
M4R: Measuring Massive Multi-Modal Understanding and Reasoning in Open Space

Shangding Gu¹, Xiaohan Wang², Donghao Ying¹, Haoyu Zhao³, Runing Yang⁴,
Boyi Li^{1,5}, Ming Jin⁴, Marco Pavone^{2,5}, Serena Yeung-Levy², Jun Wang³,
Dawn Song^{1*}, Costas Spanos^{1*}

¹UC Berkeley ²Stanford ³UCL ⁴Virginia Tech ⁵Nvidia *Equally advise

Abstract

The increasing sophistication of multi-modal models necessitates benchmarks that can rigorously evaluate their understanding and reasoning in complex, safety-pertinent, open-world scenarios. This study introduces M4R (Measuring Massive Multi-Modal Understanding and Reasoning), a large-scale benchmark uniquely designed to assess reasoning capabilities across diverse open spaces, comprehensively covering land, air, and water environments. M4R comprises approximately 2,000 videos and over 19,000 human-annotated question-answer pairs. These videos, varying in length (short, medium, long) and presenting tasks of tiered difficulty (interval-based choices and accuracy-based choices), encompass distinct operational domains: the land-based scenarios primarily focus on traffic environments, particularly traffic collisions and accident cases; the air-based scenarios center on airplane navigation; and the water-based scenarios involve ship movements. M4R systematically evaluates models on temporal-causal reasoning, spatial understanding, and intent and goal planning within these dynamic contexts. By providing a unified platform across this broad spectrum of domains, M4R aims to drive the development of more robust and generalizable AI systems. Benchmarking state-of-the-art multi-modal models on our dataset reveals that even leading models, such as ChatGPT-4o and Gemini, achieve only around a 20% success rate, highlighting the significant challenges that remain in open-space multi-modal reasoning. The code, leaderboard, and dataset are available at: <https://open-space-reasoning.github.io/>.

1 Introduction

As artificial intelligence (AI) continues to evolve, large multi-modal models have shown impressive capabilities across vision, language, and video domains. However, significant challenges remain in deploying these models for real-world, safety-critical applications such as autonomous driving, robotics, and aerial or maritime operations. While multi-modal models demonstrate remarkable performance in constrained or simulated environments, their robustness and depth of understanding in high-stakes, dynamic scenarios are still far from sufficient.

In particular, deployment in mission-critical domains requires rigorous evaluation of models' understanding and reasoning abilities under real-world conditions that involve uncertainty, physical interactions, and causal dependencies. While recent benchmarks have advanced evaluation in specific facets like temporal understanding (e.g., MVBench [22], REXTIME [7]) or domain-specific knowledge (e.g., MMMU [45], DriveLM [32]), there remains a paucity of unified platforms that assess reasoning across the combined spectrum of land, air, and water operations. To address this, our work

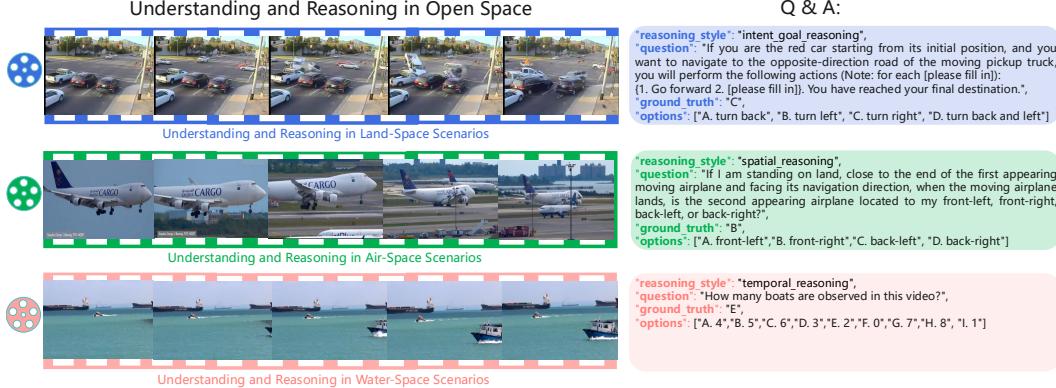


Figure 1: Examples of Multi-Modal Understanding and Reasoning in Open-Space Scenarios

defines *open space* as unstructured or semi-structured outdoor environments characterized by high variability, dynamic interactions, and minimal physical boundaries. This includes **air space** (e.g., airplane navigation), **water space** (e.g., ship and boat movements), and **land space** (e.g., road traffic involving diverse vehicle types). These settings inherently involve complex temporal dependencies, causal relationships, and real-world physical constraints, demanding advanced, robust reasoning capabilities for genuine open-world understanding.

We introduce **M4R** (Measuring Massive Multi-Modal Understanding and Reasoning), a comprehensive evaluation framework. Specifically, we present the **M4R** benchmark, which focuses on reasoning across the aforementioned land traffic, airspace, and waterway domains—settings where safety, perception, and decision-making are deeply interdependent. Unlike benchmarks focusing on isolated skills or single domains, **M4R** challenges models on several key reasoning capabilities: *temporal-causal reasoning* (understanding event sequences and causality over extended periods); *spatial understanding* (comprehending dynamic spatial relationships and multi-agent trajectories); *intent and goal planning/inference* (deducing agent intentions and goals); and *complex strategic & counterfactual reasoning* (assessing understanding of higher-order strategies, action implications, and ‘what-if’ scenarios). Several representative examples from **M4R** are illustrated in Figure 1. By systematically probing these capabilities across diverse safety-pertinent scenarios, **M4R** provides a framework for assessing progress towards AI systems that can reliably operate in the real world.

Our key contributions are summarized as follows:

- **Unified Open-World Evaluation Suite:** We introduce **M4R**, a large-scale, video-based benchmark uniquely covering land traffic, airspace, and waterway scenarios to provide a comprehensive assessment of multi-modal reasoning across these distinct yet complementary safety-critical open spaces.
- **Reasoning-Centric Evaluation:** **M4R** systematically evaluates critical reasoning facets including temporal-causal understanding, dynamic spatial awareness, intent and goal inference, and complex strategic reasoning, within dynamic and physically grounded settings.
- **Real-World Limitations and Safety Gaps:** We highlight limitations in current AI systems’ reasoning performance in open-space domains (e.g., autonomous driving, aviation, and maritime environments), and provide a challenging testbed to drive the development of safer and more robust multi-modal AI systems.

