2

3

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

32

MMR-V: What's Left Unsaid? A Benchmark for **Multimodal Deep Reasoning in Videos**

Kejian Zhu^{1,2}, Zhuoran Jin^{1,2}, Hongbang Yuan^{1,2}, Jiachun Li^{1,2}, Shangqing Tu³ Pengfei Cao^{1,2}, Yubo Chen^{1,2}, Kang Liu^{1,2}, Jun Zhao^{1,2}

¹School of Artificial Intelligence, University of Chinese Academy of Sciences ²Institute of Automation, Chinese Academy of Sciences ³Tsinghua University zhukejian2025@ia.ac.cn {zhuoran.jin, hongbang.yuan} @nlpr.ia.ac.cn {pengfei.cao, yubo.chen, kliu, jzhao} @nlpr.ia.ac.cn

Abstract

The sequential structure of videos poses a challenge to the ability of multimodal large language models (MLLMs) to locate multi-frame evidence and conduct multimodal reasoning. However, existing video benchmarks mainly focus on understanding tasks, which only require models to match frames mentioned in the question (hereafter referred to as "question frame") and perceive a few adjacent frames. To address this gap, we propose MMR-V: A Benchmark for Multimodal **Deep Reasoning in Videos.** The benchmark is characterized by the following features. (1) Long-range, multi-frame reasoning: Models are required to infer and analyze evidence frames that may be far from the question frame. (2) **Beyond perception**: Questions cannot be answered through direct perception alone but require reasoning over hidden information. (3) Reliability: All tasks are manually annotated, referencing extensive real-world user understanding to align with common perceptions. (4) Confusability: Carefully designed distractor annotation strategies to reduce model shortcuts. MMR-V consists of 317 videos and 1,257 tasks. Our experiments reveal that current models still struggle with multi-modal reasoning; even the best-performing model, o4-mini, achieves only 52.5% accuracy. Additionally, current reasoning enhancement strategies (Chain-of-Thought and scaling test-time compute) bring limited gains. Further analysis indicates that the CoT demanded for multi-modal reasoning differs from it in textual reasoning, which partly explains the limited performance gains. We hope that MMR-V can inspire further research into enhancing multi-modal reasoning capabilities.

()

Benchmark Code

https://huggingface.co/datasets/JokerJan/MMR-VBench https://github.com/GaryStack/MMR-V

Introduction

Recent models like OpenAI's o1 [1] and Deepseek-R1 [2] have significantly improved text reasoning ability through reinforcement learning. This has sparked growing interest in multimodal reasoning [3]. Models like o3 and o4-mini [4] have achieved impressive results on image reasoning tasks through tool use, integrating visual information into the reasoning process to enable deep reflection and evidence mining. However, most of these studies focus on images, with limited exploration of more challenging video reasoning tasks. Video naturally involves sequential and richer multimodal information, requiring models to perform reasoning and mine evidence over long-range, multi-frame. Since this capability is essential for real-world applications such as embodied intelligence and intelligent security monitoring [5; 6], it naturally raises an important question: can current MLLMs perform deep multimodal reasoning and mine evidence on complex videos like o3 on image tasks?

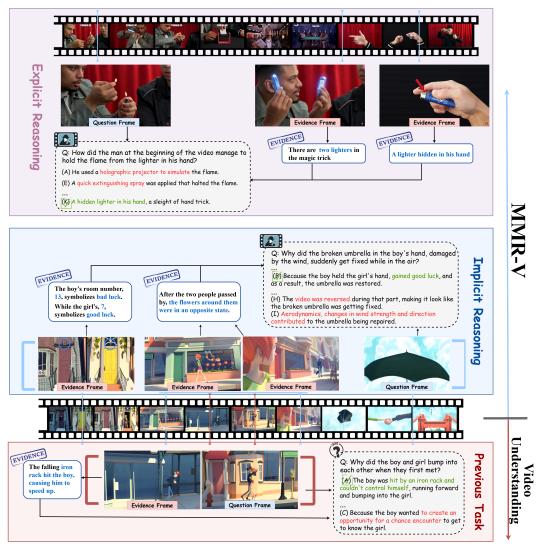


Figure 1: Examples showing the MMR-V tasks and the difference from previous video benchmarks.

However, existing video benchmarks primarily focus on perception and understanding tasks [7; 8].

These tasks often only require locating frames mentioned in the question and understanding adjacent 35 frames. For example, at the bottom of Figure 1, noticing the boy being hit by the metal frame 36 is enough to understand why he ran into the girl. Such tasks fall short in evaluating multimodal 37 reasoning abilities. We summarize their limitations as follows: (1) Limited frame context: Even for 38 long videos, existing tasks often rely on just a few adjacent frames, failing to exploit the long-range 39 40 sequential structure of the video. (2) Lack of reasoning: Many questions can be answered through direct perception. (3) Unrealistic task: Simple perception and adjacent-frame understanding tasks 41 do not meet the real-world demands for AI system strong capabilities. 42 To address these shortcomings, we propose MMR-V Bench: A Benchmark for Multi-modal Deep 43 **Reasoning in Videos.** We present two examples to illustrate the key differences with previous video 44 understanding benchmarks in Figure 1. MMR-V offers the following features: (1) Long-range, multi-frame reasoning: tasks involve multimodal reasoning over non-adjacent video frames to 46 locate and analyze multiple evidences; (2) Beyond perception: questions cannot be answered by 47 direct perception of question frame directly, requiring reasoning and the extraction of implications; 48 (3) Reliability: All tasks are annotated manually, and potential subjective bias is reduced by cross-49 referencing the most popular video comments. (4) Confusability: We employ carefully designed 50 annotation strategies to craft model-aligned distractor options, thereby ensuring confusability.

Inspired by cognitive and psychological theories [9; 10; 11], such as Kahneman's Dual Process Theory [12], we categorize the tasks in MMR-V into **implicit reasoning** and **explicit reasoning**. The 53 key distinction lies in whether the question requires reasoning beyond surface-level information to 54 infer underlying implications. Explicit reasoning is defined as questions that can be solved using 55 perceivable information from the video. For example, the task shown in Figure 1 requires noticing the 56 two lighters hidden in the hand. Implicit reasoning requires extracting and interpreting the underlying 57 subtext behind visual information. For example, in the implicit reasoning case shown in Figure 1, it requires inferring the underlying implication that the girl's room number 7 symbolizes good luck. This is more of an assessment of EQ, testing whether the model can use its deep understanding of the 60 world knowledge to make implicit and subconscious reasoning paths like humans. 61

MMR-V comprises 317 videos and 1257 tasks. The videos span six major categories, with lengths ranging from 7 to 3771 seconds, with an average of 277 seconds. Tasks are further divided into 10 categories and subcategories. Each task is in multiple-choice format with approximately ten options on average. Tasks typically require reasoning over average 12 video frames, covering about 60% of video duration. All questions and correct answers are human-annotated and reviewed. Distractors are generated using a carefully designed annotation strategy (Details in Section 3.2).

We evaluated 9 proprietary models and 11 open-source models on MMR-V. The results reveal that 68 even the best-performing model, o4-mini, achieved only 52.5% accuracy, highlighting the significant 69 challenge MMR-V poses to current multimodal large language models. Our key findings are as 70 follows. (1) Multimodal reasoning challenge: Our findings in Section 4.2 show that reasoning 71 enhancement strategies (e.g., CoT and scaling test-time compute) yield limited improvements, indicating that MMR-V presents a greater challenge to current multimodal reasoning models. Further 73 error analysis in Section 4.5 shows that the CoT demanded in multimodal reasoning differs from those 74 in textual reasoning. Current models tend to rely on textual reasoning based on visual information 75 from the question frame and few adjacent frames, lacking the multimodal reasoning needed to 76 locate and analyze evidence from long-range frames. This limitation hinders the overall reasoning 77 performance. (2) More modality will benefit: We found that for models that support all modalities, 78 adding additional audio modalities will improve the performance (Accuracy improved by 1.4%, 79 1.0%, and 1.0% for Gemini 2.0-Flash, Gemini 2.0-Flash-Thinking, and Phi-4-Multimodal-Instruct, 80 respectively). (3) Human-model gap: In human experiments, we found that although models exhibit 81 human-level performance on text reasoning tasks, there is still a significant gap between model and 82 human on multimodal, especially video, reasoning tasks. We hope MMR-V will inspire further 83 research into enhancing multimodal reasoning capabilities in AI systems.

Task Overview 85

62

63

64

67

86

87

88

89

91

92

93

94

95

96

97

98

99

100

101

102

The tasks in MMR-V require deeper multimodal reasoning. Unlike previous tasks such as math and puzzle problems [13; 14; 15], we argue that the scope of multimodal reasoning should be more broadly defined. Previous work focuses more on text-oriented reasoning based on perceived visual information. In contrast, our task requires integrating the various forms of visual evidences, such as artistic style, lighting, and depth, into the reasoning process. Even more challenging, it involves reasoning over long-range, multi-frame visual evidence. Videos have a temporal dimension, which puts a greater challenge on the ability to find clues in different frames through multimodal reasoning.

2.1 Definition for Implicit and Explicit Reasoning.

We categorize reasoning tasks in MMR-V into **Implicit Reasoning** and **Explicit Reasoning**, inspired by Kahneman's Dual Process Theory [12] and other cognitive theories [9; 10; 11]. The most obvious difference is whether or not one needs to understand the subtext beneath the surface information. Secondly, implicit reasoning for human is often achieved by experience based on world knowledge, thus consuming little attention resources. Tasks are further divided into 10 categories and 33 subcategories. Six categories are shown in Figure 2, with the first row belonging to implicit and the second row is explicit. Further explanations and examples can be found in Appendix D.

Implicit Reasoning focuses on incorporating hidden meanings behind visual information into reasoning. In these tasks, surface-level visual cues often conceal deeper layers of meaning, such as metaphor. Besides, for human, "(implicit) operates automatically and quickly, with little or no effort and no sense of voluntary control." - Dual Process Theory.

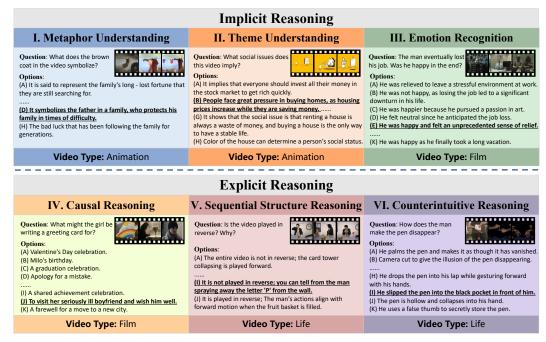


Figure 2: Overview of six tasks in MMR-V Bench.

Explicit Reasoning evaluates whether a model can perform reasoning based on multimodal details **explicitly presented** across long-range, multi-frame of a video. However, solving these tasks demands fine-grained perception and rigorous logical reasoning. "(explicit) allocates attention to the effortful mental activities that demand it, including complex computations." - Dual Process Theory.

2.2 Implicit Reasoning Tasks

105

106

107

108

109

122

- Metaphor Understanding (MU): MU tasks evaluate the ability to reason about metaphors for entities or environment. For example, the case in Figure 2 I interprets the metaphor of the brown coat.
- Theme Understanding (TU): TU assesses the ability to infer the main idea and attitude of the author through the full video. For example, the case in Figure 2 II asks what social issue the video reveals.
- Emotion Recognition (ER): ER tasks evaluate the ability to analyze character emotional states, as well as higher-level emotions such as the author's attitude and the audience's emotional response. For example, the case in Figure 2 III involves inferring whether the character feels happy at the end.
- 117 **Comment Matching (CM)**: CM task is to predict the most fitting audience comments for a video based on a criteria. For example, selecting which comment would be the most humorous after watching the video. Detailed example can be found in Appendix D.1.
- Implicit Symbol (IS): IS task is to interpret implicit symbols in the video, such as cultural elements. For example, inferring the ethnicity of the filming location. Details can be found in Appendix D.1.

2.3 Explicit Reasoning Tasks

- Causal Reasoning (CAR): CAR assesses the ability to reason about causal relationships in the video. For example, in Figure 2 IV, it involves inferring the reason why the girl is making a card.
- Sequential Structure Reasoning (SSR): SSR tasks assess reasoning about temporal structure in video editing and storytelling. In the example from Figure 2 V, the task is to infer if the video is reversed. However, the creator of this video explains the video is played normally.
- Counterintuitive Reasoning (CIR): CIR tasks evaluate the ability to analyze information that contradicts common sense, requiring detailed cross-frame analysis. In the example from Figure 2 VI, the task is to reason the principle behind the counterintuitive magic trick.

Cross-modal Transfer Reasoning (CTR): To reason and match information out of the video that shares similar meaning. For example, find the quote with same theme of the video.

Video Type and Intent (VTI): VTI tasks test the ability to infer key meta-level information such as the genre and communicative intent of the video from a global perspective. For example, the case in Appendix D.2 infers the release time by reasoning the video is set during COVID-19.

136 3 MMR-V Bench

To ensure that MMR-V effectively evaluate multimodal reasoning abilities, we follow **three principles** during construction: **P1. Multi-frame:** Questions require reference to long-range, multi-frame information, prompting the model to reason across multiple visual cues. **P2. Deep reasoning:**Answers should not be directly perceivable from the video; instead, they should demand understanding of the subtext or multimodal reasoning, reflecting a deep comprehension of the content. **P3. Realistic:**Tasks should align with real-world question-answering needs, ensuring answers are consistent with common user understanding and free from individual cognitive biases or prejudices.

144 3.1 Video Collection

We manually curated a diverse original videos from Youtube with following checklist: (1) **Avoidance of linear, descriptive content**: We excluded videos with straightforward structures, such as daily recordings or sports broadcasts, in order to ensure that the tasks require deep reasoning over multiframes (For Principle P1). (2) **Creative and thematically rich videos**: We selected videos that are intentional designed and edited by creators, often conveying well-crafted themes. This ensures that the questions require interpretation beyond surface-level visual content (For Principle P2). (3)**Alignment with real-world**: Highly Popular Videos were preferred, which are indicated by active comment sections and audience engagement. This helps avoid biases introduced by niche content and ensures alignment with general user cognition (For Principle P3). (4) **Diverse coverage**: To further promote generalizability, we ensured broad coverage across video types, topics, and durations, allowing MMR-V to reflect the diversity of real-world video content (For Principle P3). As a result, our final benchmark comprises **317 videos** spanning **six major categories**: Animation, Film, Philosophy, TV, Life, and Art. The specific categories are shown in the Appendix C. Furthermore, for problems where audio might be helpful, we ensure that the videos include audio.

