

CharacterQA: A Corpus for Multimodal Character Conversational Movie Question Answering

Anonymous ACL submission

Abstract

The rapid advancement of Large Language Models has sparked extensive exploration of their applications across various fields. Among them, the personalized conversation based on characters in movies is an attractive research area. To achieve such comprehensive conversations, the integration of extensive multimodal information, notably visual content alongside textual data, is crucial. This necessity underlines the significance of multimodal insights for enriching personalized conversations, thereby further emphasizing the urgent need for a sophisticated multimodal character conversational dataset. To this end, we introduce CharacterQA, a novel video question-answering (QA) dataset for multimodal character conversation in movies. The dataset consists of 101 selected Chinese movies, each of which is annotated with the main character profiles, the character information of the scripted conversations and their timestamps. Furthermore, a set of questions from various designed tasks and their detailed answers are annotated. Most of those questions require taking into account visual signals for logical comprehension of movie characters and plots. Subsequently, we adopt an advanced multimodal large language model MovieGPT to evaluate the CharacterQA dataset. The results yield insightful findings that are expected to drive further development of multimodal large language models in the character conversation field.

1 Introduction

In the past few years, the development of social media has greatly contributed to the user demand for personalized character conversation, and consequently inspired significant attention from researchers. Recently, breakthroughs in pre-trained large language models (LLMs) have led to a paradigm shift in the natural language processing community, which brings novel challenges for character conversation (Brown et al., 2020; OpenAI,

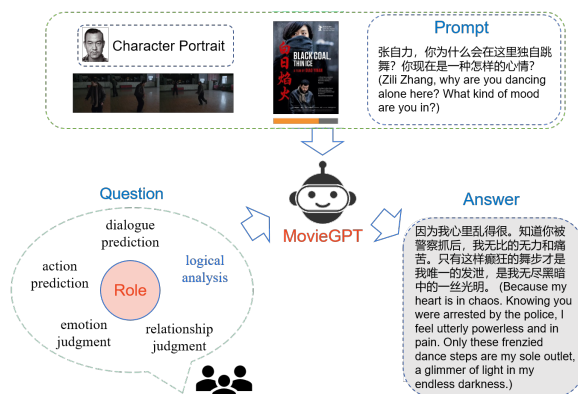


Figure 1: An example of multimodal character QA. Our CharacterQA dataset enables MovieGPT to perform character conversations by watching movies and respond to personalized questions in different scenes.

2023a; Touvron et al., 2023a; Liu et al., 2023b; Peng et al., 2023). Currently, personalized character conversation mainly focuses on text-only domains (e.g., character.ai) (Shao et al., 2023; Park et al., 2023), largely overlooking diverse multimodal applications in real-world scenarios.

As shown in Figure 1, given a video taken from a movie, a user may wish to play the character “Zhizhen Wu” and ask the character “Zili Zhang”: “Zili Zhang, why are you dancing alone here?”. If only the textual information is available, the model is unable to answer such a question since capturing the nuances of expressions and body language during dance requires understanding visual content. Conversely, compared to text-only settings, multimodal character conversations offer enhanced vividness and practicality with the visual content, facilitating much easier user interactions. Regrettably, prevailing datasets fail to equip models with such character conversation capabilities.

In light of the above considerations, we construct the first multimodal character conversation dataset CharacterQA in this paper. Diverging from existing multimodal conversation datasets primarily based

on pure English language (Tapaswi et al., 2016; Lei et al., 2018; Castro et al., 2022; Xiao et al., 2021), this dataset is derived from 101 Chinese films obtained from online platforms. Specifically, we annotate attributes such as the names, roles, and personalities of the main characters in each film, as well as their dialogues with corresponding timestamps. To evaluate character conversation ability, five tasks of varying difficulty are designed, including dialogue prediction, action prediction, relationship judgment, sentiment analysis, and logical analysis. Among them, the sentiment analysis is multiple choice questions, while the others are open-ended questions. Some questions are relatively straightforward as the answers can be found within the dialogues, while others pose greater challenges. These challenging questions require a deep understanding of the movie content, character profiles, and long conversation contexts. This depth of comprehension is necessary to capture the nuances and unique styles of the characters’ language. Furthermore, certain questions even require the ability to reason across the dialogues and the movie content based on broader commonsense knowledge related to the question.

We also developed a multimodal LLM called MovieGPT, and evaluated it alongside various LLMs using our CharacterQA dataset to evaluate their character conversation abilities. Thorough analysis of the results indicates that the CharacterQA dataset poses significant challenges for multimodal character conversation, and existing LLMs are insufficient for character portraits and visual understanding in multimodal scenarios.

This study makes several contributions. Firstly, we present a novel dataset CharacterQA for Chinese multimodal character conversation, comprising five distinct designed tasks that emphasize the understanding of character traits and multimodal content. Secondly, we introduce a multimodal pre-trained LLM MovieGPT tailored for character conversation. Thirdly, we conduct extensive evaluations on CharacterQA with MovieGPT and several mainstream LLMs, emphasizing the challenges inherent in the multimodal character conversation.

2 Related Work

Multimodal Conversational LLMs and Character-play Datasets. The success of LLMs has catalyzed advancement in multimodal conversational LLMs, such as Flamingo (Alayrac et al.,

2022), LLaVA (Liu et al., 2023a), MiniGPT-4 (Zhu et al., 2023), BLIP2 (Li et al., 2023b), and mPLUG-Owl (Ye et al., 2023d). These methods have extensively explored the visual encoders and training strategies of multimodal LLMs. However, they were not originally tailored for character-play scenarios, and previous evaluations reveal a deficiency in their capacity for robust character-playing (Shen et al., 2023; Huang et al., 2023a; Wang et al., 2023b). Concurrently, although the potential for character-playing within the LLMs has been acknowledged, the existing character-playing datasets are limited to text-only formats, lack multimodal annotations, and feature a small number of characters (often less than 200) (Tu et al., 2024; Chen et al., 2023b). These limitations highlight the importance of our proposed CharacterQA, which aims to enhance training and evaluation for character-play capabilities.

Multimodal QA Datasets. Several datasets focusing on multimodal video QA have been developed, including MovieQA (Tapaswi et al., 2016), VideoQA (Zhu et al., 2017), TV-QA (Lin et al., 2023), Life-QA (Castro et al., 2020), NextQA (Xiao et al., 2021), and Wild-QA (Castro et al., 2022). As shown in Table. 1, the existing datasets primarily focus on the visual comprehension capabilities of models, lacking the necessary annotations of the intrinsic profiles of characters essential for multimodal character conversation. Consequently, achieving an effective evaluation of this task becomes challenging. Focusing on multimodal character-based conversation, our dataset includes detailed character profiles, manually curated conversational texts, complete sets of movies, and annotations for five distinct character-based conversation tasks, as described in Section 3.

3 CharacterQA Dataset

3.1 Dataset Summary

Our CharacterQA dataset comprises 101 carefully selected Chinese movies, with an average duration of 102 minutes per movie. Among these movies, 90 are dubbed in Standard Mandarin, while 11 are dubbed in various regional Chinese dialects. The selected movies span a range of release dates, from as early as 1984 to the most recent in 2023. As shown in Figure. 2, each movie contains an average of 405 lines of dialogue, with each line annotated with its timestamp and the corresponding character. To accommodate the character conversation task,

Dataset	Lang.	Domain	Annotation	QA Type	Role Inf	Video#	QA#	Dur.(s)
MovieQA	En	Movie	Man	MC	N	6.7k	6.4k	203
VideoQA	En	Web	Aut	MC	N	109k	390k	33
TVQA	En	TV Shows	Man	MC	N	21.8k	152k	76
LifeQA	En	Daily Activaty	Man	MC	N	0.3k	2.3k	74
NExTQA	En	Daily Activity	Man	OE+MC	N	5.4k	52k	44
WildQA	En	Wild Activity	Man	OE	N	0.4k	0.9k	71
Ours	Ch	Movie	Aut+Man	OE+MC	Y	101	25k	6024

Table 1: Comparison between our dataset and representative existing datasets for videoQA. “Lang.” denotes the language of the data, “En” for English, and “Ch” for Chinese. “Annotation” indicates whether the data is annotated manually or automatically. “Aut” stands for automatical and “Man” stands for manual. “QA Type” denotes whether the answers are multiple-choice (MC) or open-ended (OE). “Dur. (s)” is the average duration of the videos in seconds. “Role Inf” is the character information.

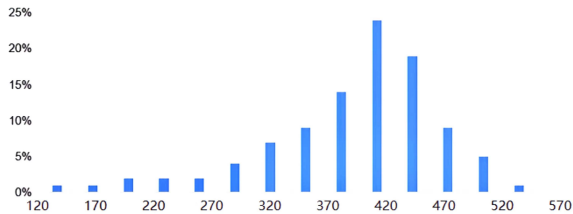


Figure 2: Illustration of dialogue counts for each movie.

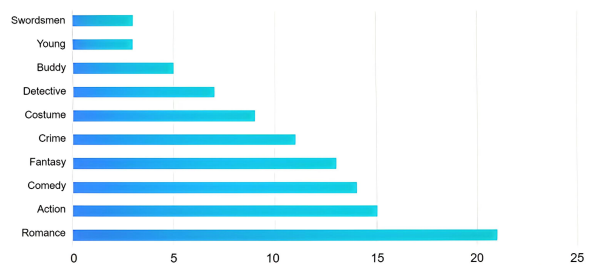


Figure 3: Illustration of genres for the selected 101 movies.

we annotate detailed character profiles delineating the names, personalities, and backgrounds of 3 to 5 main characters within each movie. Each profile is carefully crafted to summarize and portray their distinctive characteristics in about 70 words of Chinese description.