2 Related Work

2.1 General Multi-Modal Understanding Benchmarks

Recent years have witnessed growing interest in video understanding benchmarks. Foundational video question-answering (QA) efforts include MSR-VTT [43] and Next-QA [41]. More recently, MVbench [22], with its 20 diverse temporal tasks derived from static images, and MLVU [48] have expanded video QA capabilities across multiple domains. The challenge of long-form video understanding has seen contributions from benchmarks such as EgoSchema [26], Video-LLaVA [10],

MovieChat [34], and LongVideoBench [40]. Parallelly, video captioning benchmarks such as AuroraCap [6], HiCM2 [19], and LongCaptioning [39] focus on generating detailed textual descriptions.

A significant trend is the push for more rigorous temporal and causal reasoning. REXTIME [7], for instance, probes the linking of causally related events across separate video segments. For multi-domain understanding, MMWorld [15] evaluates models across diverse disciplines, requiring explanations and counterfactuals. Furthermore, LVbench [38] integrates video inputs for QA. Beyond video, reasoning from static images is explored by MME [18] (including CoT extensions), MMMU [45] (evaluating expert-level multi-discipline reasoning), and benchmarks for mathematical reasoning like Dynamath [50] and MultiModal-MATH [49]. For academic content, Video-MMLU [36] offers a large-scale lecture video benchmark.

While these diverse benchmarks significantly advance specific aspects of multi-modal understanding—be it general video comprehension, temporal analysis, long-form narrative understanding, captioning, or static image reasoning—they often do not provide a framework for unified evaluation across land, air, and maritime open-space environments, nor the specific blend of complex reasoning (including strategic and intent-based inference) that **M4R** is designed to evaluate within these contexts.

2.2 Safety-Critical Multi-Modal Understanding Benchmarks

Evaluating models in safety-critical domains, where reasoning under uncertainty is vital, is an emerging focus. Initial efforts addressed static image safety [23], model robustness against adversarial attacks (e.g., FigStep [12], JailBreakV [25]) [31, 28], or indoor robotics [44].

Autonomous driving has been a major driver of safety-critical research. Foundational datasets such as nuScenes¹ and Waymo Open Dataset², along with language-integrated efforts such as DriveLM and DriveVLM [32, 37], are closely related to **M4R**'s goals due to their real-world video and safety considerations. However, a key motivation for **M4R** was that these traditionally emphasized perception and planning, with less focus on deep safety-critical reasoning for tasks such as accident cause analysis or complex decision-making. Other specialized benchmarks tackle related issues such as video anomaly detection (e.g., VANE-Bench [11]).

While advancements continue in specialized video reasoning and domain-specific safety evaluations, existing benchmarks still largely focus on single operational domains. Critically, they often lack sufficient coverage of high-risk scenarios such as traffic collisions, ship navigation, and airplane takeoff/landing events across combined land, air, and water settings. A unified platform to consistently evaluate robust, generalizable reasoning (e.g., temporal-causal, spatial, intent, and strategic analysis) across these diverse, safety-critical open spaces also remains absent. To address this specific void, **M4R** distinctively incorporates these challenging high-risk scenarios from all three domains. The reliability of its complex reasoning evaluation is ensured as all annotations were generated by highly educated annotators (at least Master's degree). **M4R** thus provides a much-needed testbed for fostering robust, adaptable AI capable of open-world understanding.

3 Understanding and Reasoning in Open Space

3.1 Open Space Settings

We design the benchmark around three types of open-space environments: **land space**, focusing primarily on traffic accident understanding and reasoning; **air space**, centered on airplane takeoff and landing scenarios; and **water space**, which emphasizes ship navigation understanding and reasoning. Within each environment, we construct tasks that evaluate models across three key reasoning dimensions: dynamic temporal reasoning, spatial reasoning, and intent and goal reasoning. Representative examples for each reasoning type are illustrated in Figure 2.

For each reasoning style, we design tasks with varying levels of difficulty using two formats: *interval-based choices* and *accuracy-based choices*. Easy tasks provide approximately 3 coarse-grained interval choices, medium tasks offer 6 intermediate-level intervals, and hard tasks present 12 fine-grained discrete options that require an exact match with the correct answer. The number of tasks

¹<https://www.nuscenes.org/>

²<https://waymo.com/open/>



Figure 2: Examples of reasoning question settings in the M4R benchmark across three key reasoning types: *Temporal Reasoning*, which involves understanding event sequences and motion over time; *Spatial Reasoning*, which focuses on relative positioning and orientation in space; and *Intent Reasoning*, which evaluates understanding of goal-directed behaviors and decision-making in dynamic environments.



Figure 3: Land-space traffic accident scenarios for open-space video understanding and reasoning include [intersection collisions](#), [urban road accidents](#), nighttime incidents, [rural road accidents](#), [snow-covered road collisions](#), and [freeway accidents](#).

across the three difficulty levels is evenly distributed, with each comprising one-third of the total. In all cases, the model must select a single best answer, enabling the benchmark to assess performance under increasing levels of precision and ambiguity.

Land Space In our land-space setting, we include a comprehensive range of traffic scenarios, encompassing diverse collision events under varying weather conditions such as snow, rain, and sunshine, as detailed in Table 1. Specific examples of these scenarios are illustrated in Figure 3, and more detailed examples are provided in Appendix B. To enhance contextual diversity, we incorporate multiple camera perspectives—including ego-centric and third-person views—particularly for accident scenes. The dataset features incidents involving a wide variety of vehicle types, including buses, motorcycles, sedans, and several categories of trucks, across different road environments such as highways, freeways, and rural roads. The associated questions are designed to evaluate models across multiple reasoning dimensions, including temporal-causal understanding, spatial reasoning, and intent and goal planning. The original land-space video datasets are sourced from [5, 30], which primarily collected videos from YouTube and other public internet platforms.

Table 1: Overview of traffic accident scenarios in our benchmark, covering diverse road environments, weather conditions, and involved traffic participants.

Index	Categories
Road Environments:	Intersection, Highway, Freeway, Rural Road, Tunnel, Urban Road, Bridge, Parking Lot
Weather Conditions:	Snow, Rain, Sunshine, Cloudy, Foggy, Windy
Involved Participants:	Sedan, SUV, Bus, Truck, Motorcycle, Bicycle, Van, Pickup, Trailer, Pedestrian

Air Space In airspace scenarios, we primarily focus on *takeoff* and *landing* events, emphasizing the analysis of airplane navigation directions and perceptual understanding. Airplanes represent a largely unexplored domain in large multi-modal research, despite their significant real-world impact. Our benchmark investigates various aspects of airplane behavior, including differences in navigation patterns, aircraft sizes, and motion dynamics across different types of airplanes. These scenarios also incorporate videos of varying lengths and are designed to evaluate models on multiple reasoning dimensions, including spatial reasoning, temporal reasoning, and intent and goal inference. We further assess model performance across different difficulty levels using both interval-based and accuracy-based multiple-choice formats. The airspace videos are sourced from publicly available footage, including references such as ³, ⁴, and ⁵.

