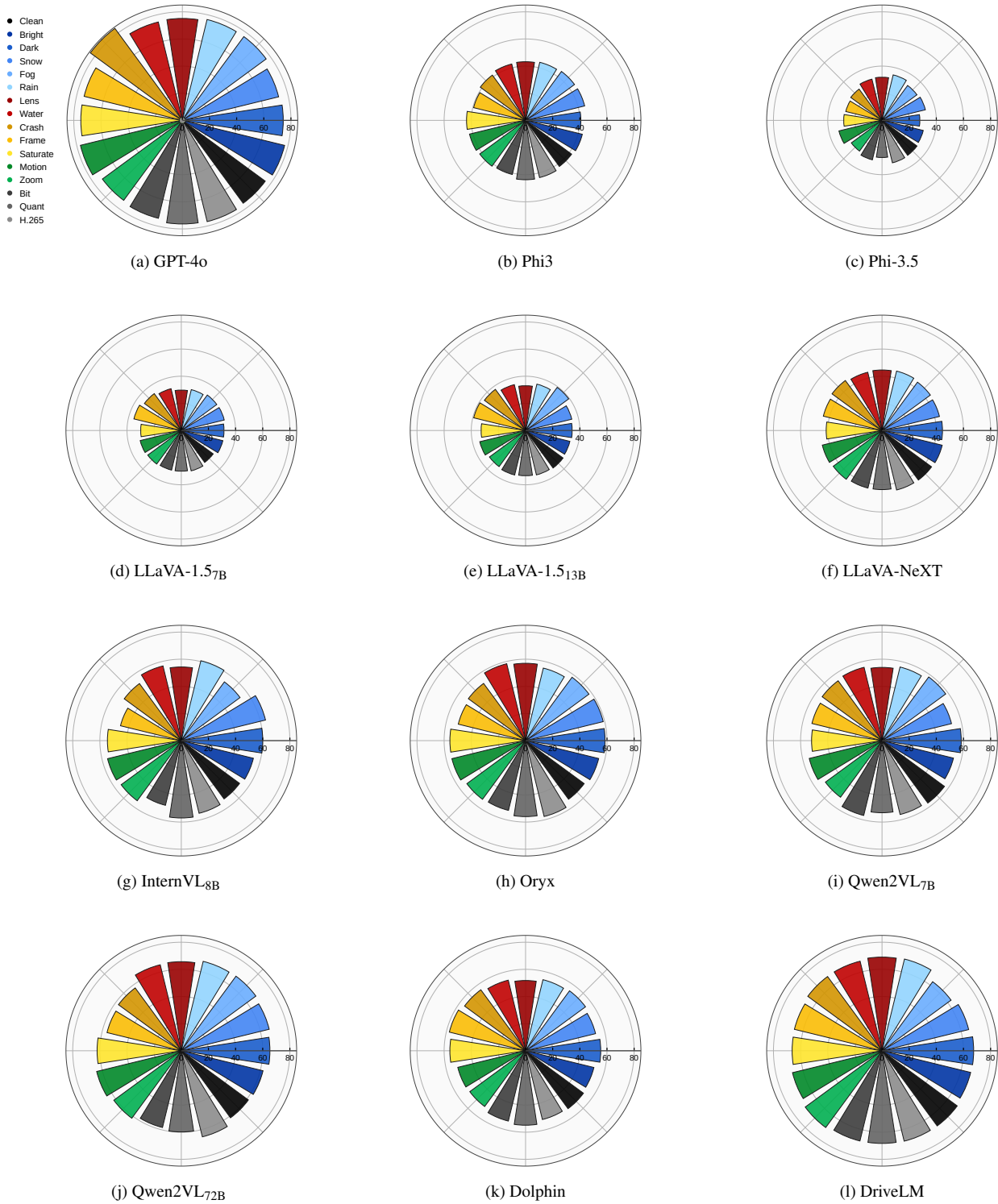


Figure G.6. **Model performance comparisons using radar graphs.** The performance for each input corruption type is averaged across all the 1,261 questions spanning four different tasks using GPT scores. The gray dashed line represents the performance of **text-only** input. We observe that VLMs have subtle performance changes under corruptions. For some models, the GPT scores under only text input are even higher than the performance when the visual information is available.

	Perception		Prediction		Planning		Behavior
	Q: What are the <u>important objects</u> in the current scene?	Q: What <u>object</u> should be noticed first when getting to the <u>next possible location</u> ?	Q: In this scenario, what are the <u>safe actions to take</u> for the ego vehicle?	Q: Predict the <u>behavior</u> of the ego vehicle in the current scene.			
	A: The <u>important objects</u> to consider for decision are: <ul style="list-style-type: none">• <u>Bus</u>, near the center of the <u>front image</u>, and is moving slowly.• <u>Pedestrian</u>, near the left side of the <u>front-left image</u>, and is walking on the sidewalk.	A: The ego vehicle should notice (1) the <u>pedestrian</u> on the crosswalk from <u>front-left image</u> , as it is walking across the road, (2) the oncoming vehicle from <u>front-right image</u> , as it is approaching from the opposite direction.	A: The recommended <u>safe actions</u> are (1) it is getting dark, ensure <u>headlights are on</u> , (2) there is a pedestrian on the left sidewalk, <u>maintain awareness</u> , (3) maintain lane discipline and <u>be cautious</u> of vehicles approaching from the opposite direction.	A: The ego vehicle is <u>steering to the right</u> side of the road. The ego vehicle is driving with <u>normal speed</u> in the current condition.			
Total Score: 74/100		Total Score: 73/100		Total Score: 87/100		Total Score: 37/100	
	Brief Summary: The answer provides a detailed analysis of <u>objects</u> in the scene, showing good <u>awareness of potential hazards</u> . However, there are some inaccuracies in the <u>motion precision</u> and <u>situational awareness</u> aspects.	Brief Summary: The answer provides a reasonable analysis of the <u>objects</u> in the scene and suggests <u>appropriate actions</u> . However, there are some inaccuracies in identifying the <u>objects</u> and <u>their states</u> , as well as minor issues with <u>clarity</u> .		Brief Summary: The answer demonstrates a high level of accuracy in predicting <u>safe actions</u> based on the scenario. The reasoning and justification are logical and align well with <u>safety considerations</u> . It could benefit from slightly <u>more conciseness</u> .		Brief Summary: The predicted answer <u>did not match</u> the correct answer. However, the explanation demonstrated a good understanding of the ego vehicle's <u>behavior</u> and provided clear reasoning and justification based on the <u>environmental cues</u> .	

Figure G.7. Examples of GPT-4o [2] responses to four tasks and the corresponding evaluation results under the **dark condition**. We observe that GPT-4o [2] is aware of the low-light environment and can identify the *bus* and *pedestrian* from the image, showing resilience.

 Perception	 Prediction	 Planning	 Behavior	
 Q: What is the <u>moving status</u> of object in the <u>front view</u> camera image?	Q: What <u>object</u> should be noticed first when getting to the <u>next possible location</u> ?	Q: What <u>actions</u> taken by the ego vehicle can lead to a <u>collision with the object</u> ?	Q: Predict the <u>behavior</u> of the <u>ego vehicle</u> in the current scene.	
 A: The object at coordinates (0.5073, 0.5778) appears to be <u>turning left</u> . The image shows a <u>vehicle</u> moving from the left side of the image towards the center, indicating a <u>leftward motion status</u> .	A: The ego vehicle should notice the <u>car</u> in the <u>front image</u> , as it is <u>moving faster than</u> the ego vehicle. The ego vehicle should <u>maintain its lane</u> and monitor the rear-view camera for any changes in the approaching car.	A: The object appears to be a <u>vehicle</u> in the center of the image. Actions that could lead to a collision include <u>sudden speeding up</u> and <u>lane change</u> . To avoid a collision, the ego should maintain a <u>steady speed</u> and <u>signal intentions</u> clearly.	A: The ego vehicle is slightly <u>steering to the left</u> side of the road. The ego vehicle is <u>driving fast</u> in the current scene.	
Total Score: 0/100	Total Score: 63/100	Total Score: 87/100	Total Score: 100/100	
 Brief Summary: The predicted answer does not align with the correct answer , and there are inaccuracies in object recognition, location, orientation, environmental conditions, and reasoning.	Brief Summary: The answer correctly identifies the object but fails to prioritize according to the instructions. The description of the object and state are generally accurate, but there are inaccuracies and lack of clarity in the recommended actions for the ego vehicle.	Brief Summary: The answer demonstrates a strong understanding of potential actions that could lead to a collision with the identified object . The reasoning provided is logical and aligns well with safety considerations .	Brief Summary: The predicted answer receives a perfect score of 100 as it exactly matches the correct answer , demonstrates a clear understanding of behavior and details , and provides logical reasoning based on environmental clues.	