3.2 Data Annotation & Quality Assurance

All tasks in MMR-V Bench are designed in a multiple-choice format. There is one correct option and several wrong options. Make sure there are carefully crafted distractors among the wrong options. To ensure the quality and plausibility of these distractors, we designed three distinct distractors annotation strategies. (1) Str. 1: We prompt a strong model GPT-40 [16] to directly answer the manually annotated question. If the model generate an incorrect answer (as verified by human annotators), that answer is retained as a high-quality distractor. If correct, we combine human-written distractors with incorrect options generated by GPT-40 as distractors. (2) Str. 2: Given the question and correct answer annotated manually, GPT-40 is prompted to generate distractors. (3) Human annotators construct distractors manually.

We conducted a test using 100 questions, using three strategies to form three test-set with 100 multiple-choice tasks. As shown in Table 1, distractors generated by strategy 1 are more confusing, significantly increasing the difficulty and quality of our tasks. It is worth noting that in the above test process, when GPT-40 directly answered 100

tasks, the accuracy rate verified by humans was

Models	Str. 1	Str. 2	Str. 3
GPT-4o	59%	70%	62%
Qwen-VL-7B	37%	51%	42%

Table 1: Performance on 100 questions annotated with different strategies (str.).

only 17%. This reflects the limitations of the current model in multimodal reasoning capabilities.

To ensure high quality, we also developed an checklist based on the construction principles and invited human annotators to verify the accuracy and difficulty of the tasks using this checklist. We invited five annotators with at least a bachelor's degree to participate in the annotation and review

process. The checklist of MMR-V is shown in the Appendix B. The overall annotation process and the annotation platform can be found in Figure 7 and Figure 8 in the Appendix B.

3.3 Data Statistics

183

197

207

208

209

210

225

226

227

228

MMR-V comprises a total of 317 videos spanning a wide 184 range of content types, and includes 1,257 multiple-choice 185 reasoning tasks. Each question is annotated with 7 to 11 186 candidate answers, with only one correct answer guaran-187 teed. As illustrated in Figure 9a, the videos are categorized 188 into six major domains, each encompassing fine-grained 189 subcategories to ensure diversity in content, style, and semantics. The reasoning tasks in our benchmark are organized across three levels of granularity, reflecting different 192 dimensions of reasoning complexity and modality. The 193 distribution of task types across these levels is shown in 194 Figure 9b. More information is shown in Table 2. 195

Table 2: Dataset Statistic of MMR-V.

Dataset	Statistic				
Task Question Count Average Option Count Average Question Words Average Option Words	1257 10 14 10				
Video Video Count Minimum Length (s) Maximum Length (s) Average Length (s)	317 7 3771 277				

4 Experiments

4.1 Settings

We conducted extensive evaluations on 9 proprietary and 11 open-source models as detailed in the Appendix E.1. Our main experiments were conducted under two settings: zero-shot and zero-shot + CoT [17], in order to examine whether reasoning enhances performance. For further analysis, we introduced the following categories of comparative models: (1) Models with different scales. (2) "Thinking" model and its base version. (e.g., Gemini-2.0-Flash and Gemini-2.0-Flash-Thinking).

Multimodal Inputs: For models supporting full-modal inputs (e.g., Gemini-2.0-flash), we further compare their performance with and without audio input to evaluate its influence on reasoning results.

Frame Selection: Since some models only support multiple images or short video clips, we standardized the number of input frames. Details of frame sampling are provided in Appendix E.

Human Experiment: To provide a meaningful upper bound for MMR-V and to examine the human-model gap, we invited participants with at least bachelor degree to conduct human experiment. We sampled 100 tasks GPT-40 answered incorrectly and 100 tasks it answered correctly for experiment.

4.2 Main Results

We report the evaluation results in Table 3. Results indicate that the MMR-V Bench poses a significant challenge to current multimodal large models. Even the best-performing model, o4-mini, achieves only 52.5% accuracy. Among open-source models, Gemma-3-27b-it performs the best, demonstrating relatively strong performance. However, there remains a gap compared to proprietary models.

Current reasoning enhancements have limitations on MMR-V. Results in Table 3 show that 216 current reasoning enhancement strategies, which are relatively effective in textual domains, such as CoT prompt reasoning and scaling test-time compute (i.e., "Thinking" models), offer only limited 217 gains on MMR-V. CoT brings only a 0.57% average gain, and "Thinking" model improves just 218 2.4%. This indicates that MMR-V presents a significant challenge to the multimodal reasoning 219 capabilities of existing models. Analysis of sampled model responses shows that visual analysis 220 accounts for only about 10% of the CoTs. This reveals that reasoning process of current model is 221 mostly text-based (reasoning on questions and options), relying on visual perception of question 222 frame, instead of integrating visual reasoning and evidence mining into CoTs. Several examples are 223 provided in Appendix H, and further analysis in Section 4.5 supports similar findings. 224

Model performance on MMR-V Bench exhibits a clear scaling law effect. Smaller models under the same architecture perform poorly on tasks that require complex reasoning. For instance, larger models like Qwen2.5-VL-72B (39.1%) and GPT-40 (44%) outperform their smaller versions Qwen2.5-VL-7B (30.1%) and GPT-40-mini (34.8%), showing relative gains of 9% and 9.2%, respectively.

Model performance across different tasks on MMR-V Bench.

		Tasks				Video Categories						
Model	Ove	erall	Imp	olicit	Exp	olicit	Art	Life	TV	Film	Ani.	Phi.
			Open-	-source	model.	s						
LLaVA-Onevision	6.5	8.8	7.0	9.6	5.4	6.6	6.5	3.4	9.5	3.8	9.8	1.2
LLaVA-Video	18.4	17.6	19.1	18.1	15.4	16.3	14.4	11.2	13.2	17.4	21.4	12.8
NVILA-8B-Video	25.5	25.3	26.2	24.2	23.9	25.9	17.3	21.3	23.5	21.6	38.0	21.8
Phi-4-multimodal-instruct	26.7	27.6	29.4	31.2	19.4	18.1	19.4	19.2	25.9	26.4	33.9	24.4
Cogvlm2-video-llama3	25.6	26.1	25.4	26.2	26.1	25.7	15.5	18.3	24.7	19.1	43.2	20.8
Qwen2.5-VL-7B	30.1	32.4	33.7	36.2	20.8	22.5	20.9	18.1	29.6	21.2	48.4	19.8
Intern3-8B	33.6	32.9	35.5	33.4	28.6	31.4	23.0	22.6	31.7	24.3	52.9	23.2
Gemma-3-12b-it	34.0	34.2	37.8	37.6	24.0	25.4	19.4	24.9	25.9	31.3	51.9	24.4
InternVL2.5-38B	39.9	39.7	43.8	43.7	29.9	29.4	30.4	28.8	30.4	37.2	<u>57.4</u>	29.1
Qwen2.5-VL-72B	39.1	40.4	41.3	42.8	33.4	34.3	28.9	28.2	29.1	36.5	55.6	37.2
Gemma-3-27b-it	<u>42.0</u>	<u>41.1</u>	46.5	<u>44.7</u>	30.3	32.0	31.7	<u>32.2</u>	<u>35.5</u>	<u>41.3</u>	56.1	33.7
			Prop	rietary	models							
GPT-4o-mini-2024-07-18	34.8	35.2	38.0	38.6	26.3	26.3	29.5	25.4	29.6	33.0	48.7	18.6
Gemini-2.0-Flash (16 frames)	42.6	44.3	44.3	45.9	38.3	40.0	30.9	32.2	40.7	40.6	58.5	24.4
Claude-3.5-Sonnet-20241022	43.3	44.2	45.0	46.1	38.9	39.1	33.8	31.1	41.3	41.3	55.8	44.4
GPT-4o-2024-11-20	44.0	46.1	46.6	46.9	37.6	44.0	38.1	37.3	34.9	41.0	61.6	32.6
Gemini-2.0-Flash-thinking	45.0	43.5	46.6	46.0	40.6	37.1	34.5	31.6	38.6	48.3	60.1	25.6
GPT-4.1-2025-04-14	46.6	48.9	49.1	51.7	40.3	41.7	43.2	35.6	43.9	46.5	57.1	34.9
Gemini-2.0-Flash (512 frames)	48.0	49.9	50.5	52.6	41.6	42.9	36.7	36.7	39.7	46.2	66.7	31.4
Gemini-2.5-Flash	51.2	50.5	52.9	52.3	46.9	45.7	45.3	39.5	50.3	47.9	65.6	34.9
o4-mini-2025-04-16	52.5	52.1	54.6	54.5	47.1	46.0	48.2	40.1	54.0	51.7	65.3	27.9
Baseline												
Best Performance of Models	52	2.5	54	4.6	47	7.1	48.2	40.1	54.0	51.7	65.6	44.4
Human	86	5.0	80	0.6	9:	1.2	57.7	92.3	90.6	92.3	90.7	70.0

Table 3: Evaluation results (%) on MMR-V. Results under CoT prompting are highlighted in gray. The random accuracy on MMR-V Bench is approximately 10%. **Bold** and <u>underlined</u> values indicate the best performance among proprietary and open-source models, respectively.

Firstly, the models performed better on implicit tasks than on explicit tasks (with an average gain of +7.9%). Through analysis of tasks and model responses, we found that in implicit tasks, video creators often embed implicit meanings throughout the entire video, resulting in abundant visual cues that can support reasoning. This reduces the requirements for multi-modal reasoning and clue localization. In contrast, explicit tasks demand finer-grained reasoning and the ability to identify specific evidence. For example, in the implicit task at the bottom of Figure 1, many frames provide clues suggesting that the girl symbolizes good luck (e.g., room number, flowers, lighting, weather, etc.). In contrast, the explicit task at the top contains only a few key frames where the hidden lighter in magician's hand can be seen.

230

231

232

233

234 235

236

237

238

239

240

241

243

244

245

246

247

248

249

250

251 252

253

254

255

Secondly, the models performed particularly poorly on Counterintuitive Reasoning (CIR), Sequential Structure Reasoning (SSR), and Comment Matching (CM) tasks. For CIR and SSR tasks, poor

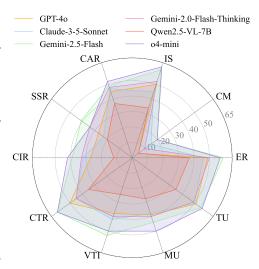


Figure 3: Performance on different tasks.

performance mainly stems from the limited ability of current models to perform multi-frame reasoning. These two tasks require the model to reason on long-range videos, rather than relying on internal knowledge. However, instead of analyzing to locate evidences in other frames, models often rely on surface-level visual perception of the question frame, followed by textual reasoning over question and options. For CM tasks, the results highlight a significant gap between model and human capabilities in implicit reasoning. While humans can infer underlying information such as humor and emotion with minimal cognitive effort [18], current models consistently fail to capture such subtleties.

		Ta	sks	Categories						
	Overall	Imp.	Exp.	Art	Life	TV	Film	Ani.	Phi.	
Gemini-2.0 +audio	42.6 44.0 ^{†1.4}	44.3 46.2 ^{†1.9}	38.3 38.3 ^{-0.0}	30.9 31.0 ^{↑0.1}	32.2 $31.6^{\downarrow 0.6}$	40.7 $42.3^{\uparrow 1.6}$	40.6 41.0 ^{10.4}	58.5 61.1 ^{2.6}	24.4 $29.1^{\uparrow 4.7}$	
Gemini-2.0-thinking +audio	45.0 $46.0^{\uparrow 1.0}$	$46.6 \\ 48.4^{\uparrow 1.8}$	40.6 $39.7^{\downarrow 0.9}$	$34.5 \\ 31.7^{\downarrow 2.8}$	31.6 $33.9^{2.3}$	38.6 $44.4^{5.8}$	$48.3 \\ 42.7^{\downarrow 5.6}$	60.1 $62.4^{\uparrow 2.3}$	25.6 $32.6^{\uparrow 7.0}$	
Phi-4-multimodal-instruct +audio	$26.7 \\ 27.7^{\uparrow 1.0}$	29.4 $31.3^{\uparrow 1.9}$	19.4 $18.1^{\downarrow 1.3}$	19.4 $15.4^{\downarrow 3.0}$	$19.2 \\ 19.7^{\uparrow 0.5}$	25.9 $24.5^{\downarrow 1.4}$	26.4 $27.8^{\uparrow 1.4}$	33.9 $37.3^{\uparrow 3.4}$	24.4 $26.7^{\uparrow 2.3}$	

Table 4: The impact of adding audio modality on the performance (accuracy %) on different tasks.

Human Performance. Humans achieved an average score of 86%, which highlights a significant human-model gap. Although studies suggest that models achieved human-level performance on text tasks [2; 19], models still lag behind on multimodal reasoning tasks. Humans can identify clues in videos easily, while models tend to focus on question frames rather than exploring other evidence frames. Specially, unlike models, humans perform slightly worse on implicit tasks, which is mainly due to the challenges posed by highly abstract implicit understanding in art and philosophy.

4.3 Influence of Frames Count

For Gemini-2.0-Flash, which supports long video inputs, we evaluated performance changes as the number of frames increases. As shown in Figure 4, accuracy improves with more frames, but the rate of improvement gradually slows. After sampling and observing the CoTs, it is found that the initial gains come from the addition of evidence frames, while the slowdown is mainly due to limited multi-frame reasoning ability of the model. Performance on implicit tasks continues to improve in later stages, as visual cues for such tasks are often dispersed throughout the video (as discussed in Section 4.2); more frames tend to provide more clues. In contrast, explicit clues are fewer and more localized.

