
FROM STEERING TO PEDALLING: DO AUTONOMOUS DRIVING VLMS GENERALIZE TO CYCLIST-ASSISTIVE SPATIAL PERCEPTION AND PLANNING?

Krishna Kanth Nakka* and Vedasri Nakka†

Munich, Bavaria, Germany

krishkanth.92@gmail.com

<https://krishnakanthnakka.github.io/CyclingVQA>

ABSTRACT

Cyclists often encounter safety-critical situations in urban traffic, highlighting the need for assistive systems that support safe and informed decision-making. Recently, vision–language models (VLMs) have demonstrated strong performance on autonomous driving benchmarks, suggesting their potential for general traffic understanding and navigation-related reasoning. However, existing evaluations are predominantly vehicle-centric and fail to assess perception and reasoning from a cyclist-centric viewpoint. To address this gap, we introduce *CyclingVQA*, a diagnostic benchmark designed to probe perception, spatio-temporal understanding, and traffic-rule-to-lane reasoning from a cyclist’s perspective. Evaluating **31+** recent VLMs spanning general-purpose, spatially enhanced, and autonomous-driving-specialized models, we find that current models demonstrate encouraging capabilities, while also revealing clear areas for improvement in cyclist-centric perception and reasoning, particularly in interpreting cyclist-specific traffic cues and associating signs with the correct navigational lanes. Notably, several driving-specialized models underperform strong generalist VLMs, indicating limited transfer from vehicle-centric training to cyclist-assistive scenarios. Finally, through systematic error analysis, we identify recurring failure modes to guide the development of more effective cyclist-assistive intelligent systems.

1 INTRODUCTION

Cycling has emerged as a cornerstone of sustainable urban mobility, offering a low-carbon and health-conscious alternative to motorized transport. This importance is reflected in its widespread adoption across Europe, where cycling accounts for a substantial fraction of daily commutes, including approximately 41% in the Netherlands and 15% in Germany (Euronews, 2023). Moreover, because cycling does not require formal licensing, it remains accessible to a broad demographic with varying levels of traffic experience. Yet this same accessibility, however, comes at a cost: cyclists are among the most vulnerable road users. In Germany alone, traffic accidents in 2023 resulted in several hundred cyclist fatalities and tens of thousands of serious injuries (Industry, 2024). Despite these risks, a significant technological disparity persists. While autonomous driving (AD) research (Zhou et al., 2024; Xie et al., 2025; Cui et al., 2024; Yang et al., 2023; Jiang et al., 2025a) has made substantial progress in vehicle-centric perception and decision-making, cyclists still lack dedicated, perception-driven assistive systems to support navigation in complex or high-risk traffic scenarios. For example, at busy intersections, cyclists may struggle with interpreting traffic signs, selecting the correct lane, and making safe crossing decisions, underscoring the value of cyclist-assistive systems with strong perceptual and reasoning capabilities.

Intuitively, cyclist-assistive intelligence requires many of the same perceptual and spatial reasoning capabilities as autonomous driving, including recognizing traffic signs, identifying lane boundaries,

*Core Contributor

†Data Curation and Annotation Support

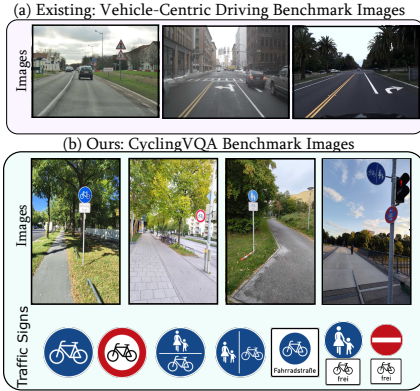


Figure 1: Comparison between vehicle-centric driving benchmarks (Corbière et al., 2025; Tian et al., 2025; Li et al., 2025b), which predominantly focus on road-level perspectives, and our cyclist-centric viewpoint, highlighting differences in camera perspective and the presence of cycling-specific traffic signage. See Appendix E for further examples from our dataset.

and reasoning over road layout. Recently, vision–language models (VLMs) (Team et al., 2025; Li et al., 2024; Liu et al., 2024a; Bai et al., 2025a) have demonstrated strong capabilities across a wide range of multimodal tasks. Alongside this progress, a substantial body of work has assessed and improved spatial reasoning in VLMs through dedicated benchmarks (Fu et al., 2024; Liu et al., 2024b; Yue et al., 2024; Lu et al., 2023; Kazemzadeh et al., 2014; Yang et al., 2025; Batra et al., 2025; Ogezi & Shi, 2025; Gan et al., 2025; Shen et al., 2025). Separately, VLMs have been adapted to autonomous driving scenarios (Li et al., 2025a; Ma et al., 2023; Huang et al., 2025; Azzolini et al., 2025; Tian et al., 2024; Ishaq et al., 2025), enabling vehicle-centric perception, reasoning, and planning tasks (Corbière et al., 2025; Ishaq et al., 2025). However, it remains unclear whether these capabilities of AD-VLMs transfer to cyclist-assistive settings, where perspective, infrastructure, and traffic rules differ. This motivates a central research question: *To what extent do existing VLMs, particularly those optimized for autonomous driving, generalize to the spatial perception and traffic understanding demands of cyclist-perspective scenarios?*

Unsurprisingly, existing traffic benchmarks for VLMs (Cao et al., 2024; Ishaq et al., 2025; Wei et al., 2025; Ghosh et al., 2025; Tian et al., 2025; Li et al., 2025b) remain predominantly vehicle-centric, focusing on motorized traffic flows from a driver’s perspective. While these benchmarks are highly relevant, a cyclist’s viewpoint, as illustrated in Figure 1, differs in several important respects from that of a car, introducing distinct navigational constraints such as bicycle-only lanes, cyclist-specific signage, and the negotiation of shared pedestrian spaces. Consequently, images derived from vehicle-centric datasets render such data ill-suited for evaluating whether VLMs can support cyclist-assistive perception and decision-making, particularly in scenarios that require fine-grained associations between traffic signs, lanes, and permitted actions.

To bridge this gap, we introduce CyclingVQA, a first effort toward evaluating cyclist-perspective spatial perception and navigation-relevant reasoning in urban traffic scenes. CyclingVQA provides a granular assessment from a cyclist’s egocentric perspective, serving as a foundational step toward cyclist-assistive intelligent systems. The benchmark comprises **2,009** multiple-choice visual question–answer pairs derived from **695** images captured in real-world urban cycling environments in Munich. As a central contribution of this work, we evaluate a diverse suite of state-of-the-art VLMs—including general-purpose models, spatially enhanced architectures (Yang et al., 2025; Batra et al., 2025; Cho et al., 2025), and autonomous-driving-specialized models (Li et al., 2025a; Azzolini et al., 2025). While several models exhibit promising zero-shot performance, they still fall short in challenging tasks. In particular, multiple driving-specialized VLMs struggle to interpret cyclist-specific traffic cues and to associate signs with the correct navigational lanes, often underperforming strong generalist baselines. We believe these empirical findings, supported by a systematic failure mode analysis, highlight important limitations in current traffic-oriented VLMs and can help guide the development of more effective cyclist-assistive systems. In summary, our contributions are as follows:

1. **Cyclist-Centric Benchmark:** We introduce CyclingVQA, a cyclist-centric benchmark comprising 2,009 multiple-choice visual question–answer pairs derived from 695 real-world egocentric images (§ 3.2).
2. **Specialized Spatial Evaluation:** We define eight evaluation tasks that probe cyclist-centric spatial perception, traffic rule compliance, and navigation-relevant reasoning in complex urban environments (§ 3.1).

-
3. **Comprehensive VLM Evaluation:** We benchmark state-of-the-art VLMs, including general-purpose, spatially enhanced (Yang et al., 2025; Batra et al., 2025; Cho et al., 2025), and autonomous-driving-focused models (Li et al., 2025a; Azzolini et al., 2025; NVIDIA, 2025; Ishaq et al., 2025), revealing substantial room for improvement in cyclist-centric reasoning (§ 4.2).
 4. **Systematic Failure Analysis:** We conduct a comprehensive analysis of recurring failure modes, providing insights and guidance for future cyclist-assistive intelligent systems (§ 4.4).

2 RELATED WORK

We provide a brief review of related benchmarks and specialist VLMs in the context of spatial intelligencedd and autonomous driving.

Spatial Reasoning and Driving Benchmarks. The rapid emergence of benchmarks over the past two years has significantly advanced the evaluation of multimodal spatial understanding in VLMs. General-purpose spatial reasoning benchmarks (Fu et al., 2024; Cheng et al., 2024; Wang et al., 2024; Jia et al., 2025; Ma et al., 2025a; Kamath et al., 2023; Stogiannidis et al., 2025) typically assess a model’s understanding of object relationships, geometry, orientation, and relative positioning through visual question answering tasks. While foundational, these benchmarks, albeit effective for general spatial reasoning, predominantly rely on generic imagery, leaving the specific perceptual and semantic challenges of cyclist-assistive traffic scenarios largely unexplored. Our benchmark extend the above spatial reasoning benchmarks by introducing a dedicated suite of tasks tailored to cyclist-centric intelligent assistance.

Expanding on this, a parallel line of research has introduced driving benchmarks such as DrivingVQA (Corbière et al., 2025) and DriveQA (Wei et al., 2025), which evaluate models from a driver’s perspective on tasks ranging from right-of-way reasoning to decision-making. Recent efforts further emphasize specialized settings, including construction-zone understanding (Ghosh et al., 2025), fine-grained perception (Li et al., 2025b), and safety-critical robustness (Xing et al., 2024). Most of these benchmarks are derived from large-scale, vehicle-centric datasets such as Waymo (Sun et al., 2020) and nuScenes (Caesar et al., 2020). However, as we mentioned, these datasets are captured exclusively from motor vehicles, they do not capture key spatial semantics of the cycling environment, including dedicated bicycle lanes, shared-use paths, and cyclist-specific traffic signals. As a result, existing benchmarks do not adequately represent the egocentric viewpoints and navigational constraints inherent to cycling, limiting their suitability for evaluating cyclist-assistive reasoning. Our work addresses this gap by introducing a benchmark explicitly designed from the cyclist’s perspective.

Specialist VLMs. The development of benchmarks has been accompanied by the emergence of specialist VLMs with enhanced spatial awareness. Spatially enhanced models (Yang et al., 2025; Batra et al., 2025; Cho et al., 2025) are often fine-tuned on grounded datasets such as SpatialVQA (Chen et al., 2024) to improve spatial grounding and relative distance understanding. Besides, autonomous-driving-focused VLMs (Huang et al., 2025; Li et al., 2025a; Ma et al., 2023; Azzolini et al., 2025) are optimized for traffic scene perception and planning in vehicle-centric settings. While these specialist models achieve strong performance within their target domains, their ability to generalize to the distinct viewpoints, constraints, and traffic semantics encountered by cyclists remains an open question. By benchmarking these specialized models alongside general-purpose models, we not only provide a rigorous assessment of their capabilities, but also establish a foundation for future cyclist-assistive intelligent systems.

3 BENCHMARK CONSTRUCTION

We first outline the design principles underlying CyclingVQA, followed by a detailed description of the benchmark tasks (§ 3.1) and the question–answer generation pipeline (§ 3.2).