Figure G.8. Examples of GPT-4o [2] responses to four tasks and the corresponding evaluation results under the **motion blur condition**. We observe that GPT-4o [2] is influenced by this type of corruption and tends to predict “driving fast” based on it. The example shows the potential of visual corruption to influence high-level driving decisions.



Figure G.9. Examples of different VLM responses under the **frame lost condition**. We observe that GPT-4o [2] responses with visible objects while LLaVA-NeXT [15] and DriveLM [19] tend to hallucinate objects that cannot be seen from the provided images.

	<p>Q: What are the important objects in the current scene? Those objects will be considered for the future reasoning and driving decision making.</p>				
	<p>Ground Truth: In this scene, there is a black SUV to the front, a black SUV to the back, a black sedan to the front, a pedestrian riding a bike to the front, and a red light to the front of the ego vehicle . . .</p>				
	GPT4-o	Qwen2-VL	LLaVA-NeXT	DriveLM-Agent	
	<p>A: Based on the images, the important objects for driving decision-making are:</p> <ul style="list-style-type: none">Vehicles, visible in the lane ahead, parked or moving slowly on the left side of the road, and visible behind.Traffic signals, visible in the distance.Road conditions, which is wet and indicating the need for cautious driving.	<p>A: In the current scene, the important objects are:</p> <ul style="list-style-type: none">Vehicles, there are multiple on the road, which are moving and need to be tracked for safe navigation.Traffic lights, which will dictate when it is safe to proceed or stop.Pedestrians, who may cross the road and need to be accounted for.	<p>A: In the current scene, the important objects include:</p> <ul style="list-style-type: none">The car in the foreground, which is the vehicle the autonomous driving assistant is controlling.The traffic lights, which are relevant for traffic rules and signals.The pedestrians, which are important for ensuring safety.The other vehicles.	<p>A: There is a black car to the back of the ego vehicle, a black sedan to the front of the ego vehicle, a black sedan to the front of the ego vehicle, and a white sedan to the front of the ego vehicle.</p>	
	<div><div>1. Action Alignment:15</div><div>2. Motion Precision: 0</div><div>3. Driving Context Appropriateness: 10</div><div>4. Situational Aware: 10</div><div>5. Concise & Clarity: 15</div><div>6. Grammar: 10</div></div> <div>Total Score: 60/100</div>	<div><div>1. Action Alignment:15</div><div>2. Motion Precision: 0</div><div>3. Driving Context Appropriateness: 10</div><div>4. Situational Aware: 10</div><div>5. Concise & Clarity: 15</div><div>6. Grammar: 10</div></div> <div>Total Score: 60/100</div>	<div><div>1. Action Alignment:15</div><div>2. Motion Precision: 0</div><div>3. Driving Context Appropriateness: 10</div><div>4. Situational Aware: 10</div><div>5. Concise & Clarity: 15</div><div>6. Grammar: 8</div></div> <div>Total Score: 58/100</div>	<div><div>1. Action Alignment:10</div><div>2. Motion Precision: 0</div><div>3. Driving Context Appropriateness: 5</div><div>4. Situational Aware: 5</div><div>5. Concise & Clarity: 10</div><div>6. Grammar: 10</div></div> <div>Total Score: 40/100</div>	
	<div><div>Brief Summary: The predicted answer provides a detailed list of important objects in the current scene, aligning well with the driving context and demonstrating good situational awareness. However, there is a lack of motion precision and failure to mention all objects mentioned in the answer.</div></div>	<div><div>Brief Summary: The predicted answer provides a detailed list of important objects in the scene, aligning well with the driving context and demonstrating good situational awareness. However, there is a lack of specific motion predictions and inaccuracies in identifying all the objects.</div></div>	<div><div>Brief Summary: The predicted answer provides a detailed list of important objects in the scene, aligning well with the driving context and demonstrating situational awareness. However, the predicted motion is missing.</div></div>	<div><div>Brief Summary: The predicted answer accurately identified some of the important objects in the scene but missed crucial details such as the presence of a pedestrian riding a bike and a red light. The motion precision was completely off, as it did not mention any motion-related information.</div></div>	
					

Figure G.10. Examples of different VLM responses under the **water splash condition**. We observe that, under severe visual corruptions, VLMs respond with ambiguous and general answers based on their learned knowledge, without referring to the visual information. Most responses include traffic signals and pedestrians, even though they are not visible in the provided images.



What is the moving status of object <c2,CAM_BACK,0.5443,0.5639>? Please select the correct answer from the following options: A. Turn left. B. Going ahead. C. Turn right.



B. Going ahead.

The object is a car that is positioned in the lane ahead of the ego vehicle and appears to be moving in the same direction as the ego vehicle.

1. Answer Correctness (50 points): 50
2. Object Recognition (10 points): 5
3. Location & Orientation (15 points): 5
4. Environmental Condition Awareness (15 points): 10
5. Clarity of Reasoning (10 points): 5

Total Score: 75/100



Brief Summary:

The answer provided accurately identifies the object as a car and correctly states its moving status as "Going ahead." The explanation includes details about the object's position relative to the ego vehicle and considers environmental factors that could impact perception. The reasoning is clear and logically structured, earning a high score of 75.



B. Going ahead.


The object is located in the center of the image and appears to be moving forward, as there are no visible indicators of turning or stopping.

1. Answer Correctness (50 points): 50
2. Object Recognition (10 points): 5
3. Location & Orientation (15 points): 5
4. Environmental Condition Awareness (15 points): 10
5. Clarity of Reasoning (10 points): 5

Total Score: 75/100

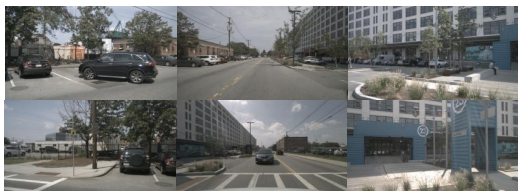
Brief Summary:

The answer provided accurately identifies the moving status of the object as "Going ahead" based on the visual description of a gray car. The explanation is clear, logically structured, and considers environmental conditions that could impact perception, earning a total score of 75.

Figure G.11. Example of Qwen2-VL_{72B} for the  **Perception** task. Under the **text-only condition** (right example), the model can fabricate answers based on the coordinate information in the question. The GPT evaluator can not reflect the fabricated answer, even with detailed rubrics, resulting in the same score as prompted with a clean image.



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 ...



First Object Noticed:

- **Object:** The black SUV in CAM_FRONT_RIGHT.
- **State:** It is parked on the right side of the road.
- **Action:** The ego vehicle should maintain its current speed and trajectory, as the parked SUV does not pose an immediate threat.

Second Object Noticed:

- **Object:** The white van in CAM_FRONT_RIGHT.
- **State:** It is parked on the right side of the road.
- **Action:** The ego vehicle should maintain its current speed and trajectory, as the parked van does not pose an immediate threat.

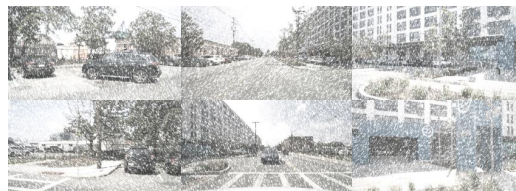
Third Object Noticed:

- **Object:** The white car in CAM_FRONT_RIGHT.
- **State:** It is parked on the right side of the road.
- **Action:** The ego vehicle should maintain its current speed and trajectory, as the parked car does not pose an immediate threat.