Based on the detailed profiles and the labeled dialogues, five different video question-answering tasks are designed, including dialogue prediction, action prediction, relationship judgment, sentiment analysis, and logical analysis. Particularly, there are 25,000 QA pairs in total for those five tasks, with 5,000 pairs per task evenly distributed across the 101 movies. Note that the selected movie clips used for task annotations are evenly distributed throughout the entire movie, facilitating subsequent model training and validation. Each QA instance corresponds to a one-minute selected movie clip, the plot of which is consecutive to avoid escalating the difficulty of video comprehension.

3.2 Dataset Details

As illustrated in Figure 3, to ensure the diversity of the dataset, we strive to gather a wide spectrum of Chinese movies spanning across different genres, including romance, action, comedy, fantasy, etc..

Note that a portion of our dataset is selected from the Movie101 dataset (Yue et al., 2023). In order to obtain dialogue transcripts that match the characters and timestamps in the movies, we use the “Tongyi Tingwu” software of Alibaba for audio-to-text transcription, which generate dialogue corresponding to timestamps. After minimal manual adjustments, we acquire the necessary dialogue content. This method forms a robust foundation for constructing the dataset of character conversation. More details of the movie and main character selection can be found in Sec. A.1.2 and Sec. A.1.3.

The second stage of dataset construction involves character profiles and five different types of VideoQA tasks. Table 2 shows the QA examples for each task. All annotations are carried out on an enterprise crowdsourcing platform. All workers are proficient in Chinese, possess a solid educational background to ensure accurate comprehension of video content, and produce corresponding descriptions precisely. Note that workers need to have completed at least 100 prior tasks on the platform with a minimum approval rate of 95%. Additionally, we conduct daily spot checks on annotations written by each worker to verify their relevance to

Task	Evaluation	QA Examples (movie: Goodbye Mr.Loser)
Dialogue prediction	BLEU,CIDEr, ROUGE-L, GPT,Human	start_time: 01:26:12 end_time: 01:26:42 Q:夏洛接下来要说什么话? (What will Xia Luo say next?) A:可我最心爱的女人被别人抢走了。(But the woman I love the most is taken away from me.)
Action prediction	BLEU,CIDEr, ROUGE-L, GPT,Human	start_time: 01:01:40 end_time: 01:02:10 Q:袁华接下来会做出什么样的行为? (What will Yuan Hua do next?) A:袁华接下来会在漫天飞雪的环境下哭泣。(Yuan Hua will cry in the midst of falling snow.)
Relationship judgment	BLEU,CIDEr, ROUGE-L,GPT, Human,Accuracy	start_time: 00:08:28 end_time: 00:09:10 Q:夏洛和马冬梅什么关系? (What's the relationship between Xia Luo and Ma Dongmei?) A:夏洛和马冬梅是夫妻关系。(Xia Luo and Ma Dongmei are husband and wife.)
Sentiment analysis	Accuracy	start_time: 01:01:40 end_time: 01:02:10 Q:袁华此时的心情如何? (What is Yuan Hua's mood at this time?) A:此时袁华的情绪是悲伤的。(At this time Yuan Hua's mood is sorrow.)
Logical analysis	BLEU,CIDEr, ROUGE-L, GPT,Human	start_time: 01:07:00 end_time: 01:08:10 Q:袁华此时作诗和之前作诗时的差别在哪? 分析原因。 (What are the differences in Yuan Hua's poetry now compared to before? Analyze reasons.) A:袁华之前作诗带有批判性, 让夏洛很没有面子, 现在作诗则极尽谄媚, 巴结夏洛。原因在于夏洛和袁华的社会地位发生了翻天覆地的变化, 袁华现在穷困潦倒, 不比之前。 (Previously,Yuan Hua's poetry bore a critical tone, which led to Xia Luo losing face. However,there has been a significant shift in his poetic style, now excessively fawning. This change can be attributed to the stark reversal in social status between Xia Luo and Yuan Hua.)

Table 2: The evaluation methods and QA examples for different tasks varied

	Optional items
Emotion	愤怒 厌恶 惊讶 悲伤 喜悦 中立 恐惧 anger disgust amazed sorrow joyful neutrality fear

Table 3: The seven options for sentiment analysis

the respective videos. We require workers to first watch the selected movie, describing the main characters' personalities using keywords and sentences, and providing concise descriptions of character profiles. Each character profile includes the name, individuality, and identity of a movie character, as shown in Table. 4. When selecting suitable movie clips, we annotate the start and end times, with the end time being one second before the answer appears. The annotation process for the dataset spans three months, involving 39 qualified workers who contribute annotations for 25,000 questions across 101 movies. Additionally, corrections are made to 40,905 dialogue data entries, and descriptions for 398 character profiles are provided. More details of the quality control of annotations is depicted in Sec.A.1.4. The design details of the five VQA tasks are as follows:

- **Dialogue prediction.** The prompt format is "What will [movie character] say next?" and the answer is the next line of dialogue for the character in the movie.

Character portrait	Name	Identity(translation)
夏洛特烦恼 Goodbye Mr.Loser	夏洛 Xia Luo	Xia Luo is a vengeful, greedy, and timid person who, in pursuit of vanity, puts on airs. However, at the same time, he is someone who has always harbored dreams of music.
夏洛特烦恼 Goodbye Mr.Loser	马冬梅 Ma Dongmei	Ma Dongmei has a straightforward and somewhat tomboyish personality, lacking a bit of feminine charm. She is unburdened by trivialities, upright, and stands up for justice. Ma Dongmei is a simple, dedicated, hardworking, and capable individual.
前任3-再见前任 The Return of the Exes	孟云 Meng Yun	Meng Yun is a career-oriented man with ambitious goals in his professional life, displaying a proactive and upwardly mobile attitude. Although he appears composed on the surface, there are fluctuations in his inner world.
我不是药神 Dying to survive	程勇 Cheng Yong	Cheng Yong is a small-time merchant who peddles Indian God Oil. He is opportunistic, selfish, engages in domestic violence against his wife, and is generally self-centered, bullying the weak and fearing the strong. However, later on, he rediscovers his inherent kindness, starts helping others, and shows a sense of responsibility.

Table 4: The names and identity descriptions of characters from several movies are displayed

- **Action prediction.** The prompt format is "What will [movie character] do next?" and the answer is the next action or expression of the character in the movie.
- **Relationship judgment.** The prompt format is "What is the relationship between [movie character A] and [movie character B]?" and the answer is the relationship between the characters in the specific scene of movies.

- 248 • **Sentiment analysis.** The prompt format is
249 “How is [movie character] feeling at this moment?”
250 and the answer describes the emotion
251 of a character based on movie clips. The re-
252 sponse should be chosen from the following
253 seven emotions: anger, disgust, joy, sorrow,
254 neutral, surprise, and fear, as shown in Ta-
255 ble. 3.
- 256 • **Logical analysis.** The prompt format is:
257 “Why is [movie character] engaging in a spe-
258 cific behavior, expression, or action?” or
259 “[movie character] changes behavior from pre-
260 vious action to current action, analyze the rea-
261 sons.” Responses must be based on the current
262 movie clip and long context, providing expla-
263 nations.

264 We choose three native Chinese speakers to
265 cross-validate the modified script dialogues. They
266 verify the accuracy of the dialogue content by
267 watching 101 Chinese movies, first confirming the
268 correctness of the dialogue and then checking the
269 alignment of the dialogue with timestamps and the
270 characters in the movies. Corrections are made for
271 any inconsistencies. For character profiling valida-
272 tion, we randomly assign 398 selected main movie
273 characters to 50 individuals, with each character
274 assessed by two people familiar with Chinese and
275 relevant movies. They provide ratings for the char-
276 acter profiles, and a consensus with satisfaction
277 levels exceeding 85% is considered a pass; other-
278 wise, it is reassessed by annotators. The validation
279 for the five video QA tasks is relatively straight-
280 forward. We randomly reassign the annotated QA
281 pairs to two additional individuals, who then assess
282 whether the answers are consistent with the movie
283 plot and characters. In cases of inconsistency, mod-
284 ifications are made.

285 3.3 Comparison with VideoQA Datasets

286 As summarized in Table. 1, most existing videoQA
287 datasets focus on English question-answering, with
288 a primary emphasis on visual understanding. The
289 QA pairs typically interpret content from specific
290 excerpts. In contrast, our dataset is dedicated to
291 conversations between characters in a multimodal
292 context, specifically within the context of movie
293 scenes. As a Chinese QA dataset, our questions
294 are designed to revolve more around the characters
295 in movies, combining video and long background
296 context to predict character actions or dialogues.

Besides, our dataset is compatible with both open-
297 domain and closed-domain QA tasks. It comprises
298 video clips and textual descriptions that are signifi-
299 cantly longer than those in existing video narration
300 datasets. The average length of video clips is 6024
301 seconds, surpassing the average length of current
302 video datasets by a considerable margin.
303

304 4 MovieGPT Model

305 4.1 Design Principles

306 As illustrated in Figure. 4, Our MovieGPT model
307 is a multimodal large language model based on
308 the transformer architecture (Vaswani et al., 2017).
309 The model utilizes the pre-trained Bloom-7B as its
310 backbone (Workshop et al., 2022). Besides, the
311 visual module contains a ViT-L/14 (Dosovitskiy
312 et al., 2021) visual encoder and a connecting layer
313 visual abstractor (Liang et al., 2022) .
314

315 To enable the model to engage in character con-
316 versation and question-answering within specific
317 scenarios, we facilitate the learning of character-
318 specific knowledge and memories by providing an
319 abundance of dialogues from various movie char-
320 acters. Additionally, detailed descriptions of char-
321 acter individuality and identity are provided to en-
322 hance the model’s understanding of the characters’
323 language styles. Given video clips from movies,
324 we task the model with comprehending the char-
325 acters in the film and engaging in conversation
326 or answering questions in specific movie scenes.
327 The character responses generated by the model
328 should adhere to several criteria of faithful repre-
329 sentation, including: (1) Lexical Consistency – the
330 model should reflect the personality of a character,
331 ensuring consistency with the character’s unique
332 conversation style; (2) Dialog Authenticity – the
333 generated responses should not only be contextu-
334 ally relevant but also align with the content and
335 plot of the movie.

