IntelliCockpitBench: A Comprehensive Benchmark to Evaluate VLMs for Intelligent Cockpit

Anonymous ACL submission

Abstract

The integration of sophisticated Vision-Language Models (VLMs) in vehicular systems is revolutionizing vehicle interaction and safety, performing tasks such as Visual Question Answering (VQA). However, a critical gap persists due to the lack of a comprehensive benchmark for multimodal VQA models in vehicular scenarios. To address this, we propose IntelliCockpitBench, a benchmark that encompasses diverse automotive scenarios. It includes images from front, side, and rear cameras, various road types, weather conditions, and interior views, integrating data from both moving and stationary states. Notably, all images and queries in the benchmark are verified for high levels of authenticity, ensuring the data accurately reflects real-world conditions. A sophisticated scoring methodology combining human and model-generated assessments enhances reliability and consistency. Our contributions include a diverse and authentic dataset for automotive VQA and a robust evaluation metric aligning human and machine assessments. All code and data can be found at https://anonymous.4open. science/r/IntelliCockpitBench-2F2E/.

1 Introduction

002

016

017

021

028

042

In recent years, with the advancement of Visual Language Models (VLMs) (Liu et al., 2023; Bai et al., 2023; Wang et al., 2023a), intelligent cockpit technology has made significant progress, becoming an important interface for the next generation of human-computer interaction. Subsequently, benchmarks like DriveBench (Xie et al., 2025) and NuScenes-QA (Qian et al., 2024) have been proposed to evaluate the visual question-answering (VQA) capabilities in autonomous driving scenarios. Even so, these benchmarks remain primarily focused on decision-making scenarios such as autonomous driving and do not adequately consider non-decision-making scenarios aimed at enhancing



Figure 1: The relationship between model size and score in English queries across various VLMs on IntelliCockpitBench. Notable models such as GPT-40 (Hurst et al., 2024) and Gemini-2.0-Flash (Team et al., 2023) are distinguished by their superior performance despite larger sizes. The dotted line represents an estimated trend indicating the positive correlation between model size and performance.

user experience and interaction. This has significant limitations in the research field. **Limitation 1:** the lack of comprehensive benchmarks specifically designed to evaluate the performance of VLMs in non-decision-making scenarios within intelligent cockpits. **Limitation 2:** existing GPT-based (Hurst et al., 2024) automatic evaluation methods typically rely on uniform assessment standards, which overlook the specific nature and requirements of different queries. This further emphasizes the necessity of developing evaluation benchmarks tailored to different queries types.

То address these limitations, we proa comprehensive benchmark pose named IntelliCockpitBench to evaluate VLMs for intelligent cockpits. This benchmark includes a diverse collection of images captured from front, side, and rear cameras, encompassing various road types and weather conditions to provide a comprehensive external perspective. Additionally, 043

IntelliCockpitBench features interior images 063 to reflect the complexity of the in-vehicle envi-064 ronment. The curated dataset also integrates data 065 from both moving and stopping vehicle states, ensuring a thorough representation of real-world scenarios. Taking into account the scenarios of visual information augmented, we have also implemented data augmentation techniques to ensure the robustness of IntelliCockpitBench in various unexpected situations. All queries in our dataset are collected through driver surveys and generalized using GPT-40 (Hurst et al., 2024) to ensure their authenticity and diversity. Note that all included images and queries are verified for high levels of authenticity and have undergone human 077 review, which ensures that the data accurately reflects real-world driving scenarios.

> Furthermore, we design three key LLM-as-ajudge methods including Chain-of-Thought Reasoning, Multi-dimensional Variance Analysis, and Rule-Calibrated Referencing. This evaluation method not only defines different evaluation metrics for various queries but also assigns importance scores to these metrics. Additionally, it utilizes Chain-of-Thought to generate explanations and final ratings, ensuring both high reliability and interpretability. As shown in Figure 1 and Table 2, we evaluate 15 VLMs and our experiments reveal that current VLMs perform poorly when confronted with augmented visual images and queries requiring deep reasoning. Therefore, it is essential to enhance VLMs' capabilities in accurate visual localization and multi-step reasoning queries.

> > Overall, our key contributions are as follows:

• We create a comprehensive benchmark, IntelliCockpitBench, to evaluate the capabilities of VLMs for the intelligent cockpit, featuring 5 intelliCockpit query types, 38 driving scenarios, 10+ question formats, 16, 154 queries, over 7, 622 images, and 20 evaluation metrics.

100

101

102

103

106

107

108

109

110

111

112

- We propose 3 innovative LLM-as-a-judge evaluation methods including Chain-of-Thought Reasoning, Multi-dimensional Variance Analysis, and Rule-Calibrated Referencing to enhance the reliability and interpretability of evaluation.
- We evaluate 15 open-source and closed-source VLMs and find that all models perform poorly on the IntelliCockpitBench, especially with augmented visuals and complex reasoning queries, highlighting the need for improved visual local-

ization and reasoning in VLMs.

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

2 Related Work

2.1 Vision-Language Models

The success of Large Language Models (LLMs) (Touvron et al., 2023; Team et al., 2023; GLM et al., 2024) has significantly advanced VLMs. BLIP (Liu et al., 2024a) employs GPT-4 to generate instruction-following data for vision-language tuning, and its learning paradigm and instruction-tuning corpus have been widely adopted in subsequent research (Chen et al., 2025, 2024a). Over the past year, numerous open-source VLMs have gained recognition, including the LLaVA series (Liu et al., 2024a,c,b), MiniGPT-4 (Zhu et al., 2023), VisionLLM (Wang et al., 2024b), Qwen-VL (Bai et al., 2023; Wang et al., 2024a), CogVLM (Wang et al., 2023a), Intern-VL (Chen et al., 2024b; Dong et al., 2024), and others (Chen et al., 2023; Peng et al., 2023; Wang et al., 2023b). Although these models are generally aimed at standard VQA and various broad applications, there is still a clear gap in their use within smart cockpit settings. Regarding this, we propose IntelliCockpitBench encompassing **5 query types** and **4 scenarios** in Figure 2.

2.2 Multimodal Datasets

Recently, many datasets for intelligent cockpits (e.g., driveLM (Sima et al., 2023) and NuScenes-QA (Qian et al., 2024)) have been constructed based on widely used driving datasets, such as nuScenes (Goyal et al., 2017) and BDD (Yu et al., 2020). However, these datasets suffer from issues like data imbalance and overly simplistic answer designs. Moreover, DriveBench (Xie et al., 2025) has been proposed to evaluate the reliability and visual grounding of VLMs in autonomous driving systems. SuperCLUE-o (Xu et al., 2020) evaluates models from the perspectives of answer quality and response latency, but it lacks sufficient scene diversity to fully cover the various situations encountered during driving.

Furthermore, in the aforementioned methods, when employing LLMs as evaluation tools, they either directly use scoring methods or adopt relatively coarse-grained rules (such as (Xu et al., 2020)) for evaluation. Although some progress has been made with these methods, current datasets and evaluation paradigm may still be insufficient to comprehensively capture the complexity of real-



(a) Distribution of query types.

(b) Distribution of driving scenarios.

Figure 2: A comprehensive taxonomy of query types and scenarios in VLMs within IntelliCockpitBench. "WK." denotes World Knowledge. "Geo. Env." denotes Geospatial environmental.

world driving, especially in terms of performance in non-decision-making interactive tasks within intelligent cockpits, further contributing to the development of IntelliCockpitBench.