Water Space We include videos from both **river** and **ocean** scenarios, featuring varying video lengths and difficulty levels. The dataset encompasses a diverse range of watercraft, including different types of boats and ships, under a broad set of navigation conditions. Despite their real-world importance, river and ocean environments remain underexplored in the context of large multi-modal models. To address this gap, we evaluate model performance across multiple reasoning styles—temporal, spatial, and intent and goal reasoning—using video-based tasks of varying durations and difficulty levels. Task difficulty is controlled through both interval-based and accuracy-based multiple-choice formats. The water-space videos are sourced from publicly available datasets, including [13, 27].

3.2 Dataset Analysis

This benchmark includes approximately 2,000 videos and 19,000 human-annotated question-answer pairs, covering a wide range of reasoning tasks. All annotations were performed by highly educated annotators, each holding at least a master’s degree in engineering-related fields such as mathematics or computer science. The dataset features a variety of video lengths, categories, and frame counts, and spans three primary open-space reasoning scenarios: **land space**, **water space**, and **air space**. An overview of the dataset’s characteristics is shown in Figure 4, which illustrates the distributions of video duration, domain coverage, and reasoning styles.

Specifically, **(a) Video Length:** A substantial portion of the videos (76.5%) are short, with durations under 10 seconds. The remaining videos are distributed across longer intervals: 10–30 seconds (3.7%), 30–60 seconds (4.6%), 60–120 seconds (4.8%), 120–300 seconds (4.4%), and over 300 seconds (6.0%). This distribution reflects a strong emphasis on short, dynamic scenarios that test rapid perception and reasoning. **(b) Video Categories:** The benchmark spans three open-space domains. Land space, which primarily involves traffic and safety-related scenarios, comprises 83.0% of the videos. Air space accounts for 10.2%, and water space makes up 6.8%. This distribution highlights both the practical importance of land-based reasoning and the inclusion of underrepresented domains such as maritime and aviation environments. **(c) Reasoning Styles:** M4R supports three major reasoning types, with a relatively balanced distribution: *spatial reasoning* (35.4%), *temporal reasoning* (34.0%), and *intent reasoning* (30.6%). This design ensures comprehensive evaluation across key dimensions essential for real-world multi-modal understanding.

Overall, the dataset provides a rich and diverse collection of real-world video scenarios across multiple modalities and time scales, offering a robust benchmark for evaluating multi-modal understanding and reasoning in open-space environments.

³<https://www.youtube.com/watch?v=i6CrbqeeksJ8>

⁴<https://www.youtube.com/watch?v=k5yvzTw08K8>

⁵<https://www.youtube.com/watch?v=Bt9tpiAmTs8>

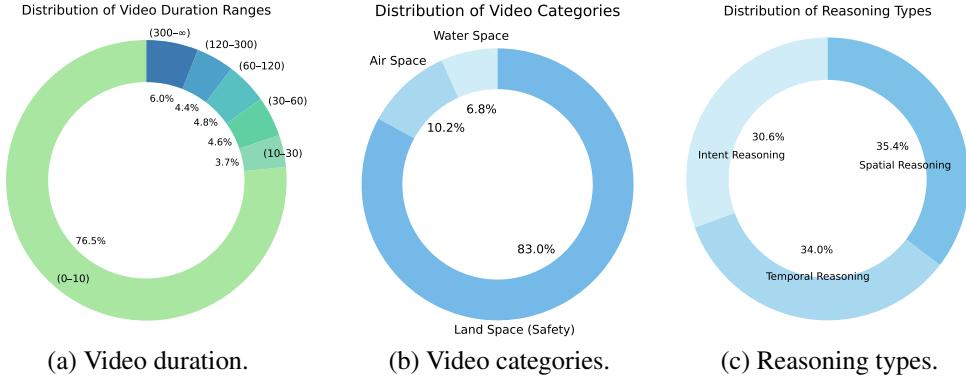


Figure 4: Distribution of video and task properties in the M4R benchmark.

3.3 Comparison with Existing Benchmarks

Table 2 provides a comparative analysis of M4R alongside existing evaluation benchmarks for MLLMs. Most benchmarks primarily focus on assessing the multimodal reasoning capabilities of MLLMs [14, 35, 48]; however, a significant limitation is the prevalent oversight of safety considerations. While a few recent benchmarks have begun to evaluate safety aspects of MLLMs [49, 23], they often do not incorporate video question-answering data. However, single-frame capture, in most cases, can introduce uncertainties in reasoning and is insufficient for adequately assessing MLLMs’ capabilities in handling safety issues. In contrast, our M4R introduces a large-scale and meticulously curated collection of video question-answer pairs that specifically focus on open-space traffic reasoning in real-world safety-related scenarios. Comprising 2,000 carefully selected videos and 19,000 reasoning question-answer pairs, the M4R features a size competitive with existing benchmarks, thus highlighting the comprehensiveness of our evaluation set.

Table 2: Benchmark comparison for multi-modal understanding and reasoning tasks.