335 4.2 Model Implementation

336 At the current stage, existing multimodal LLMs
337 can effectively process visual information. How-
338 ever, they have certain limitations in understanding
339 character conversation, particularly in Chinese con-
340 versational contexts (Yang et al., 2022; Ye et al.,
341 2023b; Zhang et al., 2023; Zheng et al., 2023). Our
342 exploration focuses on the understanding of Chi-
343 nese character conversations in multimodal scenar-
344 ios using LLMs. More training details are depicted
345 in Sec. A.2.1. To facilitate the multimodal character

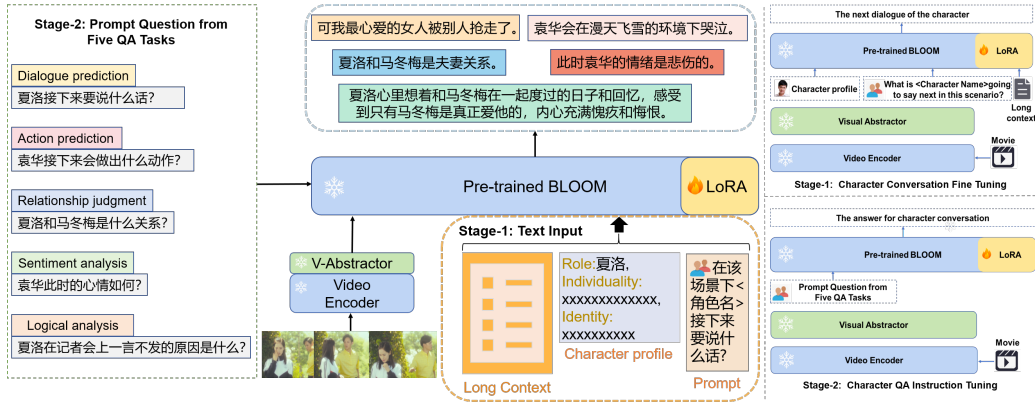


Figure 4: Illustration of the two-stage training of our MovieGPT. In the first stage, the model is fed with long context from the characters, the movie clip, and character profiles. Given the fixed prompt, the goal is to predict the next dialogue of the character. The training dataset of the second training stage is the labeled QA pairs of five designed tasks, namely “dialogue prediction, action prediction, relationship judgment, sentiment analysis, and logical analysis”, combined with the movie clips.

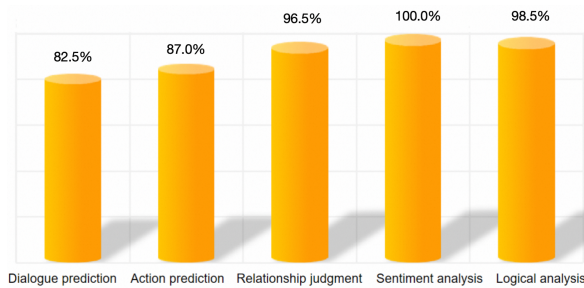


Figure 5: The accuracy of human testing is evaluated across five tasks, with 200 samples extracted for testing from a pool of 25k QA pairs for each task.

QA, our MovieGPT is trained with the following stages.

Character conversation fine-tuning At this initial stage, the pre-trained Bloom and the visual module remains in a frozen state. The model undergoes training with LoRA fine-tuning (Hu et al., 2021), with inputs comprising long context, movie clips, character profiles, and prompts in a standardized format. The expected output is the next line of dialogue of the given character.

Character QA instruction tuning. During this second stage, the parameters of Bloom-7B and visual module are still frozen. The MovieGPT is trained with LoRA fine-tuning, with inputs comprising the prompt question of five character QA tasks along with the corresponding movie clip. The expected output is the response of the character to the question.

5 EXPERIMENTS

In this section, our CharacterQA is evaluated by MovieGPT and several mainstream LLMs.

Through experiments, we investigate whether multimodal LLMs can be trained efficiently with the dataset to achieve character-based contextual understanding. Furthermore, we explore the character conversation capability of the model to characterize and interact with users across five challenging tasks, including dialogue prediction (DP), action prediction (AP), relationship judgment (RJ), sentiment analysis (SA), and logical analysis (LA).

5.1 Experimental Setup

The CharacterQA dataset is uniformly sampled to construct the training, validating, and testing sets, with 20,000, 2,500, and 2,500 QA instances, respectively. Note that the long context for each question is conversations of 30 minutes before the movie clip. To ensure efficient QA training, we only collect movie clips as video inputs, which are limited to one minute. This is because existing models encounter difficulties in encoding long videos.

5.2 Task Evaluation

To demonstrate the practicality of five tasks in real-world scenarios, we conduct manual testing with 200 randomly selected samples for each task. Specifically, the input of the model is provided to human participants for responses, which are then assessed by three native Chinese speakers. The human testing results of five tasks are in Figure 5.

To comprehensively verify the performance of character conversations, various evaluation metrics are adopted on the five VideoQA tasks, and the corresponding results are presented in Table 5. Besides, we invite five native Chinese speakers

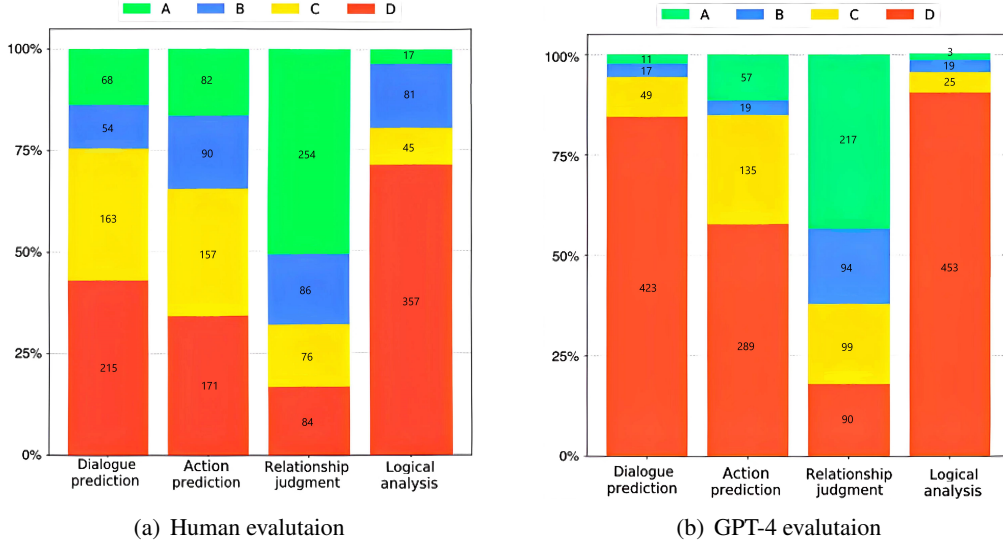


Figure 6: Human and GPT-4 evaluation of MovieGPT’s performance on four open-domain QA tasks. The scoring is done on a scale of “very consistent (A)”, “somewhat consistent (B)”, “fairly consistent (C)”, and “not consistent (D)”.

Task	BLEU-4	CIDEr	ROUGE-L	Accuracy
Dialogue prediction	15.28	21.89	7.67	-
Action prediction	29.21	41.02	0.75	-
Relationship judgment	73.62	78.30	33.89	45.76%
Sentiment analysis	-	-	-	37.29%
Logical analysis	4.45	19.51	0.00	-

Table 5: Evaluation results of our MovieGPT on five tasks with several different metrics. For sentiment analysis, only Accuracy is calculated.

to watch those movies and collectively evaluate whether the generated answers of the model align with the standard ones. Ratings are assigned across four levels: very consistent (A), somewhat consistent (B), fairly consistent (C), and not consistent (D). As shown in Figure. 6, we also employ GPT-4 to assess MovieGPT’s responses based on alignment scores (OpenAI, 2023a). Particularly, we provide the long context of dialogue and the character profile as the prompt, enabling GPT-4 to score responses of our MovieGPT on different tasks.

The results in those figures and tables illustrate that existing multimodal LLM architectures still face great challenges in five character conversational videoQA tasks. Except for sentiment analysis which is in a multi-choice form, other four tasks are open-domain questions. The results of these open-domain tasks, except for relationship judgment, deviated significantly from expectations, especially in the dialogue prediction and logical

analysis tasks. This is mainly because character dialogues and storylines in movies are full of drama and discontinuity, while existing models are unable to realize movie story reconstruction and reasoning through simple visual encoding. Notably, in comparison to human evaluations, the evaluations based on GPT-4 tend to assign lower scores to model responses. This discrepancy arises because human evaluators effectively take into account the content of the movie plot. For responses to various open-domain questions, human evaluators, even when there is some deviation from the standard answer, assign higher scores as long as the responses align with the movie context. In contrast, GPT-4 places a direct emphasis on the alignment between model responses and standard answers.