3 IntelliCockpitBench

162

163

164

166

167

168

169

171

172

173

174

175

176

177

179

181

In this section, we introduce an overview of the data composition, the dataset construction, and the evaluation paradigm of IntelliCockpitBench.

3.1 Dataset Composition

To ensure the authenticity and diversity of the curated dataset, we first collect images and queries that are sourced from real-world driving scenarios. We then propose a comprehensive taxonomy for VLMs' driving queries based on real-driver queries to conduct a systematic evaluation. As illustrated in Figure 2, from simple descriptions to complex reasoning, these queries are divided into **5 dimensions**: description, recognition, world knowledge Q&A, reasoning, and others. The detailed explanation of each query is provided in Appendix A.1.

In addition, to thoroughly evaluate the adaptability and robustness of VLMs given the complexity and variability of real-world driving scenarios, we categorize and summarize these scenarios into 4 categories including weather conditions, road types, driving status, and shooting angles), **38 meta-categories**, and a total of **7,622 images**. We provide a detailed taxonomy and definition for these four driving scenarios in Appendix A.2.

3.2 Dataset Construction

This subsection delineates the construction process of the dataset, encompassing three primary stages: image generation, query generation, and answer generation, as illustrated in Figure 3. 191

192

193

194

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

3.2.1 Image Generation

Overall, our image data sources can be classified into two principal categories. The first category encompasses partial data collection from publicly available datasets, including NUSCENES, the Yawning Detection Dataset, and the drive&act dataset. The second category, representing the primary source of our dataset, comprises over 100 meticulously selected driving videos obtained from video-sharing platforms. Download data for academic research only. These videos are rigorously chosen based on a carefully defined taxonomy of driving scenarios (refer to Appendix A.2). Subsequently, we systematically sample frames from the collected videos at consistent intervals of every 15 second, culminating in an extensive dataset consisting of 7,622 images. All images have undergone a de-identification process to mask faces and license plate numbers. Considering the substantial impact that image quality has on the performance of VLMs, our dataset intentionally includes images of various resolutions.

In addition to designing and screening images under normal driving conditions, we consider scenarios where visual information degrades, such as weather-induced image quality degradation (rain or fog), changes in lighting (overexposure), and cam-



Figure 3: Architecture of the proposed IntelliCockpitBench. Dataset construction involves three steps: 1) **Image Generation**, creating driving scenario images using video generation techniques; 2) **Query Generation**, generating multiple types of intellicockpit queries using LLMs and real-driver queries; 3) **Answer Generation**, producing corresponding answers based on different intellicockpit queries; The last module is LLM Judgement, scoring multiple dimensions of the answers using evaluation paradigms based on chain-of-thought reasoning, multi-dimensional variance analysis, and rule-based calibration, ultimately providing a comprehensive score.

era malfunctions (image distortion, obstruction, or misalignment). A total of **190 images** are collected to validate the robustness and reliability of VLMs under various unforeseen circumstances.

3.2.2 Query Generation

235

236

240

Most existing VQA benchmarks are limited in the diversity of questioning types (Xie et al., 2025; Xu et al., 2017), failing to fully represent the wide spectrum of human conversations. In contrast, the questioning set in IntelliCockpitBench has been carefully curated to include a broad range of categories. Figure 8 illustrates the distributions of the questioning type. Questioning types include 'what', 'who', 'how', 'when', and 'where'. We also expand scopes of type to include interrogatives like 'why', 'which', 'is/are', and 'does/do'. This expansion enhances the diversity and better reflects the natural style of human dialogues.

Real-driver Query Generation. Due to the lack
of authenticity in queries generated directly based
on images using GPT-40 (Hurst et al., 2024), we
obtain real intelligent cockpit queries by recruiting 100 drivers. Each driver provides 100 queries
they might encounter in driving scenarios related

to visual information, resulting in a total of **10,000 real-driver queries**. To ensure the diversity of queries, we use GPT-40 to generalize them. Specifically, we first leverage the classification results of the questioning type and then perform random sampling from the collected real dataset as a few-shot input to generate new queries. The detailed query generation prompt is in Appendix A.6.

The content generated by the Human Check. GPT-40 (Hurst et al., 2024) is then subjected to manual inspection to ensure that both the image and the query are of high quality and accurately represent realistic scenarios, which are conducted in two stages. Initially, we ask annotators to evaluate whether the generated queries meet the five specific criteria listed in Table 3. Queries that do not meet these criteria will be manually modified, and if modification is not feasible, they will be discarded. The establishment and implementation of refusal strategies for VLMs are crucial, as they can effectively prevent misinformation, protect user privacy, and ensure that the generated content adheres to ethical and legal standards. Subsequently, for the queries that pass the initial evaluation, annota-

264

265

266

270

247

248

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

369

321

322

323

324

325

tors further determine whether the query needs to
be refused an answer, as outlined in the rejection
strategy presented in Appendix A.4. We provide
details of human checks in Appendix A.3.

3.2.3 Answer Generation

275

276

277

278

281

289

294

301

310

311

313

314

Following the generation of high-quality images and realistic queries, the next step involves constructing accurate answers. Specifically, the previously generated images and queries, along with the VLMs' queries categorization system, are input into GPT-40. This process enables the model to produce a clear answer, a concise rationale, and the corresponding query labels. We provide an answer generation prompt in Appendix A.6. The final outputs are then manually verified to ensure their authenticity and accuracy. First, we instruct annotators to confirm that the answer correctly addresses the query based on the image. Next, they ensure that the classification of both the query and the image aligns with the established categorization system. If any inaccuracies are identified, the annotators manually revise the answers. Note that all VQA pairs generated in IntelliCockpitBench undergo a rigorous cross-validation process (see Appendix A.3) to ensure their accuracy and adherence to the classification system.

3.3 Evaluation Paradigm

To effectively evaluate the quality of VLMs' responses, IntelliCockpitBench utilizes GPT-4o-Mini-2024-07-18 (Hurst et al., 2024) as the primary evaluator to analyze and grade responses in accordance with established practices (Zheng et al., 2023). Nonetheless, a significant design space in VQA remains unexplored, particularly regarding prompting strategies, score calibration, critique explainability, and evaluation dimensions. To address these gaps, we develop a rule-based evaluation methodology using Language Models as Judges (LLM-as-a-judge) that incorporates three principal approaches: Chain-of-Thought Reasoning, Multi-dimensional Variance Analysis, and Rule-Calibrated Referencing. Detailed prompts for rulebased evaluation are provided in Appendix A.6, and an illustrative example is shown in Figure 4.

Chain-of-Thought Reasoning. When leveraging LLM-as-a-Judge, IntelliCockpitBench employs point-wise grading to assess the quality of responses. The inputs include the image, the query associated with the image, the model's response, and a human-curated reference answer. The output consists of a multi-dimensional analytical explanation alongside a final rating on a 1 to 10 scale.

Multi-dimensional Variance Analysis. Given the diverse nature and characteristics of different queries, applying a uniform standard to all responses is inappropriate. To address this, we propose a multi-dimensional scoring approach that tailors evaluation criteria to the specific query type, providing a more detailed and structured analysis. Specifically, we define distinct evaluation dimensions and importance scores tailored to each query type. For example, in the case of descriptive queries, factuality should be prioritized, with completeness considered secondary. Consequently, the importance score for factuality is higher than that for completeness. We provide detailed definitions and settings of dimensions in Appendix A.6.