Dataset	Safety	Traffic	Annotation	Real-World	Scenarios	# Video	# Ave. Duration (s)	Question-answering Number	Type
MovieChat-1K [35]	✗	✗	Human	✓	General	1,000	564	13,000	Open-ended
MMWorld [14]	✗	✗	Human	✓	General	1,910	107	6,627	Multiple-choice
MLVU [48]	✗	✗	Human	✓	General	1,730	930	3,102	Multiple-choice
MVBench [1]	✗	✗	Human & LLM	✓	General	4,000	16	4,000	Multiple-choice
LongVideoBench [40]	✗	✗	Human	✓	General	3,763	473	6,678	Multiple-choice
TempCompass [24]	✗	✗	Human & LLM	✓	General	410	< 30	7,540	Multiple-choice
VSI-Bench [44]	✗	✗	Human	✓	Embodied	288	50-100	5,000	Multiple-choice
Video-MMMU [16]	✗	✗	Human & LLM	✗	Professional	300	506	900	Multiple-choice
Video-MMLU [36]	✗	✗	Human & LLM	✗	Professional	1,065	109	15,746	Open-ended
DriveBench [42]	✓	✓	Human & LLM	✓	Autonomous Driving	✗	✗	19,200	Multiple-choice
DriveLM [33]	✓	✓	Human	✓	Autonomous Driving	✗	✗	15,480	Open-ended
nuScenes-QA [29]	✗	✓	Human	✓	Autonomous Driving	✗	✗	83,337	Open-ended
MSSBench [49]	✓	✗	Human & LLM	✓	General	✗	✗	1960	Open-ended
MMSBench [23]	✓	✗	LLM	✓	General	✗	✗	5040	Open-ended
M4R (ours)	✓	✓	Human	✓	General	2000	56	19,000	Multiple-choice

4 Experiments

4.1 Model Error Analysis

To demonstrate the effectiveness of our benchmark and evaluate the performance of state-of-the-art (SOTA) models, we conduct a qualitative analysis of model predictions on the M4R benchmark. As shown in Figure 5, the analysis highlights persistent challenges in spatial, temporal, and intent reasoning across open-space environments, particularly in land and air domains. Despite the strong overall performance of leading multi-modal models such as ChatGPT-4o and Gemini 2.5, the results reveal consistent failure cases in real-world scenarios. For example, both models struggle with accurately identifying spatial relationships (e.g., relative positions of vehicles), counting dynamic objects over time (e.g., cars in motion), and understanding goal-directed interactions (e.g., airplane passing events).

These failure cases underscore the limitations of current models in handling safety-critical, perception-intensive tasks. By providing richly annotated, video-based tasks that demand multi-step reasoning grounded in physics, causality, and spatial understanding, M4R serves as a rigorous diagnostic

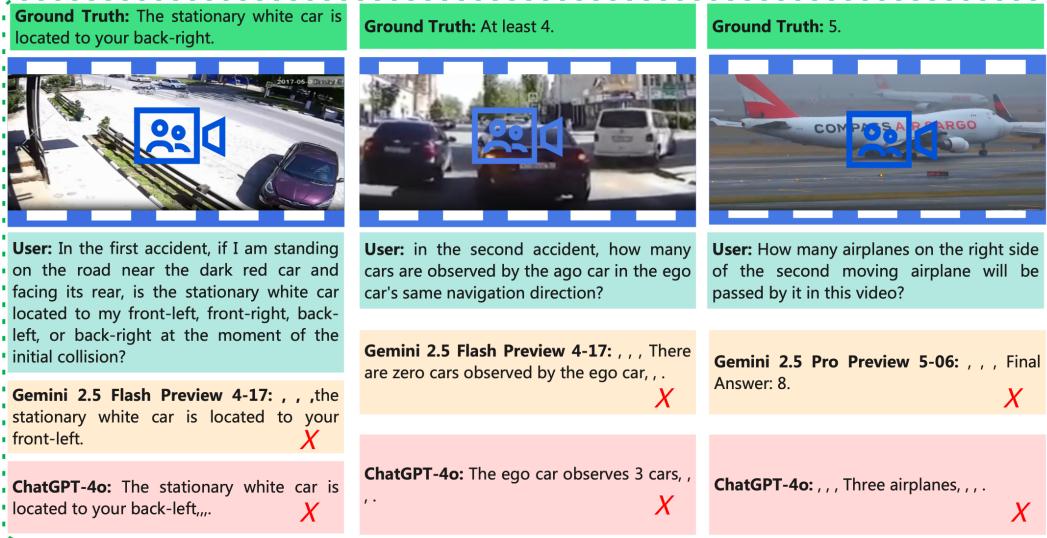


Figure 5: Qualitative error analysis of state-of-the-art multi-modal models (Gemini 2.5 and ChatGPT-4o) on the M4R benchmark. Each example illustrates a failure case in a different reasoning category: spatial reasoning (left), temporal reasoning (middle), and intent reasoning (right). Despite their capabilities, both models struggle with spatial localization, counting dynamic objects, and understanding goal-directed motion in real-world open-space scenarios.

benchmark. Our findings highlight the necessity of such benchmarks for advancing the robustness, safety, and real-world applicability of large multi-modal systems.

4.2 Comprehensive Experiments

In our experiments, we build upon the `lmms-eval` framework [46] as the foundation for our benchmark and extend it to support the specific requirements of M4R. We conduct comprehensive evaluations to assess the performance of SOTA multi-modal models across diverse open-space scenarios.

Land Space Analysis: As shown in Table 3, we present a detailed evaluation of model performance in the **Land Space** domain of M4R, categorized by reasoning type, video length, and difficulty level. InternVL2.5 [8] achieves the highest accuracy in easier settings, suggesting that simply scaling up model size does not always lead to better reasoning performance and may, in some cases, degrade specific capabilities. However, performance drops significantly across all models as tasks become harder and video contexts lengthen. Notably, both GPT-4o [17] and Gemini 1.5 Pro [9] achieve around 40% overall accuracy, reflecting competitive performance while also highlighting persistent challenges in temporal, spatial, and intent-based reasoning within complex, real-world scenarios.

Air Space Analysis: Table 4 presents model performance in the **Air Space** domain of M4R, evaluated across short, medium, and long video scenarios and categorized by temporal, spatial, and intent reasoning tasks. In the easy setting, Qwen2.5 (32B) [4] achieves the highest overall score (52.45%), outperforming GPT-4o and Gemini. However, in the medium and hard settings, Gemini 1.5 Pro outperforms all other models, achieving the top overall accuracy (38.78% in medium and 22.34% in hard), demonstrating better robustness under increasing reasoning difficulty. These results highlight the relative strengths of different models and the increasing challenge of reasoning in dynamic airspace environments as task complexity grows. Moreover, Table 5 presents model performance on the M4R benchmark in the **Water Space** domain, covering both river and ocean scenarios across varying reasoning types and difficulty levels. Gemini 1.5 Pro consistently outperforms other models across all settings.

These findings demonstrate M4R’s ability to *reveal the limitations* of existing multi-modal models, particularly in safety-critical and physically grounded domains. By highlighting domain-specific reasoning gaps, especially in underexplored high-stakes environments such as autonomous driving, ship navigation, and airspace, M4R serves as a tool for guiding the development of more robust, temporally aware, and intent-aware multi-modal systems.

Table 3: Evaluation of M4R in the **Land Space** domain using **Short**, **Medium**, and **Long** Videos, categorized by reasoning types.

