# 🚝 Mementos: A Comprehensive Benchmark for Multimodal Large Language Model Reasoning over Image Sequences

Anonymous ACL submission

#### Abstract

Multimodal Large Language Models (MLLMs) have demonstrated proficiency in handling a variety of visual-language tasks. However, current MLLM benchmarks are predominantly de-005 signed to evaluate reasoning based on static information about a single image, and the ability of modern MLLMs to extrapolate from image sequences, which is essential for understanding our ever-changing world, has been less investigated. To address this challenge, this paper introduces Mementos, a new benchmark designed to assess MLLMs' sequential image reasoning abilities. Mementos features 4,761 diverse image sequences with varying lengths. We also employ a GPT-4 assisted method to evaluate MLLM reasoning performance. Through a careful evaluation of nine 018 recent MLLMs on Mementos, including GPT-4V and Gemini, we find that they struggle to accurately describe dynamic information about given image sequences, often leading to hallucinations/misrepresentations of objects and their corresponding behaviors. Our quantitative analysis and case studies identify three key factors impacting MLLMs' sequential image reasoning: the correlation between object and behavioral hallucinations, the influence of cooccurring behaviors, and the compounding impact of behavioral hallucinations.

#### 1 Introduction

001

007

017

033

037

041

The recent emergence of Multimodal Large Language Models (MLLMs) such as GPT-4V (OpenAI, 2023b) and Gemini (Team, 2023) has shown strong visual-language understanding and generation capabilities in many areas, like image captioning and visual question answering. Despite the notable performance of existing MLLMs, they often suffer from hallucination (a phenomenon where MLLMs produce inaccurate descriptions of the given images) due to insufficient reasoning capabilities, generating inaccurate responses in visual

inference (Liu et al., 2023a; Yue et al., 2023). Thus, monitoring the reasoning capability is of great importance in understanding the ability and the limitations of MLLMs and applying MLLMs in the real world. Previous benchmarks, such as Liu et al. (2023a) and Yue et al. (2023), have primarily addressed evaluating reasoning in each individual image, relying on static and object-centric knowledge. However, they are insufficient to comprehensively assess the reasoning capabilities of MLLMs due to a lack of time-varying object behaviors or events.

042

043

044

047

048

053

054

056

058

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

079

081

To investigate the capabilities of Multi-Modal Language Models (MLLMs) in dynamic reasoning across image sequences, we present a new benchmark, Mementos. This benchmark focuses on the complex task of monitoring and deciphering the positional changes of objects within an image sequence, followed by the inference of behavioral patterns and logical connections among them. Such an endeavor requires the interpretation of the overarching context based on time-variant visual elements, posing a greater challenge than the analysis of static scenes. Concretely, Mementos consists of 4,761 image sequences with varying episode lengths, encompassing diverse scenarios from everyday life, robotics tasks, and comic-style storyboards. An episode refers to a specific event or series of events depicted in the image sequence. Each sequence is paired with a human-annotated description of the key objects and their behaviors.

To assess the reasoning capability of MLLMs on Mementos, we employ a GPT-4-assisted evaluation procedure: after an MLLM produces a description for an image sequence, we extract behavior and object keywords from both AI-generated and human-annotated descriptions using GPT-4. We then use keyword matching to assess the degree of behavioral and object hallucinations. To refine the correctness of this evaluation, we have developed behavior and object synonym graphs for each domain. These graphs facilitate more precise key-



Figure 1: Examples of hallucinations by GPT-4V in three domains on Mementos. The red box shows the description generated by GPT-4V based on the given prompt, and the human-annotated descriptions are in the blue box. Texts highlighted in yellow are hallucination parts generated by GPT-4V. This illustrates that even GPT-4V experiences severe hallucinations when reasoning from image sequences.

word matching, ensuring a thorough and nuanced analysis of the MLLMs' reasoning abilities. Besides, we also provide the comparison with human evaluation to demonstrates that the GPT-4-assisted evaluation procedure is very reliable.

We evaluated the reasoning proficiency of nine leading-edge MLLMs on Mementos, encompassing both black-box and open-source models. Our findings indicate that Mementos poses a considerable challenge to these current MLLMs. For instance, as depicted in Figure 1, GPT-4V exhibits notable behavioral and object hallucinations in various domains during image sequence reasoning. Behavioral hallucinations are defined as the MLLMs' erroneous interpretations or predictions of entity actions, while object hallucinations pertain to the inaccurate identification or creation of objects. Notably, behavioral hallucinations were more frequent

than object hallucinations, highlighting a significant deficiency in MLLMs' capability to deduce events from image sequences.

101

103

104

105

106

107

108

109

110

111

112

113

114

115

116

Furthermore, our research pinpoints three principal factors that lead to the reasoning failures of MLLMs: (1) the interconnectedness of object and behavioral hallucinations, (2) the impact of cooccurring behaviors, and (3) the cumulative effect of behavioral hallucinations. The objective of our proposed benchmark and analyses is to shed light on innovative approaches to augment the reasoning abilities of MLLMs and to reduce hallucinations in their subsequent advancements.

#### 2 **Mementos**

In this section, we introduce Mementos, a novel and challenging benchmark designed to test the reasoning capability of Multimodal Large Language 117 118Model (MLLM) under sequential image input. Ini-119tially, we detail the data gathering and annotation120methodology for Mementos, alongside an overview121of data distribution. Subsequently, we outline the122procedure and the metric employed to evaluate the123reasoning capabilities of MLLMs on Mementos.

#### 2.1 Mementos Benchmark

124

125

127

128

129

130

131

132

133

134

135

136

### 2.1.1 Dataset Composition

Mementos comprises 4,761 image sequences of varying lengths, predominantly sourced from Dailylife, Robotics, and Comics domains. Detailed statistics are provided in Table 1. This diverse collection is instrumental in evaluating the comprehensive time-varying reasoning abilities of MLLMs. Specifically, the robotics data, closely associated with embodied AI or real-world contexts, and the comic-style storyboard data, rich in stylistic and episodic diversity in image sequences, significantly enhance the benchmark's relevance and robustness.

Table 1: The number of image sequences in different categories within Mementos.

	Total	Train Set	Val set
Daily-life Robotics	3505 1101	3055 902	450 199
Comics	155	105	50

Daily-life The Daily-life image sequences in Me-137 mentos are derived from video clips in the Next-138 OA dataset, as cited in Xiao et al. (2021). These 139 sequences represent a range of everyday life sce-140 narios. We have selectively extracted videos from 141 the Next-QA Training set, specifically those with 142 frame counts ranging from 400 to 2,500. To bal-143 ance the challenge of testing MLLMs' reasoning 144 capabilities against the risk of losing critical in-145 formation, our methodology involves retaining the 146 first frame of each video. Subsequently, we sample 147 one image every 100 frames. The collected images 148 from this sampling process then form an image sequence that corresponds to the original video. This 150 approach ensures a rigorous yet feasible evalua-151 tion of MLLMs' reasoning abilities in dynamically 152 evolving everyday scenarios.

