

---

# How Good is my Video LMM? Complex Video Reasoning and Robustness Evaluation Suite for Video-LMMs

---

Muhammad Uzair Khattak<sup>1</sup> Muhammad Ferjad Naeem<sup>2</sup> Jameel Hassan<sup>1</sup>  
Muzammal Naseer<sup>1</sup> Federico Tombari<sup>3,4</sup> Fahad Shahbaz Khan<sup>1,5</sup> Salman Khan<sup>1,6</sup>

<sup>1</sup>Mohamed Bin Zayed University of AI <sup>2</sup>ETH Zurich <sup>3</sup>Google  
<sup>4</sup>TU Munich <sup>5</sup>Linköping University <sup>6</sup>Australian National University

## Abstract

1 Recent advancements in Large Language Models (LLMs) have led to the develop-  
2 ment of Video Large Multi-modal Models (Video-LMMs) that can handle a wide  
3 range of video understanding tasks. These models have the potential to be deployed  
4 in real-world applications such as robotics, AI assistants, medical surgery, and  
5 autonomous vehicles. The widespread adoption of Video-LMMs in our daily lives  
6 underscores the importance of ensuring and evaluating their robust performance  
7 in mirroring human-like reasoning and interaction capabilities in complex, real-  
8 world contexts. However, existing benchmarks for Video-LMMs primarily focus  
9 on general video comprehension abilities and neglect assessing their reasoning  
10 capabilities over complex videos in the real-world context, and robustness of these  
11 models through the lens of user prompts as text queries. In this paper, we present  
12 the Complex Video Reasoning and Robustness Evaluation Suite (CVRR-ES), a  
13 novel benchmark that comprehensively assesses the performance of Video-LMMs  
14 across 11 diverse real-world video dimensions. We evaluate 11 recent models,  
15 including both open-source and closed-source variants, and find that most of the  
16 Video-LMMs, especially open-source ones, struggle with robustness and reasoning  
17 when dealing with complex videos. Based on our analysis, we develop a training-  
18 free Dual-Step Contextual Prompting (DSCP) technique to effectively enhance  
19 the performance of existing Video-LMMs on CVRR-ES benchmark. Our findings  
20 provide valuable insights for building the next generation of human-centric AI  
21 systems with advanced robustness and reasoning capabilities. Our dataset and code  
22 are publicly available at: [mbzuai-oryx.github.io/CVRR-Evaluation-Suite/](https://mbzuai-oryx.github.io/CVRR-Evaluation-Suite/).

## 23 1 Introduction

24 Recently, Large Language Models (LLMs) [30, 38, 12] have demonstrated impressive reasoning and  
25 planning capabilities while simultaneously handling a wide range of NLP tasks [33, 2]. Consequently,  
26 their integration with the vision modality, specifically for video understanding tasks, has given rise  
27 to Video Large Multi-modal Models (Video-LMMs) [15]. These models act as visual chatbots that  
28 accept both text and video as input and handle a diverse set of tasks, including video comprehension  
29 [21], detailed video understanding [18], and action grounding [37]. As these models directly capture  
30 video data, they hold substantial potential for deployment in real-world applications such as robotics,  
31 surveillance, medical surgery, and autonomous vehicles.

32 However, as these models assume an expanding role in our everyday lives, assessing their performance  
33 in comprehending complex videos and demonstrating reliable reasoning and robustness capabilities  
34 across diverse real-world contexts becomes essential. Video-LMMs with such capabilities will be

Benchmark	Textual Robustness	Complex Reasoning	In the wild (OOD)	Contextual Dependency	Multiple Actions	Temporal Order & Fine-grained
MSVD-QA [35]	✗	✗	✗	✗	✗	✗
MSRVTT-QA [35]	✗	✗	✗	✗	✗	✗
TGIF-QA [11]	✗	✓	✗	✗	✓	✓
Activity Net-QA [36]	✗	✗	✗	✗	✗	✓
VideoChat-GPT [21]	✗	✗	✗	✓	✓	✓
MVBench [16]	✗	✓	✗	✗	✓	✓
SEED-Bench [14]	✗	✗	✗	✗	✓	✓
CVRR-ES (ours)	✓	✓	✓	✓	✓	✓

Table 1: Comparison of CVRR-ES with existing benchmarks for video question answering. The CVRR-ES benchmark represents an initial effort to assess Video-LMMs in the context of their applicability and suitability in real-world contexts.

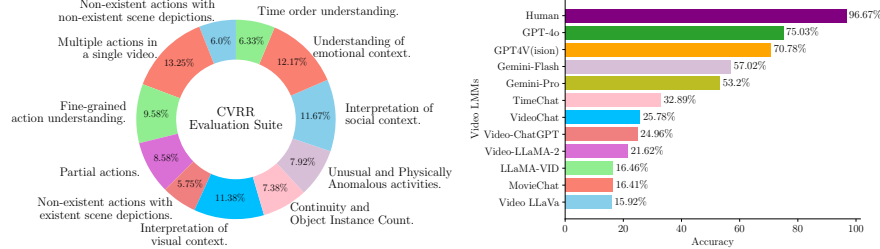


Figure 1: **Left:** CVRR-ES comprises of 11 diverse complex video evaluation dimensions encompassing a variety of complex, real-world contexts. **Right:** Overall performance of Video-LMMs on the CVRR-ES benchmark. Results for each Video-LMM are averaged across 11 video dimensions.

35 more effective when integrated into our daily lives for solving perception tasks and will be a promising  
 36 step towards building trustworthy human-centric AI-assistive systems.

37 Several attempts in literature have been made to benchmark Video-LMMs. SEED-Bench [14] curated  
 38 a MCQ-based dataset including 3 evaluation dimensions for videos. Similarly, MV-Bench [16]  
 39 constructed the Video-LMM benchmark and assembled 20 video tasks for evaluating the spatial and  
 40 temporal understanding of these models. While these methods aim at benchmarking Video-LMMs,  
 41 they predominantly evaluate video and/or temporal comprehension abilities and overlook the complex  
 42 reasoning aspects of Video-LMMs for real-world context, and their robustness towards user input text  
 43 queries; both of which are crucial to ensure their responsible engagement with humans in various real-  
 44 world situations in the wild. While some studies have explored similar areas such as hallucinations in  
 45 image-based LLMs [19, 24], no such comprehensive study exists for the case of Video-LMMs.