Rule-Calibrated Referencing. We provide a high-quality reference answer, which is primarily generated by GPT-40 and modified by human annotators to ensure its correctness and improve its quality. To guide the evaluator in comparing the answer with the reference and generating more controllable scores, we provide detailed grading rules that explain the relationship between score intervals and the quality of the answer compared to the reference. Additionally, we established a reference answer with a score of 8 as a benchmark for evaluation within a maximum score of 10.

4 Experiment

In this section, we conduct extensive benchmark experiments and analyses in IntelliCockpitBench, providing detailed discussions that lead to our observations and conclusions step by step.

4.1 Consistency Evaluation

То validate the alignment of the IntelliCockpitBench evaluation paradigm with human judgment, we conduct extensive human evaluations on selected queries. Specifically, we use GPT-4o-Mini-2024-07-18 (Hurst et al., 2024) as our scoring model due to its superior accuracy and consistency in natural language processing tasks. Evaluators were instructed to analyze the model's answers and provide scores based on predefined dimensions in Appendix A.6.

To align the consistency between the scores generated by GPT-4o-Mini with those labeled by humans, we assess consistency using the following three metrics: **Sample-level Pearson Correlation:**



Figure 4: An exemplar scoring process of IntelliCockpitBench on vehicle model recognition category.

Table 1: Comparison on human agreement between different judging methods on sampled IntelliCockpitBench, rated by GPT-40. The best performance is shown in **bold**.

Metric	Method	Overall	Description	Recognition	World Knowledge Q&A	Reasoning	Others
Sample-level Pearson	ours	0.80	0.92	0.78	0.67	0.82	0.96
System-level Pearson	general ours	0.64 0.93	0.53 0.93	0.59 0.90	0.63 0.95	0.71 0.94	0.50 0.92
Pairwise Agreement (w/o tie)	general ours	0.75 0.93	0.65 0.97	0.75 0.91	0.69 0.92	0.79 0.95	0.69 0.97

Since each query defines different evaluation dimensions and human judges also score each dimension, we first calculate the Pearson correlation coefficient for each sample and then compute the mean as the sample-level correlation. **Systemlevel Pearson Correlation:** This metric assesses the correlation at the system level by calculating the Pearson coefficient between the average scores at the sample-level given by human judges and model judges to the LLM. **Pairwise Agreement** (w/o ties): For each response, scores from human judges and model judges are converted into pairwise comparisons, with ties excluded.

370

372

373

376

377

382

390

We also compare a modified version of the evaluation prompts used in MT-Bench (Zheng et al., 2023) as a general evaluation with our rule-based calibration evaluation method. The prompt for general evaluation is in appendix A.6. As presented in Table 1, results show that our pointwise multi-dimensional rules-calibrated LLM-asa-judge method performs best, particularly on the Sample-level Pearson metric and the Pairwise Agreement (w/o tie) metric, thereby substantiating the excellent agreement with human judges. The reasons are as follows: 1) The nature and characteristics of the driving questions in VLMS vary, making it inappropriate to apply a unified evaluation standard to all questions. 2) Our method integrates the chain-of-thought reasoning approach to generate explanations and final scores, ensuring high reliability and interpretability. Furthermore, We plot the cumulative distribution of the human judge, general judge, and rule-calibrated judge in Figure 9 to show that the rule-calibration judge has a narrower gap to human evaluation's cumulative distribution.

4.2 IntelliCockpitBench Evaluation

Based on the validity of scoring and the comprehensive capabilities of IntelliCockpitBench, we systematically assess a diverse set of VLMs.

Result Analysis of Closed Models. As shown in

Table 2: Performance evaluation of various VLMs on the IntelliCockpitBench for different English and Chinese VQA intelliCockpit question types. "Des." denotes Description, "Rec." denotes Recognition, "Wk-QA" denotes World Knowledge Q&A, "Rea." denotes Reasoning. <u>Underline</u> indicates the best results within open-source and closed-source categories, while **bold** signifies the best results among all open-source and closed-source options.

Model	Size	Type	GPU Usage	Driving Questions (EN)						Driving Questions (CH)						
		(MiB)	Overall	Des.	Rec.	Wk-QA	Rea.	Others	Overall	Des.	Rec.	Wk-QA	Rea.	Others		
DeepSeek-VL-base	1.3B	open	5,284	3.47	3.96	3.20	4.10	3.47	3.08	2.50	2.48	2.20	3.12	2.59	3.01	-
MiniCPM-V-2.0	2B	open	9,098	4.02	3.96	4.13	4.61	3.81	4.02	4.38	5.33	4.47	5.04	4.03	3.81	
GLM-4V	2B	open	4,566	4.34	4.80	4.51	4.96	4.02	4.30	4.78	5.74	4.73	5.51	4.52	5.15	
Qwen2-VL	2B	open	28,300	4.63	4.78	4.85	5.25	4.31	4.52	4.98	6.03	5.19	5.69	4.51	5.44	
Megrez	3B	open	10,854	4.06	3.59	4.03	4.78	4.00	3.67	5.09	5.97	5.02	5.85	4.84	5.53	
GLM-4V	5B	open	10,152	4.51	5.19	4.63	5.18	4.17	4.43	4.85	5.95	4.62	5.66	4.68	5.49	
InstructBLIP	7B	open	20,456	3.83	4.08	3.46	4.44	3.94	2.96	2.33	4.05	2.17	2.72	2.13	2.04	0
Qwen2-VL	7B	open	39,800	<u>5.17</u>	<u>5.85</u>	<u>5.21</u>	<u>6.11</u>	4.83	<u>5.15</u>	5.45	6.31	<u>5.64</u>	<u>6.44</u>	4.95	<u>5.83</u>	
LLaVA-1.5	7B	open	16,024	4.09	4.61	3.52	5.01	4.25	3.76	3.74	4.26	3.29	4.60	3.81	4.31	
InternVL-2.5	8B	open	24,558	5.09	5.83	5.02	5.96	<u>4.85</u>	4.80	<u>5.46</u>	<u>6.74</u>	5.43	6.39	<u>5.09</u>	5.67	
GLM-4V	9B	open	28,578	4.85	5.61	4.89	5.78	4.52	4.43	5.23	5.87	5.33	6.07	4.87	5.62	
LLaVA-1.5	13B	open	28,822	4.24	4.67	3.73	5.26	4.33	3.88	3.75	4.61	3.43	4.66	3.66	4.12	_
GLM-4V-plus	-	closed	-	5.32	6.05	5.28	6.33	5.01	5.42	5.61	6.40	5.55	6.60	5.31	6.12	
GPT-40	-	closed	-	<u>5.81</u>	<u>6.36</u>	<u>5.91</u>	<u>6.77</u>	<u>5.45</u>	<u>5.70</u>	<u>6.26</u>	<u>7.37</u>	<u>6.27</u>	<u>7.26</u>	<u>5.88</u>	<u>6.27</u>	
Gemini	-	closed	-	5.34	5.86	5.38	6.29	5.02	<u>5.70</u>	5.63	6.49	5.72	6.46	5.25	6.03	

Table 2 and Table 4, main results indicate that most VLMs perform poorly on IntelliCockpitBench, achieving an average score of only 4.58. In the analysis of our experiment, we observe that the closed-source models (GLM-4V-plus, GPT-4o, and Gemini) consistently outperformed open-source models in both intellicockpit query performance metrics (EN and CH) and road type scenarios (EN). Specifically, GPT-4o demonstrates the highest overall performance in both English and Chinese driving questions, with exceptional performance in reasoning (Rea.), world knowledge Q&A (Wk-QA), and other driving questions categories, achieving scores of 6.77 and 7.26 respectively in these questions.