154RoboticsFor the Robotics data, we utilized155videos from various sub-datasets within Open X-156Embodiment (Collaboration et al., 2023). Open X-157Embodiment aggregates video datasets from multi-158ple university laboratories, showcasing a variety of159tasks performed by different robotic systems. We

meticulously selected sub-datasets from Open X-Embodiment that offer video resolutions exceeding 128x128 and exhibit a high degree of task diversity. From these chosen sub-datasets, a total of 1,101 videos were sampled. The precise number of videos sourced from each sub-dataset is detailed in Appendix A. For video sampling, our approach varied based on the length of the videos. Videos exceeding 100 frames were processed by sampling one image every n/20 frames, where n represents the total frame count. Conversely, for videos with frame counts ranging from 20 to 100, we sampled one image every 5 frames. This ensures the formation of comprehensive and representative image sequences for each video, catering to the evaluation of MLLMs in diverse and complex robotic contexts.

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

**Comics** The Comics data is composed of wordless multi-panel comics of diverse styles, curated from online sources. Unlike Daily-life and Robotics sections, where image sequences are uniformly extracted from video frames, the comics represent intentionally selected key moments within a narrative, manually illustrated by artists. This distinction sets our dataset apart from conventional video datasets. In addition to traditional comics, this category also incorporates 20 storyboards from movies reimagined in comic style. We have further deconstructed these comics into individual image sequences by taking screenshots. This approach enables a unique exploration of sequential visual reasoning, enhancing the diversity and complexity of the dataset for evaluating MLLMs.

#### 2.1.2 Dataset Annotation

For each image sequence in Mementos, we have meticulously annotated a ground truth description that captures the unfolding events. These descriptions focus on the primary objects and their respective behaviors, where *behavior* refers to a verb or verb phrase associated with the object in question.

For the Daily-life data, we initially employed GPT-4V(ision) (OpenAI, 2023a), to amalgamate and reformulate the questions and answers from the Next-QA videos into single paragraph descriptions. This method significantly expedited the manual annotation process. Following this, we conducted a thorough manual review of these automated descriptions, making necessary adjustments. This included rectifying inaccuracies, removing non-existent episodes, and adding missing details to

ensure alignment with the actual image sequences. 210 To ensure reliability, we implemented a cross-211 validation step, where a separate set of annotators 212 performed a secondary review. For the Robotics 213 and Comics categories, the annotation process was 214 entirely manual, conducted by human annotators. 215 These annotations were then subjected to a veri-216 fication process by the authors which ensures the 217 accuracy and consistency of the descriptions across 218 all categories. 219

### 2.1.3 Dataset Statistics

220

224

231

232

234

In showcasing the diversity of Mementos, we present a detailed overview of the data distribution within the Mementos validation set. Our analysis focuses on two key dimensions: the length of the image sequence and the length of the episode. The length of an image sequence is defined by the number of frames it contains, while the episode length is determined by the total number of events depicted in the sequence. A longer image sequence necessitates the MLLM to process a larger number of images, thereby challenging the model's capacity to manage sequences spanning broader timeframes. A greater episode length signifies that the image sequence encompasses more intricate scenarios.



(b) Distribution of episode length

Figure 2: Data distribution in Mementos Val set. **Image sequence length** For the image sequence length, we count the number of frames in each image sequence. As shown in Figure 2(a), the majority of image sequences are between 4 and 14 frames in length. 67.38% of image sequences contain 4 to 14 frames, yet 31.90% of sequences are composed of longer frames - more than 15 frames.

240

241

242

243

244

246

247

248

249

250

251

252

253

254

255

256

258

259

260

261

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

**Episode length** To quantify the episode length within each image sequence of Mementos, we employed GPT-4 for extracting behavior keywords, specifically verbs associated with objects, from the human-annotated descriptions. This extraction was facilitated using a pre-defined manual prompt, details of which can be found in Appendix D. Following the extraction, we calculated the length of the behavior list for each image sequence. A lengthier behavior list signifies a more extended episode within the image sequence, which inherently poses a greater challenge for the MLLM in comprehending the entire image sequence. As illustrated in Figure 2(b), a significant portion of the image sequences, particularly those from the robotics data, feature episode lengths ranging between 1 and 3. This is mainly attributed to the dominance of twoaction episodes like 'pick up and place', 'move and pull open', 'locate and push'. Meanwhile, the remaining data exhibits a normal distribution for episode lengths spanning 4 to 17.



Figure 3: GPT-4-assisted evaluation procedure. We use "O-" for objects and "B-" for behaviors.

#### 2.2 Evaluation Procedure and Metrics

In this section, we illustrate how to evaluate the descriptions generated by MLLMs, including the evaluation procedure and metrics.

**Procedure** As shown in Figure 3, we use an image sequence and a pre-designed prompt together as the input for MLLMs, and generate the description aligned with the corresponding image sequence. Next, we ask GPT-4 to extract object and behavior keywords in the AI-generated description. We then match the obtained keywords with the *synonym graph* we built, replacing the matched keywords with the root word from the synonym graph. Finally, we obtain two lists of keywords: AI-generated object list and AI-generated behavior

236 237 238 list. We note that the proposed keyword extraction
leveraging GPT-4 is surprisingly reliable and accurate, which is competitive with human extraction.
Please refer to Appendix C for more details.

**Synonym graph** The synonym graph is an unilateral digraph where each edge connects two nodes 283 representing words or phrases. For instance, given a synonym pair (pick up, lift up), an edge is directed from 'lift up' to 'pick up'. In each syn-287 onym pair, the first word, originating from the human-annotated keyword list, is referred to as 289 the root word, while the second word is from the AI-generated keyword. To construct this synonym graph, we use GPT-4 to extract object and behavior 291 keywords from all human-annotated descriptions 292 in the Val set, forming a human-annotated keyword list. Then, we generate descriptions using GPT-294 4V, LLAVA, and Gemini and use GPT-4 to extract object and behavior keywords. After that, we manually match these words with the human-annotated keyword list to identify all synonym pairs and add them as edges to the synonym graph. Given a word or phrase, this synonym graph can quickly match the corresponding root word if a synonym ex-301 ists in the human-annotated keyword list, completing the keyword replacement. For convenience in evaluation, we maintain separate synonym graphs 304 for objects and behaviors of different categories. We make all constructed synonym graphs publicly available as open-source resources.

Metrics After obtaining the AI-generated object list and behavior list, we utilize the corresponding human-annotated object list and human-annotated behavior list as the ground truth to calculate 'Recall,' 'Precision,' and 'F1 metrics' at both the ob-312 ject and behavior levels. These metrics are used to 313 measure the understanding capabilities regarding the image sequence episode. 'Recall' reflects the 315 accuracy of an MLLM's reasoning about episodes in an image sequence, while 'precision' focuses on 317 assessing the severity of hallucinations that occur when understanding the image sequence. 319

### **3** Experiments

In our experimental section, we delve into two key questions: (a) We examine the reasoning **capabilities** of current MLLMs on Mementos. Specifically, we assess the **severity** of object and behavioral hallucinations in these models. (b) We investigate the underlying **causes** of reasoning failures in MLLMs when interpreting image sequences.