5.3 Comparison Experiments

As demonstrated in Table. 6, we compare MovieGPT with several state-of-the-art LLMs (Bard (Thoppilan et al., 2022), Claude (Bai et al., 2022), GPT3.5 (Brown et al., 2020) and GPT-4) and two multimodal LLMs (NExT-GPT (Wu et al., 2023) and Video-LLaVA (Lin et al., 2023)) on five tasks to verify its character conversational abilities. Due to the lack of visual input capability in text-only LLMs, character dialogue from relevant movie clips serves as a proxy for video data. These baselines forego fine-tuning for specific tasks, opting instead for in-context learning (Li et al., 2023a). For two multimodal LLMs, they undergo the same

419
420
421
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

Model	Dialogue prediction	Action prediction	Relationship judgment	Sentiment analysis	Logical analysis
Bard	0.35	7.29	36.27	17.00%	4.21
Claude-2	1.43	14.34	63.32	30.50%	9.43
GPT-3.5	0.94	11.84	76.57	28.50%	12.04
GPT-4	2.74	21.02	87.65	35.50%	16.87
NExT-GPT	4.67	11.64	57.98	23.64%	2.31
Video-LLaVA	6.44	17.63	61.04	31.86%	3.37
MovieGPT	15.28	29.21	73.62	37.29%	4.45

Table 6: Comparison of our method with Bard, Claude-2, GPT-3.5, and GPT-4 on five tasks, namely “dialogue prediction, action prediction, relationship judgment, and logical analysis” using the BLEU-4 metric, and “sentiment analysis” using accuracy as the measure.

Model	DP	AP	RJ	SJ	LA
w/o movie clips	13.13	26.43	70.44	31.18%	3.32
w/o role profile	11.46	28.13	66.21	36.91%	2.05
w/o long context	3.47	8.25	54.23	21.66%	3.87
MovieGPT	15.28	29.21	73.62	37.29%	4.45

Table 7: Our method in comparison with others through ablation experiments under different scenarios.

two-stage training as our MovieGPT. The results show that, although MovieGPT may not exhibit comparable performance to GPT-4 in relationship judgment and logical analysis, it outperforms all baselines over the other three tasks by a large margin, highlighting the importance of the understanding of visual contents. More detailed comparisons can be found in Sec. A.2.3.

5.4 Ablation Studies

The results of three ablation variants are shown in Table. 7, where the movie clips, character profile, and long context input are removed, respectively. When the movie clips are missing, we replace them with the corresponding character dialogue. More ablation results are shown in Sec. A.2.4. It can be observed that the model obtains substantial improvements across all tasks even if only the brief movie clips are adopted, which demonstrates the important role of multimodal video semantics in the character conversation. Moreover, the absence of the long context greatly impacts the model performance, especially in dialogue prediction. This is reasonable since the long context is important for understanding character backgrounds and their expression habits. Comparatively, the influence of removing character profiles is more pronounced in logical analysis and relationship judgment, which rely more on specific information such as character identity and individuality.

6 Conclusion

We propose a Chinese multimodal character question answering dataset, comprising 101 carefully selected Chinese movies. Compared to existing datasets, our CharacterQA focuses on personalized comprehension in the Chinese multimodal conversational settings. In addition to specially annotated script conversations and character profiles, we design five videoQA tasks to evaluate character QA abilities. In addition, we built a multimodal character QA model called MovieGPT, and conducted various experiments to evaluate the multimodal character conversation capability of mainstream LLMs. The results demonstrate that the character-based QA tasks are still very challenging for current models. This requires exploring long-distance visual semantics, and mining character personality profiles needed for personalized reasoning. It also involves extending our CharacterQA to a broader range of languages and more complex problems, thereby indicating our future research directions.

7 Limitations

In this study, we investigate a multimodal LLM for character conversation through movie watching, a challenging task that requires ongoing efforts. Our work remains limited in several aspects: firstly, data constraints exist as our selected movies and character roles are limited, which are insufficient to encompass the diverse landscape of existing Chinese films. Future endeavors could benefit from a broader selection. Secondly, the foundational model—results from supervised fine-tuning—are highly influenced by factors such as the pre-training data distribution, model architecture, and scale. Subsequent work may explore trainable agents based on more potent and LLMs.

References

- 515 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc,
516 Antoine Miech, Iain Barr, Yana Hasson, Karel
517 Lenc, Arthur Mensch, Katherine Millican, Malcolm
518 Reynolds, et al. 2022. Flamingo: a visual language
519 model for few-shot learning. *Advances in Neural
520 Information Processing Systems*, 35:23716–23736.
- 521 Yuntao Bai, Saurav Kadavath, Sandipan Kundu,
522 Amanda Askell, Jackson Kernion, Andy Jones,
523 Anna Chen, Anna Goldie, Azalia Mirhoseini,
524 Cameron McKinnon, et al. 2022. Constitutional
525 ai: Harmlessness from ai feedback. *arXiv preprint
526 arXiv:2212.08073*.
- 527 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
528 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
529 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
530 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
531 Gretchen Krueger, T. J. Henighan, Rewon Child,
532 Aditya Ramesh, Daniel M. Ziegler, Jeff Wu, Clemens
533 Winter, Christopher Hesse, Mark Chen, Eric Sigler,
534 Mateusz Litwin, Scott Gray, Benjamin Chess, Jack
535 Clark, Christopher Berner, Sam McCandlish, Alec
536 Radford, Ilya Sutskever, and Dario Amodei. 2020.
537 *Language models are few-shot learners*. *ArXiv*,
538 abs/2005.14165.
- 539 Santiago Castro, Mahmoud Azab, Jonathan Stroud,
540 Cristina Noujaim, Ruoyao Wang, Jia Deng, and
541 Rada Mihalcea. 2020. Lifeqa: A real-life dataset
542 for video question answering. In *Proceedings of the
543 Twelfth Language Resources and Evaluation Confer-
544 ence*, pages 4352–4358.
- 545 Santiago Castro, Naihao Deng, Pingxuan Huang, Mi-
546 hai Burzo, and Rada Mihalcea. 2022. Wildqa: In-
547 the-wild video question answering. *arXiv preprint
548 arXiv:2209.06650*.
- 549 Ke Chen, Zhao Zhang, Weili Zeng, Richong Zhang,
550 Feng Zhu, and Rui Zhao. 2023a. *Shikra: Unleashing
551 multimodal llm’s referential dialogue magic*. *ArXiv*,
552 abs/2306.15195.
- 553 Nuo Chen, Yan Wang, Haiyun Jiang, Deng Cai, Yuhan
554 Li, Ziyang Chen, Longyue Wang, and Jia Li. 2023b.
555 Large language models meet harry potter: A dataset
556 for aligning dialogue agents with characters. In *Find-
557 ings of the Association for Computational Linguistics:
558 EMNLP 2023*, pages 8506–8520.
- 559 Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Pier-
560 giovanni, Piotr Padlewski, Daniel Salz, Sebastian
561 Goodman, Adam Grycner, Basil Mustafa, Lucas
562 Beyer, et al. 2022. Pali: A jointly-scaled mul-
563 tilingual language-image model. *arXiv preprint
564 arXiv:2209.06794*.
- 565 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,
566 Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul
567 Barham, Hyung Won Chung, Charles Sutton, Sebas-
568 tian Gehrmann, Parker Schuh, Kensen Shi, Sasha
569 Tsvyashchenko, Joshua Maynez, Abhishek Rao,
Parker Barnes, Yi Tay, Noam M. Shazeer, Vinod-
kumar Prabhakaran, Emily Reif, Nan Du, Benton C.
Hutchinson, Reiner Pope, James Bradbury, Jacob
Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin,
Toju Duke, Anselm Levskaya, Sanjay Ghemawat,
Sunipa Dev, Henryk Michalewski, Xavier García,
Vedant Misra, Kevin Robinson, Liam Fedus, Denny
Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim,
Barret Zoph, Alexander Spiridonov, Ryan Sepassi,
David Dohan, Shivani Agrawal, Mark Omernick, An-
drew M. Dai, Thanumalayan Sankaranarayanan Pil-
lai, Marie Pellat, Aitor Lewkowycz, Erica Moreira,
Rewon Child, Oleksandr Polozov, Katherine Lee,
Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark
Díaz, Orhan Firat, Michele Catasta, Jason Wei, Kath-
leen S. Meier-Hellstern, Douglas Eck, Jeff Dean, Slav
Petrov, and Noah Fiedel. 2022. *Palm: Scaling lan-
guage modeling with pathways*. *J. Mach. Learn. Res.*,
24:240:1–240:113.
- Chenhang Cui, Yiyang Zhou, Xinyu Yang, Shirley Wu,
Linjun Zhang, James Zou, and Huaxiu Yao. 2023.
*Holistic analysis of hallucination in gpt-4v(ision):
Bias and interference challenges*.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony
Meng Huat Tiong, Junqi Zhao, Weisheng Wang,
Boyang Albert Li, Pascale Fung, and Steven C. H.
Hoi. 2023. *Instructblip: Towards general-purpose
vision-language models with instruction tuning*.
ArXiv, abs/2305.06500.
- Alexey Dosovitskiy, Lucas Beyer, Alexander
Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
Thomas Unterthiner, Mostafa Dehghani, Matthias
Minderer, Georg Heigold, Sylvain Gelly, Jakob
Uszkoreit, and Neil Houlsby. 2021. *An image
is worth 16x16 words: Transformers for image
recognition at scale*.
- Danny Driess, F. Xia, Mehdi S. M. Sajjadi, Corey
Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan
Wahid, Jonathan Tompson, Quan Ho Vuong, Tianhe
Yu, Wenlong Huang, Yevgen Chebotar, Pierre Ser-
manet, Daniel Duckworth, Sergey Levine, Vincent
Vanhoucke, Karol Hausman, Marc Toussaint, Klaus
Greff, Andy Zeng, Igor Mordatch, and Peter R. Flo-
rence. 2023. *Palm-e: An embodied multimodal lan-
guage model*. In *International Conference on Ma-
chine Learning*.
- Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shi-
jie Geng, Aojun Zhou, W. Zhang, Pan Lu, Conghui
He, Xiangyu Yue, Hongsheng Li, and Yu Jiao Qiao.
2023. *Llama-adapter v2: Parameter-efficient visual
instruction model*. *ArXiv*, abs/2304.15010.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan
Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
and Weizhu Chen. 2021. *Lora: Low-rank adap-
tation of large language models*. *arXiv preprint
arXiv:2106.09685*.
- Jen-tse Huang, Wenxuan Wang, Man Ho Lam, Eric John
Li, Wenxiang Jiao, and Michael R Lyu. 2023a. Chat

