# WorldQA: Multimodal World Knowledge in Videos through Long-Chain Reasoning

Anonymous ACL submission



Figure 1: An example video from our WorldQA. To determine *where the lady went when she was absent from the video*, we rely on visual cues, auditory hints, and the application of world knowledge. This forms a reasoning chain to deduce the answer. WorldQA comprises 1007 question-answer pairs and 303 videos, spanning five types of world knowledge. On average, the reasoning chain consists of 4.45 steps.We recommand watch the video: https://www.youtube.com/watch?v=NXbJLLf9E\_I

### Abstract

Multimodal information, together with our knowledge, help us to understand the complex and dynamic world. Large language models (LLM) and large multimodal models (LMM), however, still struggle to emulate this capability. In this paper, we present WorldQA, a video understanding dataset designed to push the boundaries of multimodal world models with three appealing properties: (1) Multimodal Inputs: The dataset comprises 1007 question-answer pairs and 303 videos, necessitating the analysis of both auditory and visual data for successful interpretation. (2) World Knowledge: We identify five essential types of world knowledge for question formulation. This approach challenges models to extend their capabilities beyond mere perception. (3) Long-Chain Reasoning: Our dataset introduces an average reasoning step of 4.45, notably surpassing other videoQA datasets. Furthermore, we introduce

002

012

013

WorldRetriever, an agent designed to synthesize expert knowledge into a coherent reasoning chain, thereby facilitating accurate responses to WorldQA queries. Extensive evaluations of 13 prominent LLMs and LMMs reveal that WorldRetriever, although being the most effective model, achieved only 70% of human-level performance in multiple-choice questions. This finding highlights the necessity for further advancement in the reasoning and comprehension abilities of models. Our experiments also yield several key insights. For instance, while humans tend to perform better with increased frames, current LMMs, including WorldRetriever, show diminished performance under similar conditions. We hope that WorldQA, our methodology, and these insights could contribute to the future development of multimodal world models.

# 1 Introduction

040

041

047

063

071

077

097

Consider the scene in Fig. 1. It shows more than a woman simply drinking coffee and picking up clothes. With the background sounds of a ticking clock and a mix of radio broadcasts, along with the noise of a door opening and closing, we naturally form a story: she's just waken up and is getting ready, probably for work. Understanding this video requires combining two key human skills: perception and cognition. Perception lets us notice and recognize details, like the clock's time, and the radio's sound. Cognition, on the other hand, involves using knowledge from our own experiences, like knowing typical work hours. Together, these skills enable us to follow the video's story through a logical series of steps.

For humans, merging perception and cognition to understand video narratives is intuitive, but what about Large Multi-modal Models (LMMs)? To push the boundaries of comprehensive video understanding, we introduce WorldQA, a diagnostic benchmark dataset challenges machines to answer questions about a video by employing multimodal data and world knowledge. WorldQA is distinguished by three main features: (1) Multimodal Video Input: Success requires analyzing both auditory and visual data. (2) Emphasis on World Knowledge: Questions in the dataset necessitate engagement with broad world knowledge. We identify five knowledge types critical for question answering: societal norms, multimodal associations, selfmotivation, tool use, and social interactions, detailed in Sec. 3.1. (3) Long-Chain Reasoning: The dataset promotes integrating multimodal information and world knowledge across frames for complex reasoning. Currently, the dataset includes 1007 question-answer pairs and 303 videos, with more details in Sec. 3. An initial analysis using GPT-4 (OpenAI, 2023) shows the average reasoning steps in WorldQA to be 4.45, notably higher than the under-two average in other datasets, as demonstrated in Table 2. Evaluation protocols are discussed in the Appendix A.

We propose exploring WorldQA with **WorldRetriver**, employing large language model (LLM) agents (Shen et al., 2023; Surís et al., 2023; Chen et al., 2023; Wu et al., 2023). WorldRetriver breaks down each question into perception- or cognition-oriented tasks. These tasks are then addressed by specialized models—a **multimodal key information retriever** and a **world knowledge retriever**. The LLM integrates their outputs to form a cohesive reasoning chain, answering the question.

Our study presents a comprehensive evaluation of WorldQA, benchmarking WorldRetriver against 13 leading large language models (LLMs) and large multimodal models (LMMs), as well as comparing it to human performance. We focus on two tasks: open-ended and multiple-choice QA. WorldRetriver demonstrates superior performance in both areas, achieving 36.38% accuracy in open-ended QA and 36.59% in multiple-choice QA, surpassing even GPT-4V (Yang et al., 2023b). Key findings include: (1) While WorldRetriver generally outperforms GPT-4V, the latter excels in questions that require complex reasoning, particularly at reasoning steps 8, 9, and 10. (2) Current open-source LMMs exhibit challenges with "consistency" in multiple-choice QA, as discussed in Sec.5.3. (3) Employing GPT-4 to evaluate open-ended QA models correlates well with human judgments, explored in Sec.5.5. (4) Although human performance typically improves with additional frames, our WorldQA and current open-source LMMs show decreased performance under similar conditions, detailed in Sec. 5.5.

098

100

101

102

103

104

105

107

108

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

### 2 Related Work

## 2.1 VideoQA Datasets

Visual Question-Answering (QA) (Antol et al., 2015) is a key task in video-language research, spanning a wide range of datasets. These datasets evaluate various capabilities from basic visual perception, including activity recognition (Yu et al., 2019), concept detection (Xu et al., 2017), and counting (Jang et al., 2017), to more advanced visual reasoning such as compositional (Grunde-McLaughlin et al., 2021), causal (Xiao et al., 2021; Yi et al., 2019; Xu et al., 2021), and situated reasoning (Wu et al., 2021). Beyond using visual information for answering questions, KnowIT (Garcia et al., 2020) incorporates external knowledge from the "Big Bang Theory" in its question design, whereas Social-IQ (Zadeh et al., 2019) leverages both visual and auditory modalities, solving human-centric questions, Unlike KnowIT, which is limited to knowledge from a single TV series, WorldQA draws on a broader range of general world knowledge. Additionally, WorldQA expands beyond the human-centric focus of Social-IQ to encompass a variety of subjects including animals and machines.