### 3.1 Baseline evaluation

## 3.1.1 Models

We establish our baseline using 9 popular MLLMs. The black-box MLLMs include GPT-4V (OpenAI, 2023a) and Gemini (Team, 2023), and the opensource MLLMs are Video-LLaMA-2 (Zhang et al., 2023a), Chat-UniVi (Jin et al., 2023), LLaVA-1.5 (Liu et al., 2023c), MiniGPT4 (Zhu et al., 2023), MiniGPT5 (Zheng et al., 2023), mPLUG Owlv2 (Ye et al., 2023), and InstructBLIP (Dai et al., 2023). Considering that only a few open-source MLLMs are designed to process sequential images (Video-LLaMA-2 and Chat-UniVi), we adapt input for other models by combining all frames from an image sequence into one composite image, referred to as the combined-input (c-input) setting. For black-box MLLMs and Chat-UniVi, we conduct evaluations using both the c-input and an alternative approach where frames from the image sequence are input sequentially, termed the sequential-input (s-input) setting. For Video-LLaMA-2, we only test in s-input setting.

328

329

330

331

332

334

335

336

337

339

340

341

343

344

345

346

347

348

350

351

352

353

355

357

358

361



Figure 4: Comparison of metrics for different MLLMs.

#### 3.1.2 Evaluation results

We evaluate all MLLMs on Mementos and report the results in Figure 4. Besides, we provide the performance of each baseline method in three different domains (Daily-life, Robotics, and Comics) in Table 2. We summarize our findings as follows: **GPT-4V (s-input) and LLaVA-1.5 are the bestperforming models among black-box and opensource MLLMs, respectively.** As shown in Figure 4, except for being on par with Gemini (sinput) and LLaVA-1.5 in behavior precision, GPT-4V with s-input demonstrates the best reasoning

Domain	Input type	Model	Recall	<b>Object</b> Precision	F1	Recall	Behavior Precision	F1
Daily-life		GPT-4V	59.80%	50.96%	53.51%	36.71%	32.97%	33.59%
		Gemini	35.92%	42.06%	37.10%	18.80%	29.42%	21.64%
	Sequential	Video-LLaMA-2	31.59%	30.01%	29.37%	17.05%	28.19%	20.12%
		Chat-UniVi	40.74%	40.78%	39.13%	22.30%	31.10%	24.90%
		GPT-4V	39.45%	39.64%	38.04%	26.43%	23.59%	23.98%
		Gemini	31.17%	37.39%	32.38%	17.71%	25.65%	19.74%
		Chat-UniVi	36.19%	38.88%	36.02%	21.80%	28.52%	23.73%
	Combined	LLaVa-1.5	37.72%	47.01%	40.18%	22.17%	37.33%	26.65%
	Combined	MiniGPT4	32.25%	23.14%	25.75%	18.09%	24.16%	19.45%
		MiniGPT5	31.39%	22.62%	24.91%	18.42%	24.56%	19.85%
		mPLUG_Owl-v2	32.59%	47.17%	37.04%	17.96%	33.57%	22.13%
		InstructBLIP	31.82%	41,14%	34.28%	22.40%	30.30%	24.55%
		GPT-4V	63.94%	65.42%	62.99%	60.72%	24.43%	33.95%
	Sequential	Gemini	43.80%	46.26%	43.15%	46.43%	38.13%	39.38%
		Video-LLaMA-2	13.41%	10.33%	11.15%	17.04%	8.96%	11.23%
		Chat-UniVi	35.40%	32.57%	32.39%	32.24%	16.69%	21.14%
		GPT-4V	27.87%	31.86%	28.58%	44.72%	16.54%	23.58%
Robotics	Combined	Gemini	34.78%	41.66%	36.16%	47.29%	29.59%	34.17%
		Chat-UniVi	17.74%	18.32%	17.07%	19.81%	10.01%	12.54%
		LLaVa-1.5	36.88%	46.62%	39.31%	25.27%	14.80%	17.95%
		MiniGPT4	10.97%	7.28%	8.16%	13.40%	5.88%	7.76%
		MiniGPT5	9.75%	6.52%	7.16%	8.96%	4.53%	5.43%
		mPLUG_Owl-v2	19.75%	26.70%	21.99%	26.46%	16.59%	19.51%
		InstructBLIP	17.96%	18.65%	17.29%	31.41%	19.08%	22.69%
	Sequential	GPT-4V	49.53%	37.57%	41.71%	19.97%	17.29%	18.11%
		Gemini	38.57%	40.64%	38.53%	15.23%	19.11%	16.30%
		Video-LLaMA-2	20.26%	17.59%	18.09%	5.45%	11.07%	6.81%
		Chat-UniVi	28.04%	31.61%	28.13%	10.42%	15.74%	11.97%
		GPT-4V	29.23%	24.64%	25.90%	13.19%	13.09%	12.90%
Comics	Combined	Gemini	41.25%	45.07%	41.18%	15.37%	20.55%	16.42%
		Chat-UniVi	25.12%	28.08%	25.51%	8.85%	10.67%	9.31%
		LLaVa-1.5	29.44%	35.61%	30.97%	8.63%	13.56%	10.27%
		MiniGPT4	20.50%	13.94%	15.74%	7.95%	8.64%	7.98%
		MiniGPT5	22.94%	18.11%	19.42%	8.88%	11.92%	9.94%
		mPLUG_Owl-v2	26.82%	37.74%	29.49%	8.70%	20.85%	11.74%
		InstructBLIP	25.02%	29.15%	25.10%	8.25%	10.48%	8.97%

Table 2: Evaluation of different MLLMs on Mementos.

capability compared with all other MLLMs in understanding image sequences. Among open-source 363 models, LLaVA1.5 performs the best, nearly matching or even surpassing the black-box model Gemini in object comprehension, but its ability to infer behaviors from image sequences is weaker compared to Gemini and GPT-4V. Although Video-LLaMA-2 and Chat-UniVi are designed for video understanding, they do not show an advantage over LLaVA-1.5, especially Video-LLaMA-2, which performs notably worse compared to LLaVA-1.5. The weak-372 est models in understanding image sequences are 373 MiniGPT4 and MiniGPT5, with a significant gap in 374 every metric compared to the other baselines. It's noteworthy that under c-input setting, the performance of black-box MLLMs does not significantly 377 differ from that of open-source MLLMs. LLaVA-378 1.5 and mPLUG Owl-v2 meet or even exceed the black-box MLLMs on many metrics.

381 MLLMs possess a much stronger ability on rea-382 soning objects in image sequences than they do on reasoning behaviors. We find that all MLLM methods perform significantly better on the three metrics for objects than those for behaviors. Taking the best-performing GPT-4V as an example, it achieves over 50% on all three object metrics, with recall even reaching 60%, indicating it can effectively recognize the main objects in an image sequence. However, for behaviors, GPT-4V scores only around 30%, with the best recall metric barely exceeding 40%. Despite this, GPT-4V is still the best-performing MLLM in reasoning behaviors. This suggests that current MLLMs do not possess strong abilities to autonomously infer the behaviors from given sequential images, indicating the importance of our benchmark in highlighting the limitations in the reasoning abilities of MLLMs.