628	gpt an enfj, bard an istj: Empirical study on personalities of large language models. <i>arXiv preprint arXiv:2305.19926</i> .		
629		Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. 2023. Kosmos-2: Grounding multimodal large language models to the world . <i>ArXiv</i> , abs/2306.14824.	681 682 683 684
630			
631	Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui, Owais Khan Mohammed, Qiang Liu, Kriti Aggarwal, Zewen Chi, Johan Bjorck, Vishrav Chaudhary, Subhojit Som, Xia Song, and Furu Wei. 2023b. Language is not all you need: Aligning perception with language models . <i>ArXiv</i> , abs/2302.14045.		
632		Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. 2023. Character-llm: A trainable agent for role-playing . <i>arXiv preprint arXiv:2310.10158</i> .	685 686 687
633			
634		Tianhao Shen, Sun Li, and Deyi Xiong. 2023. Roleeval: A bilingual role evaluation benchmark for large language models . <i>arXiv preprint arXiv:2312.16132</i> .	688 689 690
635			
636			
637			
638	Jie Lei, Licheng Yu, Mohit Bansal, and Tamara L Berg. 2018. Tvqa: Localized, compositional video question answering . <i>arXiv preprint arXiv:1809.01696</i> .		
639		Makarand Tapaswi, Yukun Zhu, Rainer Stiefelwagen, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. 2016. Movieqa: Understanding stories in movies through question-answering . In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 4631–4640.	691 692 693 694 695 696
640			
641	Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. 2023a. Otter: A multi-modal model with in-context instruction tuning . <i>ArXiv</i> , abs/2305.03726.		
642			
643			
644			
645	Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. 2023b. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models . <i>ArXiv</i> , abs/2301.12597.		
646			
647			
648			
649	KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. 2023c. Videochat: Chat-centric video understanding . <i>arXiv preprint arXiv:2305.06355</i> .		
650			
651			
652			
653	Victor Weixin Liang, Yuhui Zhang, Yongchan Kwon, Serena Yeung, and James Y Zou. 2022. Mind the gap: Understanding the modality gap in multi-modal contrastive representation learning . <i>Advances in Neural Information Processing Systems</i> , 35:17612–17625.		
654			
655			
656			
657			
658	Bin Lin, Bin Zhu, Yang Ye, Munan Ning, Peng Jin, and Li Yuan. 2023. Video-llava: Learning united visual representation by alignment before projection . <i>arXiv preprint arXiv:2311.10122</i> .		
659			
660			
661			
662	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023a. Improved baselines with visual instruction tuning . <i>ArXiv</i> , abs/2310.03744.		
663			
664			
665	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. Visual instruction tuning . <i>ArXiv</i> , abs/2304.08485.		
666			
667			
668	Jiasen Lu, Christopher Clark, Rowan Zellers, Roozbeh Mottaghi, and Aniruddha Kembhavi. 2022. Unified-io: A unified model for vision, language, and multi-modal tasks . <i>ArXiv</i> , abs/2206.08916.		
669			
670			
671			
672	OpenAI. 2023a. Gpt-4 technical report . <i>ArXiv</i> , abs/2303.08774.		
673			
674	OpenAI. 2023b. Gpt-4v(ision) system card .		
675	Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior . In <i>Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology</i> , pages 1–22.		
676			
677			
678			
679			
680			
		Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. 2022. Lamda: Language models for dialog applications . <i>arXiv preprint arXiv:2201.08239</i> .	697 698 699 700 701
		Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. Llama: Open and efficient foundation language models . <i>ArXiv</i> , abs/2302.13971.	702 703 704 705 706 707 708
		Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. Llama 2: Open foundation and fine-tuned chat models . <i>ArXiv</i> , abs/2307.09288.	709 710 711 712 713 714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731
		Quan Tu, Shilong Fan, Zihang Tian, and Rui Yan. 2024. CharacterEval: A chinese benchmark for role-playing conversational agent evaluation . <i>arXiv preprint arXiv:2401.01275</i> .	732 733 734 735
		Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz	736 737

738	Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	793
739		794
740		795
741	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. <i>arXiv preprint arXiv:2308.15126</i> .	796
742		797
743		798
744		799
745		800
746	Zekun Moore Wang, Zhongyuan Peng, Haoran Que, Jiaheng Liu, Wangchunshu Zhou, Yuhan Wu, Hongcheng Guo, Ruitong Gan, Zehao Ni, Man Zhang, et al. 2023b. Rolellm: Benchmarking, eliciting, and enhancing role-playing abilities of large language models. <i>arXiv preprint arXiv:2310.00746</i> .	801
747		802
748		803
749		804
750		805
751		806
752	BigScience Workshop, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Lucchioni, François Yvon, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. <i>arXiv preprint arXiv:2211.05100</i> .	807
753		808
754		809
755		810
756		811
757		812
758	Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. 2023. Next-gpt: Any-to-any multimodal llm. <i>arXiv preprint arXiv:2309.05519</i> .	813
759		814
760		815
761	Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. 2021. Next-qa: Next phase of question-answering to explaining temporal actions. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 9777–9786.	816
762		817
763		818
764		819
765		820
766	Haiyang Xu, Qinghao Ye, Ming Yan, Yaya Shi, Jiabo Ye, Yuanhong Xu, Chenliang Li, Bin Bi, Qi Qian, Wei Wang, et al. 2023. mplug-2: A modularized multimodal foundation model across text, image and video. In <i>International Conference on Machine Learning</i> .	821
767		822
768		823
769		824
770		825
771	Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. 2022. Zero-shot video question answering via frozen bidirectional language models. In <i>NeurIPS</i> .	826
772		
773		
774		
775	Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye, Ming Yan, Yuhao Dan, Chenlin Zhao, Guohai Xu, Chenliang Li, Junfeng Tian, Qian Qi, Ji Zhang, and Fei Huang. 2023a. mplug-docowl: Modularized multimodal large language model for document understanding. <i>CoRR</i> , abs/2307.02499.	
776		
777		
778		
779		
780		
781	Qinghao Ye, Guohai Xu, Ming Yan, Haiyang Xu, Qi Qian, Ji Zhang, and Fei Huang. 2023b. Hitea: Hierarchical temporal-aware video-language pre-training. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 15405–15416.	
782		
783		
784		
785		
786		
787	Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. 2023c. mplug-owl: Modularization empowers large language models with multimodality. <i>arXiv preprint arXiv:2304.14178</i> .	
788		
789		
790		
791		
792		
	Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Anwen Hu, Haowei Liu, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. 2023d. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration.	
	Zihao Yue, Qi Zhang, Anwen Hu, Liang Zhang, Ziheng Wang, and Qin Jin. 2023. Movie101: A new movie understanding benchmark. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4669–4684, Toronto, Canada. Association for Computational Linguistics.	
	Hang Zhang, Xin Li, and Lidong Bing. 2023. Video-llama: An instruction-tuned audio-visual language model for video understanding. <i>ArXiv</i> , abs/2306.02858.	
	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena.	
	Yiyang Zhou, Chenheng Cui, Jaehong Yoon, Linjun Zhang, Zhun Deng, Chelsea Finn, Mohit Bansal, and Huaxiu Yao. 2023. Analyzing and mitigating object hallucination in large vision-language models. <i>arXiv preprint arXiv:2310.00754</i> .	
	Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. <i>ArXiv</i> , abs/2304.10592.	
	Linchao Zhu, Zhongwen Xu, Yi Yang, and Alexander G Hauptmann. 2017. Uncovering the temporal context for video question answering. <i>International Journal of Computer Vision</i> , 124:409–421.	

