EMOJI2IDIOM: BENCHMARKING CRYPTIC SYM-BOL UNDERSTANDING OF MULTIMODAL LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

Abstract

Vision and Language are two major modalities in Artificial Intelligence research. Bridging the gap between these modalities has long been a key focus in the multimodal community. Inspired by human cognition, we believe that if a model can see an image and directly associate it with its linguistic meaning, the model possesses high-level intelligence that spans vision and language. In our work, we focus on emojis in images, a widely-used "cryptic symbol", with a data form of both visual and linguistic features, i.e. emojis have the specific textual semantics while human understand the meaning from their visual information. Specifically, we first propose the novel task of translating emojis in images to corresponding idioms, thereby challenging Multimodal Large Language Models (MLLMs) to (1) understand the semantic correlation between language and emojis, and (2) reason the intricate linguistic meaning from the emojis in images. To facilitate the advancement of this task, we construct a high-quality benchmark (Emoji2Idiom) following the process of automatic model generation and human manual filtering. Based on our constructed Emoji2Idiom, we employ multiple advanced MLLMs to conduct extensive experiments and detailed analyses, demonstrating that existing MLLMs do not yet have enough capability to understand and reason the linguistic information from visual data. We believe our proposed benchmark and interesting discoveries will encourage the community to attach importance to the intelligence of MLLMs directly associating language from vision, to give MLLMs more comprehensive vision-language understanding ability¹.

034

037

000

001

003

004 005 006

012 013

014

015

016

017

018

019

021

022

025

026

027

028

029

031

032

033

1 INTRODUCTION

Multimodal Large Language Models (MLLMs) have made remarkable progress and achievements in recent years Yin et al. (2023); Wu et al. (2023); Cui et al. (2024), especially
their visual-language understanding capabilities, which have laid a solid foundation for the
widespread development of multimodal applications Chen et al. (2024); Zhang et al. (2024a).
For how to improve the visual-language understanding capabilities of MLLMs, a core challenge is how to bridge the gap between vision and language Koh et al. (2024); Peng et al. (2023); Wang et al. (2024).

Naturally, we want to know to what extent MLLMs should understand vision and language
before we can claim that "the gap between vision and language in MLLMs has been filled"?
We believe if MLLMs can behave like humans, their intelligence must have reached an extremely ideal level. We notice that when a person sees an image, he or she can often directly
associate it with the linguistic meaning behind the image. For example, when humans see
special symbols, they can directly know the words represented by these symbols. Since previous VQA-based benchmarks treat the vision and language separately, we try to transfer human analogy to MLLMs, when an image is fed to MLLMs, if MLLMs can directly

⁰⁵²

¹All our data are available in anonymous Github link https://anonymous.4open.science/r/Emoji2Idiom-0CCA.

081



Figure 1: This figure illustrates the capabilities that our Emoji2Idiom concerns. We provide the emojis in images, ground truth, and the MLLM-generated results.

associate the linguistic meaning behind the image, then we can claim that it has
 relatively advanced visual-language understanding intelligence. Images inputted
 to MLLMs here have clear and discrete semantic meanings, manifested as concrete concepts,
 specific symbols or logos, etc. The semantic information of these images can be accurately
 interpreted as linguistic concepts or tokens.

087 Inspired by the above thinking, our work aims to explore the ability of MLLMs to directly 880 understand the linguistic meaning behind images. It is exciting that many graphic symbol codes in cryptography exist as special indicators in an image so that when people see the image, they can decode the textual meaning of the image. We have noticed that emoji 090 are increasingly becoming a kind of "cryptic symbol" widely used by people from 091 all over the world and from all cultural backgrounds Mostafavi & Porter (2021). People not only use emoji to enrich their expressions and show their moods, but also directly use emoji 093 to replace the corresponding text Fischer & Herbert (2021). Emojis have specific textual 094 semantics, while human understand the meaning from their visual representations instead of directly treating emojis as characters. Thus, emojis in images are strongly coupled with their 096 corresponding linguistic meanings, contributing to becoming the basis of our benchmarks and serving as a bridge between visual and linguistic understanding. The understanding of 098 emoji not only requires MLLMs to comprehend the image information of individual emoji but also to combine its textual indications with the related contexts so that the model can further explore the deeper meanings of emoji. Therefore, understanding emojis in images is 100 a challenging vision-language task. 101

To promote the research on cryptic symbol emoji understanding of MLLMs, we propose a novel task that requires MLLMs to receive input image information of emoji sequences and generate their corresponding text information, shown in Figure 1. We introduce the visual modality of emoji and require the MLLMs not only to identify the linguistic meaning of individual emoji but also to understand the special utterances of the emoji and its associated context, contributing to generating text with special semantic and format, e.g., a word or an idiom. Specifically, our task aims to challenge the following capabilities of MLLMs:

(1) Harmonized Word Reasoning: Translating emojis into texts usually harmonizes
 with a sound-like word. Therefore, MLLMs need to have rich language knowledge, to reason
 about the harmonic words.