46 Motivated by the wide-scale applications of Video-LMMs and the lack of world-centric complex  
 47 video benchmarking efforts, we present a new benchmark, Complex Video Reasoning and Robustness  
 48 Evaluation Suite (CVRR-ES), to comprehensively assess the performance of Video-LMMs. As  
 49 shown in Tab. 1, CVRR-ES evaluates Video-LMMs on key aspects of robustness and reasoning in  
 50 videos, encompassing video domains that more accurately test models in real-world scenarios such as  
 51 videos having contextual dependency and in-the-wild aspects. CVRR-ES is an open-ended video QA  
 52 benchmark comprising 11 real-world video category dimensions (Fig. 1, left) that encompass diverse  
 53 evaluation aspects. These dimensions span from context-dependent (e.g., social, emotional, etc.)  
 54 categories to ones that often take place in the wild such as videos containing physically anomalous  
 55 activities. We comprehensively evaluate a representative set of 11 recent Video-LMMs (Fig. 1,  
 56 right) including both open-source and closed-source models on the CVRR-ES benchmark using a  
 57 LLM-assisted automatic evaluation framework [21, 4].

58 The performance of Video-LMMs on the CVRR-ES benchmark reveals that these models struggle to  
 59 correctly comprehend complex videos indicating their weak reasoning and lack of robustness to the  
 60 textual user queries (Fig. 2). For instance, state-of-the-art Video-LLaVA [18] achieves only 15.92%  
 61 performance averaged across 11 video dimensions of CVRR-ES. In contrast, closed-source models  
 62 including GPT4V(vision) [23] and Gemini-Vision-Pro [9] exhibit relatively stronger performance but  
 63 still lag behind the performance of humans. Using CVRR-ES benchmark, we extensively perform  
 64 quantitative and qualitative analysis and formulate important insights about these Video-LMMs based  
 65 on their failure cases and individual performances across the diverse video dimensions.

66 Based on our analysis, we note that standard prompting struggles in steering Video-LMMs’ focus for  
 67 complex video understanding. Additionally, their limitations in reasoning and robust video understand-  
 68 ing of real-world scenarios are dominantly driven by the quality of textual inputs (i.e., user questions).  
 69 Based on these insights, we develop a training-free Dual-Step Contextual Prompting (DSCP) tech-  
 70 nique, which effectively steers the model’s behavior during inference to elicit video-specific reasoning  
 71 and improved robustness within Video-LMMs. With DSCP, Video-LMMs substantially improve on  
 72 our benchmark, suggesting the potential of prompting methods for Video-LMMs.



Figure 2: We observe that most Video-LMMs struggle to reason over complex videos (rows 1-3) and exhibit weak robustness and rectification abilities when answering user questions that can sometimes be confusing (row 4). The QA pairs in Comprehensive Video Reasoning and Robustness Evaluation Suite (CVRR-ES) benchmark assess the performance of Video-LMMs beyond general video comprehension. (best viewed zoomed in)

73 Our main contributions are as follows: (1) We present Complex Video Robustness and Reasoning Evaluation suite (CVRR-ES), a Video Question Answering benchmark designed to assess the reasoning and robustness capabilities of Video-LMMs on 11 diverse world-centric complex video dimensions (§3). (2) We extensively evaluate both open-source and closed-source Video-LMMs on the CVRR-ES benchmark and find that most models exhibit weak performance, highlighting their limited reasoning in complex videos and lack of robustness towards user text queries (§5.1). (3) We conduct comprehensive analysis and formulate important conclusions about Video-LMMs based on their failure cases and performance on the CVRR-ES benchmark. Our findings provide key insights for building the next generation of human-centric AI systems with improved robustness and reasoning capabilities (§5.4). (4) To improve Video-LMMs’ reasoning and robustness abilities, we design a model-agnostic and training-free prompting method that effectively enhances their performance (§4).

## 84 2 Related Works

85 **Video Large Multi-modal models (Video-LMMs).** Video-LMMs [18, 17, 37] are visual chatbots capable of performing a wide range of video tasks, including video comprehension and captioning, 86 video question-answering, and action grounding. These models accept both video and textual inputs 87 and generate textual responses. From an architectural perspective, Video-LMMs combine pre-trained 88 vision backbones [25, 6, 32] with large language models [30, 38] using connector modules such 89 as MLP adapters, Q-former [5], and gated attention [1]. VideoChat [15] and VideoChat-GPT [17] 90 presented initial open-source efforts in this direction and were trained with two stages of alignment 91 and video-instruction following objectives. Recently, more advanced Video-LMMs have emerged in 92 the field, with some models focusing on improving model architectures [17], expanding to new tasks 93

94 [22], and enabling support for long videos [28, 26]. In this work, we aim to develop a comprehensive  
95 benchmarking framework to assess the reasoning and robustness capabilities of these Video-LMMs  
96 and develop a training-free prompting technique to improve their performance on these fronts.

97 **Benchmarking Video-LMMs.** With the growing number of Video-LMMs emerging in the research  
98 community, several works have presented evaluation frameworks to assess and quantify these models  
99 for benchmarking and analysis purposes. SEED-Bench [14] evaluates the visual capabilities in  
100 both image and Video-LMMs across 12 unique dimensions. MV-Bench [16] curates 20 video  
101 tasks to evaluate the spatial and temporal understanding of Video-LMMs. Video-ChatGPT [21]  
102 develops a quantitative evaluation framework to assess model understanding on five aspects of general  
103 video comprehension, such as the correctness and consistency of model captions. While these  
104 evaluation frameworks provide effective insights, their assessments do not extend beyond general  
105 video-comprehension metrics to more advanced aspects of reasoning and robustness, particularly for  
106 real-world context cases. In contrast, our work focuses on providing a complex video reasoning and  
107 robustness benchmark and offers a thorough assessment of Video-LMMs in practical applications.

108 **Training-free Prompting Techniques.** Steering model behavior at inference time using prompting  
109 has become a common paradigm in the NLP domain. Prompting [34, 31] refers to the set of  
110 instructions given as a prefix to the language model to better align model responses with human intent  
111 without the need for task-specific fine-tuning. Prompting techniques can be as simple as a single  
112 sentence (e.g., "Let's think step by step") such as zero-shot chain of thought [34] prompting, to more  
113 detailed techniques such as combining chain-of-thought prompting with few-shot learning [2] and  
114 self-consistency chain of thought prompting [31]. Surprisingly, training-free prompting techniques  
115 for Video Large Multi-modal Models (Video-LMMs) have been minimally explored. In this work,  
116 we develop a dual-step prompting technique based on principled prompt instructions specifically  
117 designed to steer the model's behavior for improved reasoning and robustness over complex videos.