A Appendix

A.1 Dataset Details

A.1.1 Visualization of our CharacterQA

<p>Movie clips: </p>	<p>Movie clips: </p>
<p>Character profile: Individuality: 白晶晶对于爱情是一切愿意牺牲。早在500年前，她与孙悟空曾有一段恋情。“我十八岁那年你说会娶我”，这也成为孙悟空的执念。她为爱生恨，但心中更多的还是爱。(Bai Jingjing follows her feelings when it comes to love. At early as 500 years ago, she had a romantic relationship with Sun Wukong. "You promised to come and marry me when I was eighteen," he once said to her. Despite harboring resentment due to unfulfilled promises, her heart still holds more love than anything else.) Identity: 五百年前，孙悟空曾与白晶晶有过一段恋情，白晶晶从至尊宝变为孙悟空，二人情同手足。随着取经二人妖身分的相行，至尊宝的归来开始与二妖开始。白晶晶在取经路上打伤至尊宝，自己也中春十三娘的毒。为了保护至尊宝，白晶晶与春十三娘决斗，白晶晶为了至尊宝而死去，心灰意冷独自求死。至尊宝和春十三娘双双救回白晶晶，为了保护至尊宝，白晶晶和春十三娘一起对孙悟空下咒，但咒语过于强大，春十三娘和白晶晶都不能对对手。为了再次救回白晶晶，至尊宝又使用月光宝盒，回到了过去。(Five hundred years ago, Sun Wukong had a romantic relationship with Bai Jingjing. Bai Jingjing recognized the Monkey King as Sun Wukong, and their affection for each other gradually grew. The Patriarch Bozhi revealed their true identities as demons, leading to conflict between Sun Wukong, Bai Jingjing, and a group of bandits. Bai Jingjing, injured by Sun Wukong during a confrontation with Princess Iron Fan, also became poisoned by her. To save Bai Jingjing, Sun Wukong went alone to find Princess Iron Fan for the antidote. Misunderstanding Sun Wukong's intentions, Bai Jingjing, feeling desolate, jumped off a cliff, apparently ending her life. Sun Wukong and Princess Iron Fan failed to locate Bai Jingjing. In an attempt to protect Sun Wukong, Bai Jingjing and Princess Iron Fan confronted the Bull Demon King together, but he proved too powerful for them both. In another attempt to revive Bai Jingjing, Sun Wukong used the magic of the Moonlight Box to return to the past.)</p>	<p>Character profile: Individuality: 宋天荫是一个村长的保长，平时吃苦耐劳，有时胆小，有时为了朋友赴汤蹈火是一个非常真的人。(Song Tianyin is the headman of a village. Usually, he is hardworking, sometimes timid, but occasionally, for the sake of friends, he is willing to face danger. He is a very kind person.) Identity: 宋天荫是孙悟空的保长，在妖怪妖石遇到时帮忙捉妖，遇到了二天荫霍小岚来捉妖，同时意外的被妖后看中，在霍小岚的蛊惑下被押往妖洞里上当妖后，晚上被妖后偷偷带走并怀孕生下妖蛋怀孕，被霍小岚发现并往妖洞里捉妖妖后并自杀霍小岚户内，他为了救妖后与霍小岚上前往妖洞里捉妖，在一家客栈住生下小妖王，后带着小妖王一起搬至深山中，在搬至小妖王时他们发现妖后与妖王小妖王产生感情，在集市时霍小岚王妖后出现的妖后妖王。(Song Tianyin is the headman of Yongning Village. When the demon queen and her followers fled to the village, he welcomed them. When the exorcist Er Qian Shi Huo Xiaolan came to capture the demons, he unexpectedly caught the eye of the demon queen. Under Huo Xiaolan's influence, he was bound and used as bait on the bed. He was secretly taken away by the demon queen at night and forcibly impregnated with demon eggs. After being rescued by Huo Xiaolan and returning to the village, he discovered that all the villagers were demons and had been captured by Ge Qianhu. In order to save the village and regain its reputation, he and Huo Xiaolan embarked on a journey to Shuntian Prefecture. While staying at an inn, he gave birth to the young demon king. Later, they hid in the deep mountains together with the young demon king. When trying to sell the young demon king, they realized they had developed a bond with him. Seeing the determination in the young demon king's bite mark on a painting at the market, they resolved to save him.)</p>
<p>Long context: 00:00:51 观音大士：孙悟空你这个畜牲，你本来答应如来佛祖保护，送你师傅唐三藏去取经，你竟然跟牛魔王要吃我的肉，你知不知道那了天大罪啊？ 00:00:51 孙悟空：少啰嗦，你逼了我三日三夜，因为你嫌女人，我才不杀你，不要以为我怕你啊。 00:31:00 白晶晶：什么不好像不像个臭猴子。 00:32:22 至尊宝：真头秃，原来晶晶姑娘喜欢我粗犷的一面，乱抓乱了点总算恢复了点死光。</p>	<p>Long context: 00:00:27 猪八戒：很久以前，人与万物共存，高空中也有妖怪，可人独占天下，于是妖魔鬼怪，奸猾奸猾奸猾深山，不敢动弹。如今，妖界内乱，老妖王被杀，新妖王为了斩除妖，大举屠杀妖魔鬼怪老妖王的痛恨归咎。 00:01:15 猪八戒：而此时妖怪已身怀妖胎，为了逃遁新妖王的杀手追杀，在余党的保护下被迫逃入人间，据说这龙胎将会颠覆人间的妖界。 00:33:49 霍小岚：你刚吃了吧，你再吃我都瞎了。 00:33:54 宋天荫：在我怀里先别吧。</p>
<p>Prompt: 白晶晶接下来要说什么？(What is Bai Jingjing going to say next?)</p>	<p>Prompt: 霍小岚与宋天荫之间是什么样的关系？(What kind of relationship is there between Huo Xiaolan and Song Tianyin?)</p>
<p>Answer: 孙悟空，你给我滚出来你个臭猴子我十八岁那年你说会娶我害我把他等了三年你给我滚出来。(Sun Wukong, get out of here, you stinky monkey! When I was eighteen, you said you would come to marry me, and you made me wait for three years. Get out of here!)</p>	<p>Answer: 霍小岚与宋天荫之间是利用关系。(The relationship between Huo Xiaolan and Song Tianyin is one of exploitation.)</p>
<p>Bard: 别想得太美，就凭你这点改变，想要我心就太天真了。(Don't think too highly of yourself. With just these minor changes, thinking you can soften my heart is too naive.)</p>	<p>Bard: 霍小岚与宋天荫之间是朋友的关系。(The relationship between Huo Xiaolan and Song Tianyin is one of friendship.)</p>
<p>Claude2: 你只在乎你在乎你山妖还是妖？我只在乎你是不是个真心的人。(Do you think I care whether you're a bandit or a scholar? I only care if you're a genuine person.)</p>	<p>Claude2: 霍小岚与宋天荫之间是伙伴关系。(Huo Xiaolan and Song Tianyin are partners.)</p>
<p>GPT3.5: 你的粗犷只是你表面的一部分，但你的内心却永远留在了那片山林的丛林里。(Your roughness is only a part of your surface, but your heart has forever remained in that jungle of bandits.)</p>	<p>GPT3.5: 霍小岚与宋天荫之间是朋友关系。(Huo Xiaolan and Song Tianyin have a friendship.)</p>
<p>GPT4.0: 孙悟空，看来你还没搞清楚，我喜欢的不是你的外表，而是你的外表或者身份。(Sun Wukong, it seems you still haven't figured it out. What I like is the real you, not your appearance or status.)</p>	<p>GPT4.0: 宋天荫与霍小岚之间是陌生人关系，互相利用。(Song Tianyin and Huo Xiaolan have a relationship of strangers, mutually exploiting each other.)</p>
<p>MovieGPT: 我十八岁那年你说会娶我。(When I was eighteen, you said you would come to marry me.)</p>	<p>MovieGPT: 宋天荫与霍小岚之间是伙伴关系，互相合作利用。(XSong Tianyin and Huo Xiaolan have a relationship of companionship, mutually cooperating and exploiting each other.)</p>

(a) An example of the dialogue prediction task.


(b) An example of the relationship judgment task.

Figure 7: Examples of the dialogue prediction and relationship judgment task. Given the labeled movie clip, character profile and long context, the answers of different LLMs to the question prompt are also illustrated.

In this section, we provide visualization results for the CharacterQA dataset. As shown in Figure 7, Figure 8, and Figure 9, each of the five tasks of dialogue prediction, action prediction, relationship judgment, sentiment analysis, and logical analysis demonstrates a character QA sample, as well as the answers of our MovieGPT and other different LLMs like Bard, Claude-2, GPT-3.5, GPT-4, to the relevant multimodal question.

As illustrated in Figure 7(a), the task of dialogue prediction poses a formidable challenge. This is primarily due to the response " 孙悟空 你给我滚出来你个臭猴子我十八岁那年你说会娶我害我把他等了三年你给我滚出来。 Sun Wukong, get out of here, you stinky monkey! When I was eighteen, you said you would come to marry me, and you made me wait for three years. Get out of here!" being imbued with intense personal emotion. Moreover, compared to preceding dialogues, this response appears particularly abrupt and necessitates a comprehensive understanding through the integration of the movie clip for an appropriate answer. Consequently, the responses generated by the majority of LLMs are overly subdued, lacked the character's emotional coloring, and strayed far from the answer. Only GPT-4's answer was close to the edge, and only our movieGPT's answer was very close to the edge. This also shows that the two-stage training is very effective.

As for the relationship judgment task shown in Figure 7(b), if one merely floats on the understanding of textual content, the answer “霍小岚与宋天荫之间是利用关系。(The relationship between Huo Xiaolan and Song Tianyin is one of exploitation.)” will be difficult to obtain (It can be seen that Bard, Claude2, GPT3.5, all these text-only LLMs answered poorly, and only by combining the visual information, GPT-4 and MovieGPT can answer the question accurately).

Movie clips: 

Character profile:

Individuality: 吕受益从一开始出现的形象就给观众留下了一个谦和、礼貌、友善的直观印象。经过一些事实发现他是有思想、有见地的人，为人和善、有城有义。 (Initially, Lv Shouyi's appearance gives us a superficial impression of someone sleazy, frail, and weak. However, through certain events, it becomes evident that he is a man of intellect and insight. He is kind-hearted and principled, displaying both compassion and a sense of justice.)

Identity: 吕受益是慢性白血病患者的重症患者，十分小家子气，每次见面都会请大家吃橘子。他闯入程勇的神油店，用一个小点子给程勇带来无限商机，也给自己带来无限生机。两人建立了深厚友谊。程勇帮吕受益筹集医药费的事情，却一拖再拖。这导致吕受益病情不断加重，且最终病情恶化到无法挽回的地步。于是低价药在病内部也悄悄销售起来。买药一事因此而起，但程勇最终决定停止卖药，将药的售卖权交给院长。院长林选贤后再度陷入买不起药价格的泥潭，病入到无可挽回的地步。最后因为不想继续连累家人自杀身亡。也是他的死让程勇下定决心。 (Lv Shouyi is a patient suffering from chronic myeloid leukemia, showing traits of being petty and narrow-minded. Every time they meet, he would treat everyone to oranges. He barges into Cheng Yong's miracle oil shop, bringing a valuable idea to Cheng Yong that brings limited business opportunities for him but infinite possibilities for himself. The two establish a deep friendship. Cheng Yong and Lv Shouyi travel to various hospitals to promote medication, but they achieve nothing. At this point, Lv Shouyi seeks help from Liu Sihai, the group leader of a patient community, to promote the medication among patients in various hospitals. Consequently, the sale of affordable medicine begins discreetly within the patient community. The whole buying medication initiative starts because of him, but ultimately, Cheng Yong decides to stop selling the medicine, handing over the selling rights to Zhang Changlin. After Zhang Changlin flees, the patients once again find themselves unable to afford expensive medicine, leading to irreparable health conditions. Lv Shouyi, unwilling to continue burdening his family, chooses to end his life, and it's his death that prompts Cheng Yong to make a change.)