Background. Our objective is to evaluate the cyclist-assistive traffic scene understanding capabilities of VLMs from a cyclist’s egocentric perspective. To ground the benchmark in real-world

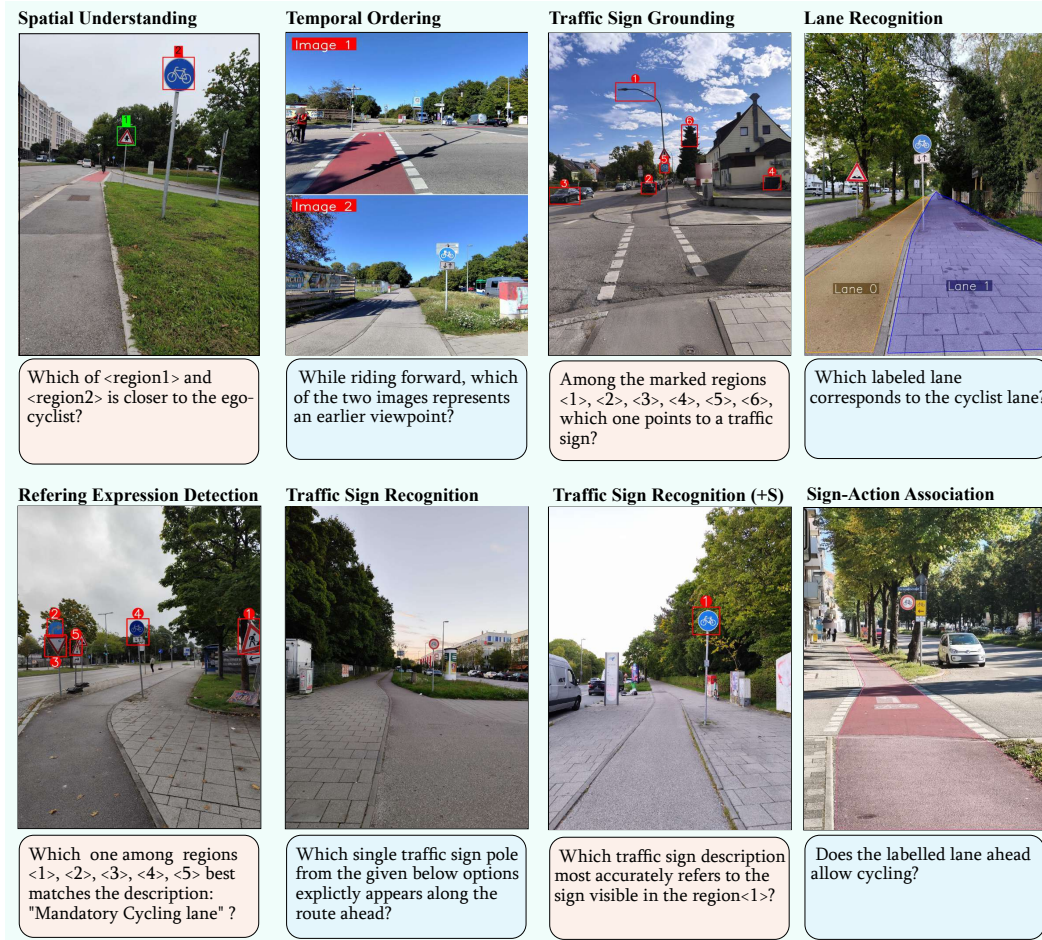


Figure 2: **Benchmark tasks.** Illustration of the eight benchmark tasks in CyclingVQA, showing example question prompts together with visual inputs augmented by lane annotations and bounding-box supervision.

traffic regulations, we draw on cycling-specific rulebooks from authoritative sources, including the Deutscher Verkehrssicherheitsrat safety booklet ([Deutscher Verkehrssicherheitsrat \(German Road Safety Council\), 2022](#)), the Vienna Convention on Road Signs and Signals ([Wikipedia contributors, 2024](#)), and European urban cycling standards ([Hiron et al., 2014](#); [Unfallforschung der Versicherer \(UDV\) et al., 2016](#); [Allgemeiner Deutscher Fahrrad-Club \(ADFC\) e. V., 2022](#)). These standards specify cyclist-specific guidance for lane usage, right-of-way, and interactions with pedestrians. Notably, cyclist-specific signage (see Table 10 in the Appendix) follows design principles analogous to standard motor vehicle signage. This is evident, for example, in the consistent use of blue backgrounds for regulatory signs ($\text{P} \rightarrow \text{B}$) and red borders for prohibitory signs ($\text{A} \rightarrow \text{R}$). This structural consistency, together with the substantial overlap in task demands between cycling and driving, motivates an investigation into whether VLMs specialized for autonomous driving can be applied to cyclist-centric traffic understanding through straightforward prompting. We examine this question through the dedicated benchmark described below.

3.1 BENCHMARK TASKS

Inspired by recent fine-grained visual question answering benchmarks ([Fu et al., 2024](#); [Danish et al., 2025](#); [Li et al., 2025b](#)), we design CyclingVQA, in a similar vein, to evaluate cyclist-centric perception and navigation-relevant reasoning. That is, the benchmark probes whether VLMs can (1) recognize and spatially localize cyclist-relevant traffic elements, and (2) reason about traffic sign-to-lane associations required for safe navigation.

| Task Name | #Im. | #Qu. | #Ch. | Anno. | Markers |
|--------------------------------|------------|-------------|------|-------|---------|
| Domain-Independent 🚴 | | | | | |
| 🚦 Traffic Sign Grounding | 410 | 410 | 6 | 👤 | 📦 Box |
| ↔️ Spatial Understanding | 153 | 179 | 2 | 👤 | 📦 Box |
| 🕒 Temporal Ordering | 110 | 111 | 2 | 👤 | – |
| Domain-Specific 🚴 | | | | | |
| 🏍️ Lane Recognition | 157 | 157 | 2 | 👤 | Lane |
| 📄 Referring Expr. Det. | 188 | 188 | 2 | 👤 | 📦 Box |
| 🔍 Traffic Sign Recog. (w/ B.B) | 259 | 259 | 7 | 🧠 + 👤 | 📦 Box |
| 📄 Traffic Sign Recog. | 465 | 465 | 7 | 🧠 + 👤 | – |
| 🚦 Sign-Action Assoc. | 225 | 240 | 2 | 🧠 + 👤 | Lane |
| Total (All Tasks) | 695 | 2009 | – | – | – |

Table 1: Summary of CyclingVQA tasks.

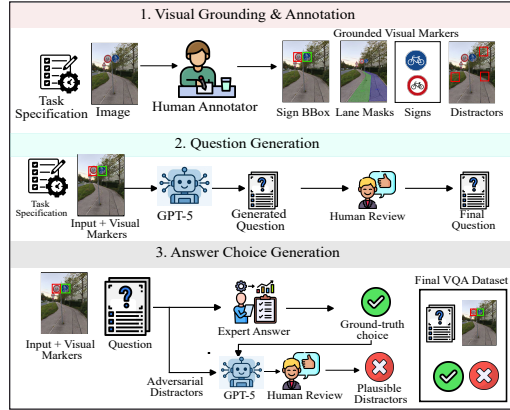


Figure 3: Overview of our annotation pipeline.

Concretely, CyclingVQA comprises eight task categories: three general visual reasoning tasks and five cyclist-specific domain tasks. Collectively, these tasks assess visual grounding, spatial and temporal understanding, and navigation-relevant decision making from a cyclist’s egocentric perspective. Each task is formulated as an image–question pair and may additionally include visual supervision, such as bounding boxes for traffic signs or pixel-level lane segmentation. Representative examples of all task categories are shown in Figure 2.

Domain-Independent Tasks. These tasks evaluate general spatial and temporal understanding and do not require cycling-specific knowledge.

- **Traffic Sign Grounding (TSG):** This task measures the ability to locate traffic signs in an image by selecting the correct bounding box.
- **Spatial Understanding (SU):** This task evaluates understanding of the relative positions of two traffic signs, such as which one is closer or whether one is to the left or right of the other.
- **Temporal Ordering (TO):** This task tests spatio-temporal reasoning by determining the correct order of two nearby frames that contain traffic-related elements.

Domain-Specific Tasks. These tasks focus on cyclist-specific perception and reasoning that are important for real-world navigation and safety.

- **Lane Recognition (LR):** This task tests lane-rule association i.e., whether the model can identify the cyclist sign and select the lane intended for cyclists among multiple lane options.
- **Referring Expression Detection (RED):** This task requires matching a textual description of a traffic sign to the correct bounding box in the image.
- **Traffic Sign Recognition with Box Supervision (TSR+S):** Given a single bounding box highlighting a traffic sign, this task evaluates whether the model can choose the correct semantic description of the sign.
- **Traffic Sign Recognition (TSR):** This task removes bounding box supervision and evaluates whether the model can recognize and interpret traffic signs directly from the full image.
- **Sign-Action Association (SAA):** This task evaluates navigation-related reasoning by requiring the model to recognize a traffic sign and decide whether a cyclist is allowed to proceed on a given user-labeled lane.

3.2 QA GENERATION

We collect **695** diverse urban street images from a cyclist’s egocentric perspective, covering a wide range of cycling-related traffic scenes and signage. All images are captured at a high resolution of approximately 3060×4080 pixels to ensure that traffic signs remain legible even at long distances,

| Model | Size | Type | Release | General Tasks | | | Domain-Specific Tasks | | | | | | Avg | Rank |
|---|------|----------|---------|---------------|-------------|-------------|-----------------------|-------------|-------------|-------------|-------------|-------------|-----|------|
| | | | | SU | TSG | TO | TSR+S | RED | TSR | LR | SAA | | | |
| Random | - | - | - | 49.2 | 16.3 | 51.4 | 12.9 | 42.6 | 13.5 | 36.9 | 51.2 | 34.3 | 31 | |
| Proprietary VLMs | | | | | | | | | | | | | | |
| Gemini-2.5-Flash(Comanici et al., 2025) | N/A | Reason | 07/2025 | 77.7 | 98.0 | 55.0 | 82.4 | 94.7 | 83.8 | 72.6 | 90.0 | 81.8 | 1 | |
| GPT-5.1(Team, 2025) | N/A | Reason | 11/2025 | 63.7 | 90.0 | 58.6 | 83.0 | 94.1 | 85.3 | 59.9 | 86.2 | 77.6 | 2 | |
| 🌐 Generalist VLMs | | | | | | | | | | | | | | |
| Eagle2.5-8B(Chen et al., 2025) | 8B | Instruct | 04/2025 | 53.6 | 87.8 | 51.4 | 50.5 | 82.4 | 39.8 | 53.5 | 82.5 | 62.7 | 12 | |
| InternVL3(Zhu et al., 2025) | 8B | Instruct | 04/2025 | 51.4 | 88.0 | 49.5 | 49.5 | 82.4 | 48.6 | 42.0 | 72.1 | 60.5 | 15 | |
| InternVL3.5-2B(Wang et al., 2025) | 2B | Instruct | 08/2025 | 62.6 | 77.6 | 53.2 | 47.5 | 74.5 | 49.0 | 42.7 | 65.8 | 59.1 | 17 | |
| InternVL3.5-8B(Wang et al., 2025) | 8B | Instruct | 08/2025 | 57.5 | 88.8 | 53.2 | 63.4 | 87.2 | 62.9 | 61.8 | 78.8 | 69.2 | 6 | |
| Molmo2-8B(Clark et al., 2025) | 8B | Instruct | 12/2025 | 56.4 | 88.0 | 51.4 | 34.2 | 73.9 | 37.5 | 48.4 | 44.2 | 54.2 | 23 | |
| Ovis2.5-2B(Lu et al., 2025) | 2B | Instruct | 08/2025 | 72.1 | 96.3 | 49.5 | 63.0 | 85.1 | 58.7 | 45.9 | 78.3 | 68.6 | 7 | |
| Ovis2.5-9B(Lu et al., 2025) | 9B | Instruct | 08/2025 | 70.9 | 97.3 | 51.4 | 81.5 | 93.6 | 72.2 | 62.4 | 79.6 | 76.1 | 4 | |
| Phi-4(Abouelenin et al., 2025) | 8B | Instruct | 02/2025 | 44.1 | 73.7 | 49.5 | 59.1 | 76.6 | 59.5 | 50.3 | 73.3 | 60.8 | 14 | |
| Qwen2.5-VL(Bai et al., 2025a) | 7B | Instruct | 02/2024 | 52.5 | 81.5 | 51.4 | 43.9 | 71.3 | 43.2 | 48.4 | 77.1 | 58.7 | 18 | |
| Qwen3-VL(Bai et al., 2025b) | 2B | Instruct | 11/2025 | 50.8 | 97.6 | 52.3 | 74.0 | 84.6 | 74.1 | 36.9 | 72.9 | 67.9 | 8 | |
| Qwen3-VL(Bai et al., 2025b) | 8B | Instruct | 11/2025 | 75.4 | 89.3 | 51.4 | 78.5 | 95.2 | 81.1 | 58.0 | 80.8 | 76.2 | 3 | |
| LLaVA-Next(Liu et al., 2024a) | 8B | Instruct | 04/2024 | 44.1 | 36.1 | 51.4 | 25.2 | 54.3 | 37.5 | 27.4 | 34.2 | 38.8 | 28 | |
| LLaVA-OneVision(Li et al., 2024) | 7B | Instruct | 06/2024 | 54.2 | 65.4 | 48.6 | 37.6 | 66.0 | 34.0 | 36.3 | 71.2 | 51.7 | 24 | |
| LLaVA-1.6(Liu et al., 2024a) | 7B | Instruct | 12/2023 | 46.9 | 23.2 | 51.4 | 26.2 | 47.9 | 32.4 | 15.9 | 59.6 | 37.9 | 29 | |
| 📍 Spatial-Aware VLMs | | | | | | | | | | | | | | |
| PerceptionLM(Cho et al., 2025) | 3B | Instruct | 04/2025 | 55.9 | 87.6 | 51.4 | 49.7 | 86.2 | 66.0 | 42.7 | 35.4 | 59.3 | 16 | |
| PerceptionLM(Cho et al., 2025) | 8B | Instruct | 04/2025 | 78.8 | 95.1 | 48.6 | 68.4 | 86.2 | 78.8 | 66.2 | 58.8 | 72.6 | 5 | |
| FoundationMotion(Bai et al., 2025a) | 7B | Instruct | 12/2025 | 49.2 | 82.4 | 51.4 | 43.2 | 69.7 | 39.4 | 59.9 | 67.5 | 57.8 | 19 | |
| SenseNova(Cai et al., 2025) | 8B | Instruct | 10/2025 | 78.8 | 70.5 | 51.4 | 48.4 | 83.0 | 49.4 | 49.7 | 68.3 | 62.4 | 13 | |
| SpatialReasoner(Ma et al., 2025b) | 7B | Reason | 04/2025 | 37.4 | 55.1 | 44.1 | 33.5 | 54.3 | 30.9 | 57.3 | 55.8 | 46.1 | 27 | |
| SpatialThinker(Batra et al., 2025) | 7B | Reason | 11/2025 | 58.1 | 95.1 | 50.5 | 57.4 | 87.2 | 47.9 | 43.3 | 71.7 | 63.9 | 11 | |
| VST(Yang et al., 2025) | 7B | Reason | 11/2025 | 78.2 | 72.4 | 51.4 | 32.5 | 58.0 | 30.5 | 51.0 | 72.9 | 55.9 | 21 | |
| 🚗 Driving-Centric VLMs | | | | | | | | | | | | | | |
| Cosmos-Reason1(Azzolini et al., 2025) | 7B | Reason | 03/2025 | 45.3 | 61.0 | 54.1 | 35.3 | 64.9 | 42.5 | 52.2 | 79.2 | 54.3 | 22 | |
| Cosmos-Reason2(NVIDIA, 2025) | 8B | Reason | 12/2025 | 52.5 | 79.3 | 53.2 | 73.5 | 87.8 | 62.9 | 56.1 | 70.4 | 67.0 | 9 | |
| DriveLMMo1(Ishaq et al., 2025) | 8B | Reason | 03/2025 | 57.5 | 76.3 | 50.5 | 43.2 | 72.9 | 46.3 | 42.7 | 70.4 | 57.5 | 20 | |
| DriveMM(Huang et al., 2025) | 7B | Instruct | 12/2024 | 54.7 | 54.1 | 51.4 | 30.1 | 60.6 | 29.3 | 45.2 | 67.5 | 49.1 | 26 | |
| ReCogDrive(Li et al., 2025a) | 8B | Instruct | 06/2025 | 48.0 | 50.2 | 51.4 | 37.0 | 54.3 | 38.6 | 53.5 | 64.6 | 49.7 | 25 | |
| Dolphins(Ma et al., 2023) | 7B | Instruct | 12/2023 | 46.4 | 15.4 | 36.9 | 14.0 | 44.1 | 16.2 | 49.7 | 70.8 | 36.7 | 30 | |