### 118 3 Complex Video Reasoning and Robustness Evaluation Suite

119 As Video-LMMs are touching new real-world applications, it is essential to ensure that they robustly  
120 handle the user inputs, comprehend the visual world, and exhibit human-like reasoning capabilities.  
121 In this work, our goal is to establish a comprehensive benchmark, Complex Video Reasoning and  
122 Robustness Evaluation Suite (CVRR-ES) to assess the *robustness* and *reasoning* capabilities of  
123 Video-LMMs over complex and contextual videos. We first provide an overview of CVRR-ES and  
124 then detail the video evaluation dimensions in Sec. 3.1. Subsequently, we discuss benchmark creation  
125 process in Sec. 3.2. We provide details on the human performance on CVRR-ES in Appendix C.

126 **Overview.** CVRR-ES encompasses evaluation dimensions that cover diverse video categories related  
127 to real-world scenarios, ranging from context-dependent (e.g., social, emotional) categories to video  
128 types that often take place in the wild (e.g., anomalous activities). Specifically, we have compiled 11  
129 video evaluation dimensions and curated 2,400 high-quality open-ended question-answer (QA) pairs,  
130 spanning 214 high-quality videos. The average video duration is 22.3 seconds, with maximum and  
131 minimum durations of 183 and 2 seconds, respectively. Fig. 2 shows some qualitative examples of  
132 collected videos for the CVRR-ES benchmark. Refer to Appendix C for additional statistical details.

#### 133 3.1 CVRR-ES Video Category definitions.

134 For curating the CVRR-ES benchmark, we carefully select 11 diverse benchmark evaluation cate-  
135 gories. As shown in Fig. 1 (left), these categories encompass a wide range of real-world complex and  
136 contextual video types. Below, we define each video evaluation dimension in detail.

137 **1) Multiple actions in a single video.** This category involves videos with 2-4 different human  
138 activities. We curate questions in this category to assess the model's ability to understand and reason  
139 about multiple actions and their interrelations in a single video.

140 **2) Fine-grained action understanding.** We collect videos that encompass fine-grained activities  
141 performed by humans, such as pushing, opening, closing, spreading, sitting, etc. This category tests  
142 the model's ability to interpret subtle and fine-grained actions through carefully crafted questions.

143 **3) Partial actions.** We observe that Video-LMMs produce content that is relevant to a video's context  
144 and likely to occur next. We collect videos with actions likely to be followed by other actions but not  
145 shown in the video e.g., cracking an egg in a kitchen suggests the next action of cooking the egg.

146 **4) Time order understanding.** Accurately recognizing the temporal sequence of activities in videos



147 is crucial for distinguishing between atomic actions, such as pushing and pulling. We collect videos  
148 of fine-grained actions occurring in a particular temporal direction and curate challenging questions.  
149 **5) Non-existent actions with existent scene depictions.** This category examines the model’s robust-  
150 ness and reasoning behavior in scenarios where we introduce non-existent activities into the video  
151 without altering the physical and spatial scenes or environmental details in it.  
152 **6) Non-existent actions with non-existent scene depictions.** In this category, we increase the  
153 difficulty of the QA task by including questions containing both non-existent activities and scenes.  
154 We alter the details of objects, attributes, and background for non-existent scene comprehension. This  
155 tests the model’s ability to correct misleading questions and avoid generating imaginary content.  
156 **7) Continuity and object instance count.** This category contains videos (real-world and simulations)  
157 designed to test the models’ ability to accurately recognize the number of instances of objects, people,  
158 etc., and distinguish between existing objects and new ones introduced later in the same video scene.  
159 **8) Unusual and physically anomalous activities.** We collect videos depicting unusual actions that  
160 seemingly defy the laws of physics, such as a person floating in the air or driving a motorbike on  
161 a running river. Assessing Video-LMMs in such scenarios is crucial, as it allows us to determine  
162 whether they can generalize to understand actions in out-of-distribution videos in practical situations.  
163 **9) Interpretation of social context.** We test Video-LMMs’ ability to understand actions influenced  
164 by social contexts, such as helping an elderly person cross the road. Video-LMMs are assessed to  
165 determine their ability to accurately infer the rationale behind actions using the social context.  
166 **10) Understanding of emotional context.** Similar to social context, humans can accurately under-  
167 stand and interpret each other’s actions by considering the emotional context. We test Video-LMMs’  
168 ability to understand actions based on emotional context, e.g., a person crying due to joy.  
169 **11) Interpretation of visual context.** This category tests the model’s ability to understand actions by  
170 leveraging the overall visual contextual cues in the video. For example, to identify the number of  
171 people present based on the presence of shadows, one must utilize the visual context of shadows.

### 172 **3.2 Building CVRR-ES Benchmark**

173 **Stage 1: Data collection and Annotation.** We first collect high-quality videos and annotate each  
174 video via human assistance. To ensure that each evaluation dimension captures relevant attributes  
175 and information, we meticulously select videos that are representative of specific characteristics  
176 associated with that dimension. Overall, 214 unique videos are selected covering 11 dimensions  
177 with around 20 videos per evaluation dimension. Around 60% of these videos are collected from  
178 public academic datasets. To introduce diversity in the benchmark distribution, we select videos from  
179 multiple datasets including Something-Something-v2 [10], CATER [8], Charades [27], ActivityNet  
180 [3], HMDB51 [13], YFCC100M [29]. The remaining 40% of videos are collected from the internet.  
181 Following the video collection process, two experienced human annotators are assigned to generate  
182 captions for each video. For videos where initial captions or metadata are available from academic  
183 datasets, the captions are generated by the annotators based on them. For videos collected from the  
184 internet, captions are entirely generated by human annotators. To ensure consistency and high quality,  
185 we provide annotation instructions to annotators, who generate captions accordingly. Personalized  
186 annotation guidelines are used for each video category. Refer to additional details in Appendix C.

187 **Stage 2: Question-Answer Generation.** The first challenge is to select an evaluation setting to assess  
188 Video-LMMs. Humans typically engage in free-form conversation to interact with each other in  
189 day-to-day life. Inspired by this, we aim to simulate a similar style of interaction with Video-LMMs  
190 by curating open-ended QA pairs to evaluate these models for robustness and reasoning. We feed  
191 detailed ground-truth video captions to GPT-3.5 LLM, which is utilized to generate open-ended  
192 questions. The QA pairs covers both the reasoning and robustness aspects as detailed below.