Moreover, our analysis of existing videoQA datasets identifies a notable limitation: most require less than two reasoning steps per question, highlighting a gap in their ability to facilitate complex reasoning. To address this, we introduce WorldQA, a dataset specifically designed to challenge models with more intricate reasoning sequences.

#### 2.2 Vision-Language Models

Recent developments in large language models (LLMs) such as GPTs (Radford et al., 2019; Brown et al., 2020), LLaMA (Touvron et al., 2023), and Vicuna (Chiang et al., 2023) have enhanced the efficacy of visionlanguage models like Flamingo (Alayrac et al., 2022; Awadalla et al., 2023) and FrozenBiLM (Yang et al., 2022), especially in zero-shot learning contexts. Researchers are now exploring instruction-tuned models, including Otter (Li et al., 2023a), InstructBLIP (Dai et al., 2023), LLaVA (Liu et al., 2023b) and others (Ye et al., 2023; Li et al., 2023b; Zhang et al., 2023; Muhammad Maaz and Khan, 2023; Gao et al., 2023), to en-

Table 1: **Dataset comparisons.** Reason. stands for reasoning. Avg. stands for average. Q/A stands for question and answer.

Dataset	Multi modal?	World knowledge?	Avg. reason steps	n. Avg. length (Q/A)
MSVD-QA	×	X	1.40	6.6/1.0
TGIF-QA	X	×	1.71	8.7/2.1
TVQA	X	×	1.91	13.4/5.3
ActivityNet-QA	. ×	×	1.62	8.7/1.2
NExT-QA	X	×	1.31	11.6/2.9
Social-IQ	1	×	1.98	11.7/11.4
WorldQA	-	1	4.45	14.2/24.3

hance the interaction abilities of these vision-language models. These models excel in intricate human-model interactions and are ideal for multi-modal chatbot applications.

# **3** WorldQA Dataset

154

155

156

157

158

159

160

163

164

166

167

168

169

170

171

172

173

174

175

178

179

181

In this section, we detail the WorldQA, aimed at creating a benchmark in video question answering for evaluating artificial intelligence (AI) in complex reasoning. The WorldQA comprises 303 videos and 1007 questions. As WorldQA is used solely for evaluation purposes, we believe the number of questions is sufficient. We first describe our thorough annotation process and multiple validation stages, then provide detailed statistics of the WorldQA.

## 3.1 Constructing WorldQA

Video Acquisition Stage The WorldQA sources data from two components: (1) PVSG(Yang et al., 2023a), which comprises 268 third-person videos from the VidOR dataset (Shang et al., 2019) along with 200 egocentric videos sourced equally from Epic-Kitchens (Damen et al., 2018) and Ego4D (Grauman et al., 2022), and (2) user-generated content from YouTube, identified using search terms such as "1 minute movie" and "short movie". In total, this initial dataset encompasses 1000 videos.

Question and Answer Creation Stage Expert annotators formulated questions to test two aspects of QA systems:

1) World Knowledge Understanding: Understanding video content goes beyond the perception level, requir-183 ing an integration of broad world knowledge. This in-184 cludes: a) Tool Use: Understanding the purposes of various tools and concepts, e.g., recognizing that a hammer is for driving nails. b) Societal Norms: Comprehending behaviors, traditions, and unwritten rules within 188 societies, e.g., the custom of handshaking in certain cul-189 tures. c) Self Motive: Identifying personal intentions and motivations, e.g., eating to satisfy hunger. d) Social Interaction: Understanding the subtleties of communication and relationships, such as interpreting social signals. e.g., recognizing that people who cannot speak 194 195 may use written notes to communicate. e) Multimodal

*Association*: Linking vision and hearing to form a complete understanding. *e.g.*, sound of alarm bells with the visual of people evacuating to infer a fire.

**2)** Long Reasoning Chain: Deducing "Who is the murderer?" in a detective movie involves complex reasoning steps. However, current videoQA, like NExT-QA (Xiao et al., 2021) and Social-IQ (Zadeh et al., 2019), often feature basic questions that require minimal reasoning. For example, typical NExT-QA questions, such as "Why are the dogs running around?" with an answer like "Chasing each other", require very few reasoning steps. Our analysis using GPT-4 found the average reasoning depth in NExT-QA to be 1.31, validating this observation.<sup>1</sup> WorldQA, aims to challenge this by including questions that demand multi-step reasoning, where at least two logical steps are necessary to arrive at the answer.

It's important to note that a single question might involve multiple types of world knowledge. Annotators were tasked with providing answers to these questions. To broaden the utility of our dataset, we used GPT-4 (OpenAI, 2023) to transform our questionanswer pairs into a multiple-choice format, as shown in Fig. 2(a).<sup>2</sup>

**Question and Answer Validation** To refine our question set, we established criteria for deletion as follows: (1) Questions that do not require world knowledge, as verified by an annotator different from the question creator; (2) Questions for which GPT-4/ChatGPT responses align with human-annotated answers, as detailed in Section 5.2; and (3) Questions that are solvable in fewer than two reasoning steps, initially verified by an annotator different from the question creator and then further filtered by querying GPT-4, as described in *Appendix* A. This process yielded a dataset of 303 videos and 1007 questions.

In improving multi-choice questions, we aimed to: (1) make all options similar in length, and (2) ensure questions can only be correctly answered by models with visual input. To achieve the second goal, we repeatedly adjusted each option until GPT-4/ChatGPT could not answer correctly. As Table 2 shows, all LLMs (Language Learning Models) scored zero on these questions. However, GPT-4 achieved 35.34 points in the NExT-QA multi-choice test, as explained in the 5.5.