Difficulty	Models	Size	Over. Avg.	Short Video Scenarios			Medium Video Scenarios			Long Video Scenarios			
				Avg.	Temporal	Spatial	Intent	Avg.	Temporal	Spatial	Intent	Avg.	Temporal
Easy	Qwen2.5 VL [4]	3B	42.00	49.33	38.0	55.0	55.0	34.67	42.0	34.0	28.0	42.00	36.0
	Qwen2.5 VL [4]	7B	40.67	51.33	55.0	42.0	57.0	36.00	32.0	42.0	34.0	34.67	34.0
	Qwen2.5 VL [4]	32B	43.22	51.00	58.0	50.0	45.0	41.33	46.0	38.0	40.0	37.33	32.0
	LLaVA OneVision [21]	7B	29.78	32.00	31.0	33.0	32.0	24.00	26.0	30.0	16.0	33.33	28.0
	LLaVA Video [47]	7B	31.44	33.00	30.0	31.0	38.0	33.33	38.0	36.0	26.0	28.00	16.0
	LLaVA Next [20]	32B	31.25	38.00	35.0	45.0	34.0	21.33	12.0	14.0	38.0	34.67	20.0
	GPT 4o [17]	-	42.17	52.35	59	47.06	51	47.16	54.9	44.9	41.67	27.00	5
Medium	Qwen2.5 VL [4]	3B	30.78	41.67	33.0	52.0	40.0	24.67	30.0	30.0	14.0	26.00	22.0
	Qwen2.5 VL [4]	7B	29.89	39.00	37.0	42.0	38.0	30.67	32.0	40.0	20.0	16.0	26.0
	Qwen2.5 VL [4]	32B	28.55	28.33	21.0	44.0	20.0	33.33	40.0	30.0	30.0	24.00	8.0
	LLaVA OneVision [21]	7B	16.67	16.00	26.0	30.0	16.0	14.67	18.0	8.0	18.0	19.33	12.0
	LLaVA Video [47]	7B	25.67	25.00	20.0	34.0	26.0	28.67	36.0	28.0	22.0	23.33	14.0
	LLaVA Next [20]	32B	20.0	27.33	16.0	49.0	17.0	10.67	14.0	10.0	8.0	22.0	16.0
	GPT 4o [17]	-	36.99	45.49	48.48	55	33	33.89	41.67	26.67	33.33	31.33	24
Hard	Qwen2.5 VL [4]	3B	22.78	23.00	17.0	33.0	19.0	26.67	38.0	26.0	16.0	18.67	10.0
	Qwen2.5 VL [4]	7B	22.89	26.00	17.0	30.0	31.0	40.0	32.0	18.0	12.67	2.0	30.0
	Qwen2.5 VL [4]	32B	22.66	19.33	11.0	34.0	13.0	35.33	46.0	24.0	36.0	13.33	4.0
	LLaVA OneVision [21]	7B	13.67	14.33	5.0	27.0	11.0	14.67	18.0	8.0	18.0	12.0	6.0
	LLaVA Video [47]	7B	19.78	19.33	12.0	35.0	11.0	24.67	26.0	30.0	18.0	15.33	10.0
	LLaVA Next [20]	32B	16.22	20.67	16.0	32.0	14.0	11.33	12.0	12.0	10.0	16.67	10.0
	GPT 4o [17]	-	24.41	26.78	34.65	34.69	11	35.70	43.14	32.14	31.82	11.00	6

Table 4: Evaluation of M4R in the **Air Space** domain using **Short**, **Medium**, and **Long** Videos, categorized by reasoning types.

Difficulty	Models	Size	Over. Avg.	Short Video Scenarios			Medium Video Scenarios			Long Video Scenarios			
				Avg.	Temporal	Spatial	Intent	Avg.	Temporal	Spatial	Intent	Avg.	Temporal
Easy	Qwen2.5 VL [4]	3B	43.67	43.33	36.00	52.00	42.00	42.67	38.00	54.00	36.00	45.00	40.00
	Qwen2.5 VL [4]	7B	39.89	33.33	28.00	18.00	54.00	38.00	48.00	16.00	50.00	48.33	55.00
	Qwen2.5 VL [4]	32B	52.45	50.00	34.00	56.00	60.00	50.67	40.00	54.00	58.00	56.67	55.00
	LLaVA OneVision [21]	7B	33.22	33.33	34.00	38.00	28.00	34.67	34.00	38.00	32.00	31.67	35.00
	LLaVA Video [47]	7B	33.22	33.33	34.00	38.00	28.00	34.67	34.00	38.00	32.00	31.67	35.00
	LLaVA Next [20]	32B	33.22	36.67	36.00	42.0	32.0	31.33	36.0	32.0	26.0	31.67	35.0
	GPT 4o [17]	-	40.67	35.33	30.00	28.00	48.00	36.67	24.00	38.00	48.00	50.00	45.00
Medium	Qwen2.5 VL [4]	3B	28.67	22.67	18.00	36.00	14.00	26.67	18.00	44.00	18.00	36.67	40.00
	Qwen2.5 VL [4]	7B	28.00	24.67	16.00	24.00	34.00	26.00	24.00	28.00	33.33	35.00	20.00
	Qwen2.5 VL [4]	32B	33.34	32.67	12.00	48.00	38.00	30.67	22.00	50.00	20.00	36.67	20.00
	LLaVA OneVision [21]	7B	23.67	23.33	20.00	34.00	16.00	22.67	20.00	32.00	16.00	25.00	20.00
	LLaVA Video [47]	7B	23.22	23.33	28.00	44.00	30.00	39.33	30.00	44.00	44.00	38.33	35.00
	LLaVA Next [20]	32B	26.11	24.67	18.0	40.0	16.0	25.33	18.0	40.0	18.0	28.33	25.0
	GPT 4o [17]	-	38.45	38.67	38.00	56.00	22.00	30.00	38.00	34.00	18.00	46.67	65.00
Hard	Qwen2.5 VL [4]	3B	15.33	16.00	8.00	32.00	8.00	13.33	6.00	26.00	8.00	16.67	10.00
	Qwen2.5 VL [4]	7B	16.55	19.33	0.00	30.00	28.00	15.33	2.00	30.00	14.00	15.00	5.00
	Qwen2.5 VL [4]	32B	16.22	20.00	6.00	36.00	18.00	15.33	4.00	24.00	18.00	13.33	0.00
	LLaVA OneVision [21]	7B	15.67	16.00	12.00	28.00	8.00	16.00	12.00	26.00	10.00	15.00	10.00
	LLaVA Video [47]	7B	14.78	16.67	14.00	28.00	8.00	12.67	6.00	22.00	10.00	15.00	5.00
	LLaVA Next [20]	32B	17.89	18.67	14.0	34.0	8.0	16.67	6.0	32.0	12.0	18.33	5.0
	GPT 4o [17]	-	22.34	26.67	24.00	26.00	30.00	18.67	20.00	22.00	14.00	21.67	10.00