383

384

387

388

390

391

392

393

394

395

396

398

399

400

401

402

403

**Reasoning capability of MLLMs varies across different domains.** From Table 2, we find that black-box models perform best in the robotics domain across the three domains, while open-source models show relatively better performance in the

daily-life domain. Analyzing each domain specifi-404 cally, it is evident that in the daily-life domain, the 405 performance of all methods, except for GPT-4V 406 (s-input), does not vary significantly. The main rea-407 son for the performance gap between open-source 408 MLLMs and black-box MLLMs is the noticeably 409 lower metrics of open-source models compared 410 to black-box models in the robotics and comics 411 domains. The recall, precision, and F1 of both ob-412 ject and behavior for black-box MLLMs are almost 413 more than double those of open-source models. We 414 speculate that one reason for this phenomenon is 415 the distribution shift between Mementos and the 416 training data of open-source MLLMs. The limita-417 tions of the training data lead to weaker reasoning 418 capability of open-source MLLMs. 419

#### 3.2 Analysis of Failure Reasoning

420

421

422

423

494

425

426

427

428

429

430

431

In this section, we will provide reasons for failure reasoning results in current MLLMs, combining specific quantitative analyses and case studies. Since behavioral hallucination is a unique phenomenon in image sequence reasoning, and the causes of object hallucination are not significantly different from those in single image reasoning, we only present the reasons leading to behavioral hallucination in this paper. Due to space limitations, please refer to the Appendix E for specific case studies. The following are our main findings:

Interplay between object and behavioral halluci-432 nations in MLLMs. A key hypothesis underpin-433 ning behavioral hallucination is that incorrect ob-434 ject identification leads to subsequent inaccuracies 435 in behavior identification. To test this, we evalu-436 ated the correlation coefficients between object and 437 behavioral hallucinations across different domains 438 for various MLLMs, as detailed in Appendix B 439 Table 4. Our findings reveal that, for most MLLMs, 440 the correlation coefficients in the three domains 441 fluctuate between 0.1 and 0.4, suggesting a weak 442 yet present correlation. This outcome supports the 443 hypothesis that object hallucination contributes to 444 behavioral hallucination to some extent. Case stud-445 ies further reveal that after an object hallucination 446 occurs, MLLMs tend to describe behaviors related 447 to the hallucinated object, even if these behaviors 448 do not exist in the image sequence. As shown in 449 450 Figure 5, after recognizing a scene as a tennis court, a MLLM might describe a person playing tennis. 451 Interestingly, in the robotics domain, there is a neg-452 ligible correlation between object and behavioral 453 hallucinations in black-box MLLMs. This diver-454

gence is likely because behaviors in robotics are predominantly linked to robotic arms, which these MLLMs generally identify correctly.



Figure 5: A sample of failure reasoning case in Dailylife domain. The failure reason is object hallucination, correlation between object hallucination and behavioral hallucination, and co-occurrence behavior. Following the object hallucination of *tennis court*, the LVLM subsequently exhibits behavioral hallucinations of *holding a tennis racket* (correlation between object hallucination and behavioral hallucination) and *appears to be playing tennis* (co-occurrence behavior).

The impact of co-occurrence on behavioral hallucination. In line with object hallucination phenomena, as noted in Li et al. (2023c) and Zhou et al. (2023a), MLLMs demonstrate a tendency to generate behaviors that are commonly paired together. This proclivity exacerbates the problem of behavioral hallucination, especially in the field of robotics. Consider the case in Figure 1 where a robotic arm is tasked with opening a drawer by grabbing its side. MLLMs might erroneously depict the sequence as the arm grabbing the handle first, followed by pulling the drawer open, since grabbing the handle is a more co-occurring behavior with 'pull open'. Despite the final outcome being accurately described, such errors in key details are unacceptable in robotics. This issue is of particular concern given the growing inclination to utilize MLLMs as reward functions in robotic training (Ma et al., 2023; Sontakke et al., 2023; Rocamonde et al., 2023; Baumli et al., 2023). Such behavioral hallucinations can critically affect the quality of the reward function, leading to potential mislearning of behaviors in robotic systems. Detailed case studies are shown in Appendix E.

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

507

509

510

511

512

513

515

516

517

519

523

482

The Snowball effect in behavioral hallucinations. The Snowball effect is a well-documented phenomenon in machine learning, referring to the progressive accumulation or intensification of errors in a system, as discussed in Asadi et al. (2019); Zhang et al. (2023b); Wang et al. (2023c); Liu et al. (2023d). Zhang et al. (2023b) notably highlight this phenomenon in Large Language Models. Experiments on Mementos reveal that the snowball effect in both behavioral and object hallucinations becomes markedly pronounced when reasoning through image sequences. The temporal nature of image sequences, consisting of a series of frames rather than a solitary image, demands that MLLMs sequentially infer the narrative. This process makes models susceptible to exacerbating hallucinations if errors occur early in the sequence. We specifically examined the trend of object and behavioral hallucination in GPT-4V and LLaVA-1.5 within the daily-life domain, correlating it with the episode length. As shown in Figure 6, there is a noticeable decrease in object and behavior recall for both MLLMs as the episode length extends. This trend suggests a heightened susceptibility to hallucinations and a pronounced snowball effect in MLLMs when processing image sequences with a greater array of objects and behaviors. Detailed case studies can be found in Appendix E.



Figure 6: The trend of changes in object and behavior recall for GPT-4V and LLaVA-1.5 in the Daily-life domain as the episode length increases.

### 4 Related work

#### 4.1 Benchmarking in MLLMs

The advent of MLLMs has prompted a reassessment of traditional benchmarks (Lin et al., 2014; Marino et al., 2019; Hudson and Manning, 2019). These benchmarks fail to sufficiently expose the hallucination issues in MLLMs. Consequently, there is a growing impetus to devise more challenging benchmarks. This trend spans various domains, from question and answering (QA) reasoning (Liu et al., 2023a; Yue et al., 2023), to optical character recognition (OCR) (Liu et al., 2023f), and extends to the study of hallucinations (Wang et al., 2023a), with benchmarks such as POPE (Li et al., 2023c) and Bingo (Cui et al., 2023). Additionally, comprehensive analyses of MLLMs, such as Mmbench (Liu et al., 2023e), Mm-vet(Yu et al., 2023b), LVLM-eHub(Xu et al., 2023), SEED(Li et al., 2023a), GAVIE(Liu et al., 2023b), and LAMM (Yin et al., 2023), are emerging.

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

555

556

557

558

559

560

561

562

564

565

566

567

568

570

571

572

Our paper presents a novel benchmark using sequences from videos or comics to study behavioral hallucinations, diverging from single-image analysis. Unlike Chen et al. (2023a)'s vision QA tasks from uniformly sampled video frames, our benchmark challenges MLLMs to describe sequences without question guidance, offering a finer evaluation of hallucinations and reasoning in MLLMs.

#### 4.2 Hallucination in MLLMs

Hallucinations in MLLMs, akin to those in Large Language Models (LLMs) (Zhang et al., 2023c; Li et al., 2023b; Zhou et al., 2024; Chen et al., 2023b), represent a significant challenge. In MLLMs, hallucinations are characterized by inconsistencies between the model's output and the visual content (Rohrbach et al., 2018; Wang et al., 2023a). Recent studies have explored various aspects of hallucination in MLLMs, covering topics such as object hallucination (Li et al., 2023c), hallucination assessment in GPT-4V (Cui et al., 2023a).

While there are methods proposed for mitigating hallucinations (Zhou et al., 2023a; Wang et al., 2023b; Leng et al., 2023; Zhou et al., 2023b; Chen et al., 2023c; Jiang et al., 2023; Huang et al., 2023; Yu et al., 2023a; Zhao et al., 2023), there is a noticeable gap in the literature regarding the study of behavioral hallucination. Moreover, the existing work does not offer a dedicated metric for evaluating behavioral hallucinations.