## 3.2 Data Statistics

**Dataset Comparison** WorldQA presents significant advancements over existing datasets, as Table 1 illustrates. Firstly, it requires complex multi-step reasoning. Using GPT-4, we evaluated the reasoning steps in each question-answer pair across datasets; WorldQA averages 4.45 steps, notably higher than others which typically involve less than two steps. This complexity is also evident in answer lengths: while answers in other VideoQA datasets average below five words, those in

<sup>&</sup>lt;sup>1</sup>The details of the prompt are shown in the *Appendix* A.

<sup>&</sup>lt;sup>2</sup>The details of the prompt are shown in the *Appendix* B.



Figure 2: (a) An example for reformating open-ended QA into multi-choice QA. (b) the distribution of different world-knowledge types. (c) the distribution of reasoning step counts.



Figure 3: **WorldRetriever**, an agent designed to synthesize expert knowledge into a coherent reasoning chain for answering questions.

WorldQA average 24.3 words. Secondly, it necessitates more than visual information for success. WorldQA encompasses audio comprehension and world knowledge, expanding its scope beyond mere video visuals for effective question resolution. To our knowledge, it represents the first VideoQA dataset that incorporates questions necessitating world knowledge.

**Knowledge Types and Reasoning Step Statistics** As illustrated in Fig. 2(b), the majority of these questions fall under the "social interaction" category. Furthermore, Fig. 2(c) demonstrates that the reasoning steps in WorldQA vary, ranging from two to ten steps.

# 4 WorldRetriver

254

260

261

263

265

269

271

273

In this section, we introduce WorldRetriver, a method that leverages LLM-as-agent for complex long-chain reasoning. As depicted in Fig. 1, human long-chain reasoning involves gathering information from various sensors and integrating world knowledge to reach a conclusion. WorldRetriver mimics this approach by using expert models for individual sub-tasks. These models perform specific functions, and then WorldRetriver combines their outputs with the original question to formulate the final answer.

274 Multimodal Key Info Retriever This component is
275 crucial for tasks that involve video information. It in276 cludes an image-language model to describe key frames
277 from videos, an audio-language model to capture audio

details, and a speech-language model for interpreting dialogues within the video. These descriptions, combined with the question, allow a pre-defined LLM to generate the *Initial Answer*.

**World Knowledge Retriever** This model transforms questions into search queries for global knowledge databases, *i.e.*, Google, and then succinctly summarizes the search results into "Retrieved Knowledge." In this context, any large language model (LLM) with a knowledge-retrieval function can be used.

**Answer Composition** In the final stage, the *Initial Answer* and the *Retrieved Knowledge*, along with the question, are fed into a pre-defined LLM to generate the final response. The prompt for this stage can be found in the *Appendix* C. Additionally, a "chain-of-thought" prompt is employed to aid the LLM in developing a logical reasoning chain.

# **5** Experiments

### 5.1 Experimental Setup

Our study assesses the performance of diverse models on WorldQA across four settings: (1) Large Multimodal Models (LMMs) for Video Processing: This category includes FrozenBiLM (Yang et al., 2022), Otter-Video(Li et al., 2023a), VideoChat(Li et al., 2023b), Video-LLaMA(Zhang et al., 2023), Video-ChatGPT(Muhammad Maaz and Khan, 2023), and mPLUG-Owl(Ye et al., 2023). These open-sourced models, which are trained with a specific number 278

281

282

306 of frames, typically combine a language model (e.g., LLama(Touvron et al., 2023) or Vicuna(Chiang et al., 307 2023)), a vision encoder (e.g., CLIP(Radford et al., 2021)), and a connector (e.g., linear layer (Liu et al., 2023a)) to convert vision embeddings into "text tokens." 310 311 Additionally, we include the recently API-accessible GPT-4V in our analysis. (2) LMMs for Image In-312 puts: This group comprises Qwen-VL(Bai et al., 2023) 313 and LLaVA-1.5(Liu et al., 2023a). (3) Large Lan-314 guage Models (LLMs): We evaluate models such as 315 Vicuna-v1.5-7B(Chiang et al., 2023) (abbreviated as Vicuna), ChatGPT(OpenAI, 2023), and GPT-4(OpenAI, 2023), focusing on scenarios where their inputs con-319 sist solely of the question or the question accompanied by options. (4) WorldRetriver: This approach 320 321 uses ChatGPT as its predefined LLM and LLaVA-v1.5-7B(Liu et al., 2023a) (abbreviated as LLaVA) for imagetext tasks, Beats(Chen et al., 2022) for audio-text, and Whisper(Radford et al., 2023) for speech-to-text con-325 version. We use LLaVA to describe images, selecting eight frames uniformly. Audio clips are extracted from 326 videos using Pydub(Robert et al., 2018) and analyzed with Beats. Our approach also integrates ReACT(Yao et al., 2022) for world knowledge retrieval. (5) Aug-330 mented LLM and Human Performance: Inspired by MathVista (Lu et al., 2023), our study evaluates human performance alongside three augmented LLMs: Augmented Vicuna/ChatGPT/GPT-4. As mentioned above, 334 the current LMM consists of a language model, a vision encoder, and a connector. We propose that LMM performance might be limited by the vision encoder and connector's ability to translate visual data into "text to-338 kens." To explore this hypothesis, we converted video information into event descriptions annotated by hu-340 mans (for detailed information on event descriptions, please refer to Appendix D). Subsequently, we prompted 341 342 the language models with these descriptions alongside 343 questions to get responses, in what we term "augmented LLM." This experiment helps us estimate the potential maximum performance for LMMs using similar lan-345 346 guage models.

Notably, except for GPT-4V and FrozenBiLM, the other LMMs use a 7B language decoder, similar in size to Vicuna. Video-LLaMA uniquely processes both audio and visual modalities.