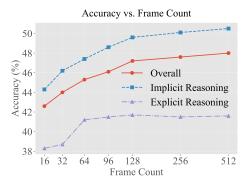


Figure 4: Accuracy with the increase of input frame counts.

4.4 Influence of Audio Input

For models that support full-modal input, we compared their performance before and after incorporating the audio modality. As shown in Table 4, overall performance improved with the addition of audio. Specifically, Gemini 2.0-Flash, Gemini 2.0-Flash-Thinking, and Phi-4-multimodal-instruct showed improvements of 1.4%, 1.0%, and 1.0%, respectively. This suggests that advancing research on fully multimodal models is a promising direction.

4.5 Error Analysis

We sampled 100 incorrect responses from GPT-40 for error analysis. The main sources of errors can be categorized as follows: (1) Lack of Visual Reasoning: the model often failed to locate the correct evidence frames and lack of long-range, multi-frame visual reasoning. (2) Implicit Misinterpretation: revealing a significant understanding gap between the model and human cognition. (3) Knowledge Insufficiency: the model lacks some intrinsic knowledge (4) Reasoning Error: during the multi-step deduction process. (5) Hallucination: the model introduced fake or unsupported information. (6) Output Formatting Issue: model refusals or formatting errors prevent answer extraction. Among error cases, Lack of Visual Reasoning accounts for the largest proportion. This indicates that

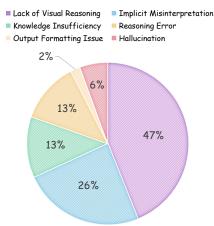


Figure 5: Error analysis of GPT-40.

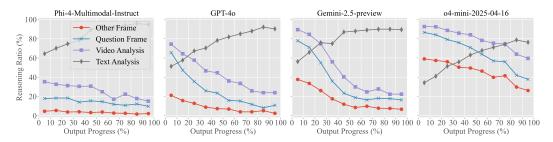


Figure 6: CoT content across different stages. The y-axis indicates the ratio of the 500 sampled CoTs that include analysis of these four types of content at each stage.

current models still lack genuine multimodal reasoning capabilities. They tend to rely on text-based reasoning after briefly perceiving frames adjacent to the question, rather than engaging in deep, long-range, multi-frame video reasoning. Most existing reasoning models remain inadequate in integrating multimodal information into the reasoning process and performing thorough analysis. In contrast, o4-mini exhibits a better reasoning paradigm, as shown in Figure 11 for comparison.

We further analyzed model CoTs by categorizing each step into video or text analysis (e.g. options), with video analysis divided into question frame and other frame analysis (details in Appendix F). We sampled 500 CoTs from models, split each into 10 equal-length segments, and used GPT-4.1 to label each segment. As shown in Figure 6, where models further to the right perform better on MMR-V, models with better performance on MMR-V show more video analysis, especially on **other frames** (red line). Notably, 40-mini stands out with strong analysis of non-question frames, highlighting the value of enhanced visual reasoning and tool use in multi-frame video reasoning tasks.

5 Related Work

Video Understanding Benchmark. Existing video benchmarks primarily focus on evaluating models' perception and intuitive understanding of visual elements in videos, such as action recognition [20, 21; 22; 23] and video description [24; 25]. Recent notable works, such as Video-MME [8], MVBench [26] and MMBench-Video [27], have extended video understanding to multiple task types and video types, enabling a more comprehensive assessment of video understanding capabilities. Additionally, benchmarks like LVBench [28] and LongVideoBench [29] have introduced long-video questionanswering tasks. However, these tasks mainly evaluate whether a model can accurately extract relevant information from long videos based on the given questions, while the subsequent steps remain largely perception-oriented. MMR-V is designed to assess whether a model can perform multi-frame, long-span, multimodal autonomous reasoning on videos based on the given questions. **Multimodal Reasoning.** Recent advancements have greatly enhanced LLM reasoning [2, 1, 30, 31]. Many top LLMs perform well on complex reasoning tasks, but their evaluation focuses on textbased reasoning [32; 33; 34; 35; 36; 37]. MLLMs still lack thorough assessment in this area. Current multimodal reasoning benchmarks mainly involve mathematical or coding tasks in image form [14; 38; 39], which primarily test visual recognition followed by text reasoning. True multimodal reasoning requires integrating details like depth, texture, and audio for complex inference. MMR-V

6 Conclusion

This paper introduces MMR-V: A Benchmark for Multimodal Deep Reasoning in Videos. All tasks are annotated by human experts and designed to evaluate abilities of multimodal reasoning. MMR-V presents a significant challenge to current models, with the best model performance still lagging 33.5% accuracy behind human. This highlights a human-model gap in interpreting and reasoning about video information. Notably, o4-mini achieves the best results on MMR-V, suggesting that integrating visual reasoning into CoT and leveraging tool use is a promising direction for tackling video reasoning tasks. We hope MMR-V will serve as a reliable evaluation benchmark for the development of MLLMs and offer valuable insights into advancing multimodal reasoning research.

Bench aims to evaluate multimodal sequential reasoning in video tasks.

References

- [1] A. Jaech, A. Kalai, A. Lerer, A. Richardson, A. El-Kishky, A. Low, A. Helyar, A. Madry, A. Beutel, A. Carney, *et al.*, "Openai o1 system card," *arXiv preprint arXiv:2412.16720*, 2024.
- [2] D. Guo, D. Yang, H. Zhang, J. Song, R. Zhang, R. Xu, Q. Zhu, S. Ma, P. Wang, X. Bi, *et al.*, "Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning," *arXiv* preprint arXiv:2501.12948, 2025.
- [3] Y. Wang, S. Wu, Y. Zhang, S. Yan, Z. Liu, J. Luo, and H. Fei, "Multimodal chain-of-thought reasoning: A comprehensive survey," *arXiv preprint arXiv:2503.12605*, 2025.
- 345 [4] OpenAI, "Openai: Introducing openai o3 and o4-mini," 2025.
- J. Hou, C. Wu, Z. Yuan, J. Tan, Q. Wang, and Y. Zhou, "Research of intelligent home security surveillance system based on zigbee," in 2008 International Symposium on Intelligent Information Technology Application Workshops, pp. 554–557, IEEE, 2008.
- J. Yang, S. Yang, A. W. Gupta, R. Han, L. Fei-Fei, and S. Xie, "Thinking in space: How multimodal large language models see, remember, and recall spaces," *arXiv preprint arXiv:2412.14171*, 2024.
- J. Zhou, Y. Shu, B. Zhao, B. Wu, S. Xiao, X. Yang, Y. Xiong, B. Zhang, T. Huang, and Z. Liu,
 "Mlvu: A comprehensive benchmark for multi-task long video understanding," arXiv preprint
 arXiv:2406.04264, 2024.
- [8] C. Fu, Y. Dai, Y. Luo, L. Li, S. Ren, R. Zhang, Z. Wang, C. Zhou, Y. Shen, M. Zhang, *et al.*, "Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis," *arXiv preprint arXiv:2405.21075*, 2024.
- [9] J. S. B. Evans, "Heuristic and analytic processes in reasoning," *British Journal of Psychology*,
 vol. 75, no. 4, pp. 451–468, 1984.
- [10] R. Sun, "The clarion cognitive architecture: Extending cognitive modeling to social simulation,"
 Cognition and multi-agent interaction, pp. 79–99, 2006.
- ³⁶² [11] M. Polanyi, *Personal knowledge*. Routledge, 2012.
- 163 [12] D. Kahneman, Thinking, fast and slow. macmillan, 2011.
- [13] P. Lu, H. Bansal, T. Xia, J. Liu, C. Li, H. Hajishirzi, H. Cheng, K.-W. Chang, M. Galley, and
 J. Gao, "Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts,"
 arXiv preprint arXiv:2310.02255, 2023.
- [14] K. Wang, J. Pan, W. Shi, Z. Lu, H. Ren, A. Zhou, M. Zhan, and H. Li, "Measuring multimodal mathematical reasoning with math-vision dataset," *Advances in Neural Information Processing Systems*, vol. 37, pp. 95095–95169, 2024.
- 570 [15] F. Zhang, L. Wu, H. Bai, G. Lin, X. Li, X. Yu, Y. Wang, B. Chen, and J. Keung, "Humaneval-v: 571 Evaluating visual understanding and reasoning abilities of large multimodal models through 572 coding tasks," *arXiv preprint arXiv:2410.12381*, 2024.
- [16] A. Hurst, A. Lerer, A. P. Goucher, A. Perelman, A. Ramesh, A. Clark, A. Ostrow, A. Welihinda, A. Hayes, A. Radford, *et al.*, "Gpt-4o system card," *arXiv preprint arXiv:2410.21276*, 2024.
- J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, *et al.*, "Chain-of-thought prompting elicits reasoning in large language models," *Advances in neural information processing systems*, vol. 35, pp. 24824–24837, 2022.
- [18] K. Krishna, Y. Chang, J. Wieting, and M. Iyyer, "RankGen: Improving text generation with large ranking models," in *Proceedings of EMNLP*, pp. 199–232, 2022.
- 80 [19] OpenAI, "Gpt-4 technical report," arXiv preprint arxiv:2303.08774, 2023.

- [20] M. U. Khattak, M. F. Naeem, J. Hassan, M. Naseer, F. Tombari, F. S. Khan, and S. Khan, "How good is my video lmm? complex video reasoning and robustness evaluation suite for video-lmms," *arXiv preprint arXiv:2405.03690*, 2024.
- [21] K. Mangalam, R. Akshulakov, and J. Malik, "Egoschema: A diagnostic benchmark for very long-form video language understanding," *Advances in Neural Information Processing Systems*, vol. 36, pp. 46212–46244, 2023.
- V. Patraucean, L. Smaira, A. Gupta, A. Recasens, L. Markeeva, D. Banarse, S. Koppula,
 M. Malinowski, Y. Yang, C. Doersch, et al., "Perception test: A diagnostic benchmark for
 multimodal video models," Advances in Neural Information Processing Systems, vol. 36,
 pp. 42748–42761, 2023.
- [23] J. Xiao, X. Shang, A. Yao, and T.-S. 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, pp. 9777–9786, 2021.
- D. Xu, Z. Zhao, J. Xiao, F. Wu, H. Zhang, X. He, and Y. Zhuang, "Video question answering via gradually refined attention over appearance and motion," in *Proceedings of the 25th ACM international conference on Multimedia*, pp. 1645–1653, 2017.
- J. Xu, T. Mei, T. Yao, and Y. 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*, pp. 5288–5296, 2016.
- [26] K. Li, Y. Wang, Y. He, Y. Li, Y. Wang, Y. Liu, Z. Wang, J. Xu, G. Chen, P. Luo, et al.,
 "Mvbench: A comprehensive multi-modal video understanding benchmark," in Proceedings
 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 22195–22206,
 2024.
- 404 [27] X. Fang, K. Mao, H. Duan, X. Zhao, Y. Li, D. Lin, and K. Chen, "Mmbench-video: A long-form multi-shot benchmark for holistic video understanding," *Advances in Neural Information Processing Systems*, vol. 37, pp. 89098–89124, 2024.
- W. Wang, Z. He, W. Hong, Y. Cheng, X. Zhang, J. Qi, X. Gu, S. Huang, B. Xu, Y. Dong, *et al.*, "Lvbench: An extreme long video understanding benchmark," *arXiv preprint arXiv:2406.08035*, 2024.
- [29] H. Wu, D. Li, B. Chen, and J. Li, "Longvideobench: A benchmark for long-context interleaved video-language understanding," *Advances in Neural Information Processing Systems*, vol. 37, pp. 28828–28857, 2024.
- [30] K. Team, A. Du, B. Gao, B. Xing, C. Jiang, C. Chen, C. Li, C. Xiao, C. Du, C. Liao, *et al.*, "Kimi k1. 5: Scaling reinforcement learning with llms," *arXiv preprint arXiv:2501.12599*, 2025.
- 415 [31] Y. Zhao, H. Yin, B. Zeng, H. Wang, T. Shi, C. Lyu, L. Wang, W. Luo, and K. Zhang, "Marco-o1: Towards open reasoning models for open-ended solutions," *arXiv preprint arXiv:2411.14405*, 2024.
- 418 [32] D. Hendrycks, C. Burns, S. Kadavath, A. Arora, S. Basart, E. Tang, D. Song, and J. Steinhardt, "Measuring mathematical problem solving with the math dataset," *arXiv preprint arXiv:2103.03874*, 2021.
- 421 [33] Y. Bai, S. Tu, J. Zhang, H. Peng, X. Wang, X. Lv, S. Cao, J. Xu, L. Hou, Y. Dong, *et al.*, "Longbench v2: Towards deeper understanding and reasoning on realistic long-context multitasks,"
 423 *arXiv* preprint arXiv:2412.15204, 2024.
- 424 [34] K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, 425 J. Hilton, R. Nakano, C. Hesse, and J. Schulman, "Training verifiers to solve math word 426 problems," *ArXiv preprint*, vol. abs/2110.14168, 2021.
- 427 [35] D. Rein, B. L. Hou, A. C. Stickland, J. Petty, R. Y. Pang, J. Dirani, J. Michael, and S. R. Bowman, "Gpqa: A graduate-level google-proof q&a benchmark," in *First Conference on Language Modeling*, 2024.