Table 2: **Evaluation of VLMs on the CyclingVQA benchmark.** We report accuracy (%) across eight tasks, observing that generalist models outperform driving-centric VLMs overall.

which is critical for early decision-making. An overview of the benchmark construction pipeline is shown in Figure 3.

Each image is first manually annotated with task-relevant traffic sign bounding boxes and labels, lane segmentations, and relative depth cues. For each task category, we then prompt GPT-5 (Team, 2025) to generate up to five closed-ended question templates (see Table 9 in the Appendix) with varied phrasings to reduce prompt sensitivity. Each question is subsequently paired with a single human-annotated ground-truth answer, while distractor options are strategically constructed to challenge the model’s reasoning. For non-binary tasks, GPT-5 (Team, 2025) is conditioned on the correct answer to generate semantically plausible distractors that require fine-grained scene analysis. Depending on the task, up to six answer choices are included to increase difficulty, and the position of the correct answer is randomized to mitigate positional bias.

Finally, and most importantly, questions are embedded within a system prompt that instructs models to provide both the selected option and a *concise explanation* (see Table 4 for the prompt). These explanations are used for downstream qualitative analysis of failure modes (§ 4.4). In total, CyclingVQA comprises **2,053** visual question–answer pairs spanning eight task categories, as summarized in Table 1. Additional statistical studies are provided in Appendix D.

4 EXPERIMENTAL RESULTS

In Section 4.2, we present the empirical results and summarize the key findings. Section 4.3 then examines the impact of Chain-of-Thought (CoT) prompting. In Section 4.4, we analyze the main failure modes, followed by a quantitative analysis of the generated tokens in Section 4.5.

4.1 EXPERIMENTAL SETUP

General-Purpose Models. We evaluate 16 open-source, general-purpose VLMs (1B–9B parameters) from the Qwen-VL (Bai et al., 2025a), Eagle (Chen et al., 2025), LLaVA-Nxt (Liu et al., 2024a), LLaVA-OneVision (Li et al., 2024), InternVL3 (Zhu et al., 2025), Phi (Abouelenin et al., 2025), Ovis (Lu et al., 2025), and Molmo (Clark et al., 2025) families, representing widely used baselines in 2025 for multimodal understanding.

Specialist Models. To assess domain-specific performance, we include 6 driving-focused VLMs: Cosmos-Reason{1,2} (Azzolini et al., 2025; NVIDIA, 2025), DriveLMM-O1 (Ishaq et al., 2025), RecogDrive (Li et al., 2025a), Dolphins (Ma et al., 2023), and DriveMM (Huang et al., 2025). Additionally, we evaluate 7 spatial-focused specialists, including PerceptionLM (Cho et al., 2025), VST (Yang et al., 2025), SpatialThinker (Batra et al., 2025), SenseNova (Cai et al., 2025), and SpatialReasoner (Ma et al., 2025b).

Proprietary Models. We also evaluate two frontier models: GPT-5.1 (Team, 2025) and Gemini-2.5-Flash (Comanici et al., 2025).

Evaluation Metric. In all cases, we report accuracy, defined as the percentage of correctly answered questions within each task category. We present both per-task accuracy and the overall mean accuracy averaged across all task categories.

Implementation. We employ Gemma2-9B (Team et al., 2024) as a parser (see Table 7) to extract discrete answer labels from generated responses, as some VLMs provide textual explanations without explicitly indicating a choice index. For inference, we allocate a maximum generation budget of 4,096 tokens per query. Additional implementation details are provided in Appendix C.

4.2 MAIN RESULTS

Table 2 summarizes performance across the eight CyclingVQA task categories. The proprietary model Gemini-2.5-Flash achieves the highest overall accuracy; several consistent trends from our evaluation are distilled below.

Generalist VLMs outperform driving-specialized models. Surprisingly, despite being fine-tuned on traffic environments, driving-specialized models consistently lag behind strong generalist VLMs on CyclingVQA. For example, specialist instruct models such as Dolphins (36.7%) and DriveMM (49.1%) perform substantially worse than generalist baselines on domain-specific tasks. Even reasoning-based driving models, including the latest Cosmos-Reason2 (67.0%) and DriveLMM-O1 (57.5%), are surpassed by generalist architectures such as Qwen3-VL-8B (76.2%) and Ovis2.5-9B (76.1%). These results, while unexpected, points to weak generalization capabilities of driving VLMs.

Instruct models outperform reasoning models. Across the six reasoning models evaluated, performance consistently trails that of instruction-following models. The strongest reasoning model, Cosmos-Reason2 (67.0%), ranks only ninth overall. In stark contrast, token-efficient instruct models dominate the leaderboard, with Qwen3-VL-8B (76.2%) ranking third and achieving the best overall performance among open-source VLMs.

Model scale is not a reliable predictor of performance. We do observe that larger models do not necessarily outperform smaller ones on CyclingVQA. Several recent compact models achieve performance comparable to—or exceeding—that of much larger older counterparts. For instance, Ovis2.5-2B (68.6%) and Qwen3-VL-2B (67.9%) outperform larger models such as Qwen2.5-7B (58.7%) and InternVL3-8B (60.5%). We hypothesize that superior training data quality in more recent models may have played an important role in this behavior.



Table 3: **Taxonomy of Failure Modes.** We characterize model errors across four recurring categories, providing a systematic overview of current VLM limitations in cyclist-centric scenarios. See Appendix E for a comprehensive qualitative analysis and additional case studies.

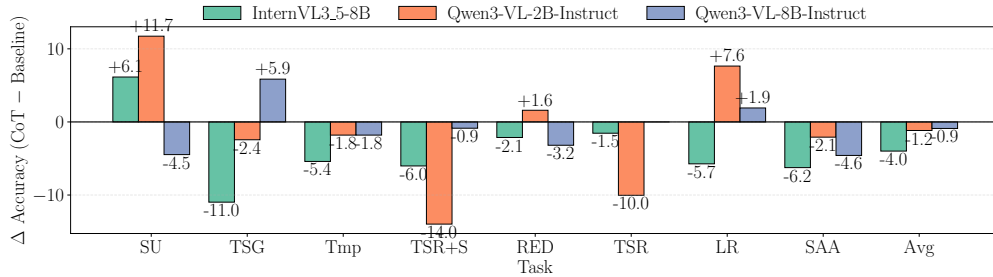


Figure 4: **CoT vs. Standard Prompting.** Overall performance degrades under CoT prompting across the three instruct models.

Driving-specialized models struggle with spatial understanding. Another interesting finding is that driving-centric models show particular weaknesses on spatial reasoning tasks. On the SU task, the best-performing driving-specialized model, DriveLMM-O1, achieves only 57.5% accuracy, barely exceeding the random baseline of 50%. In contrast, generalist models such as Qwen3-VL-8B reach 75.4%, while spatially enhanced architectures perform even better, with PerceptionLM-8B achieving the highest accuracy of 78.8%.

Temporal reasoning remains challenging. Temporal reasoning emerges as the most difficult task in our benchmark. Even frontier models such as Gemini-2.5-Flash achieve only 58.6% accuracy, again only marginally above the random baseline. This result underscores that current models lack the ability to consistently pivot their reasoning around visible traffic signs across frames. Indeed, these findings warrant deeper investigation into the limitations of current VLMs in physical and temporal understanding.

4.3 IMPACT OF CoT-PROMPTING

We examine the effect of explicit CoT prompting (See Table 6) by comparing baseline and CoT performance on three instruct models. As shown in Figure 4, CoT prompting does not lead to consistent improvements and instead results in average accuracy drops of 4.0%, 1.2%, and 0.9% for InternVL3.5-8B, Qwen3-VL-2B, and Qwen3-VL-8B, respectively. Although CoT yields occasional gains on certain spatial reasoning tasks, these improvements are outweighed by performance degradation on grounding and association tasks. Overall, the predominantly negative deltas suggest that, for cycling-specific VQA, direct instruction-following is more reliable than explicit multi-step reasoning with the current models; however, we emphasize that this finding is not conclusive and may reflect limitations of current prompting strategies rather than an inherent weakness of CoT reasoning.

4.4 QUALITATIVE ANALYSIS OF FAILURE MODES

While quantitative results are informative, they do not tell the whole story. We therefore turn to a more pressing question: what are the key failure modes that limit cyclist-assistive reasoning in current VLMs? To this end, we identify four recurring error modes related to spatial reasoning and safety-critical decision-making. Table 3 summarizes these failure modes, with additional qualitative examples in Appendix E.

1. Misinterpretation of traffic rule semantics. A frequent source of error occurs when models incorrectly interpret the semantics of traffic signs. While signs are often detected and localized correctly, their regulatory meaning is misclassified, resulting in unsafe or invalid action predictions. For instance, prohibitory signs are sometimes mistaken for permissive or mandatory indicators (Table 3(a)). These failures suggest that translating visual sign cues into cyclist-specific traffic rules remains unreliable.

2. Perception failures. Models also exhibit failures at the perception stage, particularly in visually cluttered environments or when signs are small, distant, or partially occluded (Table 3(b)). When critical signage is missed, models often fall back on contextual priors, occasionally hallucinating cues that are not supported by the visual evidence. Such behavior leads to confident but incorrect decisions that ignore safety-relevant constraints.

3. Lane-sign association errors. Even when traffic signs are correctly detected, models frequently struggle to associate them with the appropriate road lane. This issue is most pronounced at intersections or shared-use paths where multiple, potentially conflicting, signs appear in close proximity (Table 3(c)). In these cases, models may recognize all signs but fail to determine which applies to the cyclist’s trajectory, revealing weaknesses in spatial association. Moreover, models tend to rely on the physical location of traffic poles, leading to incorrect lane associations (see Table 13).

4. Incorrect attribution of relevant signage. Finally, in scenes containing multiple visible signs, models sometimes base decisions on irrelevant signage while overlooking the sign that directly governs the cyclist’s path (Table 3(d)). This failure mode suggests that models struggle to correctly rank the relevance of multiple signs, frequently conditioning their decisions on cues that do not apply to the cyclist’s intended trajectory.