1. Object Identification & Priority Order: 15/20
2. Object Category & Visual Description: 10/20
3. State of the Object: 5/15
4. Recommended Action for Ego Vehicle: 10/15
5. Logical Flow & Reasonableness: 10/20
6. Clarity and Grammar: 8/10

Total Score: 58/100

Brief Summary: The predicted answer correctly identifies the objects in the correct priority order, but there are inaccuracies in the visual descriptions and states of the objects. The recommended actions ... The logical flow ... The clarity and grammar ... Overall, the predicted answer needs improvement in providing accurate visual descriptions, states of objects, and recommended actions for the ego vehicle to align better with the correct answer.



First Object Noticed:

- **Object:** The white car in CAM_FRONT_RIGHT.
- **State:** It is parked on the right side of the road.
- **Action:** The ego vehicle should maintain its current speed and trajectory, as there is no immediate need to slow down or change direction.

Second Object Noticed:

- **Object:** The black car in CAM_FRONT_LEFT.
- **State:** It is parked on the left side of the road.
- **Action:** The ego vehicle should continue to maintain its current speed and trajectory, as there is no immediate obstacle in its path.


Third Object Noticed:

- **Object:** The building in CAM_FRONT.
- **State:** It is stationary on the right side of the road.
- **Action:** The ego vehicle should continue to maintain its current speed and trajectory, as the building does not pose an immediate threat to its path.

1. Object Identification & Priority Order: 15/20
2. Object Category & Visual Description: 18/20
3. State of the Object: 12/15
4. Recommended Action for Ego Vehicle: 12/15
5. Logical Flow & Reasonableness: 16/20
6. Clarity and Grammar: 8/10

Total Score: 81/100

Brief Summary: The predicted answer demonstrates a good understanding of the scenario and accurately identifies the objects in the correct priority order. However, there are minor inaccuracies in the description of the object states and recommended actions for the ego vehicle. The logical flow ... Overall, the predicted answer provides a reasonable response but could benefit from more precise details and clearer explanations.

Figure G.12. Example of Qwen2-VL_{72B} for the  Prediction task. Under the **snow condition** (right example), the model can still observe some visible objects in the images. Since most objects are severely occluded, the third object noticed becomes more general: the building, compared to the clean inputs, where all the objects are visible vehicles. However, the GPT score is even higher. The example illustrates that even the GPT evaluator can not accurately reflect the language quality generated by VLMs.



What actions could the ego vehicle take based on `<c1,CAM_BACK_LEFT,0.7995,0.5463>`? Why take this action and what's the probability?



The ego vehicle could take the action of **stopping** at the **stop sign** located at (0.7995, 0.5463) in the image. This action is necessary because the **stop sign indicates that the vehicle must come to a complete stop before proceeding**. The probability of taking this action is high, as it is a legal requirement to obey traffic signs and ensure safety at intersections.

1. Action Prediction Accuracy: 10/40
2. Reasoning & Justification: 15/20
3. Probability or Confidence Level: 5/15
4. Contextual & Safety Considerations: 5/15
5. Conciseness & Clarity: 8/10

Total Score: 43/100



Brief Summary: The predicted answer does not accurately align with the correct action of keeping going at the same speed. The reasoning provided is somewhat clear but lacks a direct connection to the scenario. The probability level is mentioned but not fully supported by the context. Contextual awareness and safety considerations are minimal. The answer is relatively concise and clear, earning some points in that category. Overall, the predicted answer falls short in accurately predicting the correct action and lacks strong reasoning and contextual awareness.




The ego vehicle could take the action of **slowing down or stopping**. This is because the object at (0.7995, 0.5463) in the back left view is **likely a pedestrian or a stationary object on the road**. Taking this action would help avoid a potential collision and ensure the safety of the pedestrian or the vehicle. The probability of needing to take this action is high, as the object is in close proximity to the ego vehicle's path.

1. Action Prediction Accuracy: 30/40
2. Reasoning & Justification: 15/20
3. Probability or Confidence Level: 12/15
4. Contextual & Safety Considerations: 12/15
5. Conciseness & Clarity: 8/10

Total Score: 77/100

Brief Summary: The predicted answer accurately identifies the need to slow down or stop based on the presence of a pedestrian or stationary object in the back left camera view. The reasoning provided is logical and aligns with safety considerations. The probability assessment is reasonable but could be more specific. The answer demonstrates good contextual awareness and safety considerations. However, it could be more concise and focused on key points.

Figure G.13. Example of Qwen2-VL_{72B} for  Planning task. Under the **text-only condition** (right example), the model fails to perceive the motorcycle to stop sign, giving their close distance on the image. When the image is absent, the model guesses the object to be a pedestrian or a stationary object on the road, unexpectedly leading to an answer with higher scores.



What is the moving status of object <c2,CAM_BACK,0.5443,0.5639>? Please select the correct answer from the following options: A. Turn left. B. Going ahead. C. Turn right.



B. Going ahead.

Explanation: The object is positioned in the center of the lane and appears to be moving straight ahead, as indicated by its alignment with the road markings and the absence of any visible turning action.

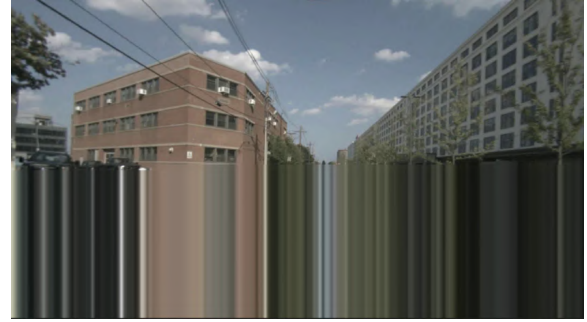
1. **Answer Correctness:** 50 points
2. **Object Recognition:** 10 points
3. **Object Location and Orientation:** 15 points
4. **Environmental Condition Awareness:** 15 points
5. **Clarity of Reasoning:** 10 points

Total Score: 100/100



Brief Summary:

The answer provided is accurate, correctly identifying the moving status of the object as "Going ahead" based on its position and alignment with the road markings. The explanation is clear, logically structured, and considers environmental conditions that could impact perception. Overall, the answer meets all criteria and earns a perfect score of 100.



B. Going ahead.


Explanation: The object is located slightly to the right of the center of the image, which suggests it is aligned with the road. Given the perspective of the image, the object appears to be moving straight along the road. Therefore, the correct answer is B.

1. **Answer Correctness:** 50 points
2. **Object Recognition:** 5 points
3. **Object Location and Orientation:** 10 points
4. **Environmental Condition Awareness:** 10 points
5. **Clarity of Reasoning:** 5 points

Total Score: 80/100

Brief Summary:

The answer provided accurately identifies the moving status of the object as "Going ahead" based on its position and alignment with the road. The reasoning is clear and logically structured, supporting the correct answer. However, there were some minor inaccuracies in object recognition and environmental condition awareness, which slightly impacted the overall score.

Figure G.14. Example of GPT-4o for  **Perception** task. Under the **bit error condition** (right example), the visual information for the target object is completely corrupted but the model can guess the answer correctly. However, the GPT evaluator here captures the nuance in the explanation and assigns a higher score to the more detailed answer with clean image inputs.



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



- **Maintain Speed:** Keep a steady speed to ensure safe distance from vehicles behind.
- **Prepare to Steer Right:** As the road curves, be ready to steer smoothly to the right.
- **Monitor Bus Movement:** Keep an eye on the bus in the front-right camera to anticipate any merging or stopping.
- **Check Mirrors:** Regularly check side mirrors for the position of the motorcycle and other vehicles.
- **Signal Intentions:** Use turn signals if planning to change lanes or make a turn to alert surrounding vehicles.



Turn left, decelerate gradually without braking, are safe actions to take for the ego vehicle in this scenario.

1. Action Prediction Accuracy: 35
2. Reasoning & Justification: 18
3. Probability or Confidence Level: 12
4. Contextual Awareness & Safety : 13
5. Conciseness and Clarity: 8

Total Score: 86/100



Brief Summary:

The predicted answer demonstrates a strong understanding of the scenario, accurately identifying key observations and recommending appropriate actions for the ego vehicle. While there are minor discrepancies and areas for improvement in conciseness and clarity, the response overall aligns well with the correct answer and provides a logical justification for the chosen actions.



- **Slow Down:** Reduce speed to assess the situation and prepare for any necessary maneuvers.
- **Stop or Maneuver:** If the object in the front cannot be safely bypassed, stop the vehicle.
- **Check Surroundings:** Use all cameras to ensure no pedestrians or vehicles are approaching before making any turns or reversing.
- **Signal Intentions:** Use turn signals if planning to maneuver around the object or change lanes.