## 5 Conclusion

In this paper, we present Mementos, a novel and challenging benchmark designed to assess the reasoning abilities of Multimodal Large Language Models (MLLMs) in interpreting image sequences. We conduct evaluations on nine most recent MLLMs using GPT-4-assisted evaluation procedure. Our findings indicate that all tested MLLMs struggle with significant behavioral and object hallucinations in generating descriptions for image sequences. Through a mix of quantitative analysis and case studies, we identify three primary factors contributing to these reasoning failures.

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

# 573 Limitations

**Domain courage** Mementos is consisted of 574 4,761 image sequences from three domains: Daily 575 life, Robotics, and Comics. It would be interesting 576 to include a broader variety of data types. This expansion could include first-person navigation ex-578 579 periences, sequential medical CT scans, and interactive gaming data. MLLMs could behave different types of hallucinations in image sequences from other domains 582

**Evaluation Process** Our evaluation process fo-583 584 cuses on the match of keywords to measure the reasoning ability of MLLMs. However, it would be possible that the MLLM generation is the same 586 as human annotations in semantics but obtains low 587 performance, since the generated tokens are not 588 covered by our synonym graph. Future work could 589 extend the evaluation method to semantic under-590 standing rather than relying predominantly on keyword matching.

Hallucination Mitigation Our work identifies two kinds of hallucination: object and behavioral hallucinations and explore the failure reason of MLLMs. We have not yet proposed a mitigation method to reduce behavioral hallucinations. Future work could utilize the three causes of reasoning failures to bolster the reasoning faculties of MLLMs, making them more adept at accurately interpreting and describing complex image sequences.

### References

596

598

603

606

607

610

611

612

613

614

615

616

617

618

619

623

- Kavosh Asadi, Dipendra Misra, Seungchan Kim, and Michel L. Littman. 2019. Combating the compounding-error problem with a multi-step model.
  - Kate Baumli, Satinder Baveja, Feryal Behbahani, Harris Chan, Gheorghe Comanici, Sebastian Flennerhag, Maxime Gazeau, Kristian Holsheimer, Dan Horgan, Michael Laskin, et al. 2023. Vision-language models as a source of rewards. *arXiv preprint arXiv:2312.09187*.
- Xiuyuan Chen, Yuan Lin, Yuchen Zhang, and Weiran Huang. 2023a. Autoeval-video: An automatic benchmark for assessing large vision language models in open-ended video question answering. *arXiv preprint arXiv:2311.14906*.
- Yuyan Chen, Qiang Fu, Yichen Yuan, Zhihao Wen, Ge Fan, Dayiheng Liu, Dongmei Zhang, Zhixu Li, and Yanghua Xiao. 2023b. Hallucination detection: Robustly discerning reliable answers in large language models. In Proceedings of the 32nd ACM International Conference on Information and Knowledge Management, pages 245–255.

- Zhiyang Chen, Yousong Zhu, Yufei Zhan, Zhaowen Li, Chaoyang Zhao, Jinqiao Wang, and Ming Tang. 2023c. Mitigating hallucination in visual language models with visual supervision. *arXiv preprint arXiv:2311.16479*.
- Open X-Embodiment Collaboration, Abhishek Padalkar, Acorn Pooley, Ajinkya Jain, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky, Anant Rai, Anikait Singh, Anthony Brohan, Antonin Raffin, Ayzaan Wahid, Ben Burgess-Limerick, Beomjoon Kim, Bernhard Schölkopf, Brian Ichter, Cewu Lu, Charles Xu, Chelsea Finn, Chenfeng Xu, Cheng Chi, Chenguang Huang, Christine Chan, Chuer Pan, Chuyuan Fu, Coline Devin, Danny Driess, Deepak Pathak, Dhruv Shah, Dieter Büchler, Dmitry Kalashnikov, Dorsa Sadigh, Edward Johns, Federico Ceola, Fei Xia, Freek Stulp, Gaoyue Zhou, Gaurav S. Sukhatme, Gautam Salhotra, Ge Yan, Giulio Schiavi, Hao Su, Hao-Shu Fang, Haochen Shi, Heni Ben Amor, Henrik I Christensen, Hiroki Furuta, Homer Walke, Hongjie Fang, Igor Mordatch, Ilija Radosavovic, Isabel Leal, Jacky Liang, Jaehyung Kim, Jan Schneider, Jasmine Hsu, Jeannette Bohg, Jeffrey Bingham, Jiajun Wu, Jialin Wu, Jianlan Luo, Jiayuan Gu, Jie Tan, Jihoon Oh, Jitendra Malik, Jonathan Tompson, Jonathan Yang, Joseph J. Lim, João Silvério, Junhyek Han, Kanishka Rao, Karl Pertsch, Karol Hausman, Keegan Go, Keerthana Gopalakrishnan, Ken Goldberg, Kendra Byrne, Kenneth Oslund, Kento Kawaharazuka, Kevin Zhang, Keyvan Majd, Krishan Rana, Krishnan Srinivasan, Lawrence Yunliang Chen, Lerrel Pinto, Liam Tan, Lionel Ott, Lisa Lee, Masayoshi Tomizuka, Maximilian Du, Michael Ahn, Mingtong Zhang, Mingyu Ding, Mohan Kumar Srirama, Mohit Sharma, Moo Jin Kim, Naoaki Kanazawa, Nicklas Hansen, Nicolas Heess, Nikhil J Joshi, Niko Suenderhauf, Norman Di Palo, Nur Muhammad Mahi Shafiullah, Oier Mees, Oliver Kroemer, Pannag R Sanketi, Paul Wohlhart, Peng Xu, Pierre Sermanet, Priya Sundaresan, Quan Vuong, Rafael Rafailov, Ran Tian, Ria Doshi, Roberto Martín-Martín, Russell Mendonca, Rutav Shah, Ryan Hoque, Ryan Julian, Samuel Bustamante, Sean Kirmani, Sergey Levine, Sherry Moore, Shikhar Bahl, Shivin Dass, Shuran Song, Sichun Xu, Siddhant Haldar, Simeon Adebola, Simon Guist, Soroush Nasiriany, Stefan Schaal, Stefan Welker, Stephen Tian, Sudeep Dasari, Suneel Belkhale, Takayuki Osa, Tatsuya Harada, Tatsuya Matsushima, Ted Xiao, Tianhe Yu, Tianli Ding, Todor Davchev, Tony Z. Zhao, Travis Armstrong, Trevor Darrell, Vidhi Jain, Vincent Vanhoucke, Wei Zhan, Wenxuan Zhou, Wolfram Burgard, Xi Chen, Xiaolong Wang, Xinghao Zhu, Xuanlin Li, Yao Lu, Yevgen Chebotar, Yifan Zhou, Yifeng Zhu, Ying Xu, Yixuan Wang, Yonatan Bisk, Yoonyoung Cho, Youngwoon Lee, Yuchen Cui, Yueh hua Wu, Yujin Tang, Yuke Zhu, Yunzhu Li, Yusuke Iwasawa, Yutaka Matsuo, Zhuo Xu, and Zichen Jeff Cui. 2023. Open X-Embodiment: Robotic learning datasets and RT-X models. https://arxiv.org/abs/2310.08864.