4.5 GENERATION VERBOSITY VS. PERFORMANCE

Generation verbosity is not strongly correlated with performance. Figure 5 reports the average number of generated tokens per datapoint for each model and task. To account for differences in tokenization across VLMs, we tokenize all generated outputs using the GPT-2 tokenizer (Radford et al., 2019). Among models explicitly optimized for long-form reasoning, Cosmos-Reason1 (Azzolini et al., 2025) produces the highest average of 546 tokens per response, yet ranks only 21st on the leaderboard in terms of overall task performance. In contrast, Qwen3-VL-8B achieves the top rank among open-source VLMs while generating just 50 tokens on average. These results suggest that increased generation verbosity does not necessarily translate to improved performance.

spatial cues, and ensure appropriate privacy safeguards. Scaling the dataset through automated labeling, while preserving annotation quality, therefore remains an important goal for future work.

Dataset Diversity. Data collection was conducted predominantly under low-traffic conditions, primarily due to privacy considerations. As a result, the dataset contains fewer instances of dense interactions involving other cyclists, pedestrians, or motor vehicles. We nevertheless include few scenarios in which cyclist intentionally navigate into pedestrian lanes to reflect unsafe behavior. Looking ahead, a broader and more varied set of such unsafe scenarios through generative models could further improve the benchmark diversity.

Geographical Coverage. All data are sourced from Munich, a city known for its highly structured cycling infrastructure ([Clean Cities Campaign, 2025](#)). We acknowledge that it may not fully capture the diversity of cycling environments worldwide, such as rural roads, regions with informal traffic norms, or cities with less developed infrastructure. Expanding geographic coverage is therefore imperative for future work.

Evaluation Framework. Our evaluation currently focuses on objective accuracy over closed-ended choices and does not include a quantitative assessment of model reasoning traces. Although the manual error analysis in Section 4.4 offers qualitative insights, future work could explore more fine-grained metrics for evaluating the correctness and usefulness of generated explanations. In addition, while our initial experiments with chain-of-thought prompting (Section 4.3) did not yield improvements, a more systematic investigation of prompt sensitivity remains an open direction.

Beyond VLMs. Finally, our study focuses exclusively on VLMs. Recent Vision–Language–Action (VLA) models have shown promising capabilities in directly predicting navigation actions ([Zhou et al., 2025](#); [Jiang et al., 2025b](#); [Hu et al., 2025](#)) in autonomous driving scenarios. Extending cyclist-assistive evaluation to such models represents an exciting avenue for future research.

REFERENCES

- Abdelrahman Abouelenin, Atabak Ashfaq, Adam Atkinson, Hany Awadalla, Nguyen Bach, Jianmin Bao, Alon Benham, Martin Cai, Vishrav Chaudhary, Congcong Chen, et al. Phi-4-mini technical report: Compact yet powerful multimodal language models via mixture-of-loras. *arXiv preprint arXiv:2503.01743*, 2025.
- Allgemeiner Deutscher Fahrrad-Club (ADFC) e. V. Traffic rules: Safe bicycle rides. Online PDF, 2022. URL https://www.adfc.de/fileadmin/user_upload/Im-Alltag/Engagement_vor_Ort/Initiativen_fuer_Gefluechtete/Download/Sicher_unterwegs_mit_dem_Fahrrad_Englisch_2022.pdf. English brochure on safe cycling rules in Germany.
- Alisson Azzolini, Junjie Bai, Hannah Brandon, Jiaxin Cao, Prithvijit Chattopadhyay, Huayu Chen, Jinju Chu, Yin Cui, Jenna Diamond, Yifan Ding, et al. Cosmos-reason1: From physical common sense to embodied reasoning. *arXiv preprint arXiv:2503.15558*, 2025.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibong Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025a.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibong Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen3-vl technical report. *arXiv preprint arXiv:2511.21631*, 2025b.
- Hunar Batra, Haoqin Tu, Hardy Chen, Yuanze Lin, Cihang Xie, and Ronald Clark. Spatialthinker: Reinforcing 3d reasoning in multimodal llms via spatial rewards. *arXiv preprint arXiv:2511.07403*, 2025.
- 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 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11621–11631, 2020.
- Zhongang Cai, Ruisi Wang, Chenyang Gu, Fanyi Pu, Junxiang Xu, Yubo Wang, Wanqi Yin, Zhitao Yang, Chen Wei, Qingping Sun, et al. Scaling spatial intelligence with multimodal foundation models. *arXiv preprint arXiv:2511.13719*, 2025.

-
- Xu Cao, Tong Zhou, Yunsheng Ma, Wenqian Ye, Can Cui, Kun Tang, Zhipeng Cao, Kaizhao Liang, Ziran Wang, James M Rehg, et al. Maplm: A real-world large-scale vision-language benchmark for map and traffic scene understanding. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 21819–21830, 2024.
- Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14455–14465, 2024.
- Guo Chen, Zhiqi Li, Shihao Wang, Jindong Jiang, Yicheng Liu, Lidong Lu, De-An Huang, Wonmin Byeon, Matthieu Le, Tuomas Rintamaki, et al. Eagle 2.5: Boosting long-context post-training for frontier vision-language models. *arXiv preprint arXiv:2504.15271*, 2025.
- An-Chieh Cheng, Hongxu Yin, Yang Fu, Qiushan Guo, Ruihan Yang, Jan Kautz, Xiaolong Wang, and Sifei Liu. Spatialrgpt: Grounded spatial reasoning in vision-language models. *Advances in Neural Information Processing Systems*, 37:135062–135093, 2024.
- Jang Hyun Cho, Andrea Madotto, Effrosyni Mavroudi, Triantafyllos Afouras, Tushar Nagarajan, Muhammad Maaz, Yale Song, Tengyu Ma, Shuming Hu, Suyog Jain, et al. Perceptionlm: Open-access data and models for detailed visual understanding. *arXiv preprint arXiv:2504.13180*, 2025.
- Christopher Clark, Jieyu Zhang, Zixian Ma, Jae Sung Park, Mohammadreza Salehi, Rohun Tripathi, Sangho Lee, Zhongzheng Ren, Chris Dongjoo Kim, YINUO Yang, Vincent Shao, Yue Yang, Weikai Huang, Ziqi Gao, Taira Anderson, Jianrui Zhang, Jitesh Jain, George Stoica, Winson Han, Ali Farhadi, and Ranjay Krishna. Molmo2: Open weights and data for vision-language models with video understanding and grounding. <https://www.datocms-assets.com/64837/1766008501-molmo2-tech-report.pdf>, 2025.
- Clean Cities Campaign. City ranking 2025: Streets for kids, cities for all – protected cycling infrastructure, 2025. URL <https://cleancitiescampaign.org/city-ranking-2025-cycling-infrastructure/>. Accessed: 2025-12-31; Ranking of European cities based on protected cycling infrastructure.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025.
- Charles Corbière, Simon Roburin, Syrielle Montariol, Antoine Bosselut, and Alexandre Alahi. Retrieval-based interleaved visual chain-of-thought in real-world driving scenarios. *arXiv preprint arXiv:2501.04671*, 2025.
- Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, Yang Zhou, Kaizhao Liang, Jintai Chen, Juanwu Lu, Zichong Yang, Kuei-Da Liao, et al. A survey on multimodal large language models for autonomous driving. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pp. 958–979, 2024.
- Muhammad Danish, Muhammad Akhtar Munir, Syed Roshan Ali Shah, Kartik Kuckreja, Fahad Shahbaz Khan, Paolo Fraccaro, Alexandre Lacoste, and Salman Khan. Geobench-vlm: Benchmarking vision-language models for geospatial tasks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 7132–7142, 2025.
- Deutscher Verkehrssicherheitsrat (German Road Safety Council). Cycling in germany: Safety brochure, 2022. Available online: https://www.germanroadsafety.de/fileadmin/user_upload/Broschueren/Fahrrad_fahren_in_Deutschland/Cycling-in-Germany_English_2022.pdf, howpublished = German Road Safety Council Publication,.
- Euronews. Cycling in europe: Which countries and cities are the most and least bicycle-friendly?, September 2023. URL <https://www.euronews.com/next/2023/09/19/cycling-in-europe-which-countries-and-cities-are-the-most-and-least-bicycle-friendly>. Accessed: 2025-03-08.

-
- Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A Smith, Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but not perceive. In *European Conference on Computer Vision*, pp. 148–166. Springer, 2024.
- Yulu Gan, Ligeng Zhu, Dandan Shan, Baifeng Shi, Hongxu Yin, Boris Ivanovic, Song Han, Trevor Darrell, Jitendra Malik, Marco Pavone, et al. Foundationmotion: Auto-labeling and reasoning about spatial movement in videos. *arXiv preprint arXiv:2512.10927*, 2025.
- Anurag Ghosh, Shen Zheng, Robert Tamburo, Khiem Vuong, Juan Alvarez-Padilla, Hailiang Zhu, Michael Cardei, Nicholas Dunn, Christoph Mertz, and Srinivasa G Narasimhan. Roadwork: A dataset and benchmark for learning to recognize, observe, analyze and drive through work zones. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 6132–6142, 2025.
- Benoît Hiron, André Isler, and François Tortel. Signs and signals for cyclists and pedestrians: comparison of rules and practices in 13 countries. In *Transport Research Arena (TRA) 5th Conference: Transport Solutions from Research to Deployment*. European Commission, CEDR, ERTRAC, WATERBORNE TP, ERRAC, IFSTTAR, 2014.
- Tianshuai Hu, Xiaolu Liu, Song Wang, Yiyao Zhu, Ao Liang, Lingdong Kong, Guoyang Zhao, Zeying Gong, Jun Cen, Zhiyu Huang, et al. Vision-language-action models for autonomous driving: Past, present, and future. *arXiv preprint arXiv:2512.16760*, 2025.
- Zhijian Huang, Chengjian Feng, Feng Yan, Baihui Xiao, Zequn Jie, Yujie Zhong, Xiaodan Liang, and Lin Ma. Robotron-drive: All-in-one large multimodal model for autonomous driving. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 8011–8021, October 2025.
- ZIV German Bicycle Industry. Single bicycle accidents, 2024. URL <https://www.ziv-zweirad.de/en/single-bicycle-accidents/>. Accessed: 2025-03-08.
- Ayesha Ishaq, Jean Lahoud, Ketan More, Omkar Thawakar, Ritesh Thawkar, Dinura Dissanayake, Noor Ahsan, Yuhao Li, Fahad Shahbaz Khan, Hisham Cholakkal, et al. DriveImm-o1: A step-by-step reasoning dataset and large multimodal model for driving scenario understanding. *arXiv preprint arXiv:2503.10621*, 2025.
- Mengdi Jia, Zekun Qi, Shaochen Zhang, Wenyao Zhang, Xinqiang Yu, Jiawei He, He Wang, and Li Yi. Omnispatial: Towards comprehensive spatial reasoning benchmark for vision language models. *arXiv preprint arXiv:2506.03135*, 2025.
- Bo Jiang, Shaoyu Chen, Qian Zhang, Wenyu Liu, and Xinggong Wang. Alphadrive: Unleashing the power of vlms in autonomous driving via reinforcement learning and reasoning. *arXiv preprint arXiv:2503.07608*, 2025a.
- Sicong Jiang, Zilin Huang, Kangan Qian, Ziang Luo, Tianze Zhu, Yang Zhong, Yihong Tang, Menglin Kong, Yunlong Wang, Siwen Jiao, et al. A survey on vision-language-action models for autonomous driving. *arXiv preprint arXiv:2506.24044*, 2025b.
- Amita Kamath, Jack Hessel, and Kai-Wei Chang. What's "up" with vision-language models? investigating their struggle with spatial reasoning. *arXiv preprint arXiv:2310.19785*, 2023.
- Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring to objects in photographs of natural scenes. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pp. 787–798, 2014.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024.
- Yongkang Li, Kaixin Xiong, Xiangyu Guo, Fang Li, Sixu Yan, Gangwei Xu, Lijun Zhou, Long Chen, Haiyang Sun, Bing Wang, et al. Recogdrive: A reinforced cognitive framework for end-to-end autonomous driving. *arXiv preprint arXiv:2506.08052*, 2025a.