- [36] Y. Wang, X. Ma, G. Zhang, Y. Ni, A. Chandra, S. Guo, W. Ren, A. Arulraj, X. He, Z. Jiang,
 et al., "Mmlu-pro: A more robust and challenging multi-task language understanding benchmark," in *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024.
- 434 [37] C. E. Jimenez, J. Yang, A. Wettig, S. Yao, K. Pei, O. Press, and K. Narasimhan, "Swe-bench: 435 Can language models resolve real-world github issues?," *arXiv preprint arXiv:2310.06770*, 436 2023.
- 437 [38] W. Shi, Z. Hu, Y. Bin, J. Liu, Y. Yang, S.-K. Ng, L. Bing, and R. K.-W. Lee, "Math-llava: Bootstrapping mathematical reasoning for multimodal large language models," *arXiv preprint* 439 *arXiv:2406.17294*, 2024.
- [39] K. Ying, F. Meng, J. Wang, Z. Li, H. Lin, Y. Yang, H. Zhang, W. Zhang, Y. Lin, S. Liu, *et al.*,
 "Mmt-bench: A comprehensive multimodal benchmark for evaluating large vision-language models towards multitask agi," *arXiv preprint arXiv:2404.16006*, 2024.
- [40] G. Lakoff and M. Johnson, *Metaphors we live by*. University of Chicago press, 2008.
- 444 [41] OpenAI, "Openai: Hello gpt-40," 2024.
- 445 [42] OpenAI, "Gpt-40 mini: advancing cost-efficient intelligence," 2024.
- 446 [43] OpenAI, "Introducing gpt-4.1 in the api." https://openai.com/index/gpt-4-1/, 2025.
- 447 [44] M. Reid, N. Savinov, D. Teplyashin, D. Lepikhin, T. Lillicrap, J.-b. Alayrac, R. Soricut,
 448 A. Lazaridou, O. Firat, J. Schrittwieser, *et al.*, "Gemini 1.5: Unlocking multimodal understand449 ing across millions of tokens of context," *arXiv preprint arXiv:2403.05530*, 2024.
- 450 [45] Google DeepMind, "Gemini 2.5: Our most intelligent ai model," March 2025.
- 451 [46] Anthropic, "Anthropic: Introducing claude 3.5 sonnet," 2024.
- [47] A. Yang, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Li, D. Liu, F. Huang, H. Wei, H. Lin,
 J. Yang, J. Tu, J. Zhang, J. Yang, J. Zhou, J. Lin, K. Dang, K. Lu, K. Bao, K. Yang,
 L. Yu, M. Li, M. Xue, P. Zhang, Q. Zhu, R. Men, R. Lin, T. Li, T. Xia, X. Ren, X. Ren, Y. Fan,
 Y. Su, Y. Zhang, Y. Wan, Y. Liu, Z. Cui, Z. Zhang, and Z. Qiu, "Qwen2.5 technical report,"
 CORR, vol. abs/2412.15115, 2024.
- [48] A. Kamath, J. Ferret, S. Pathak, N. Vieillard, R. Merhej, S. Perrin, T. Matejovicova, A. Ramé, 457 M. Rivière, L. Rouillard, T. Mesnard, G. Cideron, J. Grill, S. Ramos, E. Yvinec, M. Casbon, 458 E. Pot, I. Penchev, G. Liu, F. Visin, K. Kenealy, L. Beyer, X. Zhai, A. Tsitsulin, R. Busa-Fekete, 459 A. Feng, N. Sachdeva, B. Coleman, Y. Gao, B. Mustafa, I. Barr, E. Parisotto, D. Tian, M. Eyal, 460 C. Cherry, J. Peter, D. Sinopalnikov, S. Bhupatiraju, R. Agarwal, M. Kazemi, D. Malkin, 461 R. Kumar, D. Vilar, I. Brusilovsky, J. Luo, A. Steiner, A. Friesen, A. Sharma, A. Sharma, A. M. 462 Gilady, A. Goedeckemeyer, A. Saade, A. Kolesnikov, A. Bendebury, A. Abdagic, A. Vadi, 463 A. György, A. S. Pinto, A. Das, A. Bapna, A. Miech, A. Yang, A. Paterson, A. Shenoy, 464 A. Chakrabarti, B. Piot, B. Wu, B. Shahriari, B. Petrini, C. Chen, C. L. Lan, C. A. Choquette-465 Choo, C. Carey, C. Brick, D. Deutsch, D. Eisenbud, D. Cattle, D. Cheng, D. Paparas, D. S. 466 Sreepathihalli, D. Reid, D. Tran, D. Zelle, E. Noland, E. Huizenga, E. Kharitonov, F. Liu, 467 G. Amirkhanyan, G. Cameron, H. Hashemi, H. Klimczak-Plucinska, H. Singh, H. Mehta, H. T. 468 Lehri, H. Hazimeh, I. Ballantyne, I. Szpektor, and I. Nardini, "Gemma 3 technical report," 469 CoRR, vol. abs/2503.19786, 2025. 470
- 471 [49] J. Zhu, W. Wang, Z. Chen, Z. Liu, S. Ye, L. Gu, Y. Duan, H. Tian, W. Su, J. Shao, *et al.*,
 472 "Internvl3: Exploring advanced training and test-time recipes for open-source multimodal
 473 models," *arXiv preprint arXiv:2504.10479*, 2025.
- [50] B. Li, Y. Zhang, D. Guo, R. Zhang, F. Li, H. Zhang, K. Zhang, P. Zhang, Y. Li, Z. Liu, et al., "Llava-onevision: Easy visual task transfer," arXiv preprint arXiv:2408.03326, 2024.
- [51] B. Lin, Y. Ye, B. Zhu, J. Cui, M. Ning, P. Jin, and L. Yuan, "Video-llava: Learning united visual representation by alignment before projection," *arXiv preprint arXiv:2311.10122*, 2023.

- [52] A. Abouelenin, A. Ashfaq, A. Atkinson, H. Awadalla, N. Bach, J. Bao, A. Benhaim, M. Cai, 478 V. Chaudhary, C. Chen, D. Chen, D. Chen, J. Chen, W. Chen, Y. Chen, Y. Chen, Q. Dai, 479 X. Dai, R. Fan, M. Gao, M. Gao, A. Garg, A. Goswami, J. Hao, A. Hendy, Y. Hu, X. Jin, 480 M. Khademi, D. Kim, Y. J. Kim, G. Lee, J. Li, Y. Li, C. Liang, X. Lin, Z. Lin, M. Liu, Y. Liu, 481 G. Lopez, C. Luo, P. Madan, V. Mazalov, A. Mitra, A. Mousavi, A. Nguyen, J. Pan, D. Perez-482 Becker, J. Platin, T. Portet, K. Qiu, B. Ren, L. Ren, S. Roy, N. Shang, Y. Shen, S. Singhal, 483 S. Som, X. Song, T. Sych, P. Vaddamanu, S. Wang, Y. Wang, Z. Wang, H. Wu, H. Xu, W. Xu, 484 Y. Yang, Z. Yang, D. Yu, I. Zabir, J. Zhang, L. L. Zhang, Y. Zhang, and X. Zhou, "Phi-4-mini 485 technical report: Compact yet powerful multimodal language models via mixture-of-loras," 486 CoRR, vol. abs/2503.01743, 2025. 487
- 488 [53] W. Hong, W. Wang, M. Ding, W. Yu, Q. Lv, Y. Wang, Y. Cheng, S. Huang, J. Ji, Z. Xue,
 489 et al., "Cogvlm2: Visual language models for image and video understanding," arXiv preprint
 490 arXiv:2408.16500, 2024.
- 491 [54] Z. Liu, L. Zhu, B. Shi, Z. Zhang, Y. Lou, S. Yang, H. Xi, S. Cao, Y. Gu, D. Li, et al., "Nvila: Efficient frontier visual language models," arXiv preprint arXiv:2412.04468, 2024.

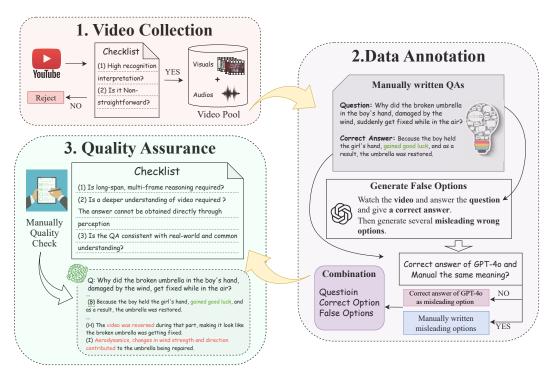


Figure 7: MMR-V Construction Pipeline.

493 A Limitations

494

495

496

497

498

500

505

506

507

508

509

Despite our efforts to improve our work, several limitations remain. (1) Scaling MMR-V is challenging due to the high cost of manual annotation and verification, as all tasks and correct answers are curated and reviewed by human annotators. (2) Although we strive to cover a wide range of video and task types, certain real-world categories (such as mystery, puzzle-solving, and gaming) are still underrepresented. (3) The majority of videos in MMR-V are in English, with only a small proportion in other languages such as Chinese, French, Thai, and German, which constrains its multilingual applicability. We will further study and try to solve this issue in the future.

501 B MMR-V Construction

502 B.1 Checklist

According to the MMR-V construction principles introduced in the main text Section 3, we wrote the following annotation checklist:

- (1) You are expected to watch the entire video before formulating any questions or answers.
- (2) Each question must require **long-distance**, **multi-frame reasoning** and cannot be answered through direct perception (ensuring compliance with Principles 1 and 2).
- (3) To ensure **consistency with real-world user** perception (Principle 3), annotators are encouraged to refer to the official interpretation of the original video author and user consensus (highly praised comments in the comment section) when writing or verifying the correct answer. This helps mitigate annotator bias and ensures that the reasoning task reflects the understanding of a wider audience.

B.2 Construction Pipeline

In this section, we present the construction process of MMR-V Bench in a macro sense. The whole process is divided into three stages: **video collection**, **data annotation**, and **quality assurance**. For **video collection**, we designed a checklist to ensure the quality and diversity of videos in the Bench. "High recognition interpretation?" ensures that the questions raised and the annotated answers based

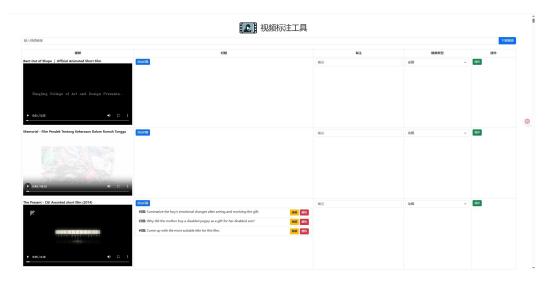


Figure 8: Annotation Platform of MMR-V.

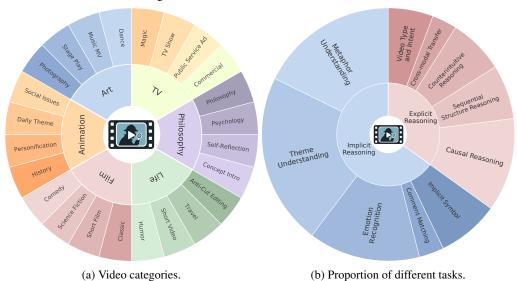


Figure 9: (a) Video categories in MMR-V Bench. (b) Proportion of different tasks in MMR-V Bench.

on the video have references that are consistent with public cognition (official interpretations or highly praised comments) to alleviate the subjective bias of the annotator. "Is it Non-straightforward?" ensures that the video is not a straightforward narrative, which is conducive to increasing the reasoning difficulty of the question. For **data annotation**, as described in section 3.2 of the main text, we use gpt-40 to assist in annotation with interference options. Let the model generate the correct answer based on the question, and manually review to ensure that the correct answer generated by the model is different from the manual annotation. If they are different, the answer generated by gpt-40 is used as the interference item, otherwise the interference item is manually written. For **quality assurance**, we designed a checklist for human reviewers to check the correctness and difficulty of the tasks. The annotation platform is shown in Figure 8.

C Diversity of MMR-V

In this section, we show the diversity of MMR-V Bench, including video diversity and task diversity.
For video, we show the six categories of videos in MMR-V in Figure 9a, including Life, Animation,
Film, Art, TV, and Philosophy. At the same time, for each category, we divide it into several

Ability Type L1	Ability Type L2	Ability Type L3				
	Matanhan Undanstanding (MU)	Structural Metaphor, Orientational Metaphor,				
	Metaphor Understanding (MU)	Ontological Metaphor, Creative Metaphor				
Implicit Reasoning	Theme Understanding (TU)	Philosophical Concepts, Social Issues, Personal Reflectio				
implicit Keasoning	Theme Orderstanding (10)	Everyday Topics, Video Naming				
	Emotion Recognition (ER)	Explicit Emotion, Implicit Emotion,				
	Emotion Recognition (ER)	Meta-emotion, Audience Emotion				
	Comment Matching (CM)	Humorous, Thought-provoking, Trending				
	Implicit Symbol (IS)	Cultural Symbols, Art Symbols, Other Symbols				
	Causal Reasoning (CAR)	Forward Reasoning, Backward Reasoning				
	Sequential Structure Reasoning (SSR)	Narrative Structure, Core Connecting Elements,				
Explicit Reasoning	Sequential Structure Reasoning (SSR)	Inference on Editing Techniques, Hallucination				
	Counterintuitive Reasoning (CIR)	Magic Deconstruction or Special Effects Editing,				
	Counterintuitive Reasoning (CIR)	Artistic Techniques, Humor and Exaggeration				
	Cross-modal Transfer Reasoning (CTR)	Video-to-Text, Video-to-Audio, Video-to-Video				
	Video Type and Intent (VTI)	Video Type, Video Intent				

Table 5: Three-level classification of tasks in MMR-V.

subcategories to better understand the classification of video categories. Secondly, in section 3.3 we show the diversity of video length, ranging from 7 seconds to 3771 seconds. For tasks, we divide them into two parts, ten categories and 33 subcategories, three levels. The division of the first and second levels, as well as the proportion of different types of tasks, can be seen in Figure 9b.

D Task Details

529

533

- The tasks in MMR-V can be divided into three levels. Level 1: Implicit Reasoning & Explicit Reasoning. Level 2: Contains ten task classes. Level 3: Contains 33 task subclasses. Next, we will introduce these tasks with some task examples.
 - D.1 Implicit Reasoning Tasks
- I. Metaphor Understanding (MU)
- For the definition of subclasses of the metaphor understanding task, we mainly refer to the book Metaphors We Live [40] By by George Lakoff and Mark Johnson, which introduces metaphor-related concepts in detail.
 - I.1. Structural Metaphor
- Task Description: There are structural similarities between the subject and object. For example, time can be compared to flowing water, both of which have the structure of flow and passing away.
- **Example Question:**

Question: What does the brown coat in the video symbolize? **Options**:

- (A) It is said to represent the family's long lost fortune that they are still searching for.
- (B) The brown coat symbolizes the lost hope of the family because it was worn during a difficult time.
- (C) It refers to a coat that has been washed and taken out to dry, likely worn by the father.",
- (D) It symbolizes the father in a family, who protects his family in times of difficulty.
- (E) It represents the fear of the outside world.
- (F) The unfulfilled dreams of the children in the family as they always saw it as a sign of something unattainable.
- (G) The brown coat in the video represents a raincoat, used to protect the clothes inside from getting wet.