193 **Reasoning QA pairs:** With Video-LMMs beginning to interact more directly with humans in our  
194 lives, it’s crucial to validate the reasoning abilities of Video-LMMs for more reliable Human-AI  
195 interaction. When evaluating the reasoning capabilities of Video-LMMs, we aim to determine whether  
196 these models can understand the input video not only by analyzing spatial content but also by grasping  
197 the underlying rationale behind the occurring activities and their relationships with the surrounding  
198 context. This involves creating questions that go beyond simple video comprehension and scene  
199 description and require the model to engage in complex logical inference, contextual understanding,  
200 and reasoning about counterfactual and hypothetical scenarios.

201 **Robustness QA pairs:** In addition to evaluating the reasoning capabilities of LLMs, it is important  
202 to assess Video-LMMs to ensure their robust and responsible performance in real-world scenarios.  
203 In the context of Video-LMMs, robustness can be evaluated from both visual (video input) and  
204 textual interfaces. Our focus in this work lies on textual interface robustness by particularly testing  
205 the model’s comprehension abilities when posed with misleading or confusing questions. This  
206 scenario mirrors realistic situations where users, based on their expertise levels, may pose irrelevant,  
207 misleading, or confusing questions. It is crucial for models to demonstrate reliability and robustness  
208 in handling such queries and avoid generating unreal or hallucinated content for input videos.  
209 We curate specific prompts for each evaluation dimension to instruct LLM in generating QA pairs.  
210 Example prompts used as an instruction to LLMs for curating QA pairs for robustness and reasoning  
211 aspects are provided in Fig. 14 in the Appendix E.

212 **Stage 3: QA Pairs Filtration.** After generating the QA pairs, we employ a manual filtration step,  
213 with human assistance to verify each generated QA pair. Approximately 30% of the QA pairs  
214 generated by GPT-3.5 are found to be noisy, containing questions that are unrelated to the video  
215 evaluation dimensions or unanswerable based on the provided ground-truth captions. Additionally,  
216 many questions contain answers within the question itself. Therefore, an exhaustive filtering process  
217 is conducted which involves QA rectification and removing those samples which are not relevant to  
218 the video or evaluation type. This process results in a final set of 2400 high-quality QA pairs for the  
219 CVRR-ES benchmark. Examples of the final QA pairs are shown in Tab. 4 in the Appendix.

220 **Stage 4: Evaluation Procedure.** Previous methods in the literature [21, 4, 19, 24] have explored  
221 using LLM models as judges for quantifying results in open-ended QA benchmarks. We adopt a  
222 similar approach and instruct LLMs to act as teachers to assess the correctness of predicted responses  
223 from Video-LMMs compared to ground-truths. We generate open-ended predictions from Video-  
224 LMMs by providing video-question pairs as inputs and then present the model predictions and their  
225 ground-truth responses to the LLM Judge using the evaluation prompt. The Judge determines whether  
226 the prediction is correct or incorrect with a binary judgment, assigns a score from 1 to 5 representing  
227 the quality of the prediction, and provides a reasoning to explain its decision. Our ablative analysis in  
228 the Appendix. E demonstrates that reasoning-constrained LLM-based evaluation aligns the most with  
229 human-based judgment. Our evaluation prompt for LLM Judge is shown in Fig. 13 in Appendix E.

230 **Quality of QA pairs.** We show examples of QA pairs from CVRR-ES benchmark in Table 4 in  
231 Appendix C. Our QA pairs are of high quality and aim to test the understanding of Video-LMMs  
232 against reasoning and robustness criteria on multiple evaluation dimensions. To quantitatively assess  
233 the quality of the benchmark, we establish a quality assessment procedure [7]. We randomly sample  
234 1120 QA pairs, which encompass all videos of the CVRR-ES benchmark, and request human experts  
235 to evaluate the quality of each QA pair by answering the following questions: (1) *"Does the QA pair  
236 correctly represent the evaluation dimension category under which it falls?"* (possible answers: "Yes",  
237 "No") (2) *"Can the question be correctly answered given only the video content?"* (possible answers:  
238 "Agree", "Disagree") and (3) *"Is the corresponding paired ground-truth answer correct? (which will  
239 be used during evaluation as ground truth)"* (possible answers: "Yes", "No"). On average, the answer  
240 of experts for the first question was "Yes" for 98.84% of the times. For the second and third questions,  
241 the averaged answer was "Agree" and "Yes" for 100% and 99.91% of the times, respectively.

## 242 4 Dual-Step Contextual Prompting for Video-LMMs.

243 Given their wide-scale potential in practical applications, new Video-LMMs are frequently introduced  
244 by the research community. Despite the availability of numerous Video-LMMs, the majority of them  
245 are trained using only positive examples and video-conversational templates that are primarily limited  
246 to tasks such as video-captioning and video question answering [15, 21, 26, 28]. This leads to highly  
247 over-affirmative behavior and a lack of self-rectification abilities in these models (Sec. 5.4).

248 Additionally, the templates have minimal focus on enhancing reasoning and robustness capabilities  
249 through reasoning instruction-tuning pairs, resulting in their weak performance against robustness  
250 and reasoning based evaluations in CVRR-ES. Consequently, enabling direct interaction of Video-  
251 LMMs with users in real-world scenarios can result in undesired responses when the user question is  
252 confusing and deceiving. Moreover, curating reasoning-based instruction fine-tuning datasets requires  
253 meticulous data curation steps, and retraining these models are computationally expensive [17, 26].

254 Alternatively, training-free prompting techniques in NLP literature have shown effectiveness in  
 255 eliciting reasoning abilities in LLMs such as chain of thought and self-consistency prompting [34, 31].  
 256 Inspired by these, we present a Dual Step Contextual Prompting (DSCP) technique, which steers  
 257 Video-LMM focus for enhanced reasoning while simultaneously encouraging the models to provide  
 258 robust and grounded answers. DSCP is a two-step prompting method that **1)** ensures that the model  
 259 comprehends the video while reasoning over crucial aspects of complex video understanding such as  
 260 contextual information and decoding the complex relationships between objects and motions, etc., and  
 261 **2)** encourages robustness by generating the response against the question while conditioning both on  
 262 video and the unbiased context retrieved in the first step. Below we discuss each step of DSCP in detail.