Slightly offset to the **left**, **change to the left lane**, and **turn left**.

1. Action Prediction Accuracy: 30
2. Reasoning & Justification: 15
3. Probability or Confidence Level: 10
4. Contextual Awareness & Safety : 12
5. Conciseness and Clarity: 8

Total Score: 75/100

Brief Summary:

The predicted answer demonstrates a good understanding of the driving scenario and provides relevant actions for the ego vehicle. The reasoning and justification are clear, although some minor details could be further elaborated. The contextual awareness and safety considerations are mostly accurate, with a few areas for improvement. The answer is concise and clear but could benefit from a more focused explanation. Overall, the predicted answer is well-structured and **aligns with the correct answer**, earning a score of 75.

Figure G.15. Example of GPT-4o for  **Planning** task. In both **low-light** (left example) and **normal** (right example) conditions, the ground truth is to turn left. However, in the left case, the vehicle is prepared to turn left. Turning right will cause a potential collision with the white sedan on the front right camera. In the right image, the road is empty, and the predicted slow-down action will not cause safety issues. However, the GPT evaluator assigns a higher score on the left case, ignoring the potential safety consequences, as there are no high-fidelity physical context prompts.

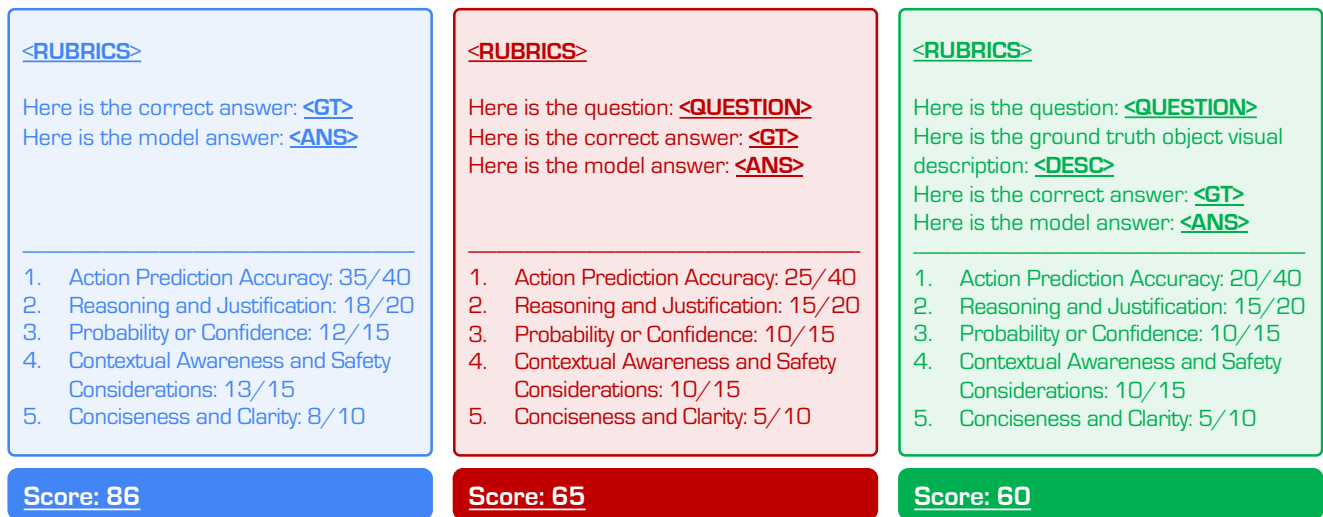
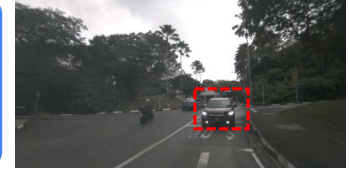


Figure G.16. Examples of **GPT evaluators with different information**. The evaluator can revise the score and give a more accurate evaluation based on more contextual information in the driving scenarios. For more prompt details, please refer to Fig. G.1 and Fig. G.2.



What actions could the ego vehicle take based on <c1,CAM BACK,0.6323,0.6314>? Why take this action and what is the probability?



Based on the image, the ego vehicle is at an intersection with a motorcycle and a car visible behind it.

1. **Maintain Current Speed:** The motorcycle and car appear to be at a safe distance behind the ego vehicle. There are no immediate obstacles or hazards directly in the path.
Probability: High, as maintaining speed is often the default action when no immediate threats are present.
2. **Prepare to Stop:** If the ego vehicle is approaching a stop sign or traffic light (not visible in the image), it may need to slow down or stop.
Probability: Moderate, depending on traffic signals or signs not visible in the image.
3. **Change Lanes:** If the ego vehicle needs to turn or avoid an obstacle ahead (not visible in the image), it might consider changing lanes.
Probability: Low, as there is no immediate indication of a need to change lanes.

Conclusion: The most likely action is to maintain current speed, given the apparent safe distance from vehicles behind and no visible immediate obstacles.



The action is to keep going at the same speed. The reason is that there is no safety issue, and the probability is high.



1. **Action Prediction Accuracy (40 points):** Award full points only for exact or highly similar action matches. Deduct points for inaccuracies or actions that do not match the correct answer, especially if they could compromise safety or appropriateness in context.
2. **Reasoning & Justification (20 points):** Deduct points if the reasoning lacks clarity, omits relevant details, or includes contradictions. The explanation should justify the action in a way that is suitable for the scenario provided.
3. **Probability or Confidence Level (15 points):** Deduct points if the probability level is missing, implausible, or does not align with the action or reasoning provided.
4. **Contextual Awareness & Safety (15 points):** Deduct points for failing to consider contextual factors that may impact the ego vehicle's decision, especially if they could lead to unsafe actions.
5. **Conciseness & Clarity (10 points):** Deduct points for verbosity, ambiguity, or lack of focus that could hinder quick comprehension. Assign 0 points if no explanation is provided.

Here is the question: <QUESTION>

Here is the ground truth object visual description: <DESC>

Here is the correct answer: <GT>

Here is the model answer: <ANS>

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. Output the number only.

Here is the correct answer: <GT>

Here is the model answer: <ANS>

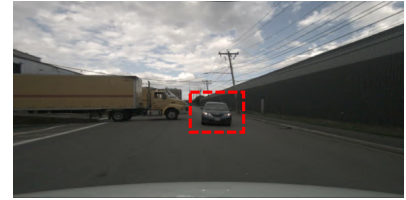
Total Score: 73/100

Total Score: 92/100

Figure G.17. Examples of the GPT evaluator with different prompts. With the same question and answer, the GPT evaluator assigns different scores based on the rubrics. Without detailed information on the physical context, the scoring point, the evaluator fails to capture critical information for driving decisions in the answer, but simply based on semantic similarities.



What actions could the ego vehicle take based on <c1,CAM_BACK,0.5073,0.5778>? Why take this action and what is the probability?



Based on the image and the coordinates provided, the ego vehicle is facing a situation where a large truck is making a turn and partially blocking the road. The coordinates (0.5073, 0.5778) likely correspond to the position of the truck, which is a critical object in this scenario.

Possible Actions:

1. Slow Down or Stop: The truck is obstructing the path, and there is an oncoming vehicle. Slowing down or stopping allows the truck to complete its turn safely and avoids potential collisions.
2. Wait for the Truck to Clear: Once the truck has cleared the path, the ego vehicle can proceed. This ensures that the road is clear and safe to continue driving.
3. Prepare to Yield: If the oncoming vehicle has the right of way, yielding ensures compliance with traffic rules and safety.

Conclusion: The most immediate and probable action is to slow down or stop to allow the truck to complete its turn safely. This action minimizes risk and ensures the safety of all vehicles involved.