### 5.2 Open-Ended QA

352

361

**Definitions and Metrics** Recent studies (Xiao et al., 2021; Zheng et al., 2023) have increasingly turned their attention to generation-based open-ended QA, where answers are not confined to a closed set. However, assessing the quality of open-ended text remains challenging. For instance, current evaluation protocols like NExT-QA may overlook semantic correlations; a response "a cute teddy bear" receives no credit if the ground truth is "a teddy bear." This issue is compounded as answer length increases.

In the field of natural language generation (NLG), recent initiatives have investigated the use of GPT-4 for

Table 2: Evaluation of Large Multimodal Models (LMMs), Language Models (LLMs), and WorldRetriver in WorldQA. We introduce two upperbounds for comparison: LLMs augmented with huamnannotated event descriptions, labeled as (Aug.) X, and the human performance. The best model in each group is highlighted in blue, while the overall top performer in all tasks is marked in red. Different types of inputs include: Q for Question, V for Video, I for Image, and  $V_d$  for Video Description. Language model param. indicates the parameter of the language model.

Model (language model param.)	Input	Open	Multi
	mput	ended ↑	choice ↑
Large Multimodal Mode	els (LM	Ms)-Video	)
FrozenBiLM (900M)	Q, V	8.21	0.32
Otter-Video (7B)	Q, V	24.22	6.11
VideoChat (7B)	Q, V	24.43	1.29
LLaMA-Adapter (7B)	Q, V	25.87	12.04
Video-LLaMA (7B)	Q, V	26.80	4.81
Video-ChatGPT (7B)	Q, V	28.51	13.25
mPLUG-Owl (7B)	Q, V	31.89	0.75
GPT-4V(ision) (-)	Q, V	35.37	32.83
Large Multimodal Models (LMMs)-Image			
Qwen-VL (7B)	Q,I	24.04	12.80
LLaVA (7B)	Q,I	31.31	0.30
Large Language Models (LLMs)			
Vicuna (7B)	Q	22.44	0.00
ChatGPT (20B)	Õ.	24.24	0.00
GPT-4 (-)	$\tilde{Q}$	28.73	0.00
LLM Agent			
Ours (ChatGPT as LLM) (20B)	<i>Q</i> , <i>V</i>	36.38	36.59
Upper Bound with Human Transcription			
(Aug.) Vicuna (7B)	$Q, V_d$	38.71	23.90
(Aug.) ChatGPT (20B)	$Q, V_d$	46.50	46.06
(Aug.) GPT-4 (-)	$Q, V_d$	48.46	56.06
Human-Level Performance			
Human	Q, V	72.43	88.79

assessing the quality of open-ended model-generated responses (Zheng et al., 2023; Xie et al., 2023). For evaluating open-ended QA, we also use GPT-4. Our scoring system for model answer A against ground truth G is: (1) A = G: Correct (1 point), (2)  $A \cap G = \emptyset$ : Incorrect (0 points), (3)  $\emptyset < A \cap G < A \cup G$ : Partially correct (0.3 points), (4)  $A \subset G$ ,  $A \neq G$ : Incomplete but correct (0.5 points), (5)  $G \subset A$ ,  $A \neq G$ : Redundant (0.5 points). GPT-4's scoring is exemplified in Fig. 5.

Main Results Our open-ended QA evaluation revealed several key insights: 1) Superiority of WorldRetriver and GPT-4V: Our method, WorldRetriver, using ChatGPT, surpasses the ChatGPT baseline by 12.14 points and GPT-4V by 1.01 points, as shown in Table. 2. In Fig. 4, WorldRetriver and GPT-4V outperform the leading open-source LMM, mPLUG-Owl, particularly in reasoning beyond five steps. While Worl-

378

379

380



Figure 4: Comparative performance of advanced LMM and our method across increasing reasoning steps.

381

387

391

401

402

406

407

408

409

411

414

417

421

dRetriver beats GPT-4V in the overall performance, the latter shows strength in complex reasoning tasks (steps 6 to 9). However, the closed-source nature of GPT-4V suggests future research opportunities to understand these differences. 2) Non-Zero Performance in Large Language Models (LLM): Although WorldQA is a videoQA dataset, it is observable that LLMs can still achieve a certain level of accuracy by using only the questions as input, which we consider reasonable. For instance, in response to the question, "What did the lady do when she left home?", an LLM might reply, "She may have gone shopping or to work." While this response is not entirely accurate, it closely approximates the actual answer, "She went to work." However, as mentioned in Sec. 3.1, questions that could be answered with 100% accuracy by LLMs were excluded. 3) Limitations of Current Video-Based LMM in Handling Multiple Frames: We found that the image-based LLaVA model outperforms most video-based models in performance, which is surprising given that video-based 400 LLMs are capable of processing multiple frames. This leads to a pivotal question: Do the tasks in WorldQA require only a minimal number of frames, or do current 403 video-based LMMs struggle to use multiple frames ef-404 405 fectively? Sec. 5.5 presents experimental evidence supporting the latter. 4) Impact of Language Model Size: When compared to LMMs with a 7B-language model, Vicuna's lower performance is expected due to its lack of vision information for answering questions. However, GPT-4 surpasses six out of the eight LMMs. This un-410 derscores the significant benefits of GPT-4's advanced language processing capabilities. It also suggests that 412 413 incorporating a languge model as powerful as GPT-4's could significantly advance the capabilities of current LMMs. 5) Potential for Improvement: By employ-415 ing the same language model, namely Vicuna, our aug-416 mented Vicuna model surpasses LLaVA by 7.4 points. This result underscores a significant opportunity for im-418 proving LMMs. Moreover, even the augmented GPT-4 419 reaches merely 67% of human performance, suggesting 420 a considerable scope for advancing current models.