Long context:

00:02:56 宾馆老板：老板。

00:02:56 木晨：欢迎光临！

00:09:00 程勇：什么事跟人家讲啊？

00:09:02 宾馆老板：你不是缺钱吗，我这就给你拉了个生意，价格可以谈的，谈啊！

Prompt:

吕受益接下来会做出怎样的动作？(What actions will Lu Shouyi take next?)

Answer

吕受益接下来会做出摘下口罩的动作。(Lu Shouyi will then proceed to take off his mask.)

Bard

吕受益会拒绝宾馆老板的提议。(Lv Shouyi refuses the proposal of the hotel owner.)

Claude2

吕受益接下来会认真考虑宾馆老板的提议。(Lv Shouyi will seriously consider the proposal from the hotel owner next.)

GPT3.5

吕受益接下来会选择暂时接受宾馆老板的提议。(Lv Shouyi will choose to temporarily accept the proposal from the hotel owner next.)


GPT4.0

吕受益接下来会做出积极主动地参与合作，但保持谨慎谨慎。(Lv Shouyi will proceed to actively engage in cooperation, but will maintain a cautious approach.)

MovieGPT

吕受益接下来会做出与程勇深入交谈的动作。(Lv Shouyi will proceed to engage in a deep conversation with Cheng Yong next.)

(a) An example of the action prediction task.

Movie clips: 

Character profile:

Individuality: 马进是一个患癌平凡但积极向上的男人，他不因别人的冷嘲热讽而自卑，也不因他人的一时抬高而迷失自己。(Ma Jin is a humble and ordinary person, but he is optimistic and upward-looking. He doesn't belittle himself because of others' ridicule, nor does he lose himself because of temporary praise from others.)

Identity: 马进在公司团建途中发现买的彩票中奖，但还未中奖及回去领奖。一场意外将车内存所有人困在一个偏僻的小岛上，并被告知他们可能是世界上最后的人。马进不相信这一说法，苦于中奖日期将近，无奈独自摸奖中奖。他本想偷偷摸摸带出岛外，但随后与队友小高离开时，马进为了拿回彩票中奖的钱财，在结账后选择告诉同事们中奖的事情，使得大家最终交钱。(During a company outing, Ma Jin discovered that the lottery ticket he brought had won. However, before he could go back to claim the prize, an accident left everyone trapped on a remote island. They were told they might be the last humans on Earth. Ma Jin didn't believe this claim, but with the prize deadline approaching, he reluctantly explored the island alone. He was manipulated by Zheng Jingting to win hearts and minds. After realizing the truth, Ma Jin left the group with his cousin Xiaoqiang. However, they missed the prize deadline. Determined to start anew on the island, Ma Jin abandoned his dreams. Just when things seemed to improve, a cruise ship passed by. He stopped his cousin from deceiving Zheng Jingting for money. After much internal struggle, he decided to tell his colleagues about the passing cruise ship, ultimately leading to their rescue.)

Long context:

00:01:39 主持人：接下来我们请著名专家史教授来分析一下这次陨石事件。

00:01:39 史教授：应该说如果这个陨石不幸掉入海洋中。

00:16:39 同事们：什么彩票专家。

00:17:03 马进：行行行行行别喊了，喊有用吗？不就这点事儿吗，事儿大了咱就解决，现在最主要的想法办法怎么回去。

Prompt:

马进此时的心情如何？(What is Ma Jin's current mood?)

Answer

惊讶 (amazed)

Bard

愤怒 (anger)

Claude2

愤怒 (anger)

GPT3.5

惊讶 (amazed)

GPT4.0

惊讶 (amazed)

MovieGPT

惊讶 (amazed)

(b) An example of the sentiment analysis task.

Figure 8: Examples of the action prediction and sentiment analysis task. Given the labeled movie clip, character profile and long context, the answers of different LLMs to the question prompt are also illustrated.

For the action prediction in Figure. 8(a), the answer “吕受益接下来会做出摘下口罩的动作。(Lu Shouyi will then proceed to take off his mask.)” necessitates the model’s ability to scrutinize the film and assimilate the visual cues pertaining to the mask for precise interpretation. Evidently, not only the conventional text-based LLMs but also GPT-4 and our MovieGPT fall short of accomplishing this task.

Regarding the task of sentiment analysis in Figure. 8(b), given the finite number of sentiment categories, this task is evaluated using the accuracy metric. It is observable that in the absence of multimodal information, the purely text-based LLMs, Bard and Claude2, still fail to provide accurate responses; whereas, the remaining three models all deliver correct answers.

Moreover, as shown in Figure 9, even with comprehension of the video content, it remains exceedingly challenging for a human to address the question “这一举动也跟窗帘旁白挂着的星星之火可以燎原呼应了。(This action also echoes with the small spark hanging by the curtains, suggesting a potential wildfire, as mentioned in the stage directions.)” This necessitates a profound understanding of the film, explaining why all LLMs uniformly responded such as “the pressure of reality, and the unreality of the dream”, which underscores the complexity of deriving nuanced interpretations from multimedia content.

A.1.2 Selection of Movies and Main Characters

Our selection was guided by a goal to ensure diversity in genres and historical span. We also focused on movies with strong narratives, clearly defined main characters, and a substantial amount of dialogue, as these elements are crucial for multimodal characterQA. Starting with a broad pool of 200 movies across various genres, we employed a meticulous review process by three annotators to identify movies meeting these criteria. Movies with weak narratives or lacking 3-5 main characters were excluded. We further filtered out movies with less than 50 lines of dialogue among the main characters. This rigorous process ultimately resulted in a curated list of 101 movies for our CharacterQA dataset.

For each movie, we chose the top 10 ranked characters from each movie’s cast list, ensuring they had significant dialogue interaction (at least 50 lines), since less conversations of other secondary characters are



Figure 9: An example of the logical analysis task. Given the labeled movie clip, character profile and long context, the answers of different LLMs to the question prompt are also illustrated.

not enough to support multimodal characterQA. We then refined this selection by focusing on characters who were central to the movie’s main storyline, resulting in selecting 3-5 movie characters for each movie. This methodical approach helped us create a robust and relevant dataset that accurately represents main characters in each movie. Finally, 398 main characters are obtained.

A.1.3 The Alignment of Dialogues and its Timestamps.

To obtain accurately matched dialogue transcripts from the movies, we utilized Alibaba’s ’Tongyi Tingwu’ software for audio-to-text transcription. This software provides dialogue texts with corresponding timestamps. However, the accuracy of the transcribed content was not always perfect. To address this, we engaged annotators to review and correct the transcriptions against the actual movie dialogues. A second round of validation by another annotator ensured the high alignment accuracy of the final dialogue texts, which further ensures the reliability of our data. Whether it’s dialects or standard Mandarin, the manually processed dialogue texts exhibit extremely high matching rates with the movie dialogues. It’s worth noting that our model does not have audio input, so the impact of dialects on performance is not as significant.

A.1.4 Quality Control for Crowd Worker Annotations

All annotations are conducted on an enterprise crowdsourcing platform by proficient Chinese workers with a solid educational background, ensuring accurate comprehension and precise description of video content. It’s important to note that workers must have completed at least 100 prior tasks on the platform with a minimum approval rate of 95%. Additionally, we perform daily spot checks on worker annotations to ensure relevance to the videos. For validation of modified script dialogues, we enlist three native Chinese speakers who verify accuracy by watching 101 Chinese movies, confirming dialogue correctness and alignment with timestamps and characters. Corrections are made for any inconsistencies. Character profiling validation involves randomly assigning 398 main movie characters to 50 individuals, each assessed by two people familiar with Chinese and relevant movies. Ratings are provided, with consensus

satisfaction levels over 85% considered a pass; otherwise, reassessment by annotators is conducted. Validation for the five video QA tasks involves reassigning annotated QA pairs to two additional individuals who assess consistency with the movie plot and characters, making modifications for any inconsistencies.

A.2 Experiments

In this section, extensive experiments will be provided about our CharacterQA dataset, including the training details, comparison of different evaluation metrics, and supplementary ablation studies.

A.2.1 Training Details

Table 8: Training hyperparameters for character conversation fine-tuning stage and character QA instruction tuning stage.

Hyperparameters	Conversation Fine-Tuning	QA Instruction Tuning
GPU type	$8 \times \text{A6000}$	$8 \times \text{A6000}$
Max token length	1024	1,024
Batch size of text instruction data	-	128
Optimizer	AdamW	AdamW
Learning rate	$2e-4$	$2e-5$
Learning rate decay	Cosine	Cosine
Adam ϵ	$1e-6$	$1e-6$
Adam β	(0.9, 0.98)	(0.9, 0.999)
Epoch	2	5
Weight decay	0.001	0.0001

Our MovieGPT is trained in two stages: the character conversation fine-tuning stage and the character QA instruction tuning stage, during which we freeze the parameters of the visual module and Bloom-7B, fine-tuning the latter with LoRA (Hu et al., 2021). LoRA allows us to indirectly train the dense layers in neural networks by optimizing the rank-decomposition matrices of dense layers during the adaptation process, while keeping the pre-trained weights unchanged.

The character conversation fine-tuning stage aims to familiarize the model with the multimodal character conversation task. In this stage, the visual module receives frames from the 60 second movie clip, while the text input contains long context, character profiles, and the prompt words “该角色在该场景下<角色名>接下来要说什么话? (What is <Character Name> going to say next in this scenario?)”. The output of the visual module and the text inputs will be fed into Bloom-7B together to generate the prediction for the next dialogue of the corresponding character. The loss function that maximizes the likelihood estimation between the prediction and groundtruth dialogue is adopted. The detailed training parameters for this stage is demonstrated in Table. 8.