4.3 Ablation Experiments

In our experiments, due to the high cost of evaluating all data points, we adopt a uniform sampling strategy to select a representative subset of tasks. Specifically, for each reasoning type, we sample 50 tasks when the total number of available tasks is fewer than 500, and 100 tasks when the number exceeds 500. The M4R benchmark spans three open-space scenarios—*land space*, *air space*, and *water space*—each with three video lengths (short, medium, long), three difficulty levels (easy, medium, hard), and three reasoning types: temporal, spatial, and intent-based reasoning.

Following this sampling strategy, we evaluate a total of 3,798 tasks, evenly distributed across the three reasoning types: 1,266 *spatial reasoning*, 1,266 *temporal-causal reasoning*, and 1,266 *intent and goal reasoning* tasks.

Table 5: Evaluation of M4R in the **Water Space** domain using **River** and **Ocean** Videos, categorized by reasoning types.

Difficulty	Models	Size	Over. Avg.	River Scenarios			Ocean Scenarios		
				Avg.	Temporal	Spatial	Intent	Avg.	Temporal
Easy	Qwen2.5 VL [4]	3B	40.21	39.74	34.62	46.15	38.46	40.67	34.00
	Qwen2.5 VL [4]	7B	31.31	34.62	38.46	19.23	46.15	28.00	36.00
	Qwen2.5 VL [4]	32B	52.77	61.54	53.85	61.54	69.23	44.00	40.00
	LLaVA OneVision [21]	7B	33.00	33.33	34.62	34.62	30.77	32.67	28.00
	LLaVA Video [47]	7B	31.03	32.05	30.77	34.62	30.77	30.00	22.00
	LLaVA Next [20]	32B	35.59	37.18	26.92	53.85	30.77	34.00	30.00
	InternVL2.5 [8]	4B	53.87	56.41	53.85	57.69	57.69	51.33	52.00
	InternVL2.5 [8]	8B	53.47	60.26	69.23	46.15	65.38	46.67	46.00
	InternVL2.5 [8]	26B	55.05	64.10	65.38	57.69	69.23	46.00	50.00
Medium	Gemini 1.5 pro [9]	-	50.69	52.56	42.31	61.54	53.85	48.81	50.00
	GPT 4o [17]	-	50.51	57.69	57.69	50.00	65.38	43.33	46.43
	GPT 4o [17]	-	50.51	57.69	57.69	50.00	65.38	43.33	46.43
	Qwen2.5 VL [4]	3B	28.08	29.49	23.08	53.85	11.54	26.67	18.00
	Qwen2.5 VL [4]	7B	24.08	29.49	19.23	30.77	38.46	18.67	18.00
	Qwen2.5 VL [4]	32B	33.31	34.62	19.23	50.00	34.62	32.00	20.00
	LLaVA OneVision [21]	7B	22.54	23.08	19.23	30.77	19.23	22.00	14.00
	LLaVA Video [47]	7B	21.92	20.51	19.23	26.92	15.38	23.33	20.00
	LLaVA Next [20]	32B	20.88	23.08	11.54	38.46	19.23	18.67	10.00
Hard	InternVL2.5 [8]	4B	44.36	48.72	23.08	65.38	57.69	40.00	28.00
	InternVL2.5 [8]	8B	41.08	46.15	34.62	61.54	42.31	36.00	34.00
	InternVL2.5 [8]	26B	41.77	44.87	30.77	57.69	46.15	38.67	24.00
	Gemini 1.5 pro [9]	-	46.31	53.84	46.15	65.38	50.00	38.78	34.00
	GPT 4o [17]	-	38.49	42.31	50.00	53.85	23.08	34.67	36.00
	Qwen2.5 VL [4]	3B	14.34	16.67	15.38	19.23	15.38	12.00	12.00
	Qwen2.5 VL [4]	7B	14.67	16.67	7.69	30.77	11.54	12.67	6.00
	Qwen2.5 VL [4]	32B	13.39	14.10	7.69	23.08	11.54	12.67	8.0
	LLaVA OneVision [21]	7B	15.67	16.67	11.54	26.92	11.54	14.67	8.00
Hard	LLaVA Video [47]	7B	14.00	16.67	15.38	23.08	11.54	11.33	8.00
	LLaVA Next [20]	32B	14.39	11.54	7.69	19.23	7.69	15.33	8.0
	InternVL2.5 [8]	4B	20.92	20.51	19.23	19.23	23.08	21.33	16.00
	InternVL2.5 [8]	8B	21.90	21.79	7.69	26.92	30.77	22.00	16.00
	InternVL2.5 [8]	26B	22.54	23.08	15.38	19.23	34.62	22.00	18.00
	Gemini 1.5 pro [9]	-	26.02	26.92	23.08	30.77	26.92	25.11	34.00
	GPT 4o [17]	-	22.10	28.20	38.46	26.92	19.23	16.00	18.00
	GPT 4o [17]	-	22.10	28.20	38.46	26.92	19.23	16.00	18.00

To assess the reliability of this sampling approach, we conduct an ablation study comparing model performance on sampled tasks versus the full set of data points in the **land space (short, easy)** setting. We use InternVL 2.5, one of the leading open-source multi-modal models, which ranks highly on several leaderboards such as ⁶ and ⁷. As shown in Table 6, performance on the sampled subset is comparable to, and in some cases slightly better than, performance on the full dataset. These results validate the effectiveness of our sampling strategy in preserving benchmark consistency while reducing evaluation cost.