#### 5.3 Multi-Choice QA

Definitions and Metrics Compared to open-ended QA, multi-choice QA tasks simplify the QA task, as the correct answer is always among the provided options. We follow the approach of MMBench (Liu et al., 2023c) and use its proposed CircularEval evaluation method to evaluate model performance. CircularEval requires the model to answer each question N times, where N is the number of choices. Each iteration involves a circular shift of the options, creating a different arrangement. This technique mitigates the effect of random guessing, which could otherwise lead to a 25% Top-1 accuracy rate in scenarios with four choices. If the model's response does not match any of the given options (e.g., A, B, C, D), ChatGPT evaluates semantic similarity to determine the most appropriate choice. For more details on CircularEval, please refer to MMBench.

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

Main Results In Table 2, WorldRetriver performs better than other models in multi-choice QA, showing a notable 3.76-point lead over GPT-4V. Also, we made sure each multi-choice question could only be correctly answered by models that can process visual input, resulting in LLMs scoring zero. Our analysis highlights important insights: 1) Consistency Challenge in Open-Source LMMs: Notably, while mPLUG-Owl and LLaVA show strong performance in open-ended tasks, their effectiveness decreases in multi-choice OA scenarios. We propose that this decline is likely due to inconsistent choice when the order of options is varied. To illustrate this issue, we introduce the "consistency rate" metric in Fig. 6, which quantifies how often a model selects the same option across different arrangements of the N options, regardless of the answer's correctness. As demonstrated in this figure, there is a significant correlation between the consistency rate and accuracy in multi-choice QA. This finding highlights the importance of a model's ability to consistently select the same answer as a key indicator of its proficiency in CircularEval. 2) Consistency Loss in Fine-Tuned Language Model: Among LMMs, LLaVA and mPLUG-Owl are unique in tuning their language models during the instruction tuning phase. However, this approach results in notably poorer consistency and, consequently, inferior performance in multi-choice QA tasks. A direct comparison between LLaVA, which uses the Vicuna, and Vicuna itself reveals a significant drop in consistency for LLaVA. This suggests that tuning language decoders during the instruction tuning stage can adversely affect the model's overall consistency. 3) Potential for Improvement: With human benchmarks at 88.79, there is substantial scope for models to match or exceed human performance in video comprehension.

#### 5.4 Ablation Studies of WorldRetriver

In this subsection, we examine the impact of distinct components within the WorldRetriver.



Figure 5: Examples of how does GPT-4 score in the open-ended question.



Figure 6: Analysis of the correlation between multichoice QA performance of models and their consistency, defined as the frequency of selecting the same response N times across varied sequences of N options.

	Open-ended $\uparrow$	Multi-choice $\uparrow$
Ours	36.38	36.59
<ul> <li>world-knowledge</li> </ul>	34.78 (-1.60)	35.46 (-1.13)
- speech-text	33.45 (-1.33)	34.12 (-1.34)
- audio-text	33.23 (-0.22)	34.45 (+0.23)
- image-text	30.98 (-2.25)	2.23 (-32.22)

As shown above, our findings indicate that the imagetext model is the most critical, while the audio-text model contributes the least to overall performance. Notably, the enhancement provided by the audio-text component is negligible. To investigate the cause, we scrutinized the output of the audio-text model, which is responsible for categorizing audio within every clips in videos. The analysis reveals that the prevalent audiotext model, Beats, seldom produces accurate classification labels for video audio content. Additionally, leveraging the capabilities of expert models—Whisper and ReACT—might be the core reason why WorldRetriver outperforms GPT-4V.

#### 5.5 Further Analysis

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493 494

495

496

497

The Relation of Reasoning Steps/Knowledge Type and Performance We investigated the impact of the number of reasoning steps on a model's ability to answer questions, specifically examining the correlation between reasoning steps and performance. In Fig. 7a, the results shown for each step are averaged from the models in LMMs-Video, LMMs-Image, WorldRetriver, and Augmented-LLM in Table 2. Our findings reveal that as the number of reasoning steps increases, performance deteriorates. Notably, by the time it reaches 10 reasoning steps, the average score is almost zero in open-ended tasks. 498

499

500

501

502

503

504

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

535

536

537

538

Additionally, we examined the average performance of models across five distinct word-knowledge types. As shown in Fig. 7b, it's notable that the poorest average performance is observed in the "multimodal association" word-knowledge. This is likely because current models are unable to handle audio/speech signals

**Does GPT-4's scoring align with human preferences?** In Sec. 5.2, a question arises: does the model that achieves the highest score in open-ended QA, as evaluated by the predefined GPT-4 scoring system, truly align with human preferences? To investigate this, we employed the Model Arena approach (Zheng et al., 2023; Xu et al., 2023), which involves rating models based on human preferences in a side-by-side comparison of two model-generated answers. The Elo rating system, employed in Model Arena, aggregates these judgments to rank models. If the rankings in Table. 2 for the openended QA task correspond with those derived from the Model Arena, it would validate the effectiveness of the GPT-4 scoring system in evaluating LLM-generated responses.

In our study, we conducted five rounds of Model Arena for each question pertaining to six video-based LLMs, huamn judges reviewed related videos for making the decision. This procedure was replicated across 4,255 comparisons.

The results, shown in Fig. 7c, primarily classify the responses as "tied (bad)", suggesting that current LMMs fail to produce high-quality responses. A comparison of the Elo rankings with the open-ended QA performance in the LLMs-video section (Table 2) shows a significant correlation. This agreement with open-ended QA scores validates our evaluation method for open-ended answers.



(a) Illustration of how performance declines as the number of reasoning steps increases, in both open-ended and multiple-choice tasks.



(c) (upper) Comparative counts of model choices from Human evaluators. (down) Elo scores assigned to various models from Human evaluators.