The action is to keep going at the same speed. The reason is that there is no safety issue, and the probability is high



<RUBRICS>

Here is the correct answer: <GT>
Here is the model answer: <ANS>

1. Action Prediction Accuracy: 35/40
2. Reasoning & Justification: 18/20
3. Probability or Confidence Level: 12/15
4. Contextual Awareness and Safety: 13/15
5. Conciseness & Clarity: 8/10

Total Score: 86/100

<RUBRICS>

Here is the question:
<QUESTION>
Here is the correct answer: <GT>
Here is the model answer: <ANS>

1. Action Prediction Accuracy: 25/40
2. Reasoning & Justification: 15/20
3. Probability or Confidence Level: 10/15
4. Contextual Awareness & Safety: 10/15
5. Conciseness & Clarity: 5/10

Total Score: 65/100


<RUBRICS>

Here is the question:
<QUESTION>
Here is the ground truth object visual description: <DESC>
Here is the correct answer: <GT>
Here is the model answer: <ANS>

1. Action Prediction Accuracy: 20/40
2. Reasoning & Justification: 15/20
3. Probability or Confidence Level: 10/15
4. Contextual Awareness & Safety: 10/15
5. Conciseness & Clarity: 5/10

Total Score: 60/100

Figure G.18. Examples of the GPT evaluators with different information. We gradually add more information about the question and the visual description of the target objects. The evaluator gives a more accurate score based on more information.

Table H. Detailed GPT score results of MCQs for the  Perception task. “Clean” represents clean image inputs. “T.O.” represents text-only evaluation. The “Corrupt” settings range from weather conditions, external disturbances, sensor failures, motion blur, and transmission errors. The benchmarked VLMs include commercial, open-sourced, and driving specialist models, respectively.



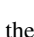
Method	 Clean	 T.O.	 Brightness	 Dark	 Snow	 Fog	 Rain	 Lens Obstacle	 Water Splash	 Camera Crash	 Frame Lost	 Saturate	 Motion Blur	 Zoom Blur	 Bit Error	 Color Quant	 H.265 Compression
 Human	47.67	-	43.33	26.67	45.00	18.33	40.00	35.33	37.33	42.00	20.00	25.00	43.00	20.00	33.67	25.33	31.33
GPT-4o	41.87	43.59	43.84	44.82	45.18	44.30	46.20	45.69	44.10	38.12	40.08	39.32	41.40	36.67	37.54	38.22	39.37
Phi-3	35.51	32.65	35.28	34.15	38.88	38.22	37.70	38.39	36.75	34.93	33.89	37.53	37.80	39.72	36.49	37.15	37.03
Phi-3.5	40.22	38.33	39.22	36.46	41.41	43.04	41.40	40.66	40.83	37.59	36.91	39.86	40.33	42.30	36.71	41.49	41.48
LLaVA-1.5 _{7B}	32.40	32.68	32.48	32.95	31.95	32.43	32.30	32.88	32.18	32.93	31.63	32.50	32.43	32.93	32.48	32.18	32.63
LLaVA-1.5 _{13B}	33.58	33.25	33.25	33.25	33.15	33.50	33.53	32.95	32.93	33.48	33.40	33.25	33.45	33.68	33.33	33.25	33.38
LLaVA-NeXT	32.98	4.20	33.85	20.43	11.62	16.58	27.33	18.24	32.30	26.80	20.83	23.50	34.00	17.75	24.50	18.03	26.24
InternVL _{8B}	46.60	52.46	43.65	44.15	43.58	46.02	42.38	41.48	43.38	45.32	49.08	43.98	41.30	41.50	38.25	44.84	42.13
Oryx	17.98	20.87	16.48	16.88	16.63	16.79	14.31	16.35	15.85	16.49	21.44	21.38	16.36	21.04	17.65	18.13	19.51
Qwen2VL _{7B}	42.64	37.76	43.08	37.29	39.72	41.67	40.87	40.69	39.89	39.75	39.17	40.85	41.32	39.62	34.28	39.90	41.20
Qwen2VL _{72B}	38.15	21.53	36.24	37.77	35.91	35.78	37.13	38.14	38.97	29.48	25.63	36.87	36.91	37.01	30.90	36.05	41.48
Dolphin	6.50	8.35	10.18	11.08	10.70	9.53	10.58	9.93	9.80	10.08	9.95	11.20	9.85	10.10	8.80	10.00	11.10
DriveLM	22.38	12.45	20.78	25.30	18.98	24.43	25.95	22.03	21.03	21.95	16.28	19.38	22.98	20.93	19.90	16.25	26.48

Table I. Detailed Accuracy score results of MCQs for the  Perception task. “Clean” represents clean image inputs. “T.O.” represents text-only evaluation. The “Corrupt” settings range from weather conditions, external disturbances, sensor failures, motion blur, and transmission errors. The benchmarked VLMs include commercial, open-sourced, and driving specialist models, respectively.



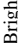
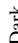
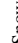















Method	 Clean	 T.O.	 Brightness	 Dark	 Snow	 Fog	 Rain	 Lens Obstacle	 Water Splash	 Camera Crash	 Frame Lost	 Saturate	 Motion Blur	 Zoom Blur	 Bit Error	 Color Quant	 H.265 Compression
 Human	93.33	-	80.00	53.33	80.00	33.33	80.00	66.67	73.33	80.00	40.00	46.67	80.00	40.00	60.00	46.67	53.33
GPT-4o	59.00	59.50	60.50	63.50	59.00	61.00	59.50	58.00	59.00	50.00	51.50	56.50	57.50	52.00	51.00	54.00	57.50
Phi3	54.50	17.50	55.00	33.00	50.00	55.00	59.00	56.00	57.50	32.50	23.50	58.00	57.50	49.00	36.00	49.50	57.50
Phi-3.5	56.50	58.50	59.00	58.00	58.00	59.00	58.50	60.00	59.50	58.50	57.00	59.50	59.00	59.50	57.50	60.50	58.50
LLaVA-1.5 _{7B}	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00
LLaVA-1.5 _{13B}	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00
LLaVA-NeXT	55.00	34.50	53.50	42.50	41.50	44.00	55.00	42.00	52.00	31.00	26.50	50.50	53.50	39.00	44.50	37.50	49.00
InternVL _{8B}	56.50	23.50	57.50	61.00	62.00	59.50	59.50	59.50	60.00	50.50	31.00	56.50	56.00	59.00	60.00	58.50	58.50
Oryx	51.00	19.00	53.50	50.50	52.00	52.50	52.50	50.00	48.00	32.50	29.00	52.00	48.00	36.00	50.50	52.00	52.00
Qwen2VL _{7B}	59.00	56.50	60.00	59.00	60.00	59.50	59.00	59.00	59.00	58.00	56.00	59.00	59.50	55.00	54.50	57.00	59.00
Qwen2VL _{72B}	60.00	23.50	58.00	60.00	58.50	59.50	61.50	60.00	59.50	39.00	29.50	60.00	59.00	55.50	55.50	58.50	61.00
Dolphin	7.00	9.00	5.50	5.50	5.50	6.00	4.50	6.50	5.50	5.50	7.50	6.50	6.00	5.50	6.00	5.50	5.50
DriveLM	40.00	23.00	37.00	45.00	35.00	46.50	45.50	40.50	37.50	39.50	29.50	34.00	42.00	36.50	36.00	30.50	47.00

Table J. Detailed GPT score results of MCQs for the  Behavior task. “Clean” represents clean image inputs. “T.O.” represents text-only evaluation. The “Corrupt” settings range from weather conditions, external disturbances, sensor failures, motion blur, and transmission errors. The benchmarked VLMs include commercial, open-sourced, and driving specialist models, respectively.