(H) The bad luck that has been following the family for generations.

CorrectAnswer: (D)

Video: father - 1 minute emotional award winning - video_url

543

544

I.2. Orientational Metaphor

Task Description: There are similarities in direction or composition between the subject and the metaphor, for example, walking up a staircase is compared to ambition.

7 Example Ouestion:

Question: Why does the dance, which is filled with artistry and beauty throughout, end with a descent?

Options:

- (A) There is a connection between the fall and the creation at some point.
- (B) It represents a dive to explore new depths, both literal and metaphorical.
- (C) It indicates the dancer's exhaustion, capturing a moment of fatigue.
- (D) It reflects the calmness of the ocean, evoking a sense of tranquility.
- (E) It highlights the theme of rebirth, symbolizing renewal and transformation.
- (F) It represents the beauty of underwater life, showcasing its unique allure.
- (G) It symbolizes being weighed down by emotions, expressing inner turmoil.
- (H) It symbolizes the end of a dream, marking a moment of conclusion.
- (I) It shows the dancer's connection to water, emphasizing fluidity and grace.
- (J) It symbolizes a return to nature and surrender to life's forces, embracing the natural flow.
- (K) It signifies the end of the dance's energy, indicating a point of culmination.

CorrectAnswer: (A)

Video: Falling - Underwater dance - video_url

548

549

I.3. Ontological Metaphor

- Task Description: This metaphor involves viewing an abstract concept as a concrete entity. Usually,
- the core concept of the entire video is turned into a concrete entity to tell the story.
- **Example Ouestion:**

Question: The scene around 1:00 metaphorically represents what aspect of communities? **Options**:

- (A) Communities can build their resilience to setbacks by working together and adapting to new challenges.
- (B) Promoting individual success in competitive environments.
- (C) Building resilience through community partnerships.
- (D) Overcoming challenges for community progress.
- (E) Celebrating the individual achievements of community members.
- (F) Developing sustainable practices for environmental harmony.
- (G) Decision-making processes of a community.
- (H) The interconnectedness of global communities.
- (I) Isolation of communities for self-sufficiency.
- (J) The role of external aid in community development.
- (K) Highlighting the diversity of cultures within a community.

CorrectAnswer: (A)

Video: Resilience: Anticipate, organise, adapt - video_url

553 554

I.4. Creative Metaphor

Task Description: This metaphor is usually carefully designed by the author for a specific video and needs to be understood in the context of the video.

Example Question:

Question: What is the pink fairy ball in the film?

Options:

- (A) It's a toy the boy picked up on the street, having no special connection to his condition.
- (B) They are the microorganisms in this world, living in every corner.
- (C) The pink fairy ball represents the boy's childhood dream of becoming a fairy.
- (D) It's a hallucination caused by lack of sleep, not related to antidepressants at all.
- (E) They are the boy's toys, which he bought to help treat his depression.
- (F) It is the effect of the antidepressants the boy is taking, which helps him see many things with vitality and positive effects.
- (G) It's an advertisement prop for a new product in the background of the scene.
- (H) The ball is a sign of the boy's wish to escape from his daily work routine.
- (I) The pink fairy ball is a symbol of the city's upcoming festival decorations

CorrectAnswer: (F)

Video: Soft Rain | Animated Short Film (2023) - video url

559

560 II. Theme Understanding (TU)

561 II.1. Philosophical Concepts

Task Description: The themes of the videos are usually about concepts and principles related to philosophy and psychology.

Example Question:

Question: What is the overall message that the animation aims to convey? **Options**:

- (A) It suggests happiness comes solely from financial achievements.
- (B) The animation emphasizes the need to avoid all responsibilities.
- (C) The animation aims to illustrate the ways to relieve stress.
- (D) It illustrates the mechanical process of water flow.
- (E) The animation encourages saving water to prevent wastage.
- (F) The animation conveys the importance of managing stress through self-care practices.
- (G) The animation highlights achieving success through hard work.
- (H) The animation suggests that ignoring stress leads to happiness.
- (I) The video underlines the significance of collective teamwork.
- (J) It depicts progress and growth through constant work.

CorrectAnswer: (C)

Video: The Stress Bucket - video_url

565

566

II.2. Social Issues

Task Description: The theme of the video is usually to reflect some problems existing in today's society and express a strong appeal of the author.

Example Question:

Question: What social reality does this video satirize? **Options**:

- (A) The rise of environmental awareness in urban settings.
- (B) The video represents the bystander effect in society.
- $(C)\ The\ economic\ disparities\ in\ urban\ vs.\ rural\ areas.$
- (D) The challenges of modern relationship dynamics.
- (E) The impact of fashion trends on daily life.
- (F) The increasing complexity of urban development planning.
- (G) The need for infrastructure improvement and road safety.
- (H) The influence of social media on public behavior.
- (I) The rapid pace of technological advancement in transportation.
- (J) The shift in societal values towards individualism.

CorrectAnswer: (B)

Video: Stone | 1 Minute Short Film | Hot Shot - video_url

II.3. Personal Reflection

- Task Description: The author hopes that the video will inspire people to reflect on and resonate with
- things in their lives.
- **Example Question:**

Question: What is the core concept that the film aims to convey? **Options**:

- (A) Romantic relationships in adolescence.
- (B) The importance of education institutions.
- (C) Overcoming supernatural challenges.
- (D) The dynamics of family disagreements.
- (E) Exploration of technological advancement.
- (F) Not to judge others too quickly.
- (G) Journey of a superhero in saving the city.
- (H) Inter-species relations on Earth.
- (I) Power struggles in political leadership.
- (J) Historical recount of a famous personality.

CorrectAnswer: (F)

Video: Award Winning SHORT FILMS Don't Judge | BATTI Hindi Heart Touching Short Movies | Content Ka Keeda - video url

575

576 II.4. Everyday Topics

- Task Description: The themes expressed in the videos are usually the sublimation of the insights and
- themes in daily life, such as praising maternal love, friendship, etc.
- **Example Question**:

Question: What is implied by the contrast between the scenes around 0:47 and 1:11? **Options**:

- (A) The contrast shows that the mother is indecisive and can't make up her mind in a crisis.
- (B) It demonstrates the father's sense of responsibility and bravery, praising paternal love.
- (C) The contrast between the beginning and the end conveys a sense of tragedy, criticizing the destruction of the ecological environment by humans.
- (D) It shows that the father wants to abandon the child when facing danger.
- (E) It shows the bravery of the bird in the background, facing authority head-on, and praises courage.
- (F) The mother still protects her child at all costs even in the face of danger, which praises maternal love.
- (G) It implies that the father is doing it for self preservation rather than out of love for the child.
- (H) It shows that even when there are many birds, they do not appear very united, and in the face of danger, they become a disorganized mess.

CorrectAnswer: (F)

Video: Mother 1 minute Sad Emotional Award Winning Iranian Short Film Animation Animated - video url

580

581

II.5. Video Naming

- Task Description: Come up with a suitable title for this video or the core content of the video (dance,
- etc.). This tests the model's control over the overall content and whether it can get the subtleties of
- the title like humans.
- **Example Question:**

Question: "Please come up with a suitable name for this dance.", **Options**:

- (A) The Dance of the Butterfly.",
- (B) The Rhythm of the Phoenix.",
- (C) The Grace of the Swan.",

```
(D) The Spirit of the Dragon.",
(E) The Charm of the Peony.",
(F) The Step of the Tiger.",
(G) The Soul of Peacock",
(H) The Beat of the Forest.",
(I) The Leap of the Deer.",
(J) The Spin of the Star.",
(K) The Waltz of the Moon."

CorrectAnswer: (G)

Video: Yang Liping - The Soul of Peacock - Peacock Dance - Traditional Dance - video_url
```

588 III. Emotion Recognition (ER)

589 III.1. Explicit Emotion

- Task Description: Analyze the emotions of the characters in the video. Explicit emotions can usually
- be directly understood through facial expressions, body movements, etc.
- **Example Question:**

```
Question: Summarize the boy's emotional changes between 6:00 and 7:00.
options:

(A) Anger - Fear - Surprise and happiness
(B) Sadness - Excitement - Helplessness
(C) Disappointment - Let - down - Sorrow
(D) Loneliness - Isolation - Solitude
```

- (E) Sadness Grief Mourning
- (F) Sadness Shock Surprise and happiness
- (G) Disappointment Astonishment Stupefaction
- (H) Loneliness Isolation Sorrow
- (I) Disappointment Excitement Helplessness

correctAnswer: (F)

Video: CGI Animated Short Film: "Crunch" by Gof Animation | CGMeetup - video_url

593

587

594 III.2. Implicit Emotion

- Task Description: Analyze the emotions of characters in the video. Implicit emotions usually need
- to be analyzed indirectly through the environment, style, etc.
- Example Question:

Question: What kind of emotional atmosphere does the stage lighting create? **options**:

- (A) Solemn and sorrowful atmosphere.
- (B) Neutral and unemotional atmosphere.
- (C) Intense and dramatic atmosphere.
- (D) Joyful and festive atmosphere.
- (E) Sadness and loss.
- (F) Confident and empowering atmosphere.
- (G) Chaotic and confusing atmosphere.
- (H) Calm and serene atmosphere.
- (I) Playful and whimsical atmosphere.
- (J) Romantic and loving atmosphere.

CorrectAnswer: (E)

Video: Stages of Grief- AVANTGARDE SHOW 2023 - video_url

III.3. Meta-emotion

Task Description: This part refers to the high-level emotions in the video, such as the emotions expressed by the author through the video, and the emotions expressed by the entire video.

Example Question:

Question: Summarize the meaning of this short film in one word.

Options: [

- (A) Creation
- (B) Transformation
- (C) Stress
- (D) Mutation
- (E) Metamorphosis
- (F) Growth
- (G) Rebirth"
- (H) Destruction
- (I) Erosion
- (J) Development
- (K) Isolation
- (L) Conversion

CorrectAnswer: (C)

Stress - Shortfilm - video url

603

604 III.4. Audience Emotion

- Task Description: Analyze the emotions that viewers are most likely to feel after watching the video.
- This is more advanced and relatively easy for humans to sense. Including the perception of humor.
- **Example Question:**

Question: What are the reasons for the high number of views on this video? **Options**:

- (A) The video features a well-known celebrity who has a large fan base, drawing a lot of attention.
- (B) The dance style is extremely unique and has never been seen before, sparking curiosity.
- (C) People are under a lot of stress and need videos that can help them unwind.
- (D) The background music is a popular hit song that many people recognize and enjoy.
- (E) The video was released during a major holiday season when people are more likely to watch videos.
- (F) The choreography is incredibly complex and impressive, showcasing the dancers' skills.
- (G) The video has a strong and inspiring message that resonates with a wide audience.
- (H) The video was featured on a popular TV show or news segment, driving more views.
- (I) The video was shared by a large number of dance schools and communities, spreading its reach.
- (J) The video was part of a viral challenge that encouraged people to share it.
- (K) The video has high-quality production values that make it stand out from other content.

CorrectAnswer: "(C)

Satisfying and Relaxing Kinetic Sand ASMR shorts - video_url

608

609 IV. Comment Matching (CM)

610 IV.1. Humorous

- Task Description: The video will spark laughter because of certain comments, making the audience
- feel funny, testing whether the model can match it correctly.
- 613 Example Question:

Question: Based on this video, which of the following comments is likely to make people laugh? **Options**:

(A) Did he just audition for a water ballet?

- (B) How many fish does it take to catch a man?
- (C) Is there a Walmart beneath the river?
- (D) The fish are holding a grudge, watch out!
- (E) Now that's what I call a splash of creativity.
- (F) I came for the fishing tips and stayed for the synchronized swimming.
- (G) That water has more personality than my neighbor!
- (H) I'm starting to think he's part fish.
- (I) I think the fish caught him instead.
- (J) That's definitely a land fish champion.
- (K) That fish will never trust humans again.

CorrectAnswer: "(C)", He DI Lao - video url

615

616 IV.2. Thought-provoking

Task Description: Some comments under the video will enhance people's thinking and test whether

618 the model can accurately understand.

619 **Example Question**:

Question: Which of the following statements can better explain the social reality expressed in this animation?

Options:

- (A) The animation showcases an idealized view of advancement within a corporate ladder.
- (B) The depiction highlights the dehumanization and mechanization of individuals in a powerful social system.
- (C) It portrays the joy of discovering one's true passions through societal pressures.
- (D) The scenes show a man achieving happiness through daily routine.
- (E) It represents personal ambition and the drive for success in individual careers.
- (F) The animation indicates the triumph of an individual's spirit in the face of adversity.
- (G) It reflects the disintegration of traditional family roles.
- (H) The animation shows the importance of family support in work-life balance.
- (I) It emphasizes the challenge of maintaining personal identity in urban settings.
- (J) We are all working for others without realizing it due to our own needs.
- (K) The animation illustrates the struggle with contemporary health issues.

CorrectAnswer: (J)

EL EMPLEO - video url

620

621 IV.3. Trending

- Task Description: It is relatively difficult to test whether the model can accurately infer and analyze
- the most popular comments under the video.
- 624 Example Question:

Question: Which of the following comments best summarizes the content conveyed by this film?

Options:

- (A) Material possessions define one's value.
- (B) Selfless acts lead to rewards that surpass material wealth.
- (C) Loneliness is a desirable state.
- (D) Personal gains are the ultimate goal of helping others.
- (E) Isolation is the path to personal growth.
- (F) True happiness is found through wealth accumulation.
- (G) Success comes from competitive behavior.
- (H) Sharing leads to financial prosperity.
- (I) He receives what money can't buy.
- (J) Adversity breeds stronger individuals.