(b) A breakdown of performance across various knowledge types, highlighting the notably lower scores in the multimodal association category.



(d) Comparative analysis of model performance in multiple-choice QA tasks under varying frame constraints. "no constraint" for models indicates their optimal performance, whereas for humans, it denotes the ability to answer questions after viewing the entire video.

Figure 7: Key findings emerged from the further analysis.

**Do More Frames Impair Performance?** In Sec. 5.3, we pose a question: do the tasks in WorldOA require only a minimal number of frames, or do current LMMs struggle with effectively using multiple frames? Our experiments suggest the latter. As depicted in Fig. 7d, a distinct pattern is observed: the performance of humans and augmented ChatGPT-which selects event descriptions based on the time period of sampled framesenhances on WorldQA when more frames are used. In contrast, Video-ChatGPT and our proposed methods exhibit peak performance at approximately 16 frames. This suggests a limitation in how current models process multiple frames.

540

541

542

545

546

547

548

550

551

555

556

557

560

561

562 563

552	Can Uni-	modal m	odel ans	wer?	
552		Social-IQ	NextQA	KnowIT	Ours(WorldQA)
555	GPT-4	29.47	35.34	36.44	0.00

In Sec. 3.2, we mention that for multi-choice questions, we repeatedly adjusted each option until GPT-4/ChatGPT was unable to provide an answer. This approach ensures that each question can only be answered by models equipped with visual input capabilities. In contrast, we noted that many existing videoQA datasets contain questions that GPT-4 can answer correctly without visual information. We explored this issue using CircularEval, and the results of these experiments are detailed above.

#### **Discussion and Conclusion** 6

In this study, we introduce WorldQA, an innovative dataset designed to assess the ability of visual-language models in understanding videos. WorldOA distinctively emphasizes the integration of multimodal information and the application of world knowledge for complex reasoning, reflecting a crucial aspect of human intelligence. Concurrently, we present WorldRetriver, a technique inspired by human approaches to video understanding. Our experiments with WorldQA reveal that current models still fall short of human-like proficiency in video understanding.

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

#### 7 Limitations

The videos in WorldQA are sourced from various platforms, including Ego4D(Grauman et al., 2022), Epic-Kitchens(Damen et al., 2018), VidOR (Shang et al., 2019), and YouTube. This diversity introduces potential biases inherent in these sources. Furthermore, there is a concern regarding the potential skew in the questionanswer pairs, possibly influenced by the annotators' perspectives. Additionally, the significant observation that WorldRetriver struggles to process multiple frames highlights a crucial area for future improvement.

## References

- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. Advances in Neural Information Processing Systems, 35:23716–23736.
- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In Proceedings of the IEEE international conference on computer vision, pages 2425–2433.
- Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. 2023. Openflamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966.*
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. Advances in neural information processing systems (NeuIPS), 33:1877–1901.
- Liangyu Chen, Bo Li, Sheng Shen, Jingkang Yang, Chunyuan Li, Kurt Keutzer, Trevor Darrell, and Ziwei Liu. 2023. Language models are visual reasoning coordinators. In *ICLR* 2023 Workshop on Mathematical and Empirical Understanding of Foundation Models.
- Sanyuan Chen, Yu Wu, Chengyi Wang, Shujie Liu, Daniel Tompkins, Zhuo Chen, and Furu Wei. 2022. Beats: Audio pre-training with acoustic tokenizers. *arXiv preprint arXiv:2212.09058*.
  - Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023. Vicuna: An opensource chatbot impressing gpt-4 with 90%\* chatgpt quality. See https://vicuna. lmsys. org (accessed 14 April 2023).
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. 2023. Instructblip: Towards generalpurpose vision-language models with instruction tuning. *arXiv preprint arXiv:2305.06500.*
- Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, et al. 2018. Scaling egocentric vision: The epic-kitchens dataset. In *Proceedings of the European conference on computer vision (ECCV)*, pages 720–736.
- Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan Lu, Conghui He, Xiangyu Yue, Hongsheng Li, and Yu Qiao. 2023. Llama-adapter v2: Parameter-efficient visual instruction model. arXiv preprint arXiv:2304.15010.

- Noa Garcia, Mayu Otani, Chenhui Chu, and Yuta Nakashima. 2020. Knowit vqa: Answering knowledge-based questions about videos. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 10826–10834.
- Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. 2022. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18995– 19012.
- Madeleine Grunde-McLaughlin, Ranjay Krishna, and Maneesh Agrawala. 2021. Agqa: A benchmark for compositional spatio-temporal reasoning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11287–11297.
- Yunseok Jang, Yale Song, Youngjae Yu, Youngjin Kim, and Gunhee Kim. 2017. Tgif-qa: Toward spatio-temporal reasoning in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2758–2766.
- Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. 2023a. Otter: A multimodal model with in-context instruction tuning. *arXiv* preprint arXiv:2305.03726.
- Kunchang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. 2023b. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023a. Improved baselines with visual instruction tuning.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*.
- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. 2023c. Mmbench: Is your multimodal model an all-around player? arXiv preprint arXiv:2307.06281.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2023. Mathvista: Evaluating math reasoning in visual contexts with gpt-4v, bard, and other large multimodal models. *arXiv preprint arXiv:2310.02255*.
- Salman Khan Muhammad Maaz, Hanoona Rasheed and Fahad Khan. 2023. Video-chatgpt: Towards detailed video understanding via large vision and language models. *ArXiv* 2306.05424.
- OpenAI. 2023. Chatgpt based on the gpt-4 architecture. Https://openai.com/research/.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning* (*ICML*), pages 8748–8763. PMLR.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *International Conference on Machine Learning*, pages 28492– 28518. PMLR.

706

707

587

588

590

592

594

596

601

609

610

611

612

613

614

615

616

617

618

619

621

622

623

624

625

626

634

636

642

803

804

805

806

807

765

766

- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- James Robert, Marc Webbie, et al. 2018. Pydub.