Table 6: Performance Comparison on **Land Space Short** (Easy): Full vs. Sample Data Points

Model	Full Data Points				Sample Data Points			
	Avg.	Temporal	Spatial	Intent	Avg.	Temporal	Spatial	Intent
InternVL2.5-26B	55.62	57.61	50.37	58.88	61.00	62.00	59.00	62.00
InternVL2.5-8B	49.26	51.89	48.57	47.31	55.67	55.00	60.00	52.00
InternVL2.5-4B	50.65	50.17	50.70	51.10	55.33	52.00	55.00	59.00

5 Conclusion

In this work, we introduce M4R, a large-scale benchmark for evaluating multi-modal understanding and reasoning in real-world open-space environments. Spanning three critical domains—land, air, and water—M4R provides richly annotated, video-based tasks designed to assess model performance across three fundamental reasoning dimensions: temporal reasoning, spatial reasoning, and intent and goal inference. The benchmark encompasses a broad range of scenarios, video lengths, and difficulty levels, enabling comprehensive evaluation in safety-critical, perception-intensive settings. Through extensive qualitative and quantitative analyses, we demonstrate that even state-of-the-art multi-modal models—both proprietary systems such as ChatGPT-4o and Gemini 2.5, and leading open-source models like Qwen and InternVL—exhibit significant limitations when reasoning over complex, dynamic physical environments. These results underscore the need for more robust, temporally-aware, and goal-sensitive multi-modal systems capable of reliable understanding in real-world scenarios. We hope that M4R will serve as a valuable resource for the research community and help advance the development of safer, more generalizable, and practically deployable multi-modal AI systems.

⁶<https://enxinsong.com/Video-MMLU-web/>

⁷https://huggingface.co/spaces/opencompass/open_vlm_leaderboard

References

- [1] Guillermo Franco Abellán, Matteo Braglia, Mario Ballardini, Fabio Finelli, and Vivian Poulin. Probing early modification of gravity with planck, act and spt. *Journal of Cosmology and Astroparticle Physics*, 2023(12):017, 2023.
- [2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [3] Anthropic. Claude 3.5 sonnet model card addendum, 2024.
- [4] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
- [5] Wentao Bao, Qi Yu, and Yu Kong. Uncertainty-based traffic accident anticipation with spatio-temporal relational learning. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 2682–2690, 2020.
- [6] Wenhao Chai, Enxin Song, Yilun Du, Chenlin Meng, Vashisht Madhavan, Omer Bar-Tal, Jenq-Neng Hwang, Saining Xie, and Christopher D Manning. Auroracap: Efficient, performant video detailed captioning and a new benchmark. *arXiv preprint arXiv:2410.03051*, 2024.
- [7] Jr-Jen Chen, Yu-Chien Liao, Hsi-Che Lin, Yu-Chu Yu, Yen-Chun Chen, and Frank Wang. Rex-time: A benchmark suite for reasoning-across-time in videos. *Advances in Neural Information Processing Systems*, 37:28662–28673, 2024.
- [8] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024.
- [9] Google DeepMind. Gemini 1.5 technical report. <https://deepmind.google/technologies/gemini/#gemini-15>, 2024. Accessed: 2025-05-12.
- [10] Chaoyou Fu, Yuhan Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. *arXiv preprint arXiv:2405.21075*, 2024.
- [11] Hanan Gani, Rohit Bharadwaj, Muzammal Naseer, Fahad Shahbaz Khan, and Salman Khan. Vane-bench: Video anomaly evaluation benchmark for conversational lmms. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 3123–3140, 2025.
- [12] Yichen Gong, Delong Ran, Jinyuan Liu, Conglei Wang, Tianshuo Cong, Anyu Wang, Sisi Duan, and Xiaoyun Wang. Figstep: Jailbreaking large vision-language models via typographic visual prompts. *arXiv preprint arXiv:2311.05608*, 2023.
- [13] Yu Guo, Ryan Wen Liu, Jingxiang Qu, Yuxu Lu, Fenghua Zhu, and Yisheng Lv. Asynchronous trajectory matching-based multimodal maritime data fusion for vessel traffic surveillance in inland waterways. *IEEE Transactions on Intelligent Transportation Systems*, 24(11):12779–12792, 2023.
- [14] Xuehai He, Weixi Feng, Kaizhi Zheng, Yujie Lu, Wanrong Zhu, Jiachen Li, Yue Fan, Jianfeng Wang, Linjie Li, Zhengyuan Yang, et al. Mmworld: Towards multi-discipline multi-faceted world model evaluation in videos. *arXiv preprint arXiv:2406.08407*, 2024.
- [15] Xuehai He, Weixi Feng, Kaizhi Zheng, Yujie Lu, Wanrong Zhu, Jiachen Li, Yue Fan, Jianfeng Wang, Linjie Li, Zhengyuan Yang, et al. Mmworld: Towards multi-discipline multi-faceted world model evaluation in videos. In *The Thirteenth International Conference on Learning Representations*, 2025.

- [16] Kairui Hu, Penghao Wu, Fanyi Pu, Wang Xiao, Yuanhan Zhang, Xiang Yue, Bo Li, and Ziwei Liu. Video-mmmu: Evaluating knowledge acquisition from multi-discipline professional videos. *arXiv preprint arXiv:2501.13826*, 2025.
- [17] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
- [18] Dongzhi Jiang, Renrui Zhang, Ziyu Guo, Yanwei Li, Yu Qi, Xinyan Chen, Liuhui Wang, Jianhan Jin, Claire Guo, Shen Yan, et al. Mme-cot: Benchmarking chain-of-thought in large multimodal models for reasoning quality, robustness, and efficiency. *arXiv preprint arXiv:2502.09621*, 2025.
- [19] Minkuk Kim, Hyeyon Bae Kim, Jinyoung Moon, Jinwoo Choi, and Seong Tae Kim. Hicm²: Hierarchical compact memory modeling for dense video captioning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 4293–4301, 2025.
- [20] Bo Li, Hao Zhang, Kaichen Zhang, Dong Guo, Yuanhan Zhang, Renrui Zhang, Feng Li, Ziwei Liu, and Chunyuan Li. Llava-next: What else influences visual instruction tuning beyond data?, May 2024.
- [21] 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.
- [22] Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, et al. Mvbench: A comprehensive multi-modal video understanding benchmark. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22195–22206, 2024.
- [23] Xin Liu, Yichen Zhu, Jindong Gu, Yunshi Lan, Chao Yang, and Yu Qiao. Mm-safetybench: A benchmark for safety evaluation of multimodal large language models. In *European Conference on Computer Vision*, pages 386–403. Springer, 2024.
- [24] Yuanxin Liu, Shicheng Li, Yi Liu, Yuxiang Wang, Shuhuai Ren, Lei Li, Sishuo Chen, Xu Sun, and Lu Hou. Tempcompass: Do video llms really understand videos? *arXiv preprint arXiv:2403.00476*, 2024.
- [25] Weidi Luo, Siyuan Ma, Xiaogeng Liu, Xiaoyu Guo, and Chaowei Xiao. Jailbreaky: A benchmark for assessing the robustness of multimodal large language models against jailbreak attacks. *arXiv preprint arXiv:2404.03027*, 2024.
- [26] Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. Egoschema: A diagnostic benchmark for very long-form video language understanding. *Advances in Neural Information Processing Systems*, 36:46212–46244, 2023.
- [27] Dilip K Prasad, Deepu Rajan, Lily Rachmawati, Eshan Rajabally, and Chai Quek. Video processing from electro-optical sensors for object detection and tracking in a maritime environment: A survey. *IEEE Transactions on Intelligent Transportation Systems*, 18(8):1993–2016, 2017.
- [28] Xiangyu Qi, Kaixuan Huang, Ashwinee Panda, Peter Henderson, Mengdi Wang, and Prateek Mittal. Visual adversarial examples jailbreak aligned large language models. In *Proceedings of the AAAI conference on artificial intelligence*, volume 38, pages 21527–21536, 2024.
- [29] Tianwen Qian, Jingjing Chen, Linhai Zhuo, Yang Jiao, and Yu-Gang Jiang. Nuscenes-qa: A multi-modal visual question answering benchmark for autonomous driving scenario. In *AAAI Conference on Artificial Intelligence*, volume 38, pages 4542–4550, 2024.
- [30] Ankit Parag Shah, Jean-Baptiste Lamare, Tuan Nguyen-Anh, and Alexander Hauptmann. Cadp: A novel dataset for cctv traffic camera based accident analysis. In *2018 15th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*, pages 1–9. IEEE, 2018.
- [31] Erfan Shayegani, Yue Dong, and Nael Abu-Ghazaleh. Jailbreak in pieces: Compositional adversarial attacks on multi-modal language models. *arXiv preprint arXiv:2307.14539*, 2023.