Method	Clean	T.O.	Brightness	Dark	Snow	Fog	Rain	Lens Obstacle	Water Splash	Camera Crash	Frame Lost	Saturate	Motion Blur	Zoom Blur	Bit Error	Color Quant	H.265 Compression
 Human	65.00	-	40.00	40.67	66.67	51.00	73.33	53.33	56.00	60.00	36.67	54.00	67.67	34.67	53.33	73.33	53.33
GPT-4o	45.40	50.03	46.27	49.20	42.54	41.79	46.55	45.30	47.78	45.28	40.44	47.89	39.91	32.25	48.40	50.35	41.07
Phi-3	45.20	40.91	45.98	44.48	47.91	45.17	47.45	44.22	44.02	43.65	44.01	43.51	42.48	41.05	43.83	46.60	44.30
Phi-3.5	36.75	39.16	37.15	38.14	37.19	39.53	38.40	36.79	37.36	36.83	37.98	39.09	37.70	38.91	38.27	38.23	36.85
LLaVA-1.5 _{7B}	13.60	14.91	12.79	12.83	15.57	12.63	14.06	13.99	12.79	14.68	13.65	13.12	13.55	13.83	13.98	13.44	13.48
LLaVA-1.5 _{13B}	32.99	32.79	33.34	33.10	33.10	31.96	32.44	32.56	32.49	31.87	31.55	31.84	33.17	31.14	33.40	31.78	33.72
LLaVA-NeXT	48.16	11.92	48.84	38.82	15.90	39.13	47.07	20.72	47.02	48.20	36.67	39.69	47.36	39.60	46.99	28.13	47.55
InternVL _{8B}	54.58	20.14	32.54	36.95	42.10	56.72	31.53	31.09	41.65	50.17	32.77	43.66	34.82	34.90	50.41	50.78	41.66
Oryx	33.92	23.94	34.19	37.77	33.02	32.89	32.56	34.16	34.83	34.51	34.82	34.05	33.95	29.61	35.25	33.27	32.33
Qwen2VL _{7B}	49.07	46.93	46.81	48.75	48.04	47.64	48.45	46.80	49.24	47.95	47.19	49.58	48.83	41.27	49.72	47.07	47.90
Qwen2VL _{72B}	51.26	39.46	52.13	51.24	51.64	49.75	53.18	52.46	50.81	51.25	47.44	51.22	48.87	35.72	52.76	49.77	48.52
Dolphin	8.81	7.11	7.17	9.54	9.02	6.48	8.05	7.95	7.10	9.29	8.94	8.02	8.02	9.42	8.37	10.07	6.32
DriveLM	42.78	27.83	47.18	36.30	40.70	39.18	40.93	43.30	40.98	39.95	38.23	40.08	45.68	38.88	41.10	33.50	39.65

Table K. Detailed Accuracy score results of MCQs for the  Behavior task. “Clean” represents clean image inputs. “T.O.” represents text-only evaluation. The “Corrupt” settings range from weather conditions, external disturbances, sensor failures, motion blur, and transmission errors. The benchmarked VLMs include commercial, open-sourced, and driving specialist models, respectively.





Method	Clean	T.O.	Brightness	Dark	Snow	Fog	Rain	Lens Obstacle	Water Splash	Camera Crash	Frame Lost	Saturate	Motion Blur	Zoom Blur	Bit Error	Color Quant	H.265 Compression
 Human	66.67	-	40.00	46.67	66.67	53.33	73.33	53.33	53.33	60.00	40.00	53.33	66.67	33.33	53.33	73.33	53.33
GPT-4o	25.50	24.00	25.50	25.00	21.50	23.50	25.00	26.00	24.00	26.50	28.50	24.50	23.50	24.00	26.00	22.50	21.50
Phi3	26.50	30.00	29.50	29.50	28.00	29.50	28.50	27.00	30.00	32.50	31.00	28.50	29.50	23.50	27.50	31.50	27.50
Phi-3.5	36.50	40.00	37.00	36.00	37.00	38.50	38.00	36.50	35.50	37.00	39.00	39.00	37.50	36.50	39.00	36.50	36.00
LLaVA1.5 _{7B}	10.00	9.50	8.50	8.00	8.00	7.50	8.00	8.50	8.00	11.00	10.00	7.50	9.00	7.50	11.00	8.50	8.00
LLaVA1.5 _{13B}	32.50	33.00	33.00	33.00	32.50	32.50	32.00	32.50	32.50	32.00	32.50	32.50	33.00	31.00	34.00	31.00	33.50
LLaVA-NeXT	23.00	15.00	23.00	24.00	22.00	23.00	22.50	24.50	25.50	27.50	24.50	26.50	24.00	23.00	24.00	21.50	25.50
InternVL _{8B}	27.50	21.50	9.00	14.50	20.50	25.50	13.00	11.50	15.00	25.00	17.50	21.00	11.50	12.00	28.00	23.50	18.50
Oryx	21.00	21.00	21.50	21.50	21.50	20.50	21.50	21.50	21.50	21.50	22.00	21.50	22.00	21.00	22.50	21.50	21.50
Qwen2VL _{7B}	30.00	23.00	29.00	28.00	25.00	28.50	27.50	25.00	28.50	31.50	28.50	33.50	26.00	21.50	27.00	28.50	30.00
Qwen2VL _{72B}	23.00	36.50	25.50	24.50	25.50	22.00	29.50	26.00	22.50	27.00	25.00	26.00	22.50	22.00	28.50	23.50	23.50
Dolphin	0.50	3.50	1.50	0.00	1.00	0.00	0.00	1.00	1.00	1.50	2.50	0.50	1.00	1.00	1.00	0.50	0.50
DriveLM	44.00	25.50	48.50	37.00	41.50	40.00	42.50	43.50	41.00	41.00	39.50	40.50	46.50	40.00	43.00	35.00	41.50

Table L. Detailed GPT score results of **open-ended questions** for the  **Perception** task. “Clean” represents clean image inputs. “T.O.” represents text-only evaluation. The “Corrupt” settings range from weather conditions, external disturbances, sensor failures, motions, and transmission errors. The benchmarked VLMs include commercial, open-sourced, and driving specialist models, respectively.

Method	Clean	T.O.	Brightness	Dark	Snow	Fog	Rain	Lens Obstacle	Water Splash	Camera Crash	Frame Lost	Saturate	Motion Blur	Zoom Blur	Bit Error	Color Quant	H.265 Compression
GPT-4o	28.87	29.37	29.51	28.15	30.19	28.89	28.63	28.49	29.42	27.12	28.16	27.96	29.82	32.79	27.25	26.44	29.89
Phi-3	10.26	10.38	10.44	11.27	10.97	10.81	10.99	10.64	11.32	10.67	10.01	11.03	11.37	12.31	11.38	10.33	10.58
Phi-3.5	14.83	18.20	14.24	13.08	15.36	16.09	18.04	16.53	14.39	16.42	14.47	12.88	14.97	12.96	14.74	13.27	18.66
LLaVA-1.5 _{7B}	14.03	11.94	13.53	13.31	13.31	13.61	13.75	13.91	13.48	13.95	12.90	12.98	13.35	13.05	13.36	12.83	14.28
LLaVA-1.5 _{13B}	13.13	11.50	12.72	13.39	12.86	13.38	13.05	13.23	13.59	15.13	14.98	12.39	13.60	13.85	12.72	13.43	12.26
LLaVA-NeXT	15.33	23.53	15.49	14.95	16.61	15.62	16.05	15.66	15.66	15.19	16.04	15.16	16.06	17.92	15.27	15.10	15.86
InternVL _{8B}	18.12	14.73	20.05	18.82	23.14	14.28	27.75	23.43	26.67	11.88	17.11	26.26	22.02	27.25	11.54	29.57	29.77
Oryx	16.07	16.07	15.72	13.46	11.98	15.54	13.73	17.75	15.86	13.35	13.29	14.45	13.74	13.90	13.17	13.40	14.66
Qwen2VL _{7B}	15.33	32.56	16.75	14.95	15.16	15.66	15.02	15.66	15.48	17.31	17.65	14.52	15.29	15.20	15.17	14.63	17.89
Qwen2VL _{72B}	22.10	13.88	20.59	17.00	14.64	18.77	20.00	24.19	20.34	17.55	15.45	19.13	16.10	18.58	15.97	18.78	16.36
Dolphin	12.68	13.67	12.07	10.34	12.37	11.46	11.73	12.33	11.04	12.34	11.39	11.47	11.34	10.17	11.40	11.35	11.65
DriveLM	11.32	5.05	11.30	10.03	11.21	9.71	10.22	10.71	11.01	11.14	10.13	9.38	10.97	9.03	11.39	10.50	10.73

Table M. Detailed GPT score results of the **open-ended questions** for  **Prediction**. “Clean” represents clean image inputs. “T.O.” represents text-only evaluation. The “Corrupt” settings range from weather conditions, external disturbances, sensor failures, motions, and transmission errors. The benchmarked VLMs include commercial, open-sourced, and driving specialist models, respectively.