CorrectAnswer: (I)

Unsung Hero - video_url

```
626 V. Implicit Symbol (IS)
```

- 627 V.1. Cultural Symbols
- Task Description: Test whether the model can infer and analyze the cultural characteristics hidden
- under the surface visual elements of the video (such as nationality, festivals, customs, religion, etc.).
- 630 Example Question:

Question: The plaque inscribed with "Dominating Three Continents" that appears in the video is most likely to be found in the architecture of which of the following religions?

Options:

- (A) Taoism
- (B) Shinto
- (C) Sikhism
- (D) Judaism
- (E) Islam
- (F) Christianity
- (G) Buddhism
- (H) Hinduism
- (I) Jainism
- (J) Zoroastrianism

CorrectAnswer: (G)

[4K] Hangzhou 2024 in the misty rain | West Lake, Lingyin Temple, Night walking in Hefang Street - video url

631

632 V.2. Art Symbols

- Task Description: Test whether the model can infer and analyze the art-related characteristics hidden
- under the surface visual elements of the video (such as dance style, artistic skills, imitation, etc.).
- 635 Example Question:

Question: What is the shadow that appears in our view at 1:40 imitating? **Options**:

- (A) The shadow is imitating a pole dancer.
- (B) The shadow is imitating a person washing a dog.
- (C) The shadow is imitating a person brushing their hair.
- (D) The shadow is imitating someone playing a violin.
- (E) The shadow is imitating two people engaged in a conversation.
- (F) The shadow is imitating someone painting a wall.
- (G) The shadow is imitating a person feeding a horse.
- (H) The shadow is imitating a person washing their car.
- (I) The shadow is imitating a dog barking at a person.
- (J) The shadow is imitating someone performing a magic trick.
- (K) The shadow is imitating a person holding an umbrella.
- (L) The shadow is imitating someone walking a large dog.

CorrectAnswer: (A)

LEAKED! Hilarious Shadow Puppets - AGT 2023 Early Release - video url

636

637 V.3. Other Symbols

- Task Description: Test whether the model can infer and analyze other special symbols (such as
- commercial advertisements, etc.) hidden under the surface visual elements of the video.
- 640 Example Question:

Question: "What do you think the chimpanzee that appears multiple times in the film symbolizes?",

Options:

(A) The chimpanzee symbolizes chaos and disruption in everyday life. (B) The chimpanzee symbolizes a childhood fear. (C) The chimpanzee symbolizes technology invading personal

space. (D) The chimpanzee symbolizes the unpredictability of fate. (E) The chimpanzee symbolizes a glue company. (F) The chimpanzee symbolizes lost opportunities. (G) The chimpanzee symbolizes an obsession with social status. (H) The chimpanzee symbolizes environmental degradation. (I) The chimpanzee symbolizes the desire for freedom. (J) The chimpanzee symbolizes misunderstanding between people. (K) The chimpanzee symbolizes reliability and trust in friendships.

CorrectAnswer: (E)

All Gorilla glue ads - video url

642

643 D.2 Explicit Reasoning Tasks

- I. Causal Reasoning (CAR)
- 645 I.1. Forward Reasoning
- Task Description: Forward reasoning can also be understood as the prediction of future events,
- 647 including prediction of outcomes, prediction of content that has not yet appeared, etc.
- **Example Question:**

Question: What is the speculated ending of the film? **Options**:

- (A) The movie concludes with an unexpected twist where the flowers reveal a hidden secret.
- (B) The ending is a cliffhanger, leaving the audience uncertain about the characters' fate.
- (C) Her boyfriend passed away due to illness, leaving the girl devastated with grief.
- (D) The film wraps up with a joyous family reunion.
- (E) The film ends with a dramatic breakup as one character leaves with a heavy heart.
- (F) The movie concludes with a comedic mishap involving the flowers.
- (G) The ending shows a tragic farewell as one character moves to a new city.
- (H) The film ends with the revelation of a long-lost sibling.
- (I) The story concludes with the characters embarking on a spontaneous road trip.
- (J) The film ends on a melancholic note, reflecting on lost opportunities.
- (K) The video closes with a heartwarming reconciliation between the main characters after exchanging heartfelt notes and gestures.

CorrectAnswer: (C)

For Milo - AWARD WINNING 1 Minute Short film (2020) - video_url

649

650 I.2. Backward Reasoning

- Task Description: Backward reasoning means finding the cause from the effect and inferring the
- reason why an event occurred.
- 653 Example Question:

Question: Why was the elderly black man warned by security at the beginning of the film? **Options**:

- (A) Mobile phones are not allowed for recording during magic shows.
- (B) He was trying to sell unauthorized merchandise.
- (C) He was recognized as a local celebrity causing disruptions.
- (D) He was accused of stealing a bicycle.
- (E) He was creating loud music disturbing the peace.
- (F) He was believed to have lost his entrance ticket.
- (G) He was inadvertently blocking the pathway.
- (H) He was associated with another person causing trouble nearby.
- (I) He was engaged in card tricks that security found suspicious.
- (J) He was loitering without a purpose.

CorrectAnswer: (A)

Now You See Me Official Opening Scene (2013) - Mark Ruffalo, Morgan Freeman Movie HD - video url

55 II. Sequential Structure Reasoning (SSR)

II.1. Narrative Structure

- Task Description: Reasoning and analyzing the narrative order of the entire video, including the
- editing order, such as sequential, flashback, and interpolation.
- **Example Question:**

656

Question: What kind of narrative sequence does the film employ? **Options**:

- (A) non-linear flashback sequence, where events are shown out of chronological order, often revealing backstory
- (B) parallel overlapping sequences, showing multiple storylines happening simultaneously with some overlap
- (C) cyclical narrative structure, repeating events or themes in a circular pattern
- (D) linear narrative sequence, following a straightforward progression from beginning to end
- (E) random jumps in the timeline, moving unpredictably between different points in time
- (F) interwoven thematic structure, weaving together different themes and ideas throughout the story
- (G) reverse chronological order, starting with the end and moving backwards in time
- (H) fragmented narrative, presenting the story in disjointed or broken segments
- (I) begins with a flashback and then proceeds in chronological order
- (J) episodic progression, advancing the story through a series of distinct episodes or chapters
- (K) multi-perspective narrative, telling the story from multiple characters' points of view **CorrectAnswer**: (I)

Identity SHORT FILM (Award Winning Inspirational Short) - video url

660

661

II.2. Core Connecting Elements

- Task Description: Videos with this type of question usually have a key connecting element that runs through the entire video. It is carefully designed by the producer and tests the model's inductive reasoning of the visual information of the entire video.
- 665 Example Question:

Question: What is the recurring element in the video, summarized in one word? **Options**:

- (A) Pareidolia
- (B) Smile
- (C) Alarm
- (D) Work
- (E) Mirror
- (F) Mundane
- (G) Routine
- (H) Suit
- (I) Coffee
- (J) Sleep
- (K) Bedroom
- (L) Portrait

CorrectAnswer: (B)

PAREIDOLIA - 1 Minute Short Film | Award Winning - video url

666

667 II.3. Inference on Editing Techniques

- Task Description: These tasks evaluate the models' deep analysis and multimodal reasoning about
- video editing strategies.
- 670 Example Question:

Question: "Please guess how many videos were needed to record the moment the man punched the punctured water ball at the beginning of the video?",

Options:

- (A) At least two separate takes would be needed.
- (B) At least one single take is needed.
- (C) Three separate takes are needed.
- (D) Four separate takes are needed.
- (E) Each scene can be captured in a single continuous take.
- (F) Five separate takes are needed.
- (G) Six separate takes are needed.
- (H) Eight separate takes are needed.
- (I) Ten separate takes are required.
- (J) Twenty separate takes are necessary.
- (K) At least ten separate takes are needed.

CorrectAnswer: (C)

Playing With Time - video_url

Note: The reasoning and analysis process of this question can refer to this disassembly video.

671

II.4. Hallucination

- Task Description: Evaluate whether the model perceives various types of hallucinations when 673
- perceiving video content. 674
- **Example Question:** 675

Question: How many dancers are there in the video? **Options**:

- $(\bar{A}) 0$
- (B) 1
- (C) 2(D)3
- (E) 4
- (F) 5
- (G)6
- (H) 7
- (I) 8 (J)9
- (K) options before are all false

CorrectAnswer: (B)

Rat dance with falling body parts - video_url

676

III. Counterintuitive Reasoning (CIR) 677

III.1. Magic Deconstruction or Special Effects Editing 678

- **Task Description:** This type of video usually creates some impossible magical effects, but some are
- magic tricks, and some are editing and special effects, which require deeply reasoning.
- **Example Question:**

Question: Starting at 4:35, how did the man achieve this magical effect in the magic trick? Options: (A) Sleight of hand technique with a hidden ring, using dexterity to conceal and reveal the ring.

- (B) Utilizing a mirror to confuse the audience, creating optical illusions through reflection.
- (C) A distraction technique with a smoke bomb, diverting attention with a sudden burst of smoke.
- (D) A special ring that retracts into a fake thumb, using a concealed mechanism to make the ring disappear.
- (E) Using a magnet hidden in the sleeve, manipulating objects with magnetic force.
- (F) A camera trick with video editing, altering footage to create the illusion of magic.
- (G) Sleight of hand technique with a hidden string, using a concealed thread to control objects.

- (H) The bottle inside the paper bag had already been altered to leave only the outer plastic skin.
- (I) Employing a twin assistant to swap the ring, using a look-alike to deceive the audience.
- (J) The use of an invisible thread, employing a nearly undetectable line to move objects.
- (K) A sound cue to mislead the audience's attention, using noise to distract from the real action.

CorrectAnswer: (H)

Level 1 to 100 Magic Tricks Anyone Can Do - video_url

683

III.2. Artistic Techniques

Task Description: This type of video usually creates some impossible scenes, but it is usually an artistic expression deliberately designed by the author.

687 Example Question:

Question: Why is the shadow on the boy's face illuminated by sunlight at 1:06?

Options: (A) Because the boy moves to a position where a strong light source is directly above him, not related to the girl.

- (B) It's just a coincidence that the angle of the sun changes suddenly at that moment, and has nothing to do with any special meaning.
- (C) The sunlight illuminates the shadow because the cameraman adjusts the lighting equipment to create a better visual effect.
- (D) The girl's appearance brings good luck, and the sunlight representing good fortune clears away the gloom of bad luck in his world.
- (E) This is because the boy has walked into a neighborhood with better weather and climate.
- (F) The sunlight lights up the shadow because there is a hidden light emitting device in the scene that is turned on at 1:06.
- (G) It's a result of the special lens filter used during filming, which makes the shadow on the boy's face appear to be lit by sunlight.
- (H) Because the boy didn't get hurt after falling and his mood improved, the sunlight is used to represent his improved mood.

CorrectAnswer: (D)

CGI Animated Short Film HD "Jinxy Jenkins & Lucky Lou" by Mike Bidinger & Michelle Kwon | CGMeetup - video_url

688

689 III.3. Humor and Exaggeration

- Task Description: A common technique in humorous videos is to use exaggerated expressions that
- seem unreasonable, but there are some clues to understand the meaning. This type of question tests
- the model's ability to reason about exaggerations and unusual techniques.
- 693 Example Question:

Question: Why does the first half of the scene look sunny but also show rain? **Options**:

- (A) It is a sunshower, when rain falls while the sun is shining.
- (B) The character is dreaming of being both wet and warm.
- (C) There are rainclouds directly above while sunlight comes from the side.
- (D) It is snow instead of rain, reflecting the sunlight.
- (E) The effect is caused by morning fog and light refraction.
- (F) It's a visual illusion caused by mist.
- (G) The character moved to a different location quickly.
- (H) A rainbow is forming which intensifies the sunlight.
- (I) Dew drops from trees reflect sunlight.
- (J) There are two unrelated weather animations merged together.
- (K) The man wet the bed, which caused the presence of water in his dream.

CorrectAnswer: (K)

It now makes sense - video_url

95 IV. Cross-modal Transfer Reasoning (CTR)

- Evaluate the ability to transfer reasoning from video to text, audio, video or image (for example,
- video-to-text: the theme of a video may have the same meaning as a famous quote) Task Description:
- Evaluate the ability to transfer reasoning from video to text (for example, the theme of a video may
- 699 have the same meaning as a famous quote)
- **Example Question:**

Question: Which of the following proverbs best explains the theme of this short film? **Options**:

- (A) When one door closes, another opens.
- (B) Opportunity knocks only once.
- (C) Time heals all wounds.
- (D) The early bird catches the worm.
- (E) Never judge a book by its cover.
- (F) All that glitters is not gold.
- (G) The grass is always greener on the other side.
- (H) Actions speak louder than words.
- (I) A stitch in time saves nine.
- (J) Beauty is in the eye of the beholder.
- (K) Absence makes the heart grow fonder.
- (L) A penny saved is a penny earned.

CorrectAnswer: (E)

Video: Award Winning SHORT FILMS Don't Judge | BATTI Hindi Heart Touching Short Movies | Content Ka Keeda - video_url

701

702 V. Video Type and Intent (VTI)

703 V.1. Video Type

Task Description: Evaluate the model's ability to analyze video types, such as commercials, science

705 fiction films, comedies, etc.

Example Question:

Question: What type of video is this most likely to be? **Options**:

- (A) A documentary about airplane technology
- (B) Advertisement for an ice-cream
- (C) A drama set on an airplane
- (D) A comedy film featuring an airline
- (E) An in-flight safety demonstration video
- (F) A travel vlog featuring aerial views
- (G) A science fiction movie on a spaceship
- (H) This is an advertisement.
- (I) A video tour of an airplane factory
- (J) A virtual reality experience of flying
- (K) A news segment on turbulence incidents

CorrectAnswer: (H)

Leo Messi vs Kobe Bryant - Legends on Board - Turkish Airlines - video_url

707

708 V.2. Video Intent

Task Description: Reasoning and analyzing the purpose and production intention of the video (e.g.

what kind of product performance is promoted in a commercial advertisement, etc.)