-
- Yue Li, Meng Tian, Zhenyu Lin, Jiangtong Zhu, Dechang Zhu, Haiqiang Liu, Yueyi Zhang, Zhiwei Xiong, and Xinhai Zhao. Fine-grained evaluation of large vision-language models in autonomous driving. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9431–9442, 2025b.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Lllanext: Improved reasoning, ocr, and world knowledge, 2024a.
- Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xu-Cheng Yin, Cheng-Lin Liu, Lianwen Jin, and Xiang Bai. Ocrbench: on the hidden mystery of ocr in large multimodal models. *Science China Information Sciences*, 67(12):220102, 2024b.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023.
- Shiyin Lu, Yang Li, Yu Xia, Yuwei Hu, Shanshan Zhao, Yanqing Ma, Zhichao Wei, Yinglun Li, Lunhao Duan, Jianshan Zhao, et al. Ovis2. 5 technical report. *arXiv preprint arXiv:2508.11737*, 2025.
- Wufei Ma, Haoyu Chen, Guofeng Zhang, Yu-Cheng Chou, Jieneng Chen, Celso de Melo, and Alan Yuille. 3dsrbench: A comprehensive 3d spatial reasoning benchmark. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 6924–6934, 2025a.
- Wufei Ma, Yu-Cheng Chou, Qihao Liu, Xingrui Wang, Celso de Melo, Jianwen Xie, and Alan Yuille. Spatialreasoner: Towards explicit and generalizable 3d spatial reasoning. *arXiv preprint arXiv:2504.20024*, 2025b.
- Yingzi Ma, Yulong Cao, Jiachen Sun, Marco Pavone, and Chaowei Xiao. Dolphins: Multimodal language model for driving, 2023.
- NVIDIA. Cosmos-reason2. <https://github.com/nvidia-cosmos/cosmos-reason2>, 2025. Accessed: 2025-01-30.
- Michael Ogezi and Freda Shi. Spare: Enhancing spatial reasoning in vision-language models with synthetic data. *arXiv preprint arXiv:2504.20648*, 2025.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Yifan Shen, Yuanzhe Liu, Jingyuan Zhu, Xu Cao, Xiaofeng Zhang, Yixiao He, Wenming Ye, James Matthew Rehg, and Ismini Lourentzou. Fine-grained preference optimization improves spatial reasoning in vlms. *arXiv preprint arXiv:2506.21656*, 2025.
- Ilias Stogiannidis, Steven McDonagh, and Sotirios A Tsaftaris. Mind the gap: Benchmarking spatial reasoning in vision-language models. *arXiv preprint arXiv:2503.19707*, 2025.
- Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 2446–2454, 2020.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Alexandre Shahriari, et al. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
- Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*, 2025.
- OpenAI Team. The gpt-5 system card. <https://cdn.openai.com/pdf/8124a3ce-ab78-4f06-96eb-49ea29ffb52f/gpt5-system-card-aug7.pdf>, 2025.

-
- Kexin Tian, Jingrui Mao, Yunlong Zhang, Jiwan Jiang, Yang Zhou, and Zhengzhong Tu. Nuscenespatialqa: A spatial understanding and reasoning benchmark for vision-language models in autonomous driving. *arXiv preprint arXiv:2504.03164*, 2025.
- Xiaoyu Tian, Junru Gu, Bailin Li, Yicheng Liu, Yang Wang, Zhiyong Zhao, Kun Zhan, Peng Jia, Xianpeng Lang, and Hang Zhao. Drivevlm: The convergence of autonomous driving and large vision-language models. *arXiv preprint arXiv:2402.12289*, 2024.
- Unfallforschung der Versicherer (UDV), Deutsche Verkehrswacht e.V., and Deutscher Verkehrssicherheitsrat (German Road Safety Council). Cycling in germany – the most important rules. Online PDF, 2016. URL https://willkommen-in-magdeburg.de/wp-content/uploads/Fahrrad_fahren_in_Deutschland_englisch.pdf. English edition; published brochure.
- Jiayu Wang, Yifei Ming, Zhenmei Shi, Vibhav Vineet, Xin Wang, Sharon Li, and Neel Joshi. Is a picture worth a thousand words? delving into spatial reasoning for vision language models. *Advances in Neural Information Processing Systems*, 37:75392–75421, 2024.
- Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang Liu, Linglin Jing, Shenglong Ye, Jie Shao, et al. Internvl3. 5: Advancing open-source multimodal models in versatility, reasoning, and efficiency. *arXiv preprint arXiv:2508.18265*, 2025.
- Maolin Wei, Wanzhou Liu, and Eshed Ohn-Bar. Driveqa: Passing the driving knowledge test. *arXiv preprint arXiv:2508.21824*, 2025.
- Wikipedia contributors. Vienna convention on road signs and signals, 2024. URL https://en.wikipedia.org/wiki/Vienna_Convention_on_Road_Signs_and_Signals. Accessed: 2025-03-08.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 38–45, Online, October 2020. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/2020.emnlp-demos.6>.
- Shaoyuan Xie, Lingdong Kong, Yuhao Dong, Chonghao Sima, Wenwei Zhang, Qi Alfred Chen, Ziwei Liu, and Liang Pan. Are vlms ready for autonomous driving? an empirical study from the reliability, data, and metric perspectives. *arXiv preprint arXiv:2501.04003*, 2025.
- Shuo Xing, Hongyuan Hua, Xiangbo Gao, Shenzhe Zhu, Renjie Li, Kexin Tian, Xiaopeng Li, Heng Huang, Tianbao Yang, Zhangyang Wang, et al. Autotruster: Benchmarking trustworthiness in large vision language models for autonomous driving. *arXiv preprint arXiv:2412.15206*, 2024.
- Rui Yang, Ziyu Zhu, Yanwei Li, Jingjia Huang, Shen Yan, Siyuan Zhou, Zhe Liu, Xiangtai Li, Shuangye Li, Wenqian Wang, et al. Visual spatial tuning. *arXiv preprint arXiv:2511.05491*, 2025.
- Zhenjie Yang, Xiaosong Jia, Hongyang Li, and Junchi Yan. Llm4drive: A survey of large language models for autonomous driving. *arXiv preprint arXiv:2311.01043*, 2023.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9556–9567, 2024.
- Xingcheng Zhou, Mingyu Liu, Ekim Yurtsever, Bare Luka Zagar, Walter Zimmer, Hu Cao, and Alois C Knoll. Vision language models in autonomous driving: A survey and outlook. *IEEE Transactions on Intelligent Vehicles*, 2024.
- Xingcheng Zhou, Xuyuan Han, Feng Yang, Yunpu Ma, Volker Tresp, and Alois Knoll. Opendrivevla: Towards end-to-end autonomous driving with large vision language action model. *arXiv preprint arXiv:2503.23463*, 2025.

Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, et al. InternV13: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv preprint arXiv:2504.10479*, 2025.

A APPENDIX SUMMARY

In the following sections, we provide additional details, experiments, and visualizations to supplement our main findings.

Prompt Details (Section B): We provide the system prompts used in our benchmarks, including the various prompt templates and judge templates.

Implementation Details (Section C): We outline the experimental implementation details, such as exact model cards and model sources.

Additional Quantitative Results (Section 4.5): We present further experimental insights regarding the tokens generated by different models.

Additional Qualitative Results (Section E): We provide qualitative results through additional case studies and detailed output analysis.

B PROMPT DETAILS

This section describes the prompts and parsing strategies used across our experiments.

Prompting for Instruct Models: Table 4 lists the prompts used for instruct models (e.g., OVis2.5-8B (Lu et al., 2025)). The prompt instructs the models to return a choice followed by a brief justification. We find that this format is computationally efficient and facilitates interpretation of model decisions; however, some instruct models occasionally omit the justification despite these explicit instructions.

Prompting for Reasoning Models: Table 5 shows the prompts used for reasoning-oriented models (e.g., Cosmos-Reason1 (Azzolini et al., 2025)). For these models, we do not explicitly request intermediate reasoning, as it is intrinsic to their generation behavior.

CoT Prompting Ablation: As an ablation, we also evaluate Chain-of-Thought (CoT) prompting for instruct models (Table 6), which requires a step-by-step analysis before providing the final answer.

Response Parsing Strategy: To reduce parsing errors arising due to inconsistent formatting outputs, we use a dedicated choice-parsing model. Given the question, candidate options, and the raw model-generated output, the parsing model extracts a standardized final choice and the corresponding rationale. The parsing prompt is provided in Table 7.

Standard Prompt for Instruct Models

```
Answer the following multiple-choice question by selecting exactly one option.
```

```
Question: {question}
```

```
Options: {options}
```

```
Provide only the selected option letter, followed by a brief reasoning
```

Table 4: The instruction template used for instruct-tuned models.

Standard Prompt for Reasoning Models

```
Answer the following multiple-choice question by selecting exactly one option.
```

```
Question: {question}  
Options: {options}
```

Table 5: The instruction template used for reasoning-focused models.

CoT Prompt for Instruct Models

```
Answer the following multiple-choice question by selecting exactly one option.
```

```
Question: {question}  
Options: {options}
```

```
Let's think step by step first and then return the choice.
```

Table 6: The instruction template used for the Chain-of-Thought (CoT) ablation.

Parsing Prompt for Choice Extraction

```
You are a choice-parsing language model.  
Your task is to infer the final predicted choice from the generated text.
```

```
Rules:
```

- Valid choices: A, B, C, D, E, F, G, H
- If a clear final choice is stated, extract it.
- If multiple choices appear, select the FINAL one.
- If no choice is explicitly stated, infer the most confident option.
- Return ONLY a valid JSON object.

```
Question: {question}  
Options: {options}  
Generated Text: {generated_text}
```

```
{  
  "predicted_choice": <A/B/.../null>,  
  "predicted_reasoning": "Explanation"  
}
```

Table 7: The instruction template used for the choice-parsing model.

C IMPLEMENTATION

Target VLMs. We use the opensource models present on the HuggingFace (Wolf et al., 2020) or respective official GitHub sources for our experiments. We report details about the model card in Table 8. For reasoning models (eg., Cosmos-Reason2 (Azzolini et al., 2025)), we use a budget of 4096 tokens during generation and use greedy-decoding.

| Model | Scale | Hugging Face Model ID |
|------------------|-------|--------------------------------|
| Qwen3-VL | 2B | Qwen/Qwen3-VL-2B-Instruct |
| Qwen3-VL | 8B | Qwen/Qwen3-VL-8B-Instruct |
| Ovis2.5 | 2B | AIDC-AI/Ovis2.5-2B |
| Ovis2.5 | 9B | AIDC-AI/Ovis2.5-9B |
| InternVL3.5 | 2B | OpenGVLab/InternVL3.5-2B |
| InternVL3.5 | 8B | OpenGVLab/InternVL3.5-8B |
| Eagle2.5 | 8B | NEXT-GPT/Eagle2.5-8B |
| Phi-4 | 6B | microsoft/phi-4 |
| InternVL3 | 8B | OpenGVLab/InternVL3-8B |
| Molmo2 | 8B | allenai/Molmo-7B |
| FoundationMotion | 7B | WoWolf/Qwen2.5vl-7b-fm-tuned |
| LLaVA-OV | 7B | llava-hf/llava-onevision-7b-ov |
| LLaVA-Next | 8B | llava-hf/llava-next-8b |
| LLaVA-1.6 | 7B | liuhaotian/llava-v1.6-7b |
| PerceptionLM | 2B | facebook/Perception-LM-2B |
| PerceptionLM | 8B | facebook/Perception-LM-8B |
| Qwen2.5-VL | 7B | Qwen/Qwen2.5-VL-7B |
| SenseNova | 8B | sensenova/SenseNova-SI-1.1 |
| SpatialThinker | 7B | OX-PIXL/SpatialThinker-7B |
| VST | 7B | rayruiyang/VST-7B-RL |
| SpatialReasoner | 7B | ccvl/SpatialReasoner |
| Cosmos-Reason2 | 8B | nvidia/Cosmos-Reason2-8B |
| Cosmos-Reason1 | 7B | nvidia/Cosmos-Reason1-7B |
| DriveLMM-O1 | 7B | ayshaishaq/DriveLMMo1 |
| DriveMM | 7B | DriveMM/DriveMM |
| ReCogDrive | 7B | owl110/ReCogDrive-VLM |
| Dolphins | 7B | github.com/SaFo-Lab/Dolphins |

Table 8: **Target VLMs** evaluated in our study.