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

Result Analysis of Open-sourced Models. 426 Qwen2-VL (7B) and InternVL-2.5 (8B) are the 427 top performers. Qwen2-VL achieves the highest 428 scores in both overall intellicockpit query perfor-429 mance in Chinese (CN) with a score of 5.45 and in 430 the special roads category for English road types, 431 scoring 5.48. Meanwhile, InternVL-2.5 demon-432 strates strong performance across various English 433 434 query, achieving an overall score of 5.09, including high scores in the reasoning (5.96) and urban roads 435 categories (4.85). Notably, the larger open-source 436 models (sizes 8B and 13B) do not consistently out-437 perform smaller models (sizes 2B to 7B), suggest-438

ing that model architecture and training data might play more crucial roles than mere parameter size in determining query-specific performance. We follow the default open-source code to evaluate and show the model's GPU usage as a reference.

In particular, we observe that InstructBLIP, with a parameter size of 10B, performed worse on this dataset compared to smaller models (5B parameters and below). This may be due to the shorter training duration of InstructBLIP. Additionally, Instruct-BLIP score lower on the world knowledge questionanswering queries, likely because the model is exposed to less driving scenario-related data during training.

Result Analysis of Query Types. Moreover, models of all sizes seem to outperform in Wk-QA questions compared to other categories of questions. This might be attributed to the fact that Wk-QA questions primarily evaluate the knowledge capacity of the models, and the answers to such questions are typically more singular. But for reasoning problems, especially in driving decision-making, the accuracy is notably low. This not only requires the model to have strong visual localization capabilities but also demands robust reasoning abilities. As illustrated in Figure 6, we provide the failure cases generated by advanced GPT-40 for better understanding.

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466



Figure 5: Results of data augmentation generated by GPT-40.

4.3 Data Augmentation

467

468 469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

Description. In real driving scenarios, the clarity of images can often not be guaranteed due to various reasons such as lighting brightness, shooting distortion, radar imaging (no color), low-pixel cameras, vehicle movement, camera occlusion, and exposure. To evaluate the robustness of VLMs in these scenarios, we employ data augmentation techniques including Clear (reduced brightness), Distorted (distortion), Grayscale (removal of image color), Low Resolution, Motion Blur, SnowEffect, and Overexposed to construct abnormal image data. We select a total of 190 images from the entire dataset, with the original images, questions, and GPT-4o's responses serving as the control group, and the augmented images, questions, and GPT-40's responses as the experimental group.

Result Analysis. The experimental results are 484 shown in Figure 5, the key findings are: 1) SnowEf-485 486 fect (simulating lens obstruction) have the greatest impact on the model's performance, with the score 487 dropping from 5.67 to 3.30 (-2.37). This indicates 488 that the model's recognition ability significantly de-489 creases when the lens is partially obstructed. 2) The 490 effects of Overexposed at 4.63 (-0.4), Grayscale 491 at 4.83 (-0.57), Clear (reduced brightness) at 4.80 492 (-1.2), and Low Resolution at 5.27 (-0.46) show 493 that the model is quite sensitive to changes in light-494 ing, color, and resolution. 3) Under Motion Blur 495 at 5.27 (-0.14) and Distorted (image distortion) at 496 4.87 (-0.26), the model still maintain good robust-497 ness, showing less impact. These results provide 498 499 important references for future improvements of the model. For example, optimizing the model in terms of occlusion, lighting variations, color, and resolution to enhance the overall robustness and adaptability of the model. 503



Figure 6: Bad cases generated by GPT-4 across five query categories. Each category presents a question and the model's generated answer is compared against the ground truth. Visual elements within each image are highlighted to indicate relevant information. Correct model responses are marked with a check, and incorrect responses are marked with a cross.

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

529

530

531

532

533

534

535

536

537

4.4 Case Study

To gain a deeper understanding of VLMs' performance and robustness, we conduct case studies and choose a specific category for an in-depth case analysis focusing on reasoning questions, with a detailed examination of the scenario depicted in Figure 6 (d). **Reasoning query:** This requires the model to accurately identify the image content based on instructions and make correct conclusions based on the scenario's knowledge. Analysis: However, GPT-40 incorrectly identified the number of cairns on the right-hand side of the road. The model's response of "eight cairns" deviated significantly from the actual count of "four cairns". This error indicates a need for improvement in the model's reasoning capabilities, particularly in object counting when the objects are similar in appearance and evenly spaced. Potential improvement: Providing more diverse and extensive training data is essential for fine-tuning VLMs, specifically targeting scenarios that require precise counting and complex visual differentiation.

5 Conclusion

In this paper, we present IntelliCockpitBench, a comprehensive benchmark designed specifically to evaluate VLMs for the intelligent cockpit. This benchmark addresses a significant gap in multimodal VQA research by incorporating a diverse and representative dataset that includes various image perspectives and four driving scenarios. We propose three innovative evaluation methods and use them to evaluate 15 VLMs. Experimental results demonstrate that GPT-40 performs well but all models struggle with complex reasoning tasks.

538

549

551 552

553

554

555

556

557

558

559

560

563

565

574

575

577

578

581

582

584

588

Limitations

Although the IntelliCockpitBench dataset includes a diverse range of scenarios, there are still some scenarios that are not fully covered, such as passenger drowsiness status. These can be included in future releases. In addition, our current dataset includes only two modes: image and text. Given that other modes (e.g., voice) are also widely used in the context of car scenes, automated driving, and intelligent driving, we will consider incorporating these additional modes in future updates.

References

- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv* preprint arXiv:2308.12966.
- Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. 2023. Shikra: Unleashing multimodal llm's referential dialogue magic. *arXiv preprint arXiv:2306.15195*.
- Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2025. Sharegpt4v: Improving large multi-modal models with better captions. In *European Conference on Computer Vision*, pages 370–387. Springer.
- Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. 2024a. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *Science China Information Sciences*, 67(12):220101.
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. 2024b. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24185–24198.
- Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Xilin Wei, Songyang Zhang, Haodong Duan, Maosong Cao, et al. 2024. Internlm-xcomposer2: Mastering freeform text-image composition and comprehension in vision-language large model. arXiv preprint arXiv:2401.16420.
- Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, et al. 2024. Chatglm: A family of large language models from glm-130b to glm-4 all tools. *arXiv preprint arXiv:2406.12793*.

Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913.

589

590

592

593

595

596

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-40 system card. *arXiv preprint arXiv:2410.21276*.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024a. Improved baselines with visual instruction tuning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 26296–26306.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024b. Llavanext: Improved reasoning, ocr, and world knowledge.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. In *NeurIPS*.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024c. Visual instruction tuning. *Advances in neural information processing systems*, 36.
- Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. 2023. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*.
- Tianwen Qian, Jingjing Chen, Linhai Zhuo, Yang Jiao, and Yu-Gang Jiang. 2024. Nuscenes-qa: A multimodal visual question answering benchmark for autonomous driving scenario. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 4542–4550.
- Chonghao Sima, Katrin Renz, Kashyap Chitta, Li Chen, Hanxue Zhang, Chengen Xie, Ping Luo, Andreas Geiger, and Hongyang Li. 2023. Drivelm: Driving with graph visual question answering.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu,

742

743

744

745

746

747

748

699

700

Chang Zhou, Jingren Zhou, and Junyang Lin. 2024a. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.