Chenhang Cui, Yiyang Zhou, Xinyu Yang, Shirley Wu,

- 687
- 701 703
- 704 705 706 707 711 712 713 714 715 716 717
- 718 721 725 727 728 729 730

- 731 733 734
- 735 736 737

740 741

- Linjun Zhang, James Zou, and Huaxiu Yao. 2023. Holistic analysis of hallucination in gpt-4v (ision): Bias and interference challenges. arXiv preprint arXiv:2311.03287.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. 2023. Instructblip: Towards general-purpose vision-language models with instruction tuning.
- Qidong Huang, Xiaoyi Dong, Pan Zhang, Bin Wang, Conghui He, Jiaqi Wang, Dahua Lin, Weiming Zhang, and Nenghai Yu. 2023. Opera: Alleviating hallucination in multi-modal large language models via over-trust penalty and retrospection-allocation. arXiv preprint arXiv:2311.17911.
- Drew A Hudson and Christopher D Manning. 2019. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 6700-6709.
- Chaoya Jiang, Haiyang Xu, Mengfan Dong, Jiaxing Chen, Wei Ye, Ming Yan, Qinghao Ye, Ji Zhang, Fei Huang, and Shikun Zhang. 2023. Hallucination augmented contrastive learning for multimodal large language model. arXiv preprint arXiv:2312.06968.
- Peng Jin, Ryuichi Takanobu, Caiwan Zhang, Xiaochun Cao, and Li Yuan. 2023. Chat-univi: Unified visual representation empowers large language models with image and video understanding. arXiv preprint arXiv:2311.08046.
- Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian Lu, Chunyan Miao, and Lidong Bing. 2023. Mitigating object hallucinations in large visionlanguage models through visual contrastive decoding. arXiv preprint arXiv:2311.16922.
- Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. 2023a. Seed-bench: Benchmarking multimodal llms with generative comprehension. arXiv preprint arXiv:2307.16125.
- Junyi Li, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. 2023b. Halueval: A largescale hallucination evaluation benchmark for large language models. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 6449-6464.
- Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023c. Evaluating object hallucination in large vision-language models. arXiv preprint arXiv:2305.10355.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In Computer Vision-ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V13, pages 740-755. Springer.

Fuxiao Liu, Tianrui Guan, Zongxia Li, Lichang Chen, Yaser Yacoob, Dinesh Manocha, and Tianyi Zhou. 2023a. Hallusionbench: You see what you think? or you think what you see? an image-context reasoning benchmark challenging for gpt-4v (ision), llava-1.5, and other multi-modality models. arXiv preprint arXiv:2310.14566.

742

743

744

745

746

747

749

750

751

753

754

755

756

758

759

760

761

762

764

765

766

767

768

770

771

772

773

774

775

776

777

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

- Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. 2023b. Aligning large multi-modal model with robust instruction tuning. arXiv preprint arXiv:2306.14565.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023c. Improved baselines with visual instruction tuning.
- Xiaoyu Liu, Jiaxin Yuan, Bang An, Yuancheng Xu, Yifan Yang, and Furong Huang. 2023d. Cdisentanglement: Discovering causally-independent generative factors under an inductive bias of confounder. arXiv preprint arXiv:2310.17325.
- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. 2023e. Mmbench: Is your multi-modal model an all-around player? arXiv preprint arXiv:2307.06281.
- Yuliang Liu, Zhang Li, Hongliang Li, Wenwen Yu, Mingxin Huang, Dezhi Peng, Mingyu Liu, Mingrui Chen, Chunyuan Li, Lianwen Jin, et al. 2023f. On the hidden mystery of ocr in large multimodal models. arXiv preprint arXiv:2305.07895.
- Yecheng Jason Ma, William Liang, Vaidehi Som, Vikash Kumar, Amy Zhang, Osbert Bastani, and Dinesh Jayaraman. 2023. Liv: Language-image representations and rewards for robotic control. arXiv preprint arXiv:2306.00958.
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. Ok-vqa: A visual question answering benchmark requiring external knowledge. In Proceedings of the IEEE/cvf conference on computer vision and pattern recognition, pages 3195-3204.

OpenAI. 2023a. Gpt-4 technical report.

OpenAI. 2023b. Gpt-4v(ision) system card.

- Juan Rocamonde, Victoriano Montesinos, Elvis Nava, Ethan Perez, and David Lindner. 2023. Vision-language models are zero-shot reward models for reinforcement learning. arXiv preprint arXiv:2310.12921.
- Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. 2018. Object hallucination in image captioning. arXiv preprint arXiv:1809.02156.
- Sumedh Anand Sontakke, Jesse Zhang, Séb Arnold, Karl Pertsch, Erdem Biyik, Dorsa Sadigh, Chelsea

795

847

848

Finn, and Laurent Itti. 2023. Roboclip: One demonstration is enough to learn robot policies. In Thirtyseventh Conference on Neural Information Processing Systems.

Gemini Team. 2023. Gemini: A family of highly capable multimodal models.

- Junyang Wang, Yiyang Zhou, Guohai Xu, Pengcheng Shi, Chenlin Zhao, Haiyang Xu, Qinghao Ye, Ming Yan, Ji Zhang, Jihua Zhu, et al. 2023a. Evaluation and analysis of hallucination in large vision-language models. arXiv preprint arXiv:2308.15126.
- Lei Wang, Jiabang He, Shenshen Li, Ning Liu, and Ee-Peng Lim. 2023b. Mitigating fine-grained hallucination by fine-tuning large vision-language models with caption rewrites. arXiv preprint arXiv:2312.01701.

Xiyao Wang, Ruijie Zheng, Yanchao Sun, Ruonan Jia, Wichayaporn Wongkamjan, Huazhe Xu, and Furong Huang. 2023c. Coplanner: Plan to roll out conservatively but to explore optimistically for model-based rl.

Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. 2021. Next-qa: Next phase of questionanswering to explaining temporal actions. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 9777-9786.

Peng Xu, Wenqi Shao, Kaipeng Zhang, Peng Gao, Shuo Liu, Meng Lei, Fanqing Meng, Siyuan Huang, Yu Qiao, and Ping Luo. 2023. Lvlm-ehub: A comprehensive evaluation benchmark for large visionlanguage models. arXiv preprint arXiv:2306.09265.

Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Haowei Liu, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. 2023. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration. arXiv preprint arXiv:2311.04257.

Zhenfei Yin, Jiong Wang, Jianjian Cao, Zhelun Shi, Dingning Liu, Mukai Li, Lu Sheng, Lei Bai, Xiaoshui Huang, Zhiyong Wang, et al. 2023. Lamm: Language-assisted multi-modal instruction-tuning dataset, framework, and benchmark. arXiv preprint arXiv:2306.06687.

- Oifan Yu, Juncheng Li, Longhui Wei, Liang Pang, Wentao Ye, Bosheng Qin, Siliang Tang, Qi Tian, and Yueting Zhuang. 2023a. Hallucidoctor: Mitigating hallucinatory toxicity in visual instruction data. arXiv *preprint arXiv:2311.13614*.
- Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2023b. Mm-vet: Evaluating large multimodal models for integrated capabilities. arXiv preprint arXiv:2308.02490.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. 2023. Mmmu:

A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. arXiv preprint arXiv:2311.16502.