The character QA instruction tuning stage aims to enable the model’s abilities of the specific task. Particularly, the labeled data of our five designed tasks (dialogue prediction, action prediction, relationship judgment, logical analysis, and sentiment analysis) is adopted to train our MovieGPT, i.e., frames from 60 second movie clips are fed into the frozen visual module, the specific task-related question is the text input. Given the output visual feature and the question, Bloom-7B produces the answer to the question, with the loss function again focusing on maximizing the likelihood estimation between the predicted response and the correct answer. The detailed training parameters for this stage is demonstrated in Table. 8.

The underlying principle of the model training of our MovieGPT is as follows:

1. The first training stage is designed to enable the model to learn personalized representations of characters within the movie, which allows the model to grasp the intricacies of the plot and facilitate accurate dialogue predictions.

2. In the second stage, the model can actively engage in question-answering tasks tailored to various personalized scenarios, leveraging its understanding gained from analyzing movie clips based on provided instructions.

932 A.2.2 Evaluation Details

933 To comprehensively verify the character conversation performance, apart from the evaluation metrics, we
 934 also adopt human annotators to evaluate the open-domain VideoQA tasks in our CharacterQA dataset.
 935 For human evaluation of the responses of different models, we adopted a common method, engaging
 936 several expert annotators for evaluation to maintain opinion alignment. Concretely, groups of five native
 937 Chinese speakers independently reviewed the relevant movie and assessed answer consistency. In cases of
 938 disagreement, group discussion started, and three additional evaluators will assess answer consistency.
 939 This iterative process continued until unanimous agreement was reached within the group, thereby
 940 maximizing the precision of our consistency ratings.

941 A.2.3 Comparison of Different Evaluation Metrics

CIDEr Score	DP	AP	RJ	SA	LA
Bard	0.79	11.49	39.57	-	15.43
Claude-2	2.46	23.13	66.34	-	38.15
GPT-3.5	2.61	16.76	83.23	-	44.13
GPT-4	6.75	33.51	91.22	-	51.96
MovieGPT	21.89	41.02	78.30	-	19.51

Table 9: The CIDEr score of different models on five tasks. “DP” stands for dialogue prediction, “AP” stands for action prediction, “RJ” stands for relationship judgment, “SA” stands for sentiment analysis, “LA” stands for logical analysis.

ROUGE-L Score	DP	AP	RJ	SA	LA
Bard	0.19	0.00	12.34	-	0.00
Claude-2	0.57	0.04	28.09	-	0.01
GPT-3.5	0.53	0.17	37.16	-	0.01
GPT-4	1.84	0.33	48.26	-	0.02
MovieGPT	7.67	0.75	33.89	-	0.00

Table 10: The ROUGE-L score of different models on five tasks.

942 In the paper, we have furnished a comprehensive comparison of BLEU scores across various tasks.
 943 Given the substantial challenges inherent in evaluating open-domain responses, where multiple correct
 944 answers are possible, especially for our multimodal character QA, we also further present the CIDEr,
 945 ROUGE-L and human evaluation results of different models across five tasks in Table. 9, Table. 10 and
 946 Table. 11, respectively. The inputs for the text-only LLMs during the experiments include the long context
 947 preceding the movie clips, the text dialogue of the movie clips, character profiles, and the question of the
 948 corresponding task for the current clip.

949 Note that the BLEU, CIDEr, and ROUGE-L metrics measures the qualities of answers from different
 950 perspective. For example, ROUGE-L calculates the longest common subsequence between the answer (C)
 951 and the groundtruth sentence (S), as shown in Equ. 1:

$$R_{LCS} = \frac{LCS(C, S)}{\text{len}(S)}$$

$$P_{LCS} = \frac{LCS(C, S)}{\text{len}(C)} \quad (1)$$

$$ROUGE - L = F_{LCS} = \frac{(1 + \beta^2) R_{LCS} P_{LCS}}{R_{LCS} + \beta^2 P_{LCS}}$$

953 where R_{LCS} represents recall, while P_{LCS} represents precision, and F_{LCS} is ROUGE-L. Typically, β is
 954 set to a large number, so F_{LCS} almost only considers R_{LCS} (i.e., recall). Note that when β is large, F_{LCS}

will focus more on R_{LCS} than P_{LCS} . If β is very large, the P_{LCS} term can be disregarded. Obviously, this metrics is very strict for open-domain QA since it requires responses and answers to be identical in the longest possible sequence. Most tasks begin with a common sequence for ROUGE-L scores, e.g., “The character will next” for action prediction, “The character will next say” for dialogue prediction, and “The relationship between character 1 and character a is” for relationship prediction; while logical analysis scored lower due to the lack of a common sequence.

As shown in Table. 9 and Table. 10, the CIDEr and ROUGE-L scores were in good agreement with the BLEU scores in terms of overall trends, although there were some differences in the spread between the different model effects. Furthermore, as previously noted, the ROUGE-L scores for all models approached zero, attributable to the inherent challenges of this open-domain task and the constraints of the metric itself.

For the more accurate human evaluations shown in Table. 11, different models obtaining A or B scores also showed the same results. This suggests that metrics such as BLEU-4 and CIDEr, despite their limitations, are still valuable in assessing answer quality. Furthermore, our MovieGPT still achieves the best results over all tasks, further confirming its effectiveness.

Model	Dialogue Prediction				Action Prediction				Relationship Judgment				Logical Analysis			
	A	B	C	D	A	B	C	D	A	B	C	D	A	B	C	D
Bard	0	3	14	483	11	74	41	374	137	137	129	97	11	43	55	391
Claude2	2	7	32	459	38	109	72	281	220	121	122	37	21	62	78	339
GPT-3.5	4	8	40	448	21	102	74	303	289	111	81	19	37	83	72	308
GPT-4	9	11	51	429	63	146	53	238	317	109	63	11	59	71	86	284
MovieGPT	68	54	163	215	82	90	157	171	254	86	76	84	17	81	45	357

Table 11: The human evaluation results of different models on four open-domain tasks, where ratings are assigned across four levels: very consistent (A), somewhat consistent (B), fairly consistent (C), and not consistent (D).

A.2.4 Supplementary Ablation Studies

Movie Clips	Role Profile	Long Context	DP	AP	RJ	SA	LA
✓			3.12	7.88	46.27	20.06%	1.97
	✓		0.71	5.18	29.36	13.21%	1.22
		✓	9.87	14.21	60.43	27.64%	2.01
	✓	✓	13.13	26.43	70.44	31.18%	3.32
✓		✓	11.46	28.13	66.21	36.91%	2.05
✓	✓		15.28	29.21	73.62	37.29%	4.45
✓	✓	✓	15.28	29.21	73.62	37.29%	4.45

Table 12: Our method in comparison with others through ablation experiments under different scenarios.

In this section, we have conducted extensive ablation studies on our CharacterQA dataset. For all results in this section, the BLEU score is adopted for dialogue prediction, action prediction, relationship judgment, and logical analysis, and the Accuracy metric is adopted for sentiment analysis.

Models	Training Setting	DP	AP	RJ	SA	LA
NExT-GPT	In-context Learning	0.29	5.34	21.05	11.63%	1.46
Video-LLaVA	In-context Learning	0.45	6.69	46.38	26.86%	1.87
NExT-GPT	Two-stage Training	4.67	11.64	57.98	23.64%	2.31
Video-LLaVA	Two-stage Training	6.44	17.63	61.04	31.86%	3.37
MovieGPT	Two-stage Training	15.28	29.21	73.62	37.29%	4.45

Table 13: The comparison between our MovieGPT and other multimodal LLMs on five tasks.

Moreover, to underscore the complexity of the CharacterQA dataset and affirm the efficacy of MovieGPT, we embarked on comparative experiments with other well-regarded multimodal LLMs. Our selection was constrained by the scarce availability of open-source multimodal LLMs capable of processing both Chinese language and video inputs. For instance, in our preliminary evaluations, Video-LLaMA (Zhang et al., 2023) exhibited proficiency in handling video temporal information but fell short in accommodating Chinese conversational contexts. Consequently, NExT-GPT (Wu et al., 2023) and Video-LLaVA (Lin et al., 2023), two popular open-source multimodal platforms, were chosen for evaluation.

Specifically, the results of both NExT-GPT and Video-LLaVA, when performing in-context learning only and when performing the identical two-stage training as our MovieGPT, are shown in the Table. 13. As we can see, with only in-context learning, Video-LLaVA and NExT-GPT exhibit suboptimal performance across all tasks, even underperforming some text-only LLMs. This can be attributed to their inadequate comprehension of video content, which may exacerbate errors in the absence of training tailored to multimodal character conversation. After two stages of training, both NExT-GPT and Video-LLaVA show significant improvement over their results with in-context learning only. However, their results are still inferior to our MovieGPT over all tasks.

Length of movie clip	DP	AP	RJ	SA	LA
10s	14.68	28.12	73.06	35.22%	3.48
60s	15.28	29.21	73.62	37.29%	4.45
300s	16.97	30.58	76.29	38.62%	6.67

Table 14: The performance of our MovieGPT with the movie clip of different lengths.

Furthermore, we attempt ablation experiments with different video lengths in Table. 14, which shows marked improvement when the duration of the movie clip increases to 300 seconds, confirming the necessity of incorporating "visual context" is crucial for multimodal characterQA. However, much longer videos will face a dilemma of high frame extraction computational costs, and less frame extraction will result in more temporal information loss. The visual processing capabilities of existing multimodal LLMs for long videos are also limited (Zhang et al., 2023; Liu et al., 2023b) (For the common multimodal LLMs like Video-LLaMA (Zhang et al., 2023) and mPLUG-owl (Ye et al., 2023c), the lengths of video input are usually less than 3 minutes due to their inferior visual capacities.). Designing a model capable of rapidly processing longer movie clips to enhance multimodal character QA performance is a future research direction for us.