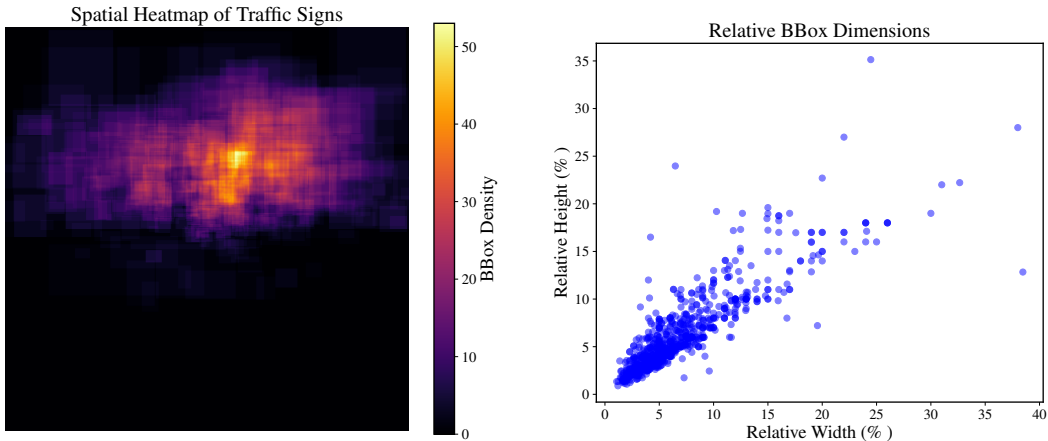


Figure 6: **Traffic sign bounding-box statistics.** Left: spatial heatmap of traffic-sign locations. Right: distribution of relative bounding-box widths and heights.

D ADDITIONAL DETAILS OF THE BENCHMARK DATASET

We visualize the spatial heatmap of traffic sign locations in Figure 6(a) and the relative area distribution of traffic signs in Figure 6(b). We also provide representative sample question templates used in our benchmark in Table 9.

E ADDITIONAL QUALITATIVE RESULTS

We present several case studies showing recurring errors for different tasks. These are shown in Tables 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, and 27. Note that, the generated text shown in these Figures are automatically paraphrased for brevity using GPT-5 (Team, 2025). We refer the reader to the project website for the full evaluation logs.

Questions

Traffic Sign Localization

- "In this image, the regions <1>, <2>, <3>, <4>, <5>, <6>, <7> are highlighted. Which highlighted region corresponds to a traffic sign?"
- "From the cyclist's perspective, several regions are marked in the image as <1>, <2>, <3>, <4>, <5>, <6>, <7>. Which region actually contains a traffic sign?"
- "From the cyclist's perspective, given the marked regions <1>, <2>, <3>, <4>, <5>, <6>, <7>, <8>, which one localizes a real traffic sign?"
- "From the cyclist's viewpoint, among the marked regions <1>, <2>, <3>, <4>, <5>, <6>, <7>, which region correctly points to a traffic sign?"
- "Looking at the highlighted regions <1>, <2>, <3>, <4>, <5>, <6>, which one contains the traffic sign?"

Temporal Ordering

- "From the cyclist's perspective while riding forward, which of the two images represents an earlier viewpoint?"
- "As the cyclist progressed along the road, which of these two views corresponds to the earlier moment?"
- "While the cyclist was moving forward, two snapshots were captured moments apart. Which image shows the earlier point in the journey?"
- "During continuous forward motion, the cyclist took these two photos. Which one was captured first in time?"
- "These two images were captured sequentially during the cyclist's forward travel. Which image was taken first?"

Spatial Scene Understanding

- "In the image, captured from the cyclist's perspective, focus on the regions marked as <region1> (highlighted by a green bounding box) and <region2> (highlighted by a red bounding box). Describe the spatial relationship between the object in <region1> and the object in <region2>."
- "In the image, captured from the cyclist's perspective, focus on the regions marked as <region1> (highlighted by a green bounding box) and <region2> (highlighted by a red bounding box). How is the object in <region1> positioned or related to the object in <region2>?"
- "In the image, captured from the cyclist's perspective, focus on the regions marked as <region1> (highlighted by a green bounding box) and <region2> (highlighted by a red bounding box). What is the relationship between the object in <region1> and the object in <region2>?"

Bicycle Lane Prediction

- "From the cyclist's perspective while riding forward, which of the two images represents an earlier viewpoint?"
- "As the cyclist progressed along the road, which of these two views corresponds to the earlier moment?"
- "While the cyclist was moving forward, two snapshots were captured moments apart. Which image shows the earlier point in the journey?"
- "During continuous forward motion, the cyclist took these two photos. Which one was captured first in time?"
- "These two images were captured sequentially during the cyclist's forward travel. Which image was taken first?"

Sign-Action Association

- "The image shows a forward-facing view from a cyclist's perspective. Based on the shaded path (in pink) shown ahead, should the cyclist proceed forward on this lane? The shading is externally drawn on the image as an overlay."
- "The image shows a forward-facing view from a cyclist's perspective. Based on the shaded route (in pink) ahead, should the cyclist enter this street? The shading is externally drawn on the image as an overlay."
- "The image shows a forward-facing view from a cyclist's perspective. From the cyclist's forward view, can one ride into the shaded street area (in pink)? The shading is externally drawn on the image as an overlay."
- "The image shows a forward-facing view from a cyclist's perspective. From the cyclist's forward view, is cycling allowed along the highlighted lane (in pink)? The highlight is externally drawn on the image as an overlay."

Traffic Sign Recognition (w/ B.B supervision)

- "In the image, focus on the region marked as <1> (highlighted by a red bounding box). Which traffic sign description most accurately refers to the sign visible in that region?"
- "In the image, focus on the region marked as <1> (highlighted by a red bounding box). Which specific traffic sign description correctly matches the sign visible in that region?"
- "In the image, focus on the region marked as <1> (highlighted by a red bounding box). What is the correct description of the traffic sign visible in that region?"
- "In the image, focus on the region marked as <1> (highlighted by a red bounding box). Which official traffic sign description best represents the sign that appears in that region?"
- "In the image, focus on the region marked as <1> (highlighted by a red bounding box). Which road sign description corresponds to the sign present in that region?"

Traffic Sign Recognition

- "In the image captured from the cyclist's forward viewpoint, which single traffic sign pole from the given below options explicitly appears along the route ahead?"

Referring Expression Detection

- "From the cyclist's perspective, focus on the bounding box regions marked <1>, <2>, <3>. Which region is being referred to as {"Regulatory sign indicating a mandatory bicycle-only lane ?"},"

Table 9: Sample question templates part of our benchmark.









| Traffic Sign(s) | Description |
|---|--|
|  | Sign 237 – Mandatory cycle path. Indicates a dedicated lane or path reserved exclusively for cyclists. |
|  | Sign 239 – Pedestrian path. Cycling is not permitted on this path. |
|  | Sign 240 – Shared path for pedestrians and cyclists. Both users share the same lane. |
|  | Sign 241-30 – Segregated pedestrian and cycle path. The path is divided into separate lanes for each. |
|  | Sign 254 – No cycling permitted. Indicates that cyclists are prohibited from entering or using this section. |
|  | Sign 244.1 – Beginning of a cycle street. Indicates a roadway primarily designated for bicycles; motor vehicles are permitted only under specific conditions. |
|  | Sign 239 + Zusatzzeichen 1022-10 – Pedestrian path with the “Cyclists allowed” supplementary plate, indicating that bicycles may also use the pedestrian walkway. |
|  | Sign 267 + Zusatzzeichen 1022-10 – Entry restricted with the “Cyclists allowed” supplementary plate, indicating that bicycles may enter the street. |

Table 10: **Traffic signs relevant to cyclists.** Examples adapted from the Vienna Convention on Road Signs and Signals ([Wikipedia contributors, 2024](#)) and the German traffic regulations (*Straßenverkehrs-Ordnung, StVO*).



Question

In the image, focus on the region marked as <1> (highlighted by a red bounding box). Which road sign description corresponds to the sign present in that region?

Choices

Options:

- A. "Regulatory sign prohibiting pedestrian entry to this path"
- B. "Regulatory sign indicating a shared path for pedestrians and cyclists"
- C. "Warning sign indicating presence of children nearby, typically near schools or playgrounds"
- D. "Regulatory sign marking the end of a cycle-priority street"
- E. "Warning sign indicating narrowing of the road on the right side"
- F. "Regulatory sign prohibiting cycling"
- G. "Regulatory sign indicating a pedestrian-only path (cyclist not permitted)"

Ground Truth: G — Regulatory sign: Pedestrian-only path (Cyclists prohibited).

DRIVELMM-O1 (8B)

Reasoning: The highlighted region <1> contains an icon depicting a person holding a child. This typically indicates a warning related to the presence of pedestrians, specifically near schools.

Answer: C. Warning sign indicating presence of children nearby.

QWEN3-VL-8B

Reasoning: The sign shows a blue circle with a white pictogram of an adult and child. The presence of a "Schule" (school) sign nearby supports the interpretation of a standard European warning sign.

Answer: C. Warning sign indicating presence of children nearby.

COSMOS-R2 (8B)

Reasoning : "Blue circles usually indicate mandatory regulations rather than warnings. Option C is a warning, but the color is blue... perhaps it's a regulatory sign marking a shared path for pedestrians and cyclists, even if symbols are simplified."

Answer: B. Regulatory sign indicating a shared path for pedestrians and cyclists.

INTERNVL3

Reasoning: "Blue circles usually indicate mandatory regulations rather than warnings. Option C is a warning, but the color is blue... perhaps it's a regulatory sign marking a shared path for pedestrians and cyclists, even if symbols are simplified."

Answer: B. Regulatory sign indicating a shared path for pedestrians and cyclists.

Table 11: **Qualitative Failure Case Analysis (TSR+S)**. A regulatory semantic failure in which models correctly ground the sign but misinterpret the icon (pedestrian-only vs. school warning). Note that generated text is paraphrased for brevity.



Question

In the image, focus on the region marked as <1> (highlighted by a red bounding box). Which road sign description corresponds to the sign present in that region?

Choices

- A. "Regulatory sign indicating a pedestrian-only path (cyclist not permitted)"
- B. "Combined sign indicating a parking area for the bicycles"
- C. "Warning sign indicating cyclists or bicycle traffic ahead"
- D. "Regulatory sign indicating a shared path for pedestrians and cyclists"
- E. "Regulatory sign indicating parallel separated path—one for cyclists and one for pedestrians"
- F. "Regulatory sign prohibiting cycling"
- G. "Regulatory sign indicating a mandatory bicycle-only lane"

GT: G — Regulatory sign indicating a mandatory bicycle-only lane

DRIVELMM-O1 (8B)

Reasoning: The image shows a paved path. The highlighted region shows a road sign with an icon of a bicycle and downward arrows, which typically indicate that the path is shared between cyclists and pedestrians without separation. The icons and the lack of a prohibition symbol suggest shared use.

Answer: E. Regulatory sign indicating parallel separated paths—one for cyclists and one for pedestrians.

QWEN3-VL-8B

Reasoning: The sign shows a blue circle with a white bicycle symbol above a rectangular sign with two arrows, which typically indicates separate paths for different users, often cyclists and pedestrians, in the same direction or in opposite directions.

Answer: E

COSMOS-R2 (8B)

Reasoning: The top circular blue sign with a white bicycle symbol typically means the path is designated for cyclists only. The rectangular sign below showing two arrows pointing in opposite directions usually indicates separate lanes. The combination suggests separate paths; the presence of opposing arrows implies separation rather than sharing.

Answer: E

INTERNVL3

Reasoning: The sign in the highlighted region shows a bicycle symbol with arrows indicating movement in both directions, which typically indicates a shared path for pedestrians and cyclists.

Answer: D. Regulatory sign indicating a shared path for pedestrians and cyclists.

Table 12: **Qualitative Failure Case Analysis (TSR+S):** Here, some models incorrectly relies on the double arrows to base its reasoning. Note that generated text is paraphrased for brevity.



Question

The lanes are shaded and labeled numerically. Which labeled lane corresponds to the cyclist lane?

Choices

- A. Lane 1
- B. Lane 0

GT: B — Lane 0

DRIVELMM-O1 (8B)

Reasoning: Identify the relevant road signs in the image. A sign is shown indicating a bike with an arrow pointing down, which indicates a lane for bicycles. The signs are above the lanes labeled "Lane 0" and "Lane 1." Since the sign indicates "Lane 1" as having a cyclists lane, this corresponds to the intended direction of the bicycles.

Answer: A. Lane 1

QWEN3-VL-8B

Reasoning: The blue bicycle symbol on the sign above Lane 1 indicates it is designated for cyclists, and the lane is paved with bricks typically used for bike paths. Lane 0 is a wider, smoother surface likely for pedestrians or general traffic, not specifically for cyclists.

Answer: A. Lane 1

COSMOS-R2 (8B)

Reasoning: Lane 0 is a smooth, light brown asphalt path with scattered leaves. Lane 1 is made of gray rectangular paving stones. The key detail is the presence of a blue circular sign with a white bicycle symbol above Lane 1. In many countries, such signs indicate that the lane is designated for cyclists. The positioning of the sign directly over Lane 1 reinforces this association.