711 Example Question:

Question: Which year do you think this video was most likely released?

Options:

(A) 2018

```
(B) 2017
(C) 2016
(D) 2015
(E) 2023
(F) 2019
(G) 2023
(H) 2020
(I) 2014
(J) 2013
CorrectAnswer: (H)
Lockdown | One Minute Short Film Challenge | Film Riot - video_url
```

Evaluation Details

715 E.1 Baselines

713

735

The baselines include closed-source models: (1) GPT series: GPT-4o [41], GPT-4o-mini [42], and GPT-4.1 [43]; (2) Gemini series: Gemini-2.0-flash, Gemini-2.0-flash-thinking-01-21 [44], and Gemini-2.5-flash [45]; (3) Claude-3-5-Sonnet-20241022 [46]; (4) o4-mini [4]; open-source models: (1) Qwen series: Qwen2.5-VL (7B/72B-Instruct) [47]; (2) Gemma series: Gemma-3 (12B/27B) [48]; (3) InternVL series: Intern3-VL (8B/38B) [49]; (4) LLava series: LLava-Onevision-7B [50], Video-LLava-7B [51]; (5) Phi-4-multimodal-Instruct [52]; (6) Other video models: Cogvlm2-video-llama3-chat [53], NVILA-8B-Video [54]. All local experiments are conducted on 4×A100 80GB GPUs.

723 E.2 Frame Selection

We followed the official configurations of models that support multi-image input, as well as settings 724 in previous works [8; 28], to define the number of input frames for each model. Specifically, we fixed 725 the number of frames per model and sampled them evenly across the video duration. We sampled 726 8 frames for LLaVA-OneVision, Video-LLaVA, and NVILA-8B-Video, Phi-4-multimodal-instruct; 727 16 frames were sampled for Qwen2.5-VL-7B, Qwen2.5-Omni-7B, CogVLM2-Video-LLaMA3-728 Chat, InternVL-8B, Gemma-3-it-12B and Gemini-2.0-Flash-Thinking; 32 frames were sampled 729 for Qwen2.5-VL-72B, InternVL-38B, Gemma-3-it-27B, GPT-4o, GPT-4o-Mini, o4-mini, Gemini-2.5-Flash-preview and Claude-3.5-Sonnet. Exceptionally, since Gemini-2.0-Flash supports long video and multimodal context inputs, we sampled one frame per second across each video, with a 732 maximum cap of 512 frames to ensure API stability. Additionally, to enable a fair comparison with 733 Gemini-2.0-Flash-Thinking, we also tested a version of Gemini-2.0-Flash with 16 frames. 734

F Details of CoT Analysis Experiment

This section presents the CoT analysis experiments discussed in Section 4.5. We use a representative model CoT to illustrate 4 categories of analysis in Figure 10. Specifically, Text Analysis refers to the examination of textual information such as the question and options; Video Analysis focuses on the content of the video; Question Frame targets the specific frame referenced in the question—for instance, the frame where the magician controls two flames; and Other Frame pertains to frames outside the scope of the question. In Figure 10, yellow, red, and blue represent text, question frame, and other frame analysis respectively. Red, blue, and green all represent video analysis.

Prompt for CoT Annotation

You will be given a model's textual reply to a video-based question along with the video frames. Your task is to determine four boolean labels for each chunk of the reply:

- 1. 'other frame desc': Does this chunk describe visual information from frames other than question frame?
- 2. 'question frame desc': Does this chunk correctly describe visual information from the question frame specified in the question?
- 3. 'video analysis': Does this chunk perform analysis of the video content?
- 4. 'text analysis': Does this chunk perform analysis of the text (e.g., question text, options) rather than visual content
- The question frame refers to the specific frame(s) referenced by the question prompt.
- Other-frame descriptions are visual details not present in the question frame but from other frames.
- Video analysis includes describing trends, motions, or visual inference beyond plain description.
- Text analysis includes reasoning over question text, options, or external text context.

```
Respond strictly in JSON:
{
"other frame desc": true or false,
"question frame desc": true or false,
"video analysis": true or false,
"text analysis": true or false
}
Question: {question}
Reply Chunk: {chunk}
Whole CoT Reply: {CoT}
```

Table 6: CoT analysis prompt.

CoT Analysis Example









Question: How does the man at the beginning of the video pick up and casually control the flame on the lighter?

Options:

- (A) He used a holographic projector to simulate the flame.
- (B) He used a special flame-retardant chemical on his hand to create the illusion.
- (C) He possessed an innate immunity to fire.
- (D) He practiced yoga meditation to withstand any flame heat.
- (E) A quick extinguishing spray was applied that halted the flame.
- (F) He surrounded the flame with an invisible film.
- (G) He mastered the art of fire manipulation.
- (H) The flame was made of non-flammable gas.
- (I) He applied a hidden cooling technology under his sleeve.
- (J) The flame was actually an LED light.
- (K) A hidden lighter in his hand, a sleight of hand trick.

My thinking process to arrive at the answer is as follows:\n\n1. **Understand the Question:** The question is specifically asking about *how* the man at the beginning of the video picks up and *casually controls* the flame on a lighter. The key words are \"how\" and \"casually controls,\" implying a method or technique rather than a supernatural ability or complex technology. $\n\$ 2. **Analyze the Video**I noticed that the \u201cfloating\u201d flame always stayed just above his palm in a way that perfectly hid the source, and that when he distracted us he slipped a small butane lighter into position. In the close\u2010up reveal later, you can see how the second lighter remains concealed until he produces the flame on his fingers. It\u2019s all done by classic palm concealment and sleight of hand with a hidden lighter. \n\n2. **Analyze the Options:** I reviewed each option to see which one is most plausible and aligns with the context of a magic trick video. I considered the practicality and likelihood of each explanation.\n\n * **(A) Holographic Projector, (J) LED Light, (H) Non-flammable gas:** These options suggest the flame isn't real. While possible for illusions, they are less likely for a *casual* control and might be visually detectable as fake in a video focusing on close-up magic.\n\n * **(C) Innate immunity to fire, (D) Yoga meditation, (G) Mastered fire manipulation:** These options are fantastical or mystical and don't explain a practical method. They are less likely to be the intended explanation in a video about magic *tricks*.\n\n 4. **Prioritize Simplicity and Plausibility:** In the context of magic tricks, simpler explanations are often better. Option (K) - sleight of hand - is the simplest and most direct explanation that fits the description of \"casual control.\" It's a wellknown method for creating illusions with small objects like lighters.\n\n5. **Eliminate Less Likely Options:** I eliminated the options that were:\n * Too complex or technologically advanced (A, I).\n * Unrealistic or fantastical (C, D, G, H).\n * Less direct or slightly less probable for \"casual control\" in a close-up magic context compared to sleight of hand (B, E, F, J).\n\n6. **Select the Best Fit:** Option (K) - \"A hidden lighter in his hand, a sleight of hand trick\" emerges as the most likely and best-fitting answer because it is:\n * A common and effective magic technique.\n * A simple and direct explanation.\n * Consistent with the idea of \"casual control.\"\n * Plausible in the context of a magic trick video.\n\n\nTherefore, based on the question, the options, and the likely context of a magic trick video, option (K) is the most accurate answer.\n\n**Final Answer:** [[K]]

Figure 10: CoT example of experiments in Section 4.5. Yellow, red, and blue represent text, question frame, and other frame analysis respectively. Red, blue, and green all represent video analysis.

743 G Evaluation Prompt

We evaluated two settings in the main experiment, zero shot and zero shot + CoT. The prompts used

745 are as follows.

Prompt for Zero-Shot Setting

[[INSTRUCTION]]

Please select the best answer to the following multiple-choice question based on the video.

Only one option is the most accurate answer in relation to the question and the video.

What is the correct answer to this question

{Question}

Options:

{Options}

[[END OF INSTRUCTION]]

[[OUTPUT FORMAT]]

Format your answer as follows:

Please directly output the answer letter without thinking and explanation.

If the correct option letters (A, B, C, D...) for the multiple-choice question is X, give the final correct option number in the following format: "[[X]]"

[[END OF OUTPUT FORMAT]]

Table 7: Evaluation prompt for the Zero-Shot Setting.

Prompt for Zero-Shot + CoT Setting

[[INSTRUCTION]]

Please select the best answer to the following multiple-choice question based on the video.

Only one option is the most accurate answer in relation to the question and the video.

What is the correct answer to this Question:

{Question}

Options:

{Options}

Let's think step by step.

[[END OF INSTRUCTION]]

[[OUTPUT FORMAT]]

Format your answer as follows:

Your thinking process.

If the correct option letters (A, B, C, D...) for the multiple-choice question is X, give the final correct option number in the following format: "[[X]]"

[[END OF OUTPUT FORMAT]]

Table 8: Evaluation prompt for the Zero-Shot + CoT Setting.

746 H Case Study

747 In this section, we present reasoning processes and results from selected models on the MMR-V

benchmark. Through these case studies, we aim to better illustrate the current shortcomings of models

749 in multimodal reasoning tasks and provide insights that may inspire future research and advancements

750 in this area.

Firstly, there is a comparison between a good CoT and a poor CoT in 11. Yellow highlights indicate

752 text-based analysis, while green highlights denote video-based analysis. As shown, the upper CoT

engages in deep multimodal reasoning grounded in the video content, ultimately leading to the correct

answer. In contrast, the lower CoT performs only shallow perception of the video and relies heavily

on deep reasoning over the textual content, resulting in a text-dominant reasoning process that yields

756 an incorrect answer.

Note: The reasoning and analysis process of example Figure 15 can refer to this disassembly video.

Good MCoT vs Bad MCoT



Question: What does the object being chased by the people refer to? Options:

- (A) Difficulties in life.
- (B) His fully automatic house that can walk autonomously.
- (C) The house that this person is saving money to buy.
- (D) The old house that the family decided to renovate.
- (E) The villa that the rich man bought as a vacation home.
- (F) The mansion that the celebrity owns in the city center.
- (6) The cottage that the couple wants to rent for their honeymoon.
- (H) The enchanted house.

o4-mini-CoT: The animation shows people in hot pursuit of a walking \u2010house \u2014literally chasing after a \u201chome.\u201d In the course of the film we see the house slipping through rain, snow, across seas and canyons, always just out of reach, until in the end it stands on the far side of a great divide between a crowd of people and a cluster of buildings. This is a metaphor for the endless struggle to afford and finally own a home. Thus the chased object refers to the house the person is saving up to buy. [[C]]

Gemini-2.0-Flash-Thinking-CoT: The**Thinking Process:**\n\n1. **Identify the question:** The question asks about the meaning of \"the object being chased by the people\".\n\n2. **Watch the video and observe the object:** The video clearly depicts a house with legs walking and being chased by people. \n\n3. **Analyze the options and compare them to the video:**\n * **(A) Difficulties in life:** While metaphorical interpretations are possible, the video literally shows a house being chased. This option is too abstract and doesn't directly relate to the visual content.\n * **(B) His fully automatic house that can walk autonomously:** This option directly describes what is seen in the video $\u2013$ a house that can walk on its own. This is a strong contender. \n * **(C) The house that this person is saving money to buy:** There is no visual or narrative element in the short video clip about saving money or buying a house. This option is not supported by the video.\n ***(D) The old house that the family decided to renovate:** No renovation or family context is present in the clip. This option is irrelevant.\n ***(E) The villa that the rich man bought as a vacation home:** The video does not suggest wealth, vacation, or a villa specifically. This option is not supported.\n ***(F) The mansion that the celebrity owns in the city center:** No celebrity, mansion, or city center is depicted. This option is irrelevant.\n ***(G) The cottage that the couple wants to rent for their honeymoon:** There is no mention of a couple, honeymoon, or cottage rental in the video clip. This option is irrelevant.\n * **(H) The enchanted house:** While the walking house could be considered \"enchanted,\" option (B) is more precise and descriptive of the *mechanism* of the house's unusual behavior (autonomously walking). \"Enchanted\" is a more general and less accurate description compared to \"fully automatic house that can walk autonomously.\"\n\n4. **Select the best option:** Option (B) is the most direct, accurate, and visually supported answer based on the video. It literally describes the object being chased.\n\n**Final Answer:**\n\n[[B]]

Figure 11: A comparison of CoTs from two models on the same task. Yellow and green indicate text and video analysis, respectively. As shown, o4-mini's reasoning paradigm demonstrates a deeper analysis of the video content.

Lack of Visual Reasoning









Question: Why did the broken umbrella in the boy's hand, damaged by the wind, suddenly get fixed while the boy and the girl were in the air?

Options:

- (A) Because the boy initially used too much force, causing the umbrella to be blown apart, and it was repaired after he adjusted his effort.
- (B) Because the boy held the girl's hand, gained good luck, and as a result, the umbrella was restored.
- (C) Because the wind in the air happened to die down, the umbrella returned to normal.
- (D) Because the girl held the boy's hand, making him feel extremely happy, which in turn allowed the umbrella to be repaired smoothly.
- (E) It's just a visual illusion in the video, and the umbrella wasn't actually blown apart in the first place.
- (F) The umbrella was designed to be a self repairing model, and it activated its repair function automatically.
- (G) It was a pre arranged magic trick, and the umbrella was never truly damaged.
- (H) The video was reversed during that part, making it look like the broken umbrella was getting fixed.
- (I) Aerodynamics, changes in wind strength and direction contributed to the umbrella being repaired.
- (J) It was a computer generated special effect added during post production to create a magical look."

Video Type: Animation

Task Type L1: Explicit Reasoning.