263 **Step 1: Video reasoning.** We prompt Video-LMMs to  
 264 interpret video from a reasoning perspective using ten  
 265 principled instructions (Fig. 3, in blue) to direct the mod-  
 266 els to understand general video content, reason over the  
 267 rationale behind actions and their relationships with the  
 268 context, and consider factors like contextual priors, the  
 269 temporal order of actions, instance count, and attributes.  
 270 The prompting technique also includes instructions to  
 271 ensure conciseness and factuality to mitigate hallucina-  
 272 tions. Given a Video-LMM  $\mathcal{F}$  and input video  $\mathcal{V}$ , we  
 273 retrieve contextual reasoning information  $I_{\text{context}}$  by pro-  
 274 viding principled reasoning prompt  $P_{\text{reason}}$  along with  
 275 the video to the LMM,  $I_{\text{context}} = \mathcal{F}(P_{\text{reason}}|\mathcal{V})$ . This  
 276 contextual information is then used in the second step of  
 277 DSCP to generate a grounded response to user question.

278 **Step 2: Context conditioned question answering.** To address the challenges of over-affirmative  
 279 behavior and hallucinations in Video-LMMs when prompted with confusing or misleading questions,  
 280 we propose an additional inference step. We note that Video-LMMs often possess factual knowledge  
 281 about the video content but become distracted and hallucinate when prompted with confusing or  
 282 misleading questions (Appendix D). Our DSCP technique conditions the model to first comprehend  
 283 the video without attending to the user question and, therefore eliminates its influence. This complex  
 284 video comprehension information,  $I_{\text{context}}$  (formulated in step 1) is then used to condition the model  
 285 on both the video and  $I_{\text{context}}$ . Finally, we pose the user question using prompt  $P_{\text{user}}$  which combines  
 286 the user question and the contextual reasoning information (Fig. 3, in green). The final response is  
 287  $\mathcal{F}(P_{\text{user}}|\mathcal{V})$ , where  $P_{\text{user}} = [\text{question}; I_{\text{context}}]$ . Here  $[\ ; ]$  denotes the text prompt concatenation.  
 288 The factual content generated in step 1 guides the model towards a robust response in step 2, pro-  
 289 ducing factual and correct responses even with noisy or misleading user questions. We show the  
 290 qualitative results of DSCP technique in Fig. 11 in Appendix D. This approach leads to responses  
 291 that are better grounded in the actual video content and are robust against lower-quality user queries.  
 292 The DSCP technique effectively enhances the performance of Video-LMMs on CVRR-ES (Sec. 5.2).

## 293 5 Evaluation Experiments on CVRR-ES.

294 **Video-LMMs.** Among the open-source models, we evaluate 7 recent Video-LMMs, including  
 295 Video-LLaVA [18], TimeChat [26], MovieChat [28], LLaMA-ViD [17], VideoChat [15] Video-  
 296 ChatGPT [21], and Video-LLaMA-2 [37]. For evaluating closed-source models, we use Gemini-Pro,  
 297 Gemini-Flash, [9], GPT-4V and recent GPT-4o [23]. Refer to Appendix B for additional details.

### 298 5.1 Main Experiments on CVRR-ES.

299 Tab. 2 shows the evaluation results of Video-LMMs on CVRR-ES. Below, we discuss main results.  
 300 **Open Source Video-LMMs struggles on CVRR-ES benchmark.** All open-source LMMs show in-  
 301 ferior performance across the different evaluation dimensions of CVRR-ES. Interestingly, some of the  
 302 earlier developed open-source Video-LMMs, like Video-LLaMA, VideoChat, and Video-ChatGPT,  
 303 exhibit higher performance compared to more recent models such as Video-LLaVA, MovieChat, and  
 304 LLaMA-ViD. Overall, TimeChat achieves the highest performance of 32.89% averaged across the 11  
 305 evaluation dimensions among open-source LMMs, followed by VideoChat with a score of 25.78%.  
 306 **Humans rank highest in CVRR-ES benchmark.** Human evaluation achieves the highest perfor-

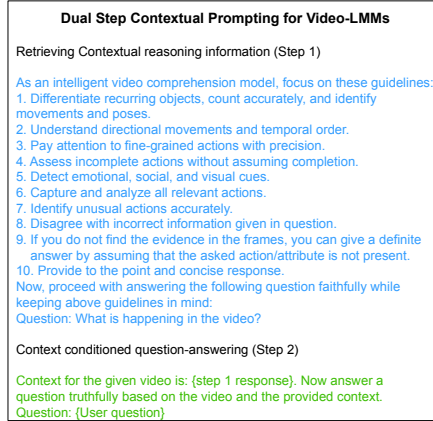


Figure 3: Principled prompt instructions in DSCP for Video-LMMs.

Table 2: Evaluation results of Video LLMs across various video-evaluation categories on the CVRR-ES benchmark. We present results for both open-source and closed-source models and human evaluation.

Benchmark Category	Video-LLaMA-2	VideoChat	Video-ChatGPT	Video-LLaVA	MovieChat	LLaMA-VID	TimeChat	Gemini-Y Pro	Gemini-V Flash	GPT4V	GPT4o	Human
Multiple Actions in single video.	16.98	23.90	27.67	15.72	12.58	17.92	28.30	43.08	44.65	57.55	62.89	<b>93.40</b>
Fine-grained action understanding.	29.57	33.48	26.96	25.22	23.48	26.09	39.13	51.61	64.78	77.39	80.43	<b>95.65</b>
Partial actions.	24.76	33.01	22.82	13.59	21.36	14.56	49.51	67.48	62.14	73.79	77.67	<b>98.54</b>
Time order understanding.	16.45	31.58	27.63	21.05	16.45	19.74	34.21	45.39	55.26	57.89	71.05	<b>97.37</b>
Non-existent actions with existent scene.	10.14	15.22	23.19	5.07	5.07	2.90	23.19	57.25	60.14	71.01	83.33	<b>97.10</b>
Non-existent actions with non-existent scene.	13.19	14.58	17.36	3.47	11.81	6.94	13.89	49.64	56.30	75.00	70.14	<b>100.00</b>
Continuity and Object instance Count.	28.25	24.29	28.41	21.47	19.77	24.86	34.46	36.16	43.50	62.71	62.71	<b>96.49</b>
Unusual and Physically Anomalous activities.	18.95	18.42	18.95	15.79	17.89	16.32	27.37	60.00	60.53	74.74	78.42	<b>96.84</b>
Interpretation of social context.	25.00	31.07	32.50	18.93	17.14	13.93	39.29	64.29	69.64	79.64	83.57	<b>97.51</b>
Understanding of emotional context.	21.92	23.63	21.23	15.07	13.70	14.73	27.40	47.26	52.74	66.44	70.89	<b>95.55</b>
Interpretation of visual context.	32.60	34.43	27.84	19.78	21.25	23.08	45.05	63.00	57.51	82.42	84.25	<b>94.87</b>
<b>Average</b>	<b>21.62</b>	<b>25.78</b>	<b>24.96</b>	<b>15.92</b>	<b>16.41</b>	<b>16.46</b>	<b>32.89</b>	<b>53.20</b>	<b>57.02</b>	<b>70.78</b>	<b>75.03</b>	<b>96.67</b>