708

710

711

712

713

714

715

716

718

721

723

724

725

727

730

731

732

733

734 735

736

740

741

742

743

744

745

746

747

750

751

754

755

756

763

- Xindi Shang, Donglin Di, Junbin Xiao, Yu Cao, Xun Yang, and Tat-Seng Chua. 2019. Annotating objects and relations in user-generated videos. In *Proceedings of the 2019 on International Conference on Multimedia Retrieval*, pages 279–287. ACM.
- Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. 2023. Hugginggpt: Solving ai tasks with chatgpt and its friends in huggingface. *arXiv preprint arXiv:2303.17580*.
- Dídac Surís, Sachit Menon, and Carl Vondrick. 2023. Vipergpt: Visual inference via python execution for reasoning. arXiv preprint arXiv:2303.08128.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971.
- Bo Wu, Shoubin Yu, Zhenfang Chen, Joshua B Tenenbaum, and Chuang Gan. 2021. Star: A benchmark for situated reasoning in real-world videos. In *Thirty-fifth Conference* on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2).
- Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. 2023. Visual chatgpt: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*.
- Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. 2021. Next-qa: Next phase of question-answering to explaining temporal actions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9777–9786.
- Binzhu Xie, Sicheng Zhang, Zitang Zhou, Bo Li, Yuanhan Zhang, Jack Hessel, Jingkang Yang, and Ziwei Liu. 2023. Funqa: Towards surprising video comprehension. arXiv preprint arXiv:2306.14899.
- Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. 2017. Video question answering via gradually refined attention over appearance and motion. In *Proceedings of the 25th ACM international conference on Multimedia*, pages 1645–1653.
- Li Xu, He Huang, and Jun Liu. 2021. Sutd-trafficqa: A question answering benchmark and an efficient network for video reasoning over traffic events. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9878–9888.
- 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 vision-language models. *arXiv preprint arXiv:2306.09265*.
- Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. 2022. Zero-shot video question answering via frozen bidirectional language models. In *NeurIPS*.

- Jingkang Yang, Wenxuan Peng, Xiangtai Li, Zujin Guo, Liangyu Chen, Bo Li, Zheng Ma, Kaiyang Zhou, Wayne Zhang, Chen Change Loy, et al. 2023a. Panoptic video scene graph generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18675–18685.
- Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. 2023b. The dawn of lmms: Preliminary explorations with gpt-4v (ision). *arXiv preprint arXiv:2309.17421*, 9.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2022. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*.
- Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, Chaoya Jiang, Chenliang Li, Yuanhong Xu, Hehong Chen, Junfeng Tian, Qian Qi, Ji Zhang, and Fei Huang. 2023. mplug-owl: Modularization empowers large language models with multimodality. *Preprint*, arXiv:2304.14178.
- Kexin Yi, Chuang Gan, Yunzhu Li, Pushmeet Kohli, Jiajun Wu, Antonio Torralba, and Joshua B Tenenbaum. 2019. Clevrer: Collision events for video representation and reasoning. arXiv preprint arXiv:1910.01442.
- Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao. 2019. Activitynet-qa: A dataset for understanding complex web videos via question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 9127–9134.
- Amir Zadeh, Michael Chan, Paul Pu Liang, Edmund Tong, and Louis-Philippe Morency. 2019. Social-iq: A question answering benchmark for artificial social intelligence. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 8807–8817.
- Hang Zhang, Xin Li, and Lidong Bing. 2023. Video-Ilama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-ajudge with mt-bench and chatbot arena. *arXiv preprint arXiv:2306.05685*.

808

#### 809 810

811

812

817

818

819

820

822

# A Prompt for Reasoning Hop

The prompt for getting the required reasoning hop for each question is shown in Table. 4. An example of a question-answer pair with ten reasoning steps is shown in Fig. 9.

# **B** Prompt for Multi-choice QA

The prompt for getting the required reasoning hop for each question is shown in Table. 5.

# C Prompt for Answer Composition

# System Message

As an AI visual assistant, your task involves synthesizing an answer to a question about a video by integrating responses from two expert models: Model A and Model B. Model A provides the initial answer to the question, while Model B focuses on contributing additional external knowledge to solve the question. Let's begin this task: Question: {question} Response from Model A: {Model\_A\_response} Response from Model B: {Model\_B\_response} Step-by-step reasoning: <reasoning process> Answer: <answer> Let's think step by step.

Table 3: System message for creating answer composition. {XXX} is the placeholder for specific information, *i.e.*, question, response from Model A and Model B.

# **D** Dataset Construction

In each video, we ask annotators to describe the events occurring, along with the timestamps. An example of such an event description is shown in Fig. 8:

# E Case Study

More examples of question-answer pairs are displayed in Fig. 8-16. Additionally, the full videos for each example can be viewed at https://www.youtube.com/ watch?v=NXbJLLf9E\_I.

## System Message

As an AI visual assistant, your objective is to determine the number of reasoning steps—each representing a logical or causal link—between a given question and paired answer. In this context, a reasoning step elucidates the process through which the answer relates to the question. The goal is to quantify the depth of reasoning, whereby a direct causation between the answer and question constitutes one step, while the presence of intermediate steps to establish a connection warrants a higher count.

Let's begin this task, you reasoning steps should be as fewer as possible. In your reasoning step, do not generate anything unmentioned in the question and answer.

You format should be : Reasoning steps: <reasoning steps> Number of Reasoning steps: <number of Reasoning steps> Question: {question} Answer: {answer}

### **In-context Examples**

Augustion: why did the kid drink water? Answer: the kid is thirsty Reasoning steps: 1. When people feel thirsty, they want to drink water Number of Reasoning steps:1

Question: Why did the kid touch the cup? Answer: the kid is thirsty Reasoning steps:

1. When people feel thirsty, they want to drink water.

2. Cups are usually used to collect water. Number of Reasoning steps:2

Question: Why does the tank turn red?
Answer: It stands to reason that there was a piranha in the tank, The piranha bit the man in the tank.
The person in the water tank bleeds because of this.
Blood turns the water tank red.
Reasoning steps:

There was a piranha in the tank.
Piranha bit the man in the tank.