647

649

650

658

664

666

667

670

673

674

675

676

677

678

679

683

690

692

697

- Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. 2023a. Cogvlm: Visual expert for pretrained language models. arXiv preprint arXiv:2311.03079.
- Weiyun Wang, Min Shi, Qingyun Li, Wenhai Wang, Zhenhang Huang, Linjie Xing, Zhe Chen, Hao Li, Xizhou Zhu, Zhiguo Cao, et al. 2023b. The all-seeing project: Towards panoptic visual recognition and understanding of the open world. arXiv preprint arXiv:2308.01907.
- Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, et al. 2024b. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. Advances in Neural Information Processing Systems, 36.
- Shaoyuan Xie, Lingdong Kong, Yuhao Dong, Chonghao Sima, Wenwei Zhang, Qi Alfred Chen, Ziwei Liu, and Liang Pan. 2025. Are vlms ready for autonomous driving? an empirical study from the reliability, data, and metric perspectives. *arXiv preprint arXiv:2501.04003*.
- Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. 2017.
 Video question answering via gradually refined attention over appearance and motion. In *Proceedings of the 25th ACM international conference on Multime-dia*, pages 1645–1653.
- Liang Xu, Hai Hu, Xuanwei Zhang, Lu Li, Chenjie Cao, Yudong Li, Yechen Xu, Kai Sun, Dian Yu, Cong Yu, et al. 2020. Clue: A chinese language understanding evaluation benchmark. *arXiv preprint arXiv:2004.05986*.
- Fisher Yu, Haofeng Chen, Xin Wang, Wenqi Xian, Yingying Chen, Fangchen Liu, Vashisht Madhavan, and Trevor Darrell. 2020. Bdd100k: A diverse driving dataset for heterogeneous multitask learning. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 2636– 2645.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*.

A Appendix

A.1 Taxonomy of VLMs' Driving Questions

Our IntelliCockpitBench covers five mainstream query types, including description, recognition, world knowledge Q&A, reasoning, and others, examples are shown in Figure 7.

Description. Simple queries that require basic descriptions or presentation of information, e.g., "What's the view from the front?".

Recognition. Moderately complex queries that involve pattern recognition and basic reasoning. The subcategories include vehicle model recognition, information extraction, object recognition, emotion recognition, behavior recognition.

World Knowledge Q&A. These queries demand the application of domain-specific knowledge and common sense, combined with intermediate reasoning skills. The subcategories consist of traffic laws and regulations, geospatial environmental information, socio-cultural knowledge, general knowledge.

Reasoning. Queries at this level represent the highest complexity, necessitating advanced logical reasoning and refined cognitive skills. The subcategories include quantitative statistics, distance measurement, angle measurement, area and volume, intent recognition/ probabilistic reasoning, driving decisions.

Others. These queries combine multiple types of reasoning and require the synthesis of diverse skills. The subcategories include: creation, translation, others.

A.2 Taxonomy of Driving scenarios

We have classified the data based on driving scenarios into 4 categories, as shown in Figure 2. Taking road type as an example, from densely populated urban streets to isolated rural roads, the distinct visual attributes of these varied driving environments serve as a robust can be used to assess the adaptability and generalizability of VLMs.

Weather Conditions. Our dataset covers a spectrum of weather conditions such as Clear, Cloudy, Overcast/Nighttime, Light Rain, Heavy Rain, Snowy, Foggy, Dusty/Stormy, and Others. Each condition presents unique visual features and challenges, ensuring that VLMs can handle a wide range of environmental scenarios, thus enhancing their robustness.

Road Types. These images cover various types of roads, including Urban Roads, Rural Roads,



Figure 7: Examples of Various VLMs' Driving Questions on IntelliCockpitBench.



Figure 8: Distribution of questioning type.

Highways, Special Roads, Parking Lots or Private Roads, and Others Roads. The specific classifications are as follows:

749

750

752

754

756

Urban Roads: Residential Area Roads, Commercial Area Roads, Ring Roads/Express Loops,
Urban Arterial Roads. Rural Roads: Small
Village Roads, Rural Multi-lane Roads, Farm
Roads, Forest or Hill Roads. Highways: National/Provincial Roads, Intercity Highways, Urban Highways. Special Roads: Mountain Roads,
Coastal Roads, Desert Roads, Forest Roads, High
Mountain Ice and Snow Roads. Parking Lots or



Figure 9: Cumulative distribution of judging by human, general, and rule-calibrated on sampled IntelliCockpitBench along their ratings.

Private Roads: Parking Lots, Private/Exclusive Roads. **Other Roads**: Construction Zones, Tunnels, Bridges, Flooded Roads/Waterlogged Sections, Other Roads.

This diversity ensures that VLMs can understand and respond accurately in distinct driving environments, ranging from congested city streets to remote rural roads.

Driving Status. Images are categorized based on the vehicle's driving status, either Moving or **Stopping**. This distinction is crucial because it affects the context and relevance of visual informa-

837

tion, enabling VLMs to adapt to both dynamic and static conditions.

Shooting Angles. To capture the complete environment of the vehicle, images are taken from different angles: Inside the Vehicle and outside (Front of the Vehicle, Side of the Vehicle, Rear of the Vehicle). This multiangle approach allows VLMs to process and understand perspectives from various points of view, improving their situational awareness.

A.3 Human Check Details

773

774

775

776

778

779

790

791

794

795

798

799

805

807

810

811

813

814

816

We conduct a high-standard human check of the generated VQA pairs. Specifically, a total of 12 data annotators participate in this process, with each annotator labeling approximately 150 items per day, resulting in a total of 16, 154 items over the course of 108 person-days. Quality control identifies 4, 000 items that require rework, which takes an additional 27 person-days, bringing the entire query to 135 person-days. Additionally, two senior annotators conduct quality checks, inspecting 20% of each batch of 200 items. Any batch with an accuracy rate below 95% is sent back for rework, and this process takes another 24 person-days.

A.4 Rejection Strategy

In the construction of VQA pairs, we have developed a comprehensive refusal strategy to ensure information security, answer accuracy, and query quality. We refuse to answer for the following situations.

- The image with poor quality, including those that are difficult to see due to being too far away, too dark at night, blurry due to shooting, or distorted images.
- The image from cameras other than the front/rear cameras or the left/right side mirrors (such as those depicting the trunk or underneath the vehicle).
- The image does not contain sufficient information to answer the user's query.
- The query is a declarative sentence.
 - The query that involves user privacy.

We present examples of refusal queries in IntelliCockpitBench in Figure 10.

817 A.5 AI Assistants In Writing

We use AI Assistants (e.g., ChatGPT) in our research to help us improve writing.

A.6 Prompts and Details of Methods

In our evaluation paradigm, we select different dimensions for various categories to ensure a more comprehensive and accurate assessment. The detailed selections of the dimensions are described in Table 6 and the detailed definitions of these dimensions are provided in Table 7.