Prompting Method	VideoChat	Video-LLaVA	MovieChat	LLaMA-VID	TimeChat
Standard prompting	25.78	15.92	16.41	16.46	32.89
Chain of Thought (CoT) prompting	22.44	25.87	15.89	29.68	<b>39.57</b>
DSCP (Stage 1)	38.07	32.12	28.05	25.13	33.04
DSCP (Both stages)	<b>47.92</b>	<b>37.93</b>	<b>35.87</b>	<b>46.85</b>	39.45

Table 3: **Prompting methods.** DSCP stage 1 uses only principled instructions of step 1 and DSCP (Both stages) uses complete dual-step technique.

mance on the CVRR-ES benchmark, with over 95% accuracy across all evaluation dimensions. These results suggest that the CVRR-ES QA pairs are reasonable and suitable for benchmarking.

**Closed source models perform competitively on CVRR-ES.** As shown in Tab. 2, both Gemini and GPT variants improve over open-source models and achieve high gains across all evaluation dimensions. The competitive results of GPT4o and Gemini-Flash on complex video evaluation dimensions such as partial actions, non-existent action/scene depiction, and context-dependent categories show that these models have a more sophisticated understanding of the complex visual contents of videos and have strong capabilities to rectify misleading and confusing user questions. Overall, GPT4o improves over Gemini-Flash by 18.01% and provides the highest average accuracy of 75.03%.

## 5.2 Effectiveness of DSCP method for improving Video-LMMs performance

We next integrate DSCP technique with Video-LMMs and present results for CVRR-ES in Fig. 4. DSCP improves the model’s performance compared with models that use standard prompting (i.e., using only the question itself). These results also suggest that prompting techniques in Video-LMMs can better guide models for improved reasoning and robustness. With DSCP, initially low-performing Video-LMMs like Video-LLaVa, MovieChat, and LLaMA-Vid show much better relative gains and become competitive with other models. The highest relative gain of 184% is achieved by LLaMA-ViD, which moves from 7th place in the leaderboard to 2nd among the open-source models after using the DSCP technique. We observe similar overall positive trends of using DSCP with closed-source model Gemini, which improves on the benchmark by an absolute overall gain of 5.02%. We provide more detailed results comparisons in Appendix D.

## 5.3 Different prompting techniques.

We now study the contribution of each step of DSCP and compare it with chain-of-thought (CoT) prompting [34]. Results for the top 5 performing open Video-LMMs are shown in Tab. 3. CoT prompting improves over standard prompting in 3 out of 5 Video-LMMs, suggesting that prompting

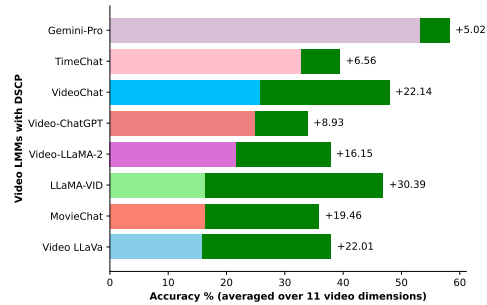


Figure 4: Video-LMMs with DSCP technique effectively improves their performance (gains are shown in green) on CVRR-ES benchmark.



337 techniques from NLP literature can also guide multi-modal Video-LMMs to enhance reasoning and  
338 robustness. Next, we ablate on the first step of DSCP prompting, which uses principled instructions  
339 of DSCP step 1 as a prefix alongside the actual user question. DSCP step 1 notably improves model  
340 performance on all Video-LMMs, suggesting the effectiveness of the principled prompt instructions  
341 designed specifically for Video models. DSCP with both steps, which additionally uses the initial  
342 context in the second step, shows additional gains and achieves highest results on 4 out of 5 models.

#### 343 5.4 Main findings and Qualitative Results

344 We now present key insights that can guide the development of the next generation of robust and  
345 reliable Video-LMMs. We show qualitative results and additional analysis in the Appendix A.

346 **Models excelling at standard VQA benchmarks struggle on CVRR-ES.** Our analysis in Sec.  
347 5.1 reveals that the latest open-source Video-LMMs, like Video-LLaVA, MovieChat, and LLaMA-  
348 VID, perform less effectively on CVRR-ES compared to Video-LMMs that were introduced earlier  
349 in the community, such as VideoChat and Video-ChatGPT. Interestingly, the same recent models  
350 demonstrate superior performance on general video comprehension benchmarks. This suggests  
351 that current VQA benchmarks, like ActivityNet-QA [36] and MSRVT [35], do not adequately  
352 correlate with the complex video reasoning and robustness scenarios highlighted in our benchmark.  
353 Consequently, this also indicates that most newer Video-LMMs are heavily trained to excel on the  
354 general video benchmarks while reducing their generalizability, reasoning, and robustness capabilities.  
355 **Over-affirmative behavior of open-source Video-LMMs.** We observe that open-source models  
356 exhibit positive and over-affirmative responses. Open-source Video-LMMs consistently respond with  
357 "Yes" even when faced with confusing questions that describe non-existent actions and objects (Fig.  
358 5 in Appendix A). This highlights the vulnerability of these models when interacting with users in  
359 real-world scenarios. In our CVRR-ES benchmark, open-source models are notably vulnerable to  
360 evaluation dimensions of "*Non-existent actions with the existent scene*" and "*Non-existent actions with  
361 the non-existent scene*" compared to closed models. These models lack negation and self-rectification  
362 capabilities, especially when users provide misleading or confusing questions. We conjecture that  
363 such behavior arises due to the absence of negative instruction tuning pairs during training.

364 **Tendency towards activity completion.** Most open-source Video-LMMs have shown lower results  
365 on the evaluation dimension of partial actions, which focuses on incomplete or atomic actions. We  
366 note that most open-source models tend to complete actions, even when only part of the action is  
367 provided in the video (Fig. 6 in Appendix A). Upon examining the fine-tuning strategies [21, 20], we  
368 find that almost all models are trained on end-to-end actions-based instruction-tuning data, causing  
369 them to generate complete action descriptions at inference. This tendency highlights the vulnerability  
370 of Video-LMMs after deployment, as real-world scenarios often involve atomic, sub-atomic, and  
371 general actions alike. To improve the performance of Video-LMMs, it is crucial to incorporate diverse  
372 action types during the training phase, including partial and incomplete actions.