Answer: A

INTERNVL3.5-8B

Reasoning: The cyclist lane is indicated by the blue sign with a bicycle symbol, and it is positioned next to Lane 1, suggesting it is the designated lane for cyclists.

Answer: A. Lane 1

Table 13: **Failure case study for the Bicycle Lane Prediction task.** In this example, the model relies on the physical location of the signboard to infer the lane. Note that generated text is paraphrased for brevity.



Question

The lanes are shaded and labeled numerically. Which labeled lane corresponds to the cyclist lane?

Choices

- A. Lane 1
- B. Lane 0

GT: B — Lane 0

DRIVE LMM-O1 (8B)

Reasoning: In the image, there is a road sign visible that indicates the presence of a cyclist, accompanied by a bicycle symbol. The sign is located along the path, and the lane corresponding to this mark is clearly shaded as "Lane 1". Therefore, the lane intended for cycling activity is labeled as "Lane 1".

Answer: A. Lane 1

QWEN3-VL-8B

Reasoning: The blue circular sign with a bicycle symbol indicates a designated bike lane, and Lane 1 is shaded blue, matching the color associated with the bike lane in the image.

Answer: A. Lane 1

COSMOS-R2 (8B)

Reasoning: Lane 0 has light brown tiles; Lane 1 has purple tiles. The presence of a blue circular sign with a bicycle icon above Lane 1 suggests Lane 1 is designated for bicycles. Additionally, the color coding (purple for bikes) and the positioning of the sign directly adjacent to Lane 1 reinforces this association.

Answer: A

INTERNVL3.5-8B

Reasoning: Lane 1 is marked in purple with a blue line, which usually indicates a bike lane, and it also has a bicycle symbol on the sign nearby.

Answer: A. Lane 1

Table 14: **Failure case study for the Bicycle Lane Prediction task.** In this example, the model relies on the physical location of the signboard to infer the lane. Note that generated text is paraphrased for brevity.



Question

In the image, captured from the cyclist's perspective, focus on the regions marked as <region1> (green box) and <region2> (red box). Describe the spatial relationship between the object in <region1> and the object in <region2>.

Choices

- A. The object in [region2] is closer to the cyclist than the object in [region1].
- B. The object in [region1] is closer to the cyclist than the object in [region2].

GT: A — The object in [region2] is closer to the cyclist than the object in [region1].

INTERNVL3.5-8B

Reasoning: The green bounding box (region1) is positioned in front of the red bounding box (region2), indicating that the object in region1 is closer to the cyclist.

Answer: B. The object in [region1] is closer to the cyclist than the object in [region2].

COSMOS-R2 (8B)

Reasoning: Traffic Sign 1 (left) is near the edge of the pavement, while Traffic Sign 2 (right) is mounted higher up. In the cyclist's line of sight, the left-side sign would be closer because it's on the near side of the road, whereas the higher sign might be positioned further back.

Answer: B

SPATIALTHINKER-7B

Reasoning: The cyclist would see the bicycle prohibition sign (region1) first as they approach the intersection, and then the pedestrian crossing sign (region2) further ahead. The arrangement shows region1 is in the foreground.

Answer: B

VST-7B-RL

Reasoning: Using the concept of depth in 3D space, region1 is positioned such that it lies between the cyclist and region2. The spatial arrangement shows region1 is in the foreground relative to region2.

Answer: B

Table 15: **Failure case study for the Spatial Understanding (SU) task.** Here, models struggle with depth perception of the traffic elements. Note that generated text is paraphrased for brevity.



Question

In the image, captured from the cyclist's perspective, focus on the regions marked as <region1> (green box) and <region2> (red box). Describe the spatial relationship between the object in <region1> and the object in <region2>.

Choices

- A. The object in [region2] is closer to the cyclist than the object in [region1].
- B. The object in [region1] is closer to the cyclist than the object in [region2].

GT: A — The object in [region2] is closer to the cyclist than the object in [region1].

INTERNVL3.5-8B

Reasoning: The cyclist's perspective shows the green sign in region1 in front of the blue sign in region2, indicating that the green sign is closer to the cyclist.

Answer: B. The object in [region1] is closer to the cyclist than the object in [region2].

COSMOS-R2 (8B)

Reasoning: Objects higher on the pole appear further away because the cyclist looks upwards; objects lower down seem closer. The rectangular sign is beneath the blue one, so its position relative to the cyclist's line of sight makes it appear nearer.

Answer: B

SPATIALTHINKER-7B

Reasoning: The image shows a road intersection. Region 1 is highlighted in green and is closer to the cyclist, while Region 2 is highlighted in red and is further away. Therefore, the object in Region 1 is closer.

Answer: B

VST-7B-RL

Reasoning: Analyzing depth in 3D space, the object in region1 is positioned such that it lies between the cyclist and the object in region2 in terms of depth. Thus, region1 is in the foreground relative to region2.

Answer: B

Table 16: **Failure case study for the Spatial Understanding (SU) task.** This example highlights a systematic failure where all models incorrectly perceive depth by equating vertical height or 2D image coordinates with 3D proximity. Note that generated text is paraphrased for brevity.



Question

The image shows a forward-facing view from a cyclist's perspective. From the cyclist's perspective, does the shaded lane (in pink) ahead allow cycling? The shading is externally drawn on the image as an overlay.

Choices

- A. Yes, continue riding forward on the labelled lane
- B. No, cycling is not allowed on this lane

GT: B — No, cycling is not allowed on this lane

DRIVELMM-O1 (8B)

Reasoning: There is a clear "Bike Route Only" sign on the far right, indicating this lane is specifically for bicycles. The pink shaded lane corresponds with the bike route as depicted by the sign. The warning triangle below it indicates caution, not prohibition.

Answer: A. Yes, continue riding forward on the labelled lane

COSMOS-R2 (8B)

Reasoning: There's a red circle with a black bicycle symbol. That usually means "bicycle lane" or "cycle path," suggesting cyclists are permitted. Even if there's a shared area, the presence of the bicycle-specific sign implies cycling is allowed.

Answer: A

QWEN3-VL-8B

Reasoning: The red circle with a bicycle symbol indicates a "No Cycling" zone. Therefore, cycling is not permitted on the shaded lane.

Answer: B. No, cycling is not allowed on this lane

INTERNVL3.5-8B

Reasoning: The traffic sign shows a symbol of a bicycle with a red circle, which means that cycling is prohibited in this lane. The lane also has a "Straßen schaden" (road damage) warning.

Answer: B. No, cycling is not allowed on this lane

Table 17: **Failure case study for the Sign-Action Association (SAA) task.** This example demonstrates the lack of traffic rule understanding error, where some models correctly identify the bicycle pictogram but misinterpret the red circular border as a permissive indicator rather than a prohibition. Note that generated text is paraphrased for brevity.



Question

The image shows a forward-facing view from a cyclist's perspective. From the cyclist's forward view, is cycling allowed along the highlighted lane (in pink)? The shading is externally drawn on the image as an overlay.

Choices

- A. No, cycling is not allowed on this lane
- B. Yes, continue riding forward on the labelled lane

GT: A — No, cycling is not allowed on this lane

DRIVELMM-O1 (8B)

Reasoning: The lane is pink, indicating a designated cycling path. There is a "P" sign visible for parking, but no specific prohibition shown. The lack of obstructions or signs prohibiting biking directly implies that it is likely permitted.

Answer: B. Yes, continue riding forward on the labelled lane.

COSMOS-R2 (8B)

Reasoning: The scene includes a bicycle lane marked by dashed lines and a red bike symbol on a pole. Such symbols indicate designated paths. Since the lane is specifically marked for bicycles, it suggests cycling is permitted.

Answer: B

QWEN3-VL-8B

Reasoning: The bicycle symbol on the traffic sign (red circle, white background, black icon) typically indicates that cycling is permitted and encouraged. The lane is physically separated by a grass strip, confirming it is a bike lane.

Answer: B. Yes, continue riding forward on the labelled lane

INTERNVL3.5-8B

Reasoning: The highlighted pink lane is a designated bicycle path, evident from the sign showing a bicycle symbol. The forward-facing view from the cyclist's perspective shows they are in the correct lane.

Answer: B. Yes, continue riding forward on the labelled lane

Table 18: **Failure case study for the Sign-Action Association (SAA) task.** This case illustrates a high-frequency failure where all models misinterpret a red circular prohibitory sign as a permissive one. Note that generated text is paraphrased for brevity.



Question

The image shows a forward-facing view from a cyclist's perspective. From the cyclist's forward view, is cycling allowed along the highlighted lane (in pink)? The shading is externally drawn on the image as an overlay.

Choices

- A. No, cycling is not allowed on this lane
- B. Yes, continue riding forward on the labelled lane

GT: A — No, cycling is not allowed on this lane

DRIVELMM-O1 (8B)

Reasoning: Pay close attention to the highlighted lane in pink. Notice the sign on the left indicating parking. The presence of this sign, alongside the specific regulatory markings, suggest restricted use where cycling is not allowed.

Answer: A. No, cycling is not allowed on this lane.

COSMOS-R2 (8B)

Reasoning: There's a red circular sign with a black bicycle symbol. That usually indicates that cycling is permitted there. Since the sign explicitly marks the lane for bicycles, the correct answer should be that the cyclist can continue using it.

Answer: B

QWEN3-VL-8B

Reasoning: The pink highlighted lane is a designated bike path. The signpost shows a bicycle symbol within a red circle, which typically indicates a bike route. Since the cyclist is on the path and the sign suggests it is for cyclists, cycling is allowed.

Answer: B. Yes, continue riding forward on the labelled lane

INTERNVL3.5-8B

Reasoning: The image shows a no-bicycle sign above the parking sign, indicating that cycling is not permitted on this lane. The red circular border serves as the primary prohibitory indicator.

Answer: A. No, cycling is not allowed on this lane

Table 19: **Failure case study for the Sign-Action Association (SAA) task.** This example traffic rule misunderstanding error. In particular, models correctly identify the bicycle icon and the red circle but fail to associate the red circular frame with a prohibition, misinterpreting it as a permissive marker. Note that generated text is paraphrased for brevity.



Question

The image shows a forward-facing view from a cyclist's perspective. Considering the shaded lane (in pink) ahead, should the cyclist stay on this route or avoid it? The shading is externally drawn on the image as an overlay.

Choices

- A. No, cycling is not allowed on this lane
- B. Yes, continue riding forward on the labelled lane

GT: A — No, cycling is not allowed on this lane

DRIVELMM-O1 (8B)

Reasoning: The image shows a lane designated for cyclists by the arrow on the road. The pink shaded area indicates this lane. There is a "no entry" sign for vehicles nearby, indicating this area is reserved for cycling.

Answer: B. Yes, continue riding forward on the labelled lane.

COSMOS-R2 (8B)

Reasoning: There's a bicycle lane marked by arrows pointing left. The pink area is part of that lane. The presence of the bike symbol suggests it's allowed. Continuing on that lane would be appropriate unless obstacles exist.

Answer: B

QWEN3-VL-8B

Reasoning: The visual cues, including the lane markings and the bicycle pictogram on the sign, suggest a designated route for cyclists. Forward movement on the pink shaded lane is therefore permitted.

Answer: B. Yes, continue riding forward on the labelled lane

INTERNVL3.5-8B

Reasoning: The presence of a no-entry sign for cyclists (red circle with a white bike inside) indicates that cycling is prohibited in this lane. This regulatory sign overrides other visual indicators.

Answer: A. No, cycling is not allowed on this lane

Table 20: **Failure case study for the Sign-Action Association (SAA) task.** This case highlights a critical failure in directional and regulatory sign interpretation. Most models misinterpreted the lane markings and the prohibitory sign, incorrectly concluding that the path was permissive for forward travel. Note that generated text is paraphrased for brevity.



Question

The image shows a forward-facing view from a cyclist's perspective. From the cyclist's perspective, does the shaded lane (in pink) ahead allow cycling? The shading is externally drawn on the image as an overlay.

Choices

A. Yes, continue riding forward on the labelled lane
 B. No, cycling is not allowed on this lane

GT: B — No, cycling is not allowed on this lane

DRIVELMM-O1 (8B)

Reasoning: The image shows a sign indicating a bicycle lane. Currently, the light is green indicating it is safe to proceed. The lane is explicitly colored and separated from the road, further confirming its intent for bikers.

Answer: A. Yes, continue riding forward on the labelled lane.

COSMOS-R2 (8B)

Reasoning: The paved area specifically mentions it's designated for cyclists. Since it's labeled as a cyclist's paved area, that suggests it's meant for cycling. You can continue riding forward on that lane.

Answer: A

QWEN3-VL-8B

Reasoning: The red circle with a bicycle symbol indicates a "No Bicycles" zone, meaning cycling is prohibited in this area regardless of other lane markings.