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

- Hang Zhang, Xin Li, and Lidong Bing. 2023a. Videollama: An instruction-tuned audio-visual language model for video understanding. arXiv preprint arXiv:2306.02858.
- Muru Zhang, Ofir Press, William Merrill, Alisa Liu, and Noah A. Smith. 2023b. How language model hallucinations can snowball.
- Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, et al. 2023c. Siren's song in the ai ocean: A survey on hallucination in large language models. arXiv preprint arXiv:2309.01219.
- Zhiyuan Zhao, Bin Wang, Linke Ouyang, Xiaoyi Dong, Jiaqi Wang, and Conghui He. 2023. Beyond hallucinations: Enhancing lvlms through hallucinationaware direct preference optimization. arXiv preprint arXiv:2311.16839.
- Kaizhi Zheng, Xuehai He, and Xin Eric Wang. 2023. Minigpt-5: Interleaved vision-and-language generation via generative vokens.
- Yiyang Zhou, Chenhang Cui, Jaehong Yoon, Linjun Zhang, Zhun Deng, Chelsea Finn, Mohit Bansal, and Huaxiu Yao. 2023a. Analyzing and mitigating object hallucination in large vision-language models. arXiv preprint arXiv:2310.00754.
- Yuhang Zhou, Suraj Maharjan, and Beiye Liu. 2023b. Scalable prompt generation for semi-supervised learning with language models. arXiv preprint arXiv:2302.09236.
- Yuhang Zhou, Paiheng Xu, Xiaoyu Liu, Bang An, Wei Ai, and Furong Huang. 2024. Explore spurious correlations at the concept level in language models for text classification.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. arXiv preprint arXiv:2304.10592.

#### **Details of Open X-Embodiment Data** Α Selection

In this section, we provide the names of all subsets selected from Open X-Embodiment dataset and the corresponding sampling video numbers. For detailed information, please refer to Table 3.

#### **Correlation Coefficients between** B **Object and Behavioral Hallucinations**

In this section, we provide detailed correlation coefficients between object and behavioral hallucinations in Table 4.

Sub-dataset name	Number of videos selected
fractal20220817_data	400
kuka	50
bridge	300
jaco_play	50
berkeley_autolab_ur5	50
toto	10
columbia_cairlab_pusht_real	5
stanford_hydra_dataset_converted_externally_to_rlds	5
ucsd_kitchen_dataset_converted_externally_to_rlds	50
bc_z	50
utokyo_pr2_opening_fridge_converted_externally_to_rlds	5
utokyo_pr2_tabletop_manipulation_converted_externally_to_rlds	10
utokyo_xarm_pick_and_place_converted_externally_to_rlds	1
utokyo_xarm_bimanual_converted_externally_to_rlds	5
dlr_sara_pour_converted_externally_to_rlds	5
dlr_edan_shared_control_converted_externally_to_rlds	100
asu_table_top_converted_externally_to_rlds	20
utaustin_mutex	30
berkeley_fanuc_manipulation	30

## **C** Human Evaluation

In this section, to verify the reliability of the GPT-4 assisted evaluation procedure, we compare the results of GPT-4 assisted evaluation with those of human evaluation. We randomly select 200 image sequences from the entire Val set and manually extract object and behavior keyword lists for each image sequence's AI-generated description and human-annotated description. Then, we calculate six metrics and compare them with the metrics obtained using keyword lists extracted by GPT-4. We choose the four MLLMs that performed best in reasoning on Mementos as representatives: GPT-4V (s-input), Gemini (s-input), Chat-UniVi (s-input), and LLaVA-1.5. The evaluation results are shown in Table 5.

After comparison, we find that there is not a significant gap between the results of GPT-4 assisted evaluation and human evaluation, with the absolute value of the difference mostly ranging between 1% to 4%. For most metrics, the GPT-4 assisted evaluation tends to overestimate the performance of MLLMs, meaning the evaluation results are higher than those of human evaluation. However, the relative ranking among different MLLMs remains essentially unchanged. Overall, the GPT-4 assisted evaluation is quite reliable.

## **D** Prompt Details

In this section, we provide all the prompts used in our paper, including those used to merge questions and answers from Daily-life videos into a single description, prompts for MLLMs to generate descriptions corresponding to image sequences, and prompts for extracting object and behavior keywords from both human-annotated and AIgenerated descriptions. The detailed prompts are showm in Table 6. 928

929

930

931

932

933

934

935

936

937

903

904

905

921

923

924

925

927

Domain	Input type	Model	Recall	Precision	F1
Daily-life		GPT-4V	0.120	0.188	0.132
	Sequential	Gemini	0.165	0.179	0.146
		Video-LLaMA-2	0.197	0.067	0.125
		Chat-UniVi	0.138	0.178	0.137
		GPT-4V	0.242	0.182	0.199
		Gemini	0.158	0.179	0.152
		Chat-UniVi	0.127	0.184	0.172
	Combined	LLaVa-1.5	0.112	0.134	0.106
	Combined	MiniGPT4	0.135	0.145	0.115
		MiniGPT5	0.126	0.188	0.146
		mPLUG_Owl-v2	0.106	0.113	0.069
		InstructBLIP	0.133	0.125	0.127
		GPT-4V	-0.012	0.022	0.011
	Sequential	Gemini	0.027	0.144	0.101
		Video-LLaMA-2	0.107	0.107	0.109
		Chat-UniVi	0.038	0.121	0.089
	 	GPT-4V	0.041	-0.022	0.008
Robotics		Gemini	-0.049	-0.086	-0.106
		Chat-UniVi	0.189	0.242	0.207
	Combined	LLaVa-1.5	0.135	0.123	0.157
	Combined	MiniGPT4	0.186	0.316	0.233
		MiniGPT5	0.056	0.027	0.045
		mPLUG_Owl-v2	0.244	0.163	0.231
		InstructBLIP	0.227	0.235	0.253
		GPT-4V	0.045	0.225	0.158
	Sequential	Gemini	0.176	0.081	0.144
		Video-LLaMA-2	0.261	0.280	0.299
		Chat-UniVi	0.239	0.331	0.221
		GPT-4V	0.343	0.539	0.471
Comics		Gemini	0.187	0.121	0.167
		Chat-UniVi	0.293	0.113	0.279
	Combined	LLaVa-1.5	0.062	0.101	0.088
		MiniGPT4	0.199	0.134	0.213
		MiniGPT5	0.324	0.366	0.339
		mPLUG_Owl-v2	0.231	-0.043	0.157
		InstructBLIP	0.288	0.005	0.262

Table 4: Correlation coefficient between behavioral hallucination and object hallucination of different MLLMs on Mementos.

#### **Case Study** Ε

938

939

940

941

942

943

944

945

In this section, we present failure reasoning cases of different domains (Figure 7-22), with specific reasons for failure detailed in the captions of each figure.

#### **Status of Exemption from Institutional** F **Review Board**

Before starting any segments of the study involving 946 human evaluation, the research team completed and submitted a "Human Subjects Research Determina-947 tion" form to the appropriate Institutional Review 948 Board (IRB). We obtained a determination letter 949 from the IRB before any human study activities 950 commenced, indicating that our project proposal 951 had been granted 'Exempt' status. This classifica-952 tion implies that the proposed research was deemed 953 'Not Human Subjects Research'. 954



Figure 7: A sample of failure reasoning case in Daily-life domain, we highlight the hallucination parts in yellow. Failure reason: co-occurrence behavior and Snowball.