373 **Video-LMMs struggles in understanding the emotional and social context.** For more reliable  
374 interaction with humans in practical scenarios, Video-LMMs models should comprehend the video  
375 scenes with social and contextual reasoning capabilities similar to humans. The lower performance of  
376 Video-LMMs on social and emotional contextual dimensions in CVRR-ES highlights their limitations  
377 and lack of understanding of scenes based on contextual cues (Fig. 9 in Appendix A).

## 378 6 Conclusion

380 Given the expanding role of Video-LMMs in practical world-centric applications, it is crucial to ensure  
381 that these models perform robustly and exhibit human-like reasoning and interaction capabilities  
382 across various complex and real-world contexts. In this work, we present the CVRR-ES benchmark for  
383 Video-LMMs, aiming to evaluate Video-LMMs on these very fronts. Through extensive evaluations,  
384 we find that Video-LMMs, especially open-source ones, exhibit limited robustness and reasoning  
385 capabilities over complex videos involving real-world contexts. Based on our analysis, we formulate  
386 a training-free prompting technique that effectively improves the performance of Video-LMMs across  
387 various evaluation dimensions of the CVRR-ES benchmark. Furthermore, we analyze and investigate  
388 the failure cases of Video-LMMs on the CVRR-ES benchmark and deduce several important findings.  
389 We hope that the CVRR-ES benchmark, accompanied by our extensive analysis, will contribute  
390 towards building the next generation of advanced world-centric video understanding models.

## References

- 391
- 392 [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson,  
393 Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual  
394 language model for few-shot learning. 2022. 3
- 395 [2] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
396 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
397 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. 1,  
398 4
- 399 [3] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet:  
400 A large-scale video benchmark for human activity understanding. In *Proceedings of the ieee  
401 conference on computer vision and pattern recognition*, pages 961–970, 2015. 5, 24
- 402 [4] Rizhao Cai, Zirui Song, Dayan Guan, Zhenhao Chen, Xing Luo, Chenyu Yi, and Alex Kot.  
403 Benchmm: Benchmarking cross-style visual capability of large multimodal models. *arXiv  
404 preprint arXiv:2312.02896*, 2023. 2, 6
- 405 [5] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng  
406 Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose  
407 vision-language models with instruction tuning. *arXiv:2305.06500*, 2023. 3
- 408 [6] Yuxin Fang, Quan Sun, Xinggang Wang, Tiejun Huang, Xinlong Wang, and Yue Cao. Eva-02:  
409 A visual representation for neon genesis. *arXiv:2303.11331*, 2023. 3
- 410 [7] Kanishk Gandhi, Jan-Philipp Fränken, Tobias Gerstenberg, and Noah Goodman. Understanding  
411 social reasoning in language models with language models. *Advances in Neural Information  
412 Processing Systems*, 36, 2024. 6
- 413 [8] Rohit Girdhar and Deva Ramanan. CATER: A diagnostic dataset for Compositional Actions  
414 and TEmporal Reasoning. In *ICLR, 2020*. 5, 24
- 415 [9] Google. Gemini, 2023. 2, 7
- 416 [10] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne  
417 Westphal, Heuna Kim, Valentin Haebel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag,  
418 et al. The "something something" video database for learning and evaluating visual common  
419 sense. In *ICCV, 2017*. 5, 24
- 420 [11] Yunseok Jang, Yale Song, Youngjae Yu, Youngjin Kim, and Gunhee Kim. Tgif-qa: Toward  
421 spatio-temporal reasoning in visual question answering. In *Proceedings of the IEEE conference  
422 on computer vision and pattern recognition*, pages 2758–2766, 2017. 2
- 423 [12] Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris  
424 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand,  
425 et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024. 1
- 426 [13] Hildegard Kuehne, Hueihan Jhuang, Estíbaliz Garrote, Tomaso Poggio, and Thomas Serre.  
427 Hmdb: a large video database for human motion recognition. In *2011 International conference  
428 on computer vision*, pages 2556–2563. IEEE, 2011. 5, 24
- 429 [14] Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-  
430 bench: Benchmarking multimodal llms with generative comprehension. *arXiv preprint  
431 arXiv:2307.16125*, 2023. 2, 4
- 432 [15] KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang,  
433 and Yu Qiao. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*,  
434 2023. 1, 3, 6, 7

- 435 [16] Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo  
436 Chen, Ping Luo, et al. Mvbench: A comprehensive multi-modal video understanding benchmark.  
437 *arXiv preprint arXiv:2311.17005*, 2023. 2, 4
- 438 [17] Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An image is worth 2 tokens in large  
439 language models. *arXiv preprint arXiv:2311.17043*, 2023. 3, 6, 7
- 440 [18] Bin Lin, Bin Zhu, Yang Ye, Munan Ning, Peng Jin, and Li Yuan. Video-llava: Learning united  
441 visual representation by alignment before projection. *arXiv preprint arXiv:2311.10122*, 2023.  
442 1, 2, 3, 7
- 443 [19] Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Aligning  
444 large multi-modal model with robust instruction tuning. *arXiv preprint arXiv:2306.14565*, 2023.  
445 2, 6
- 446 [20] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. 2023.  
447 9, 22
- 448 [21] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt:  
449 Towards detailed video understanding via large vision and language models. *arXiv preprint*  
450 *arXiv:2306.05424*, 2023. 1, 2, 4, 6, 7, 9, 22
- 451 [22] Shehan Munasinghe, Rusiru Thushara, Muhammad Maaz, Hanoona Abdul Rasheed, Salman  
452 Khan, Mubarak Shah, and Fahad Khan. Pg-video-llava: Pixel grounding large video-language  
453 models. *arXiv preprint arXiv:2311.13435*, 2023. 4
- 454 [23] OpenAI. GPT-4V(ision) System Card, 2023. 2, 7
- 455 [24] Yusu Qian, Haotian Zhang, Yinfei Yang, and Zhe Gan. How easy is it to fool your multimodal  
456 llms? an empirical analysis on deceptive prompts. *arXiv preprint arXiv:2402.13220*, 2024. 2, 6
- 457 [25] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
458 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
459 models from natural language supervision. 2021. 3
- 460 [26] Shuhuai Ren, Linli Yao, Shicheng Li, Xu Sun, and Lu Hou. Timechat: A time-sensitive multi-  
461 modal large language model for long video understanding. *arXiv preprint arXiv:2312.02051*,  
462 2023. 4, 6, 7
- 463 [27] Gunnar A Sigurdsson, Gül Varol, Xiaolong Wang, Ali Farhadi, Ivan Laptev, and Abhinav Gupta.  
464 Hollywood in homes: Crowdsourcing data collection for activity understanding. In *Computer*  
465 *Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14,*  
466 *2016, Proceedings, Part I 14*, pages 510–526. Springer, 2016. 5, 24
- 467 [28] Enxin Song, Wenhao Chai, Guan hong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu,  
468 Xun Guo, Tian Ye, Yan Lu, Jenq-Neng Hwang, et al. Moviechat: From dense token to sparse  
469 memory for long video understanding. *arXiv preprint arXiv:2307.16449*, 2023. 4, 6, 7, 22
- 470 [29] Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland,  
471 Damian Borth, and Li-Jia Li. Yfcc100m: The new data in multimedia research. *Communications*  
472 *of the ACM*, 59(2):64–73, 2016. 5, 24
- 473 [30] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timo-  
474 thée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez,  
475 Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation  
476 language models. *arXiv:2302.13971*, 2023. 1, 3
- 477 [31] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha  
478 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language  
479 models. *arXiv preprint arXiv:2203.11171*, 2022. 4, 7