He person in the water tank bleeds because of this.

Number of Reasoning steps:4

Table 4: System message and in-context exemplars for reasoning step prompting.

#### System Message

As an AI visual assistant, your task involves first analysing video content, including various textual descriptions, and then rewriting the given question into multi-choice questions (with four options and only correct answer) related to the video.

##Video Textual Description##

1. Event Descriptions: Descriptions of specific video segments identified by their timestamps (start and end times).

##Guidelines##

While executing this task, please adhere to the following guidelines:

1. The answer to each question should be in the form of a single letter: A, B, C, or D.

2. All the options you provide should be roughly the same length.

3. The choices you present should be formulated in a way that makes them tricky to differentiate, thus creating some confusion for the individual answering the question.

4. You should rewriting the given question even if you think it do not seem to match the content described.

5. The correct option should follow the answer of the given question.

6. Starting your response without saying anything unrelated to the output format.

##Format## Your output format should be like: Question: <question> Option A.XXX **B.XXX** C.XXX D.XXX Answer: <answer> Let's begin this task, you should rewrite all the below questions (in the (2)) into multi-choice question, based on these video event descriptions(in the (1)). (1) Event Descriptions: {event\_descriptions} (2) The questions for rewriting: Question: {question} Answer: {answer}

Table 5: System message for creating multi-choice question. {XXX} is the placeholder for specific information, *i.e.*, event description, question and answer.



Figure 8: Example 1.



Figure 9: Example 2.



Figure 10: Example 3.



Figure 11: Example 4.

t=45 t=185 t=215 t=275 t=325
t=405 t=535 t=555 t=595 t=595
1 35 35 35 X 35 X
Questions
Question: Why did the man apologize to the toilet paper and show a pained and remorseful expression?
G I Answer: Because the toilet paper with a smiling face brought joy to his duil life, so he considered it as a good friend. However, when
ne nad diarrinea, ne nad to sacrifice this intimate friend, leading to his feeling of pain and seif-reproach.
GP1-4 generated reasoning step:
1. The tonet paper with a smilling race brought joy to the mark soull life.
2. The man bad discharge the tollet paper as a good mend.
A the man had to carrifice the toilet namer as a result of his diarrhea
5. Sarrificing the intimate friend caused the main to feel main and self-renroach
Bequired Knowledge Type:
Self Motive: Social Interaction
Question: Why did the man adult feel distressed when he found the bag in the restroom empty?
GT Answer: The man was distressed when he found the bag in the restroom empty because he had run out of toilet paper. The only
remaining toilet paper was the one with a smiling face that he had developed an emotional attachment to, and he was reluctant to use it.
GP1-4 generated reasoning step:
1. The man had run out of toilet paper.
2. The only remaining toilet paper was the one with a smilling face.
<ol> <li>The main had beveruped an endutional attachment to the tonet paper with a similing face.</li> <li>The map was rejustant to use that failet paper with a similar face.</li> </ol>
4. The man was relaciant to use the tone paper with a similar access 5. The man field distressed when he found the bag in the restroom emity.
Required Knowledge Type:
Self Motive: Social Interaction
CDestion: what historical significance does this tim bear?
G I Answer: From the onset, the film states it is "Quarantine Day 22", highlighting the efforts of an individual in self-isolation during the
COVID-19 pandemic to regulate their emotions.
Gr 1-4 generated reasoning step:
1. The nime is set during the COVID-19 pandemic.
2. The mini recurses on an individual in Set-Solation.
Required Knowledge Type:
Social Interaction
Social interaction

Figure 12: Example 5.

t=25 $t=145$ $t=165$ $t=215$ $t=215$ $t=255$
Questions
Question: Why did the man repeatedly scratch his head? GT Answer: Because the man encountered questions he didn't know how to solve, while it was also close to the end of the exam, causing him to feel nervous and frustrated about not finishing the questions. GPT-4 generated reasoning step:
<ol> <li>The man encountered questions he didn't know how to solve</li> <li>It was close to the end of the exam</li> <li>Feeling nervous and frustrated about not finishing the questions</li> <li>Scratching his head repeatedly</li> <li>Required Knowledge Type:</li> <li>Self Motive</li> </ol>
Question: Why does the man have wandering eyes?         GT Answer:         Because it's nearing the end of the exam, and the man hasn't finished answering the questions yet leading to inner anxiety.         GPT-4 generated reasoning step:         1. It's nearing the end of the exam.         2. The man hasn't finished answering the questions yet.         3. The situation of not finishing the exam leads to inner anxiety.         4. Inner anxiety causes wandering eyes         Required Knowledge Type:         Societal Norm; Self Motive
Question: Why is the man observing what his classmates are doing?         GT Answer:         Because it's nearing the end of the exam, and the man wants to know how his classmates are progressing with their answers.         GPT-4 generated reasoning step:         1. It's nearing the end of the exam.         2. The man wants to know how his classmates are progressing with their answers         Required Knowledge Type:         Societal Norm; Self Motive;
Question: The man finally choose to turn back but then the bell rings. Why he seems so frustrated?         GT Answer:         Because when he finally decides that he want to cheat, the time is up and leaves him no time to cheat.         GPT-4 generated reasoning step:         1. It's nearing the end of the exam.         2. The man wants to know how his classmates are progressing with their answers         Required Knowledge Type:
Self Motive; self Motive; Multi-modal Association

Figure 13: Example 6.



Figure 14: Example 7.



Figure 15: Example 8.



Figure 16: Example 9.