Task Type L2: Counterintuitive Reasoning

Task Type L3: Artistic Techniques

Gemini-2.0-Flash-Thinking: [[Thinking Process]]\nThe user is asking for the best answer to the question \"Why did the broken umbrella in the boy's hand, damaged by the wind, suddenly get fixed while the boy and the girl were in the air?\" based on the provided video and options.\n\nAfter watching the video, it's clear that:\n1. The umbrella is initially broken by the wind on the ground.\n2. The boy and girl are lifted into the air by pianos attached to balloons.\n3. While in the air, the broken umbrella appears to become fixed.\n\nNow let's evaluate each option against what's observed in the video and consider the most logical explanation within the context of the animation:\n\n* **(B) Because the boy held the girl's hand, gained good luck, and as a result, the umbrella was restored.** - This is a whimsical, magical explanation. While the video has a fantastical element, \"good luck\" causing physical repair is not the most direct or likely intended answer.\n\n(I) Aerodynamics, changes in wind strength and direction contributed to the umbrella being repaired.** This is a more detailed and scientifically grounded explanation. As they ascend, changes in wind strength, direction, and the aerodynamic forces acting on the umbrella could allow it to reshape itself. This is very plausible.\n\n* - While the video *is* computer-generated animation, this option describes *how* it's shown, not *why* it happens within the video's narrative. It doesn't explain the in-universe reason for the repair.\n\n\nComparing options (C) and (I), option (I) is more comprehensive as it includes \"aerodynamics\" in addition to changes in wind, which is a more complete physics-based explanation relevant to an umbrella in motion in the air. Option (C) is a simplified version of (I). \n\nOption (I) is the most scientifically plausible and best explanation from the given options for why the umbrella appears to get fixed within the context of the animation. $\ln[I]$

Human Error Annotation: The model only analyzes visual information in the first small orange part of the analysis, and only analyzes the frame corresponding to the question. After that, most of the CoT is doing text analysis (comparing the pros and cons of each option)

Figure 12: Error Case: Lack of Visual Reasoning.

Lack of Visual Reasoning













Question: How did the man manage to insert the top card into the deck and then reveal the top card again, only to find it back in its original position?

Options

"(A) The card deck was equipped with a hidden compartment.",

- (B) He had two identical cards, and he switched them.",
- (C) The card deck was magnetic, allowing precise control.",
- (D) The cards were fitted with a special adhesive.",
- (E) The man used camera tricks to manipulate what was shown.",
- (F) He used invisible threads to control the card's position.",
- (G) The cards had sensors that automatically rearranged them.",
- (H) He inserted the card in a fake cut, not the real deck.",
- (I) The man used sleight of hand to secretly return the card to the top.",
- (J) The deck was rigged with springs to pop the card back up.",
- (K) When showed the cards, he actually turned over two cards."

Video Type: TV

Task Type L1: Explicit Reasoning.

Task Type L2: Counterintuitive Reasoning

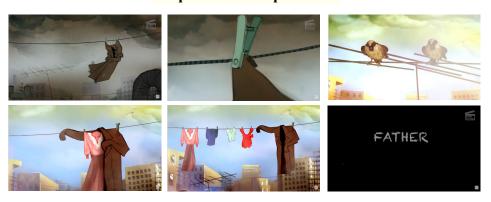
Task Type L3: Magic Deconstruction

Gemini-2.0-Flash-Thinking: My thinking process to arrive at the correct answer, option (I), is as follows:\n\n1. **Understand the Question: ** The question is asking for the *method* used by the man to make the top card reappear on top of the deck after being inserted into the deck. This is a classic magic trick scenario.\n\n2. **Watch the Video Carefully:** I watched the video multiple times, paying close attention to the man's hand movements, especially when he \"inserts\" the card and then reveals it again.\n\n3. **Analyze the Actions in the Video:**\n * He shows the top card (2 of Clubs).\n * He *appears* to insert the 2 of Clubs into the deck.\n * He immediately reveals the top card again, and it's still the 2 of Clubs.\n\n4. **Evaluate Each Option Against the Video and Common Magic Principles:** I went through each option, considering if it's plausible based on what I saw and general knowledge of magic tricks. n\n ***(A) Hidden Compartment: ** Possible in theory, but decks with hidden compartments are less common for simple card tricks. It's also not visually indicated in the video. Less likely for this type of quick trick.\n***(I) Sleight of Hand:** This is the most probable and general explanation for card magic. Sleight of hand encompasses various techniques to manipulate cards secretly. This option directly addresses the core of card magic and fits the observed actions. *This is the most likely answer.*\n ***(J) Rigged Springs:** Mechanically complex and unnecessary. Springs would be hard to conceal and control for this specific effect. Unlikely.\n ***(K) Turned Over Two Cards:** Doesn't explain how the card returns to the top after insertion. Also, visually, it looks like he's only turning over one card. Unlikely and irrelevant to the question.\n\n5. **Select the Best Fit:** After evaluating all options, option (I) \"The man used sleight of hand to secretly return the card to the top\" stands out as the most general, accurate, and likely explanation. It encompasses the core principle of card magic and directly addresses the question of how the card returned to the top. Option (B) is a *type* of sleight of hand, but (I) is broader and more encompassing as a general explanation. Option (H) is also related to sleight of hand but option (I) is more direct to the question.\n\n6. **Final Answer:** Therefore, option (I) is the best answer. $\n\[[I]]$

Human Error Annotation: The model only analyzes the actions in the frames adjacent to the question frame in the video. In fact, at the end of the video, the magician explains that he turned over two cards at the beginning and showed the second card. The **model did not observe the evidence frame of the decryption part through visual reasoning**.

Figure 13: Error Case: Lack of Visual Reasoning.

Implicit Misinterpretation



Question: What is the main message that this video wants to convey?

Options

- (A) The video is mainly about showing how easy it is to solve family problems without much effort.
- (B) It praises the mother's strength and kindness, and how she teaches her children to interact with others.
- (C) It praises the sacrifice of the brown coat, ensuring that the other clothes do not get wet.
- $(D)\ It\ aims\ to\ illustrate\ that\ family\ members\ should\ always\ rely\ on\ external\ help\ when\ in\ trouble.$
- (E) The video tries to convey that family is not important compared to personal success.
- (F) The video praises the father's sense of responsibility and commitment.
- The main point is that fathers always make mistakes and children have to fix them.
- (H) It aims to show that fathers should focus more on their careers than on family.

Video Type: Animation

Task Type L1: Implicit Reasoning. Task Type L2: Theme Understanding Task Type L3: Everyday Topics

GPT-40: [Analyze the best option for question] \nThe video showcases a brown coat protecting other clothes from rain on a clothesline. The actions and symbolism suggest the brown coat's sacrifice, highlighting the theme of protection or sacrifice for others, akin to a familial role. \n\n[Justification for your final choice based on the thinking process.]\nThrough the coat's protective measures, the video illustrates the theme of sacrifice by the brown coat to keep other clothes dry. This aligns closely with option C, where the sacrifice of the brown coat ensures that the other clothes do not get wet.\n\nGive the final correct option number in the following format: \"[[C]]\"

Human Error Annotation: The model only analyzes the actions in the frames adjacent to the question frame in the video. In fact, at the end of the video, the magician explains that he turned over two cards at the beginning and showed the second card. The model did not observe the evidence frame of the decryption part through visual reasoning.

Figure 14: Error Case: Implicit Misinterpretation.

Reasoning Error







Question: Please guess how many videos were needed to record the moment the man punched the punctured water ball at the beginning of the video?

Options:

- (A) At least two separate takes would be needed.
- (B) At least one single take is needed.
- (C) Three separate takes are needed.
- (D) Four separate takes are needed.
- (E) Each scene can be captured in a single continuous take.
- (F) Five separate takes are needed.
- (G) Six separate takes are needed.
- (H) Eight separate takes are needed.
- (I) Ten separate takes are required.
- (J) Twenty separate takes are necessary.
- (K) At least ten separate takes are needed.

Video Type: TV

Task Type L1: Explicit Reasoning.

Task Type L2: Counterintuitive Reasoning Task Type L3: Magic Deconstruction

Gemini-2.0-Flash-Thinking: [[Thinking Process]]\nThe question asks to guess the number of videos needed to record the water balloon punch scene at the beginning of the video. This is a question about the *process* of video making, specifically related to achieving a good take.\n\nBased on common video production knowledge and the nature of capturing a dynamic event like a water balloon bursting in slow motion, it is highly unlikely that such a scene would be captured perfectly in a single take. There are many variables such as timing of the punch, the balloon's burst pattern, camera focus, and framing.\n\nConsidering the options:\n* **Options suggesting a single take (B) and (E)** are less likely. Capturing a dynamic and potentially messy action like a water balloon burst perfectly on the first try is improbable, especially for a visually appealing slow-motion shot.\n***Options suggesting multiple takes (A), (C), (D), (F), (G), (H), (I), (J), (K)** are more reasonable. Filmmaking often involves multiple takes to ensure the best possible shot.\n* **Option (A) 'At least two separate takes would be needed.'** This is the most conservative option indicating that more than one take is probable. It's a very reasonable minimum guess.\n* **Options with higher numbers of takes (C), (D), (F), (G), (H), (I), (J), (K)** are also possible, but they are more specific guesses. Without explicit information from the video, these are harder to justify as the *best* guess compared to the minimal \"at least two\".\n* **Option (K) 'At least ten separate takes are needed.*** and **(J) 'Twenty separate takes are necessary.'** are likely overestimations for a relatively simple shot like punching a water balloon. While possible, they are less probable as the *best* guess without more information.\n\nGiven that we are asked to *guess* and choose the *best* answer, the most reasonable and safe guess, based on general video production practices and the nature of the scene, is that **at least two takes would be needed**. This accounts for the likelihood of needing to adjust something or re-perform the action for a better result.\n\nTherefore, Option (A) is the most accurate and reasonable answer from the provided options.\n\n[[Correct Option]]\n[[A]]

Human Error Annotation: The model accurately found the frame corresponding to the problem and perceived the evidence frame required for reasoning, so it is not a Lack of Visual Reasoning. However, the model made mistakes in the orange part during reasoning.

Figure 15: Error Case: Reasoning Error.

Correct Case









Question: The entire video is viewed through a black circle. What is this?

Options:

- (A) It's a circular hole in a piece of paper the viewer is looking through.
- (B) It's a decorative frame added for aesthetic purposes.
- (C) It is the scope of a gun.
- (1) It is the view through a telescope.
- (E) It is a special pair of observation glasses used to observe the behavior of wild animals.
- (F) It is an astronomical telescope.
- (G) It is the lens of a telephoto DSLR camera.
- (H) It is the view through a kaleidoscope.
- (I) It's the lens of a vintage camera.
- (J) It's a magic hoop that gives special vision.

Video Type: Animation

Task Type L1: Explicit Reasoning. Task Type L2: Causal Reasoning Task Type L3": Backward Reasoning

claude-3-5-sonnet: Let me think about this step by step:\n\n1) Looking at the video, I notice the view is consistently shown through a circular frame with crosshairs in it.\n\n2) The crosshairs (reticle) are a key identifying feature - they have the characteristic \"+\" shape with a circle.\n\n3) Throughout the video, the view is used to target/aim at birds in the trees.\n\n4) The combination of the circular view, crosshairs, and targeting function strongly indicates this is a scope used for aiming.\n\n5) Among the options, C mentions \"the scope of a gun\" - this matches perfectly with all the visual evidence.\n\n6) While other options like telescopes or cameras might have circular views, they typically don't include targeting crosshairs like we see here.\n\n7) The clear targeting/hunting context and reticle design confirms this is specifically a gun scope rather than any other optical device.\n\n[[C]]

Human Analysis: This is a relatively good answer. The model's thinking process before answering the question is mainly centered around the visual information of the video, rather than the entire CoT being spent analyzing the textual reasoning of the options.

Figure 16: MMR-V Construction Pipeline.

58 NeurIPS Paper Checklist

766

767

768 769

784

785

786

787

788

789

790

791

793

794

795

796

797

798

799

800

801

802

803

804

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and follow the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer [Yes], [No], or [NA].
- [NA] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

The checklist answers are an integral part of your paper submission. They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a proper justification is given (e.g., "error bars are not reported because it would be too computationally expensive" or "we were unable to find the license for the dataset we used"). In general, answering "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification please point to the section(s) where related material for the question can be found.

783 IMPORTANT, please:

- Delete this instruction block, but keep the section heading "NeurIPS Paper Checklist",
- · Keep the checklist subsection headings, questions/answers and guidelines below.
- Do not modify the questions and only use the provided macros for your answers.

1. Claims

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

Answer: [Yes]

Justification: In this paper, we propose MMR-V, a multimodal deep reasoning benchmark in videos, and conduct extensive experiments and analysis on various models.

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?

805 Answer: [Yes]

Justification: See Appendix A

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 work is entirely empirical and application-driven, without introducing any theoretical results, assumptions, or formal proofs.

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: The paper provides a comprehensive overview of the key details involved in reproducing the main experiments, as presented in Appendix E, G and Section 4.1. All code and data used in the paper have also been publicly released.

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: We provide open access to all the data and code. For detailed information, please refer to: https://huggingface.co/datasets/JokerJan/MMR-VBench, https://github.com/GaryStack/MMR-V

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]

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

945

946

947

948

949

950

951

952

953

954 955

956

958 959

960

961

962

963

964

965

Justification: Please see Appendix E, F, G and Section 4.1

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: [No]

Justification: Given the high computational cost associated with evaluating models on the video benchmark, we were unable to conduct error analysis across multiple random seeds.

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 Appendix E

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: We have read the NeurIPS Code of Ethics and confirm that our research fully complies with its principles.

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: During the video collection process for MMR-V, we strictly adhered to community guidelines of YouTube, avoiding content with negative social impact or restricted material. Furthermore, all dataset content were rechecked during the annotation stage to ensure compliance and appropriateness.

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)?

1020 Answer: [NA]

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1039

1040

1041

1042 1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1058

1059

1060

1062

1063 1064

1065

1066

1067

1068

1069

1070

Justification: This work does not involve the release of high-risk models or potentially unsafe datasets.

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: We properly cited all assets and adhered to their licenses and terms of use.

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: We release new assets with proper documentation and usage details.

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: [Yes]

Justification: See Appendix B

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: [Yes]

Justification: MMR-V data was collected and annotated by internal annotators, and all necessary consent and disclosures were obtained.

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 rigorousness, or originality of the research, declaration is not required.

Answer: [Yes]

Justification: This study evaluates the performance of MLLM on video reasoning tasks, so MLLM is widely used for evaluation. Secondly, MLLM is also used for auxiliary annotation and analysis of model CoT.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.