For queries with relatively fixed answers (e.g., Quantitative Statistics, Vehicle Model Recognition), we set the temperature to 0.1, ensuring deterministic and reproducible outputs. For queries requiring creativity and diversity (e.g., description), we use a higher temperature (e.g., 0.7) to encourage longer and more varied generations.

The following are all the prompts used in our experiments, including query generation prompt, answer generation prompt, rule-based evaluation prompt, and general evaluation prompt.



Figure 10: Examples of refusal VQA pairs in IntelliCockpitBench.

Table 3: Criteria for determining whether a query is discarded, if the answer is no, then the query is discarded.

Specific Criteria

- 1. Whether it matches human expression habits.
- 2. Whether it is consistent with the questions typically asked in vehicle scenarios.
- 3. Whether it is reasonable and legal.
- 4. Whether the question is accurate and relevant.
- 5. Whether the question aligns with the visual content ("in the picture"),
- or if it necessitates discarding due to similarity to existing expressions.

Model	Size	Type	Road Types (EN)							
		-510	Highways	PLPR.	Rural Roads	Special Roads	Other Roads	Urban Roads		
DeepSeek-VL-base	1.3B	open	3.36	3.25	3.65	3.89	3.74	3.18		
MiniCPM-V-2.0	2B	open	3.99	3.76	4.13	4.25	4.31	3.84		
GLM-4V	2B	open	4.31	4.27	4.40	4.57	4.37	4.20		
Qwen2-VL	2B	open	4.61	4.67	4.67	4.86	4.71	4.46		
Megrez	3B	open	4.11	3.96	4.32	4.45	4.27	3.72		
GLM-4V	5B	open	4.57	4.37	4.56	4.74	4.70	4.32		
InstructBLIP	7B	open	3.88	3.64	4.05	4.36	3.91	3.45		
Qwne2-VL	7B	open	5.17	5.24	5.23	5.48	5.28	4.94		
LLaVA-1.5	7B	open	4.04	3.81	4.23	4.73	4.43	3.66		
InternVL-2.5	8B	open	5.03	4.98	5.18	5.51	5.04	4.85		
GLM-4V	9B	open	4.89	4.73	5.03	5.19	4.82	4.62		
LLaVA-1.5	13B	open	4.14	3.93	4.37	4.87	4.52	3.83		
GLM-4V-plus	-	closed	5.30	5.23	5.46	5.73	5.34	5.04		
GPT-4o	-	closed	5.86	5.90	5.78	6.03	5.80	5.67		
Gemini	-	closed	5.33	5.30	5.52	5.63	5.20	5.15		

Table 4: Performance evaluation of various VLMs on the IntelliCockpitBench for different English road types. The best performance is shown in **bold**. 'PLPR." denotes Parking Lots and Private Roads. The best performance is shown in bold.

Table 5: Performance evaluation of various VLMs on the IntelliCockpitBench for different English weather conditions. The best performance is shown in **bold**.

Model Size Tu				Weather Condition (EN)								
Model	Size Type		Clear	Cloudy	Dust/Sandstorm Weather	Foggy	Light Rain	Moderate or Heavy Rain	Overcast or Night	Snowy	Unknown	
DeepSeek-VL-base	1.3B	open	3.20	3.35	3.65	4.30	3.61	3.99	3.24	4.05	3.62	
MiniCPM-V-2.0	2B	open	3.76	4.09	4.35	4.58	4.11	4.23	3.92	4.48	4.13	
GLM-4V	2B	open	4.13	4.37	4.62	5.07	4.44	4.70	4.14	4.79	4.33	
Qwen2-VL	2B	open	4.42	4.59	4.69	5.22	4.75	5.02	4.39	5.10	4.73	
Megrez	3B	open	3.79	4.01	4.29	4.75	4.31	4.39	3.84	4.63	4.16	
GLM-4V	5B	open	4.31	4.39	4.89	5.19	4.63	4.95	4.30	4.93	4.60	
InstructBLIP	7B	open	3.48	3.64	4.56	4.80	3.90	4.43	3.54	4.79	3.83	
Qwne2-VL	7B	open	4.91	5.00	5.43	5.97	5.37	5.57	4.98	5.73	5.25	
LLaVA-1.5	7B	open	3.66	3.73	4.34	5.36	4.25	4.76	3.82	5.11	4.31	
InternVL-2.5	8B	open	4.87	4.90	5.81	5.86	5.21	5.35	4.92	5.71	5.01	
GLM-4V	9B	open	4.56	4.67	5.70	5.68	5.10	5.53	4.60	5.58	4.78	
LLaVA-1.5	13B	open	3.83	3.99	4.70	5.39	4.37	4.78	3.96	5.24	4.34	
GLM-4V-plus	-	closed	5.03	5.24	6.17	6.11	5.45	5.81	5.17	5.96	5.25	
GPT-40	-	closed	5.60	5.66	6.49	6.14	6.04	6.19	5.75	6.25	5.82	
Gemini	-	closed	5.15	5.14	6.01	6.04	5.53	5.54	5.22	5.82	5.24	

Query Generation Prompt

```
def make_example(df: pd.DataFrame, col):
    if col=="Description":
       questions = df[df['label2'].isna() | (df['label2'].str.strip() == '')['english'].tolist()
    else:
       questions = df[df['label2'] == col]['english'].tolist()
    return "e.g.: " + " ".join([f"[idx+1]. [value]" for idx, value in enumerate(questions)])
prompt = f"""You are a driver operating a vehicle. Based on the content of images captured by
the onboard cameras, generate 'questions' and their 'perspectives' from the **question perspective**,
**question requirement**, and **question classification**. Please ensure:**Usage Scenario and Goal**:
The goal is to generate a variety of questions for vehicle multimodal scenarios that align with
human habits and cover diverse perspectives. Avoid meaningless or illogical questions. The questions
must specifically require the use of multimodal models (combining visual and other data) and should
not be solvable by language models alone or language models paired with tools like maps, weather
applications, or navigation systems. Additionally, avoid using terms such as "in the image" or "in the
background," as these are not typical questions asked by drivers.
1.Safety Check: Verify that the AI assistant's response is safe and harmless. If the response is found
to be unsafe or harmful in any way, it will be given a score of 0.
2.Positional Accuracy Check: If the query involves specific positions or orientations, check whether the
AI's answer correctly describes the location or object. Incorrect descriptions will result in a lower
score.
3.Answer Comparison: Compare the AI's response to a reference answer. Identify any discrepancies
and provide explanations for the shortcomings.
4.Multi-Dimensional Scoring: Assess the response across various dimensions, Each dimension will be
```

4.Multi-Dimensional Scoring: Assess the response across various dimensions, Each dimension will be scored on a scale of 1-10.Overall Score Calculation: Combine the scores from different dimensions, applying a weighted average based on the importance of each dimension, to compute the overall score. 5.Strict Score Adjustment: Adjust the overall score according to specific rules to ensure a stringent evaluation. This step is critical to maintain the integrity and accuracy of the scoring process. **Question Perspectives**

- **Why** - **What** - **Where** - **When** - **Who/Which** - **How** - **How much/How many** - **How feel** - **Can/Have** - **Is/Do/Others**

```
Query Generation Prompt
```

Question Classification System :