Method	Clean	T.O.	Brightness	Dark	Snow	Fog	Rain	Lens Obstacle	Water Splash	Camera Crash	Frame Lost	Saturate	Motion Blur	Zoom Blur	Bit Error	Color Quant	H.265 Compression
GPT-4o	51.30	49.05	52.15	50.28	47.97	49.66	51.39	50.49	53.30	46.62	45.95	49.18	51.90	49.75	48.59	47.56	54.36
Phi-3	40.11	22.61	45.21	26.28	35.54	44.05	34.87	40.90	33.56	44.26	38.74	41.39	35.08	36.46	32.25	32.82	37.57
Phi-3.5	45.13	4.92	46.57	9.67	48.02	45.08	52.16	42.75	47.02	47.18	24.66	23.26	49.16	28.38	41.39	18.02	49.87
LLaVA-1.5 _{7B}	22.02	14.64	24.79	20.95	15.97	15.30	18.98	16.28	24.16	6.11	13.90	20.30	25.20	10.56	11.10	15.61	23.92
LLaVA-1.5 _{13B}	36.98	23.98	36.00	35.59	40.51	39.23	38.90	38.11	38.92	36.25	36.54	37.57	38.10	39.74	34.28	37.16	36.31
LLaVA-NeXT	35.07	28.36	37.15	35.31	37.59	37.62	35.44	37.00	35.87	36.25	30.10	40.56	34.66	39.36	31.74	34.07	35.66
InternVL _{8B}	45.52	48.89	45.73	40.71	35.75	38.43	33.18	38.69	40.71	45.00	45.75	30.03	39.52	36.55	40.12	28.40	30.31
Oryx	48.13	12.77	49.52	44.33	47.67	45.77	47.20	45.90	50.30	42.18	44.59	42.26	48.21	52.54	45.56	43.77	49.61
Qwen2VL _{7B}	37.89	37.77	40.82	35.90	38.92	44.15	40.15	41.89	41.57	36.61	35.87	41.25	40.89	39.23	36.69	40.52	38.84
Qwen2VL _{72B}	49.35	5.57	43.89	43.25	43.74	44.49	45.57	41.61	47.46	40.89	34.75	39.38	48.21	51.15	40.49	40.82	46.67
Dolphin	32.66	39.98	29.85	32.31	24.64	29.92	31.38	33.41	31.79	29.05	30.93	30.49	31.59	26.38	30.13	25.64	30.62
DriveLM	44.33	4.70	46.82	43.90	42.33	35.84	44.13	44.00	42.59	46.25	33.56	29.69	42.15	19.00	38.20	44.33	42.87

Table N. Detailed ROUGE-L score results of **open-ended questions** for the  **Predicion** task. “Clean” represents clean image inputs. “T.O.” represents text-only evaluation. The “Corrupt” settings range from weather conditions, external disturbances, sensor failures, motion blur, and transmission errors. The benchmarked VLMs include commercial, open-sourced, and driving specialist models, respectively.



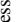















Method	 Clean	 T.O.	 Brightness	 Dark	 Snow	 Fog	 Rain	 Lens Obstacle	 Water Splash	 Camera Crash	 Frame Lost	 Saturate	 Motion Blur	 Zoom Blur	 Bit Error	 Color Quant	 H.265 Compression
GPT-4o	19.74	18.58	19.77	19.58	19.67	19.71	19.69	20.22	19.94	19.89	19.66	19.70	19.39	19.48	19.47	19.90	19.67
Phi3	17.76	14.71	17.21	21.55	25.81	15.28	23.77	18.56	19.27	17.59	17.42	16.90	16.75	14.36	25.23	26.88	22.55
Phi-3.5	19.36	18.76	18.37	5.78	17.28	15.74	18.24	18.38	17.98	17.10	12.63	10.84	17.85	12.46	17.12	8.50	17.03
LLaVA1.5 _{7B}	21.18	23.21	21.75	22.64	22.21	22.17	22.19	19.61	22.05	22.53	22.36	22.65	21.71	20.28	22.46	22.57	21.17
LLaVA1.5 _{13B}	24.12	24.79	24.03	24.04	24.09	24.37	24.14	24.18	24.35	23.91	23.91	23.84	24.04	24.59	23.93	24.31	24.12
qwen2-7b	25.49	24.15	25.64	24.80	25.21	24.82	24.86	25.17	25.64	25.62	25.99	24.74	25.11	23.73	25.09	25.31	25.35
Qwen2VL _{72B}	23.42	16.10	20.05	18.46	17.97	19.08	19.59	18.68	18.51	22.05	21.17	17.97	18.98	19.03	18.89	19.35	18.13
LLaVA1.6-7b	16.09	13.93	16.29	16.99	17.02	17.14	16.54	16.27	17.08	16.25	16.59	18.05	16.69	16.93	16.58	17.33	16.46
InternVL _{8B}	13.92	13.25	20.57	15.14	14.67	13.64	15.15	16.39	15.45	13.92	16.12	14.47	16.29	14.79	14.02	14.26	14.07
Oryx	21.03	15.43	18.79	17.45	16.99	18.04	18.19	17.81	17.68	20.38	19.56	16.54	18.21	17.84	18.20	17.92	17.03
Qwen2VL _{7B}	25.49	23.15	25.64	24.80	25.21	24.82	24.86	25.17	25.64	25.62	25.99	24.74	25.11	23.73	25.09	25.31	25.35
Qwen2VL _{72B}	23.42	16.10	20.05	18.46	17.97	19.08	19.59	18.68	18.51	22.05	21.17	17.97	18.98	19.03	18.89	19.35	18.13
Dolphin	25.18	23.64	24.71	24.84	25.95	24.98	26.19	25.09	25.18	24.52	24.75	24.48	24.99	24.63	24.87	24.34	24.95
DriveLM	40.00	23.00	37.00	45.00	35.00	46.50	45.50	40.50	37.50	39.50	29.50	34.00	42.00	36.50	36.00	30.50	47.00

Table O. Detailed GPT score results of the **open-ended questions** for the  **Planning**. “Clean” represents clean image inputs. “T.O.” represents text-only evaluation. The “Corrupt” settings range from weather conditions, external disturbances, sensor failures, motions, and transmission errors. The benchmarked VLMs include commercial, open-sourced, and driving specialist models, respectively.



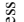















Method	 Clean	 T.O.	 Brightness	 Dark	 Snow	 Fog	 Rain	 Lens Obstacle	 Water Splash	 Camera Crash	 Frame Lost	 Saturate	 Motion Blur	 Zoom Blur	 Bit Error	 Color Quant	 H.265 Compression
GPT-4o	75.75	73.21	77.56	74.33	73.00	76.58	76.53	75.14	74.65	84.08	74.54	74.84	77.34	73.09	74.30	76.36	76.82
Phi-3	60.03	46.88	61.48	59.14	64.59	64.81	63.47	61.83	62.31	58.32	54.63	63.04	61.61	58.93	58.11	64.04	63.40
Phi-3.5	31.91	46.30	30.57	28.64	32.88	28.38	34.70	30.88	29.53	23.16	24.53	25.47	32.73	23.00	28.24	23.42	29.31
LLaVA-1.5 _{7B}	29.15	32.45	31.52	31.49	32.58	32.42	31.00	29.81	31.72	33.70	35.95	29.93	31.20	30.65	30.05	30.00	30.61
LLaVA-1.5 _{13B}	34.26	38.85	33.63	34.36	35.06	39.43	35.18	33.01	34.80	37.61	39.32	32.77	34.74	32.99	33.75	33.33	33.67
LLaVA-NeXT	45.27	27.58	45.64	44.54	43.55	44.17	45.08	44.62	44.21	45.69	44.51	41.17	45.30	44.57	43.67	43.57	45.13
InternVL _{8B}	53.27	34.56	54.70	60.02	63.89	53.69	60.64	54.31	56.68	52.14	46.12	54.44	55.94	54.93	49.08	57.02	55.13
Oryx	53.57	48.26	55.46	58.60	58.91	57.93	55.04	57.01	58.35	52.28	51.05	55.76	55.88	53.63	52.94	56.00	57.57
Qwen2VL _{7B}	57.04	41.66	54.19	58.37	52.70	58.18	55.85	53.98	55.64	54.71	53.18	51.93	55.22	51.73	56.79	53.16	56.04
Qwen2VL _{72B}	61.30	53.35	62.01	65.31	66.69	68.15	67.83	65.67	65.26	57.42	56.34	62.06	64.20	61.66	58.23	59.86	65.31
Dolphin	52.91	60.98	51.85	55.39	53.09	54.78	53.92	51.79	53.57	55.73	57.81	55.78	51.42	50.95	53.35	54.89	52.17
DriveLM	68.71	65.24	67.25	67.52	65.72	63.08	69.60	69.04	67.97	67.85	66.47	66.25	67.93	70.17	68.46	68.30	68.59

Table P. Detailed ROUGE-L score results of **open-ended questions** for the  **Planning** task. “Clean” represents clean image inputs. “T.O.” represents text-only evaluation. The “Corrupt” settings range from weather conditions, external disturbances, sensor failures, motion blur, and transmission errors. The benchmarked VLMs include commercial, open-sourced, and driving specialist models, respectively.