- [32] Chonghao Sima, Katrin Renz, Kashyap Chitta, Li Chen, Hanxue Zhang, Chengen Xie, Jens Beißwenger, Ping Luo, Andreas Geiger, and Hongyang Li. Drivelm: Driving with graph visual question answering. In *European Conference on Computer Vision*, pages 256–274. Springer, 2024.
- [33] Chonghao Sima, Katrin Renz, Kashyap Chitta, Li Chen, Hanxue Zhang, Chengen Xie, Jens Beißwenger, Ping Luo, Andreas Geiger, and Hongyang Li. Drivelm: Driving with graph visual question answering. In *European Conference on Computer Vision*, pages 256–274, 2024.
- [34] Enxin Song, Wenhao Chai, Guanhong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Haozhe Chi, Xun Guo, Tian Ye, Yanting Zhang, et al. Moviechat: From dense token to sparse memory for long video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18221–18232, 2024.
- [35] Enxin Song, Wenhao Chai, Guanhong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Xun Guo, Tian Ye, Yan Lu, Jenq-Neng Hwang, et al. Moviechat: From dense token to sparse memory for long video understanding. *arXiv preprint arXiv:2307.16449*, 2023.
- [36] Enxin Song, Wenhao Chai, Weili Xu, Jianwen Xie, Yuxuan Liu, and Gaoang Wang. Video-mmlu: A massive multi-discipline lecture understanding benchmark. *arXiv preprint arXiv:2504.14693*, 2025.
- [37] 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. In *8th Annual Conference on Robot Learning*, 2025.
- [38] Weihan Wang, Zehai He, Wenyi Hong, Yean Cheng, Xiaohan Zhang, Ji Qi, Xiaotao Gu, Shiyu Huang, Bin Xu, Yuxiao Dong, et al. Lvbench: An extreme long video understanding benchmark. *arXiv preprint arXiv:2406.08035*, 2024.
- [39] Hongchen Wei, Zhihong Tan, Yaosi Hu, Chang Wen Chen, and Zhenzhong Chen. Longcaptioning: Unlocking the power of long video caption generation in large multimodal models. *arXiv preprint arXiv:2502.15393*, 2025.
- [40] Haoning Wu, Dongxu Li, Bei Chen, and Junnan Li. Longvideobench: A benchmark for long-context interleaved video-language understanding. *Advances in Neural Information Processing Systems*, 37:28828–28857, 2024.
- [41] Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. Next-qa: Next phase of question-answering to explaining temporal actions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9777–9786, 2021.
- [42] 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.
- [43] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5288–5296, 2016.
- [44] Jihan Yang, Shusheng Yang, Anjali W Gupta, Rilyn Han, Li Fei-Fei, and Saining Xie. Thinking in space: How multimodal large language models see, remember, and recall spaces. *arXiv preprint arXiv:2412.14171*, 2024.
- [45] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruqi 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*, pages 9556–9567, 2024.
- [46] Kaichen Zhang, Bo Li, Peiyuan Zhang, Fanyi Pu, Joshua Adrian Cahyono, Kairui Hu, Shuai Liu, Yuanhan Zhang, Jingkang Yang, Chunyuan Li, and Ziwei Liu. Lmms-eval: Reality check on the evaluation of large multimodal models, 2024.

- [47] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Video instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*, 2024.
- [48] Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Shitao Xiao, Xi Yang, Yongping Xiong, Bo Zhang, Tiejun Huang, and Zheng Liu. Mlvu: A comprehensive benchmark for multi-task long video understanding. *arXiv preprint arXiv:2406.04264*, 2024.
- [49] Kaiwen Zhou, Chengzhi Liu, Xuandong Zhao, Anderson Compalas, Dawn Song, and Xin Eric Wang. Multimodal situational safety. *arXiv preprint arXiv:2410.06172*, 2024.
- [50] Chengke Zou, Xingang Guo, Rui Yang, Junyu Zhang, Bin Hu, and Huan Zhang. Dynamath: A dynamic visual benchmark for evaluating mathematical reasoning robustness of vision language models. *arXiv preprint arXiv:2411.00836*, 2024.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: See the abstract and introduction.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: See the end of the main paper.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: This is a dataset and benchmark paper.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Yes, see our code link.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: See the code link.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: See our code link and our paper.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: See the experiment section.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: See the experiment settings.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes].

Justification: conducted the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: See the end of the paper.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to

generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.

- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: Our study poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The benchmark used some existing data points and mentioned them properly.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: The paper provides a new dataset.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorosity, or originality of the research, declaration is not required.

Answer: [NA]