- 1. Descriptive:[make_example(df, 'Description')]
- 2. Recognition:
 - **Vehicle Model Recognition**: [make_example(df, 'Vehicle Model Recognition')]
 - **Information Extraction**: [make_example(df, 'Information Extraction')]
 - **Object Recognition**: [make_example(df, 'Object Recognition')]
 - **Emotion Recognition**: [make_example(df, 'Emotion Recognition')]
 - **Human Activity Recognition**: [make_example(df, 'Human Activity Recognition')]

```
3. World Knowledge Q&A:
```

- **Traffic Laws and Regulations**: [make_example(df, 'Traffic Laws and Regulations')]
- **Geospatial Environmental Information**: [make_example(df,'Geospatial Environmental Information')]
- **Socio-cultural Knowledge**: [make_example(df,'Socio-cultural Knowledge')]
- **General Knowledge**: [make_example(df,'General Knowledge')]
- 4. Reasoning:

```
- **Quantitative Statistics**: [make_example(df,'Quantitative Statistics')]
```

- **Distance Measurement**: [make_example(df, 'Distance Measurement')]
- **Angle Measurement**: [make_example(df, 'Angle Measurement')]
- **Area and Volume**: [make_example(df, 'Area and Volume')]

- **Probabilistic Reasoning/ Intent Recognition**: [make_example(df,'Probabilistic Reasoning/ Intent Recognition')]

- **Driving Decisions**: [make_example(df, 'Driving Decisions')]

5. Others:

- **Creation**: [make_example(df, 'Creation')]
- **Translation**: [make_example(df, 'Translation')]
- **Others**: [make_example(df,'Others')]

Query Generation Prompt

Question Requirements
(a) Relevance
- Definition: Is the question relevant to the given image?
(b) Answerability
- Definition: Can the question be clearly answered?
(c) Innovativeness
- Definition: Is the question novel and not easily repetitive?
(d) Authenticity
- Definition: Is the question typical of an in-car scenario, consistent with human preferences?
(e) Simplicity
- Definition: Is the question concise, avoiding unnecessary complexity?
Output Format:
[[["Question":"Generated Question 1","Perspective":"Question Perspective 1"]],
[["Question":"Generated Question n", "Perspective":"Question Perspective n"]],]
Begin generating questions, ensuring diverse perspectives, and output only in the specified 'Output
Format' without any extra text!!!

Answer Generation Prompt

You are an in-car intelligent agent.Based on the content of images captured by the onboard camera and the given question, generate matching 'primary tags' and 'secondary tags' from the **Question Classification System**, and provide the 'answer' to the question along with the 'reason' for the answer. Ensure the following:

1. **Clarity**: Descriptions must be clear and concise. 2. **Consistency**: The generated primary and secondary tags must strictly correspond to the relevant categories in the classification system without cross-category questions. 3. **Conciseness**: Ensure questions and explanations are short and easy to process for quick comprehension during real-time operations. 4. **Relevance**: If the question is unclear or does not require the capabilities of the in-car multimodal model (i.e., it can be answered solely by the language model or by using tools like 'weather software', 'maps' for precise location, 'navigation', etc.), please directly generate "Sorry, I can't answer" in the 'Answer' field of the **Output Format**. 5. **Context Relevance**: If the question contains phrases such as 'in the picture', 'in the background', etc., which are not typical of questions a driver would ask while driving, please directly generate "Sorry, I can't answer" in the 'Answer' field of the **Output Format**. **Output Format**.

[question]

Question Classification System

1. Description

2. Recognition: - **Vehicle Model Recognition**: e.g., What is the vehicle model in the far left foreground? - **Information Extraction**: e.g., What is the content of the yellow billboard on the top right? - **Object Recognition**: e.g., What is on the ground on the left? - **Emotion Recognition**: e.g., Is that person on the road crying? Why is that man laughing? - **Human Activity Recognition**: e.g., What is that person doing? Why is he crawling on the road?

3. World Knowledge Q&A: - **Traffic Laws and Regulations**: e.g., What is the meaning of the sign ahead? Can I turn left at this intersection? - **Geospatial Environmental Information**: e.g., Where is this place? Is this a commercial or residential area? What building is in front? What is the current weather? - **Socio-cultural Knowledge**: e.g., How is this left-turn signal represented in other countries? - **General Knowledge**: e.g., Is the building on the street a restaurant or a hotel?

4. Reasoning: - **Quantitative Statistics**: e.g., How many black cars are in the left foreground lane? How many lanes are there on the road ahead? How many floors does the white building on the right have? - **Distance Measurement**: e.g., How far is the bus stop from me? How far is the man in black from the mall? How far is the car from the crosswalk? - **Angle Measurement**: e.g., What is the approximate distance between the black car ahead and the pedestrian? - **Area and Volume**: e.g., What is the ground area of the object on the right ahead? - **Probabilistic Reasoning/ Intent Recognition**: e.g., What is that person standing in the middle of the road trying to do? Is there an accident ahead? Why is this car signaling a left turn? - **Driving Decisions**: e.g., Based on the sign, which lane should be chosen to head to a specific address? Please evaluate the road conditions ahead; how should I operate to avoid danger in the situation ahead? How to get to a specific address?

5. Others: - **Creation**: e.g., Please write a poem based on the road conditions. - **Translation**: e.g., Please translate the content of the advertisement ahead into English. - **Others**: Questions not included in the above categories Output Format:

[[["Primary Tag": "Primary Tag of the Question", "Secondary Tag": "Secondary Tag of the Question", "Answer": "Answer to the Question"]]]

Please begin generating and output only in the specified 'Output Format' without any extra text.

Rule-based Evaluation Prompt

You are an assistant skilled at evaluating text quality.

Please act as an impartial judge and assess the quality of an AI assistant's response to a user's query. Since the primary category of the responses you are evaluating is **[category]** and the secondary category is **[subcategory]**, you need to evaluate the responses from the following dimensions:

[dimensions] We will provide you with the user's uploaded image, the user's question based on the image, a high-quality reference answer, and the AI assistant's answer that you need to evaluate. When performing your evaluation, you must reference the input image, not just the reference answer, and you need to compare the image with the reference answer and the AI assistant's answer to determine which one is more reasonable. When you begin your evaluation, you need to follow these steps:

1. Safety Check Determine if the AI assistant's answer is safe and harmless, meaning that the response should not incite dangerous or harmful behavior, nor should it disseminate harmful information. If the AI assistant's answer does not meet the safety and harmlessness criteria, each dimension's score must be 0.

2. Positional Accuracy Check If the question specifies a particular location, then you need to check the corresponding location's object in the image to confirm whether the AI assistant's response aligns with the object at the specified location in the image. The reference answer certainly describes the object at the corresponding location. If the AI assistant's answer correctly describes the content in the image but the described location doesn't match the specified location in the guestion, then the scores for all evaluation dimensions should be lowered.

3. Answer Comparison Compare the AI assistant's answer with the reference answer and, in conjunction with the input image, point out the deficiencies in the AI assistant's answer, providing further explanations.

4. Multi-Dimensional Scoring Evaluate the AI assistant's answer from different dimensions, giving a score between 1 and 10 for each dimension after evaluation. You must score all given dimensions.

5. Overall Score Calculation Finally, provide an overall score between 1 and 10 for the AI assistant's answer, based on the evaluations of each dimension. Each evaluation dimension has an importance score ranging from 1 to 3, with higher scores indicating greater importance. When calculating the overall score, please weight each dimension's scores according to their importance scores.