Method	● Clean	● T.O.	● Brightness	● Dark	● Snow	● Fog	● Rain	● Lens Obstacle	● Water Splash	● Camera Crash	● Frame Lost	● Saturate	● Motion Blur	● Zoom Blur	● Bit Error	● Color Quant	● H.265 Compression
GPT-4o	6.54	6.42	6.54	6.34	6.54	6.54	6.52	6.40	6.52	6.47	6.35	6.39	6.39	6.40	6.35	6.37	6.67
Phi3	9.39	10.05	9.56	9.26	9.51	9.62	9.47	9.65	9.60	10.01	10.01	9.90	9.40	9.42	9.07	10.02	9.36
Phi-3.5	6.19	8.27	6.52	7.97	6.30	6.10	6.12	6.57	6.59	7.18	7.59	5.98	6.17	7.87	6.64	6.04	6.37
LLaVA1.5 _{7B}	8.76	9.98	8.93	8.99	9.18	8.96	8.76	9.00	9.04	9.06	9.24	8.88	8.94	9.27	9.25	9.16	8.85
LLaVA1.5 _{13B}	8.06	8.68	8.14	8.09	8.10	8.06	8.10	8.03	7.99	8.68	8.64	8.43	8.10	8.04	8.39	8.33	8.24
LLaVA-NeXT	5.69	7.75	5.72	5.66	5.50	5.49	5.60	5.67	5.60	5.62	5.88	5.59	5.64	5.45	5.64	5.64	5.58
InternVL _{8B}	4.06	9.45	10.37	11.65	12.07	4.05	10.81	10.97	9.88	4.16	9.29	11.02	10.04	9.23	4.03	10.24	9.65
Oryx	21.03	15.43	18.79	17.45	16.99	18.04	18.19	17.81	17.68	20.38	19.56	16.54	18.21	17.84	18.20	17.92	17.03
Qwen2VL _{7B}	9.61	8.33	9.31	8.70	8.14	8.33	8.69	8.32	8.82	9.15	8.81	8.69	8.64	9.05	9.13	8.77	7.92
Qwen2VL _{72B}	12.26	13.17	12.13	11.60	11.45	10.97	11.62	11.50	11.67	12.01	11.43	12.13	11.85	11.55	12.06	11.84	11.93
Dolphin	12.90	14.82	12.90	13.01	12.61	12.82	12.45	12.81	12.74	13.56	13.91	12.82	12.86	12.89	13.22	13.29	12.69
DriveLM	53.12	46.83	51.55	51.26	49.26	46.74	52.02	53.38	52.77	52.44	50.91	48.66	53.28	52.20	49.88	52.79	51.43

References

- [1] Marah Abidin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 Technical Report: A Highly Capable Language Model Locally on Your Phone. *arXiv preprint arXiv:2404.14219*, 2024.
- [2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 Technical Report. *arXiv preprint arXiv:2303.08774*, 2023.
- [3] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen Technical Report. *arXiv preprint arXiv:2309.16609*, 2023.
- [4] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuScenes: A Multimodal Dataset for Autonomous Driving. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11621–11631, 2020.
- [5] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. InternVL: Scaling Up Vision Foundation Models and Aligning for Generic Visual-Linguistic Tasks. *arXiv preprint arXiv:2312.14238*, 2023.
- [6] Yuhao Dong, Zuyan Liu, Hai-Long Sun, Jingkan Yang, Winston Hu, Yongming Rao, and Ziwei Liu. Insight-V: Exploring Long-Chain Visual Reasoning with Multimodal Large Language Models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9062–9072, 2025.
- [7] Jinkyu Kim, Anna Rohrbach, Trevor Darrell, John Canny, and Zeynep Akata. Textual Explanations for Self-Driving Vehicles. In *European Conference on Computer Vision*, pages 563–578. Springer, 2018.
- [8] Lingdong Kong, Youquan Liu, Xin Li, Runnan Chen, Wenwei Zhang, Jiawei Ren, Liang Pan, Kai Chen, and Ziwei Liu. Robo3D: Towards Robust and Reliable 3D Perception Against Corruptions. In *IEEE/CVF International Conference on Computer Vision*, pages 19994–20006, 2023.
- [9] Lingdong Kong, Shaoyuan Xie, Hanjiang Hu, Lai Xing Ng, Benoit Cottureau, and Wei Tsang Ooi. RoboDepth: Robust Out-of-Distribution Depth Estimation Under Corruptions. *Advances in Neural Information Processing Systems*, 36:21298–21342, 2023.
- [10] Boyun Li, Xiao Liu, Peng Hu, Zhongqin Wu, Jiancheng Lv, and Xi Peng. All-In-One Image Restoration for Unknown Corruption. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17452–17462, 2022.
- [11] 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 *European Conference on Computer Vision*, pages 1–18. Springer, 2022.
- [12] Zhiqi Li, Zhiding Yu, Shiyi Lan, Jiahao Li, Jan Kautz, Tong Lu, and Jose M Alvarez. Is Ego Status All You Need for Open-Loop End-to-End Autonomous Driving? In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14864–14873, 2024.
- [13] Chin-Yew Lin. ROUGE: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out*, pages 74–81, 2004.
- [14] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved Baselines with Visual Instruction Tuning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306, 2024.
- [15] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. LLaVA-NeXT: Improved Reasoning, OCR, and World Knowledge, 2024.
- [16] Yingzi Ma, Yulong Cao, Jiachen Sun, Marco Pavone, and Chaowei Xiao. Dolphins: Multimodal Language Model for Driving. In *European Conference on Computer Vision*, pages 403–420. Springer, 2024.
- [17] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. BLEU: A Method for Automatic Evaluation of Machine Translation. In *Annual Meeting of the Association for Computational Linguistics*, pages 311–318, 2002.
- [18] 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. In *AAAI Conference on Artificial Intelligence*, pages 4542–4550, 2024.
- [19] Chonghao Sima, Katrin Renz, Kashyap Chitta, Li Chen, Hanxue Zhang, Chengen Xie, Jens Beißwenger, Ping Luo, Andreas Geiger, and Hongyang Li. DriveLM: Driving with Graph Visual Question Answering. In *European Conference on Computer Vision*, pages 256–274. Springer, 2024.
- [20] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-VL: Enhancing Vision-Language Model’s Perception of the World at Any Resolution. *arXiv preprint arXiv:2409.12191*, 2024.
- [21] Yue Wang, Vitor Campagnolo Guizilini, Tianyuan Zhang, Yilun Wang, Hang Zhao, and Justin Solomon. DETR3D: 3D Object Detection from Multi-View Images via 3D-to-2D Queries. In *Conference on Robot Learning*, pages 180–191. PMLR, 2022.
- [22] Shaoyuan Xie, Lingdong Kong, Wenwei Zhang, Jiawei Ren, Liang Pan, Kai Chen, and Ziwei Liu. Benchmarking and Improving Bird’s Eye View Perception Robustness in Autonomous Driving. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 47(5):3878–3894, 2025.

- [23] Yiran Xu, Xiaoyin Yang, Lihang Gong, Hsuan-Chu Lin, Tz-Ying Wu, Yunsheng Li, and Nuno Vasconcelos. Explainable Object-Induced Action Decision for Autonomous Vehicles. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9523–9532, 2020.
- [24] Chenyu Yi, Siyuan Yang, Haoliang Li, Yap-peng Tan, and Alex Kot. Benchmarking the Robustness of Spatial-Temporal Models Against Corruptions. In *Advances in Neural Information Processing Systems*, 2021.