Answer: B. No, cycling is not allowed on this lane

INTERNVL3.5-8B

Reasoning: The lane is marked with a traffic sign indicating it is a no-bicycle zone. Prohibitory signs take precedence over the existence of a paved path.

Answer: B. No, cycling is not allowed on this lane

Table 21: **Failure case study for the Sign-Action Association (SAA) task.** This example highlights a traffic rule misunderstanding error where some models incorrectly rely on adjacent traffic lights while ignoring relevant regulatory "No Bicycles" sign. Note that generated text is paraphrased for brevity.



Question

The image shows a forward-facing view from a cyclist's perspective. From the cyclist's forward view, is cycling allowed along the highlighted lane (in pink)? The shading is externally drawn on the image as an overlay.

Choices

- A. Yes, continue riding forward on the labelled lane
- B. No, cycling is not allowed on this lane

GT: B — No, cycling is not allowed on this lane

DRIVELMM-O1 (8B)

Reasoning: The pink lane suggests it is specifically designated for cyclists. The arrangement of markings and indicators confirms that this is a lane intended for cyclists and indicates that it is safe to proceed.

Answer: A. Yes, continue riding forward on the labelled lane.

COSMOS-R2 (8B)

Reasoning: There's a pink lane with white arrows and a blue circular bicycle sign. Even though there's a red circle nearby, the blue sign over the lane overrides it. The primary indicator for the highlighted lane allows cycling.

Answer: A

QWEN3-VL-8B

Reasoning: The blue circular sign indicates permitted cycling. The "no cycling" sign is positioned above the pink lane, suggesting it applies to the area beyond the section, not the highlighted lane itself. The pink color confirms it is allowed.

Answer: A. Yes, continue riding forward on the labelled lane

INTERNVL3.5-8B

Reasoning: The signs indicate that cycling is prohibited in this lane due to the red circle with a line through the bicycle symbol. This prohibitory status remains the primary regulatory instruction.

Answer: B. No, cycling is not allowed on this lane

Table 22: **Failure case study for the Sign-Action Association (SAA) task.** This case highlights a Lane-misassociation error where models struggle to resolve contradictory signals (prohibitory vs. permissive) and fail to assign to correct lanes. Note that generated text is paraphrased for brevity.



Question

The image shows a forward-facing view from a cyclist's perspective. Looking at the shaded lane (in pink) in the image, is the cyclist allowed to keep following it? The shading is externally drawn on the image as an overlay.

Choices

A. Yes, continue riding forward on the shaded lane
 B. No, cycling is not allowed on this lane

GT: B — No, cycling is not allowed on this lane

DRIVELMM-O1 (8B)

Reasoning: The blue circular sign with a bicycle symbol indicates this is a bike lane. The pink shaded lane corresponds with the bike lane marked by the sign. Since the lane is marked for cycle use, the cyclist is permitted to continue.

Answer: A. Yes, continue riding forward on the shaded lane

COSMOS-R2 (8B)

Reasoning: There's a blue circular sign with a bicycle symbol and another with pedestrians. Usually, these shared spaces allow cyclists to proceed unless a red prohibition sign is present. The shaded lane aligns with the bicycle indicator.

Answer: A

QWEN3-VL-8B

Reasoning: The sign directly above the bicycle indicator (pedestrian and child) indicates a pedestrian zone. The bidirectional arrows on the auxiliary plate suggest this is a shared or pedestrian-priority area. The shaded path is likely restricted.

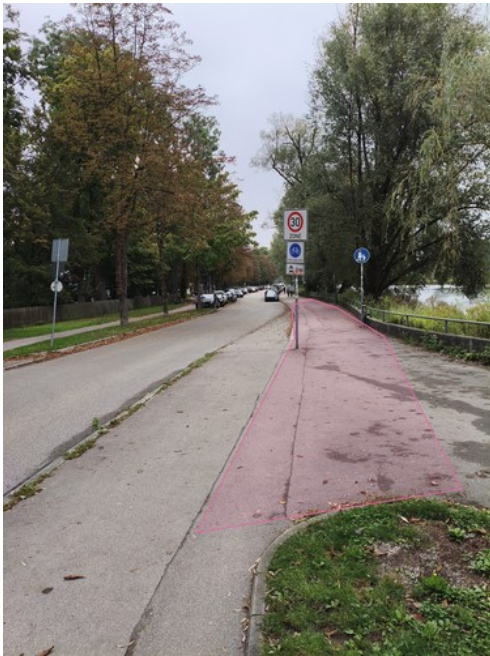
Answer: B. No, cycling is not allowed on this lane

INTERNVL3.5-8B

Reasoning: The sign indicates that the lane is for both cyclists and pedestrians. Therefore, the cyclist is allowed to continue following the pink shaded lane forward as it is a shared space.

Answer: A. Yes, continue riding forward on the shaded lane

Table 23: **Failure case study for the Sign-Action Association (SAA) task.** Models fail to spatially perceive when multiple signs with different instructions are present next to each other, leading to confusion about which regulation applies to the shaded lane. Note that generated text is paraphrased for brevity.



Question

The image shows a forward-facing view from a cyclist's perspective. Based on the shaded path (in pink) shown ahead, should the cyclist proceed forward on this lane? The shading is externally drawn on the image as an overlay.

Choices

A. Yes, continue riding forward on the shaded lane
 B. No, cycling is not allowed on this lane

GT: B — No, cycling is not allowed on this lane

DRIVELMM-O1 (8B)

Reasoning: The image shows a lane specifically designated for bicycles, indicated by the symbol on the sign above the lane. The pink mark suggests it is earmarked for cycling purposes. The cyclist's forward view implies the lane ahead is available.

Answer: A. Yes, continue riding forward on the shaded lane

COSMOS-R2 (8B)

Reasoning: A blue sign with a white bicycle symbol indicates a bike lane. The presence of a "30 ZONE" sign suggests restricted car use. Since the shaded path is separated by a railing and labeled for cyclists, the cyclist should continue using it.

Answer: A

RECOGDRIVE-VLM

Reasoning: The pink shade indicates a designated cycling path specifically for cyclists. This path is separate from the main road, ensuring a safe route. Therefore, the cyclist should proceed as it is clearly intended for cycling activities.

Answer: A. Yes, continue riding forward on the shaded lane.

SPATIALTHINKER-7B

Reasoning: The sign indicates a "Fahrradstraße" (bike lane) and a "30 ZONE" limit. The shaded pink path appears to be a designated bike lane, and the cyclist should follow it for safety and adherence to rules.

Answer: A. Yes, continue riding forward on the shaded lane

Table 24: **Failure case study for the Sign-Action Association (SAA) task.** In this models fails to focus on the pedestrian lane sign on the right side. Note that generated text is paraphrased for brevity.



Question
As the cyclist progressed along the road, which of these two views corresponds to the earlier moment?

Choices
A. Image 1
B. Image 2
GT: B — Image 2 corresponds to the earlier moment.

DRIVELMM-O1 (8B)
Reasoning: Image 1 shows a clear route with fewer obstructions. Image 2 shows a more complex intersection with traffic signals, suggesting a more "developed" or later progression of the cyclist through the environment.
Answer: A. Image 1

INTERNVL3.5-8B
Reasoning: The blue bicycle sign is clearly visible in Image 1. Typically, such instructional signs are encountered earlier along a road as a cyclist progresses. Image 2 appears to be further along after passing this marker.
Answer: A. Image 1

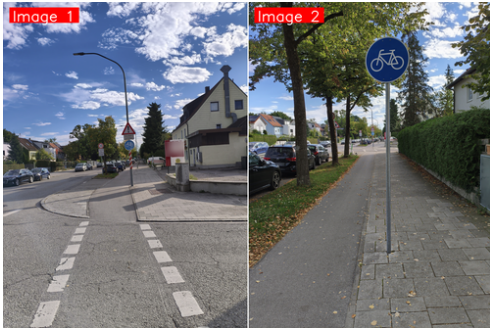
COSMOS-R2 (8B)
Reasoning: Bus stops are usually encountered before major intersections. The person waiting at the shelter in Image 1 suggests an earlier point in the journey before the cyclist reaches the traffic control elements seen in Image 2.
Answer: A

OVIS2.5-9B
Reasoning: Prominent signage like the bicycle lane indicator in Image 1 marks an approaching zone. In Image 2, the cyclist has progressed to the intersection, approaching the traffic lights, which represents a later stage of the trip.
Answer: A. Image 1

QWEN3-VL-8B
Reasoning: Image 1 provides a direct, close-up view of the bus stop, suggesting the cyclist is just arriving at this point. Image 2 shows a wider perspective with more road ahead, implying the cyclist is further along.
Answer: A. Image 1

VST-7B-RL
Reasoning: Image 1 captures a more central/early section of the road characterized by the bus stop. Image 2 captures a later section with different signage and a more open road environment.
Answer: A

Table 25: **Failure case study for the Temporal Ordering (TO) task.** Models misinterpret spatial cues and fail to accurately sequence the images based on the cyclist’s progression along the road. Note that generated text is paraphrased for brevity.



Question

These two images were captured sequentially during the cyclist's forward travel. Which image was taken first?

Choices

- A. Image 1
- B. Image 2

GT: B — Image 2 was taken first.

DRIVELMM-O1 (8B)

Reasoning: Image 1 shows a broad view of a curved road and intersection, suggesting the cyclist is approaching the path. Image 2 is a close-up of a sign, implying it was captured after moving past the broader scene.

Answer: A. Image 1

COSMOS-R2 (8B)

Reasoning: The cyclist would first encounter the road intersection and speed signs (Image 1) before reaching the specific bike path indicated by the blue sign (Image 2) further along the journey.

Answer: A

OVIS2.5-9B

Reasoning: Image 1 shows the cyclist at a street intersection, while Image 2 shows a sidewalk path. Since travel usually moves from main road intersections to residential paths, Image 1 was captured first.

Answer: A. Image 1

QWEN3-VL-8B

Reasoning: Image 1 shows the cyclist approaching an intersection with road damage signs. Image 2 shows a straight path past that intersection. Given forward travel, the intersection view occurs first.

Answer: A. Image 1

SPATIALTHINKER-7B

Reasoning: Image 1 shows the street from a distance, while Image 2 shows a close-up of a sign on the sidewalk. Forward travel implies getting closer to the sign, so the distant view is first.

Answer: A (A) Image 1

VST-7B-RL

Reasoning: The road with vehicles and buildings (Image 1) represents the initial visual for a forward-moving cyclist before they transition to the specific sidewalk area (Image 2).

Answer: A

Table 26: **Failure case study for the Temporal Ordering (TO) task.** The models consistently misinterpreted the spatial progression of the cyclist, incorrectly identifying a more distant, complex intersection as the "earlier" moment and the immediate path sign as the "later" moment. Note that generated text is paraphrased for brevity.



Question

During continuous forward motion, the cyclist took these two photos. Which one was captured first in time?

Choices

- A. Image 1
- B. Image 2

GT: B — Image 2 was captured first.

DRIVELMM-O1 (8B)

Reasoning: Image 2 depicts the forward-facing road marking that the cyclist travels over, representing the path ahead. Image 1 captures specific side details and navigation signs near the curb. Given the forward motion, the road-path view in Image 2 occurs first.

Answer: B. Image 2

INTERNVL3.5-8B

Reasoning: Visual cues such as shadows and the positions of objects suggest that Image 1 was captured when the cyclist was closer to the buildings. Image 2 shows them further away, indicating it must have been taken later in the sequence.

Answer: A. Image 1

COSMOS-R2 (8B)

Reasoning: Image 1 captures the initial set of signs and a parked car. As the cyclist progresses, new elements like graffiti and broader building views appear in Image 2, suggesting the cyclist has moved past the initial parked car location.

Answer: A

OVIS2.5-9B

Reasoning: Image 1 shows a closer view of the building and signs, while Image 2 shows a more distant perspective. Since the cyclist is moving forward, the photo taken earlier would be the one closer to the initial landmarks.

Answer: A. Image 1

QWEN3-VL-8B

Reasoning: Image 1 shows the cyclist in closer proximity to the signpost. Image 2 shows the cyclist further down the street with a bike symbol on the road ahead. The earlier photo would naturally be closer to the starting reference point.

Answer: A. Image 1

SPATIALTHINKER-7B

Reasoning: In Image 1, the cyclist is closer to the building and immediate surroundings. In Image 2, the view is more open and shows a wider street. Therefore, the closer view in Image 1 must have been captured first.

Answer: A. Image 1

Table 27: **Failure case study for the Temporal Ordering (TO) task.** Here, most models assume that being near a building or signpost represents the start of a sequence, rather than correctly identifying the traffic elements present in the scene to base its reasoning. Note that generated text is paraphrased for brevity.