6. Strict Score Adjustment Your scoring needs to be as strict as possible. After scoring each dimension and calculating the total score, you need to adjust the scores for each dimension and the total score based on the following rules: Factuality, User Satisfaction, and Visual Location are the most important dimensions. If any of these dimensions perform poorly, the scores for other dimensions should be lowered accordingly. If the response contains irrelevant issues or has significant factual errors or generates harmful content, the total score must be 1 to 2. If the response has no major errors and is generally harmless but of low quality and fails to meet user needs, the total score is 3 to 4. If the response generally meets user requirements but performs poorly in some dimensions, indicating moderate quality, the total score can be 5 to 6. If the response quality is close to the reference answer and performs well in all dimensions, the total score is 7 to 8. Only when the response quality significantly exceeds the reference answer, fully resolving the user's issues and needs and performing near-perfectly in all dimensions can it score 9 to 10. As an example, the reference answer can be scored 8.

Remember, you must conduct evaluation and explanation before scoring. After explaining each dimension, you need to add the score for that dimension. At the end of your response, return all your scores in the following dictionary format (including brackets), ensuring your scores are whole numbers:

"Dimension One": [Score, Importance Score], "Dimension Two": [Score, Importance Score], ..., "Overall Score": Score. User's Question: [question]

[Reference Answer Start] [reference] [Reference Answer End] [Assistant's Answer Start] [answer] [Assistant's Answer End]

General Evaluation Prompt

You are an assistant skilled at evaluating text quality. Please act as an impartial judge and assess the quality of the AI assistant's responses to user queries. Your evaluation should take into account factors such as correctness (high priority), helpfulness, relevance, depth, innovativeness, and level of detail. You will be provided with a high-quality reference answer and the assistant's response to be evaluated. When you begin your assessment, compare the assistant's response to the reference answer, identify errors in the assistant's response, and provide a brief explanation. Please be as objective as possible. After providing an explanation, you must rate the response strictly in the following format on a scale of 1 to 10: "[[Rating]]," for example, "Rating: [[5]]."

User's Query: [Question]

[Reference Answer Start][Reference Answer][Reference Answer End] [Assistant's Response Start][Model Answer][Assistant's Response End]

Category	Query Type	Evaluation Dimension	Reply Temperature
		["Factuality", 3], ["User Satisfaction", 3], ["Visual Location", 3],	
Description	Description	["Clarity", 1], ["Naturalness", 1], ["Richness", 2],	0.7
		["Completeness", 2]	
Descerition	Vahiala Madal Daga mitian	["Factuality", 3], ["User Satisfaction", 3], ["Visual Location", 3],	0.1
Recognition	venicle Model Recognition	["Clarity", 1], ["Completeness", 2]	0.1
	T. Constant T. Constant	["Factuality", 3], ["User Satisfaction", 3], ["Visual Location", 3],	
	information Extraction	["Clarity", 1], ["Completeness", 2]	
	Object Deservition	["Factuality", 3], ["User Satisfaction", 3], ["Visual Location", 3],	
	Object Recognition	["Clarity", 1], ["Completeness", 2]	
		["Factuality", 3], ["User Satisfaction", 3], ["Visual Location", 3],	
	Emotion Recognition	["Clarity", 1]	
	Daharian Daaramitian	["Factuality", 3], ["User Satisfaction", 3], ["Visual Location", 3],	
	Benavior Recognition	["Clarity", 1]	
W. 11 K 1. 1 0.8 A	To College and Decision	["Factuality", 3], ["User Satisfaction", 3], ["Visual Location", 3],	0.1
world Knowledge Q&A	Traine Laws and Regulations	["Clarity", 1], ["Completeness", 1], ["Responsibility", 2]	0.1
	Geospatial Environmental	["Factuality", 3], ["User Satisfaction", 3], ["Visual Location", 3],	
	Information	["Clarity", 1], ["Completeness", 1], ["Responsibility", 2]	
	Sacio aultural Knowladza	["Factuality", 3], ["User Satisfaction", 3], ["Visual Location", 3],	
	Socio-cultural Knowledge	["Clarity", 1], ["Completeness", 1], ["Responsibility", 2]	
	Conoral Knowledge	["Factuality", 3], ["User Satisfaction", 3], ["Visual Location", 3],	
	General Knowledge	["Clarity", 1], ["Completeness", 1], ["Responsibility", 2]	
Reasoning	Quantitative Statistics	["Factuality", 3], ["User Satisfaction", 3], ["Visual Location", 3],	0.1
Reasoning	Quantitative Statistics	["Clarity", 1]	0.1
	Distance Measurement	["Factuality", 3], ["User Satisfaction", 3], ["Visual Location", 3],	
	Distance weasurement	["Clarity", 1]	
	Angle Measurement	["Factuality", 3], ["User Satisfaction", 3], ["Visual Location", 3],	
	The weasternent	["Clarity", 1]	
	Area and Volume	["Factuality", 3], ["User Satisfaction", 3], ["Visual Location", 3],	
	Thea and volume	["Clarity", 1]	
	Intent Recognition	["Factuality", 3], ["User Satisfaction", 3], ["Visual Location", 3],	
	/ Probabilistic Reasoning	["Clarity", 1], ["Responsibility", 2], ["Logical Coherence", 2]	
		["Factuality", 3], ["User Satisfaction", 3], ["Visual Location", 3],	
	Driving Decisions	["Clarity", 1], ["Responsibility", 2], ["Logical Coherence", 2],	
		["Completeness", 2]	
Others	Creation	["Factuality", 3], ["User Satisfaction", 3], ["Visual Location", 3],	0.7
Oulors	Crownon	["Clarity", 1], ["Creativity", 2]	0.7
	Translation	["Factuality", 3], ["User Satisfaction", 3], ["Visual Location", 3],	
	mansharion	["Clarity", 1], ["Completeness", 2]	
	Others	["Factuality", 3], ["User Satisfaction", 3], ["Visual Location", 3],	
	Guidio	["Clarity", 1], ["Completeness", 2]	

Table 6: Judging dimensions and VLM reply generation temperatures of IntelliCockpitBench on different categories. ["Factuality", 3] represents a Factuality importance score of 3.

Dimension	Definition
Factuality	Whether the information provided in the response is accurate and based on reliable facts and data,
	or derived from the content in the provided images, and whether it helps answer the user's question.
User Satisfaction	Whether the response meets the user's question and needs,
User Saustaction	and provides a comprehensive and appropriate answer to the question.
Visual Location	Whether the response accurately perceives the specific orientation in the image
VISUAI LOCATION	when the user's question involves specific spatial orientation.
Clarity	Whether the response is clear and understandable, and whether it uses concise language
	and structure so that the user can easily understand it.
Naturalness	Whether the content of the response is fluent and smooth,
	consistent with everyday language norms and colloquial expressions.
Diahnaga	Whether the response includes rich info, depth, context, diversity, detailed explanations
Kichness	and examples to meet user needs and provide a comprehensive understanding.
Completeness	Whether the response provides sufficient information and details to meet the user's needs,
Completeness	and whether it avoids omitting important aspects.
Daanansihility	Whether the recommendations or information provided in the response are practical and responsible,
Responsibility	and whether they consider potential risks and consequences and comply with safety standards.
Logical Coherence	Whether the response maintains overall consistency and logical coherence between different sections,
	avoiding self-contradiction.
Creativity	Whether the response is innovative or unique, providing novel insights or solutions.

Table 7: The definition of different dimensions.