Figure 8: A sample of failure reasoning case in Daily-life domain, we highlight the hallucination parts in yellow. Failure reason: object hallucination, and correlation between object hallucination and behavioral hallucination.



Figure 9: A sample of failure reasoning case in Daily-life domain, we highlight the hallucination parts in yellow. Failure reason: lack of common sense and Snowball.



Figure 10: A sample of failure reasoning case in Daily-life domain, we highlight the hallucination parts in yellow. Failure reason: object hallucination, correlation between object hallucination and behavioral hallucination, and co-occurrence behavior.



Figure 11: A sample of failure reasoning case in Daily-life domain, we highlight the hallucination parts in yellow. Failure reason: Snowball. In this case, we observe that in addition to the significant behavioral hallucinations caused by Snowball effect mentioned in Section 3.2, another result of Snowball is that LVLMs may not fully describe all episodes in an image sequence. That is, after a behavioral hallucination occurs, the LVLM might assume the episode has ended and stop describing. For instance, in this case, the LVLM stopped describing after mentioning the child reaching the living room and the adult leaving, without continuing to describe the child pushing the box back along the hallway.



Figure 12: A sample of failure reasoning case in Robotics domain, we highlight the hallucination parts in yellow. Failure reason: co-occurrence behavior.



Figure 13: A sample of failure reasoning case in Robotics domain, we highlight the hallucination parts in yellow. Failure reason: Snowball. This case effectively demonstrates the lack of LVLM's reasoning ability in image sequence comprehension. In the first image, the robotic arm indeed appears to be moving towards the cube, but from the second image, the arm lowers and moves towards the disc-shaped object. The LVLM failed to infer this behavior from the first two images and based its subsequent description solely on the understanding in the first image, leading to a Snowball effect.



Figure 14: A sample of failure reasoning case in Robotics domain, we highlight the hallucination parts in yellow. Failure reason: co-occurrence behavior and Snowball. This case also reflects another outcome of the Snowball effect that we mentioned in Figure 11. After assuming that the robotic arm is cooking, the LVLM do not continue to describe the behavior of the robotic arm moving the pot from the right stove to the left.



Figure 15: A sample of failure reasoning case in Robotics domain, we highlight the hallucination parts in yellow. Failure reason: Snowball.



Figure 16: A sample of failure reasoning case in Robotics domain, we highlight the hallucination parts in yellow. Failure reason: object hallucination, correlation between object hallucination and behavioral hallucination, and Snowball.



Figure 17: A sample of failure reasoning case in Robotics domain, we highlight the hallucination parts in yellow. Failure reason: object hallucination, correlation between object hallucination and behavioral hallucination, and Snowball.



Figure 18: A sample of failure reasoning case in Comics domain, we highlight the hallucination parts in yellow. Failure reason: object hallucination, correlation between object hallucination and behavioral hallucination, and Snowball.



Figure 19: A sample of failure reasoning case in Comics domain, we highlight the hallucination parts in yellow. Failure reason: object hallucination, correlation between object hallucination and behavioral hallucination, and Snowball.

[Prompt]: Write a description for the given image sequence in a single paragraph, what is happening in this episode

 Image: the paragraph is the provided of the paragraph is the provided of the paragraph is the paragraph is

Figure 20: A sample of failure reasoning case in Comics domain, we highlight the hallucination parts in yellow. Failure reason: Snowball.



Figure 21: A sample of failure reasoning case in Comics domain, we highlight the hallucination parts in yellow. Failure reason: object hallucination, correlation between object hallucination and behavioral hallucination, and Snowball.

Table 5: Human evaluation.

Model	Eval type	Recall	<b>Object</b> Precision	F1	Recall	Behavior Precision	F1
GPT-4V (s-input)	GPT-4	60.91%	51.04%	54.13%	38.02%	33.05%	34.12%
	Human	57.69%	49.54%	52.01%	35.26%	31.60%	32.67%
Gemini (s-input)	GPT-4	37.54%	39.43%	36.88%	23.38%	34.19%	24.02%
	Human	35.82%	38.11%	37.09%	20.46%	33.72%	22.99%
ChatUnivi (s-input)	GPT-4	40.32%	42.04%	39.52%	24.95%	28.06%	27.15%
	Human	37.65%	38.59%	36.46%	25.73%	27.40%	26.64%
LLaVA-1.5 (c-input)	GPT-4	35.77%	44.18%	38.09%	24.47%	38.79%	28.59%
	Human	36.84%	41.37%	39.77%	22.95%	39.82%	29.18%



Figure 22: A sample of failure reasoning case in Comics domain, we highlight the hallucination parts in yellow. Failure reason: object hallucination, correlation between object hallucination and behavioral hallucination, and Snowball.

## Prompt

Task: Rewrite questions and answers into a single paragraph

Image: <Image sequence>

Text: <Write a description for this image based on the following questions and answers in one paragraph. Please remember that some objects or actions in the following questions and answers may not be included in the images. Please do not include the excluded items in your description. Here are the questions and answers: Question: {Question 1} Answer: {Answer 1} Question: {Question 2} Answer: {Answer 2} ... Question: {Question n} Answer: {Answer n}>

Task: Generate description for the given image sequence

Image: <Image sequence>

Text: <Write a description for the given image sequence in a single paragraph, what is happening in this episode?>

Task: Extract object and behavior keywords

Text: <I will provide you two paragraphs. The first paragraph is human-composed and the second paragraph is generated by AI models. I want to evaluate the hallucination in the second paragraph. Please extract the object and action words or phrases from the following text. The objects should have a tangible meaning and consist of no more than two words; non-tangible objects should not be extracted. The action words or phrases should only relate to the extracted objects. Also, you must convert the corresponding actions to their complete root form. Then, for the final answer, please examine 4 lists and must transfer the synonyms in 4 lists into the same word. Please directly output the final object and action lists in two paragraphs, respectively as in the form in the example below without any justifications or intermediate steps.

Here is an example:

1. The sequence of images captures a dog's cautious interaction with a metal toy inside a house. The dog appears wary and maintains a distance from the unfamiliar object, barking to express its disapproval and possibly intimidation. As the toy moves, the dog's reaction is to bark and lean backward, showing a clear sign of being unsettled by the toy's motion. When the toy momentarily ceases movement, the dog also stops, remaining alert and attentive. At the end of the image, when the toy comes to a halt, the dog looks up, still processing the strange encounter with the inanimate object.

2. The image is a collage of multiple pictures featuring two dogs playing with a toy alligator. The dogs are in various positions, with some of them standing on the toy alligator, while others are interacting with it in different ways. The collage captures the dogs' playfulness and excitement as they engage with the toy alligator.

The lists are

Object list 1: [dog, toy, house]

Action list 1: [interaction, bark, express intimidation, move, lean backward, stop, look up]

Object list 2: [dog, toy]

Action list 2: [play, stand, interaction]

Here is the paragraphs:

# 1. {Human-annotated description}

2. {AI-generated description}

The lists are:>