- 480 [32] Yi Wang, Kunchang Li, Yizhuo Li, Yinan He, Bingkun Huang, Zhiyu Zhao, Hongjie Zhang,  
481 Jilan Xu, Yi Liu, Zun Wang, et al. Internvideo: General video foundation models via generative  
482 and discriminative learning. *arXiv preprint arXiv:2212.03191*, 2022. 3
- 483 [33] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani  
484 Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large  
485 language models. *arXiv preprint arXiv:2206.07682*, 2022. 1
- 486 [34] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le,  
487 Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models.  
488 *Advances in neural information processing systems*, 35:24824–24837, 2022. 4, 7, 8
- 489 [35] Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang.  
490 Video question answering via gradually refined attention over appearance and motion. In  
491 *Proceedings of the 25th ACM international conference on Multimedia*, pages 1645–1653, 2017.  
492 2, 9
- 493 [36] Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao.  
494 Activitynet-qa: A dataset for understanding complex web videos via question answering.  
495 In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 9127–9134,  
496 2019. 2, 9, 14
- 497 [37] Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language  
498 model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023. 1, 3, 7
- 499 [38] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
500 Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.  
501 Judging llm-as-a-judge with mt-bench and chatbot arena. *arXiv:2306.05685*, 2023. 1, 3



502 **Checklist**

503 1. For all authors...

504 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
505 contributions and scope? [Yes]

506 **Justification:** Yes, we have ensured that the main claims in the abstract and introduction  
507 accurately reflect the paper’s contributions and scope.

508 (b) Did you describe the limitations of your work? [Yes]

509 **Justification:** We have discussed the limitations of our work in the Appendix. **F**.

510 (c) Did you discuss any potential negative societal impacts of your work? [N/A]

511 **Justification:** This is a dataset paper aimed at studying and benchmarking the reasoning  
512 of Video-LMMs in real-world context and robustness from the lens of user text queries.  
513 Therefore, to the best of our knowledge, there are no potential negative societal impacts  
514 of our work.

515 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
516 them? [Yes]

517 **Justification:** Yes we have read the ethics review guidelines and ensured that our paper  
518 conforms to them.

519 2. If you are including theoretical results...

520 (a) Did you state the full set of assumptions of all theoretical results [N/A]

521 **Justification:** There is no theoretical result in this paper that requires a full set of  
522 assumptions and correct proof.

523 (b) Did you include complete proofs of all theoretical results? [N/A]

524 **Justification:** There is no theoretical result in this paper that requires a full set of  
525 assumptions and correct proof.

526 3. If you ran experiments (e.g. for benchmarks)...

527 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
528 mental results (either in the supplemental material or as a URL)? [Yes]

529 **Justification:** We have attached the code, link to data, and all instructions to reproduce  
530 the main experimental results in the supplemental material.

531 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
532 were chosen)? [Yes]

533 **Justification:** We have provided implementation details in the Appendix. **B**.

534 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
535 ments multiple times)? [No] .

536 **Justification:** We did not have enough compute resources to completely re-run all the  
537 experiments for different seeds and report error bars for different runs. We are currently  
538 re-running the error bar experiments, and we plan to include all the experiments with  
539 different seeds in the final version.

540 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
541 of GPUs, internal cluster, or cloud provider)? [Yes]

542 **Justification:** We have provided details on the compute resources in the Appendix. **B**.

543 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

544 (a) If your work uses existing assets, did you cite the creators? [Yes]

545 **Justification:** We have cited the creators of datasets used in our benchmark in the main  
546 paper in Sec. 3.2.

547 (b) Did you mention the license of the assets? [Yes]

548 **Justification:** Our dataset is released for educational and research purposes under  
549 the CC-BY-4.0 license. We have mentioned the license of assets in the files in our  
550 supplemental material as well as on our GitHub dataset hosting platform.

- 551 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
552 **Justification:** Yes we have included the assets in the supplemental material and also on  
553 the public URL. Our assets can be publically accessed at [mbzuai-oryx.github.io/CVRR-](https://mbzuai-oryx.github.io/CVRR-Evaluation-Suite/)  
554 [Evaluation-Suite/](https://mbzuai-oryx.github.io/CVRR-Evaluation-Suite/).
- 555 (d) Did you discuss whether and how consent was obtained from people whose data you're  
556 using/curating? [Yes]  
557 **Justification:** We collected most of the videos from academic datasets while respecting  
558 their license information. The videos obtained from the web from YouTube are subject  
559 to the copyright of the original owners and are used only for research and academic  
560 purposes, consistant with previous works and benchmarks such as ActivityNet [36] etc.
- 561 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
562 information or offensive content? [Yes]  
563 **Justification:** In our initial analysis using the subset (50%) of our CVRR-ES dataset, we  
564 noted that no video contained specific personally identifiable information or offensive  
565 content.
- 566 5. If you used crowdsourcing or conducted research with human subjects...
- 567 (a) Did you include the full text of instructions given to participants and screenshots, if  
568 applicable? [Yes]  
569 **Justification:** The instructions to humans for the benchmark quality assessment are  
570 provided in Appendix C.
- 571 (b) Did you describe any potential participant risks, with links to Institutional Review  
572 Board (IRB) approvals, if applicable? [N/A]  
573 **Justification:** Not applicable.
- 574 (c) Did you include the estimated hourly wage paid to participants and the total amount  
575 spent on participant compensation? [N/A]  
576 **Justification:** The annotation process was carried out by the authors of this manuscript.  
577 As a result, the aspect of compensation for human subjects does not apply in this case.