 (2) Abstract visual understanding of image: Emoji symbols often have strong indicative meanings, which requires MLLMs to deeply understand the abstract visual characteristics of emoji, rather than just the visual shape.

(3) Many-to-one or one-to-many mapping generation problem: According to our
observation, it is common for multiple emojis to correspond to one word, or one emoji to
correspond to multiple words. This requires the MLLMs to make correct predictions based
on the origin emoji understanding and to realize the complex reasoning via context.

119 Furthermore, we construct the Emoji2Idiom benchmark to support the task of translating cryptic symbol emojis in images to corresponding texts. To enrich the diversity of the bench-120 mark, we set up emoji-to-Chinese idiom, emoji-to-English word, and emoji-to-English idiom 121 tasks, taking into account the language and semantic diversity. After automatic filtering and 122 manual filtering by human experts from raw data, we obtain a high-quality dataset. From 123 the above challenges of understanding and the design of diverse benchmark tasks, we hope 124 that MLLM can not only realize the complex understanding of real-world text substitution 125 expressions using emoji, but also generalize to other cryptographic symbols understanding. 126 We hope to realize a generalized unified visual-verbal understanding benchmark instead of 127 the traditional VQA-based benchmarks, which treat visual and verbal information sepa-128 rately. Based on our constructed Emoji2Idiom, we employ multiple advanced MLLMs to 129 conduct extensive experiments and detailed analyses, demonstrating that existing MLLMs 130 do not yet have enough capability to understand and reason the linguistic information from 131 visual data. Our contributions are summarized as follows:

- 132
- 133 134 135

136

137

138

139

140 141

143

1. We first propose the task of translating a sequence of emojis in images to corresponding texts, aiming to guide MLLMs to perform high-level vision-language understanding like humans.

- 2. We build the high-quality Emoji2Idiom benchmark, which is a new data resource that can facilitate MLLMs to better understand cryptic symbol in images.
- 3. We conduct experiment of advanced MLLMs on Emoji2Idiom and provide some detailed analysis, interesting discoveries, and valuable insights for the community to further improve the visual-language understanding capabilities of MLLMs.
- 142 2 Related Work

144 Language Model Based Cryptic Understanding Emoji can be represented by UTF-145 8 Abel (2019), and many treat emoji as text and encode them as vectors Eisner et al. 146 (2016). Leveraging the emoji Unicode library, numerous studies have explored emoji-text translation, including translation text into emoji Monti et al. (2016); Leonardi (2022); Klein 147 et al. (2024), and bidirectional translationDanesi (2022). Beyond this, emoji-based sentiment 148 analysis has become a significant area of emoji research Gibson et al. (2018); Chen et al. 149 (2019; 2018); Liu et al. (2021). However, to the best of our knowledge, our Emoji2Idiom is 150 the first to apply the visual representation and textual semantics of emojis. 151

152 MLLMs Benchmark Earlier unified MLLM benchmarks collect a substantial number 153 of images and generate corresponding QA pairs to evaluate MLLMs Fu et al. (2023a), 154 with a focus on uniformity and objectivity, as seen in SEEDBENCH Li et al. (2024b) 155 and SEEDBENCH-2 Li et al. (2023a). Recent benchmarks have started to assess different 156 capabilities from different dimensions, including visual comprehension Fu et al. (2023b); Li 157 et al. (2024a); Tong et al. (2024); Cai et al. (2023), reasoning ability Zhang et al. (2024c); 158 Roberts et al. (2023), in-context learning capability Shukor et al. (2023); Liu et al. (2023), 159 hallucination challenge Cui et al. (2023); Liu et al. (2023), and multiple domains (math, physics, music, medical, etc.)Lu et al. (2024b); Li et al. (2023c); Yue et al. (2024); Zhang 160 et al. (2024b). However, most benchmarks are based on the VQA annotations and natural 161 scenario image, rather than directly associating an abstract image with its linguistics.



Figure 2: This figure illustrates the data collection pipeline, which is divided into two stages, raw data collection, and data annotation and filtering.

3 THE EMOJI2IDIOM BENCHMARK

182 3.1 TASK DEFINITION183

163

164

165

166

168

170 171

172

173

174

175 176 177

178

179 180

181

187

195

197

207

210

211

212

213

Given an image of a sequence of emoji $I_i^{\text{emoji}} = \{\text{emoji}_1, \text{emoji}_2, \cdots, \text{emoji}_n\}$, Emoji2Idiom task aims to translate emojis in images to corresponding idiom text by model F:

$$\Gamma \text{ext} = F(I_i^{\text{emoji}}). \tag{1}$$

These emoji sequences correspond to texts with specific formats and semantics, representing a Chinese idiom word, an English word, or an English idiom sentence. It requires not only understanding the direct corresponding text of a single emoji, but also inferring complex linguistic meaning based on the surrounding emoji context, containing some harmonic characters, multiple emoji mapping to a single character, or one emoji mapping to multiple words. These more complex emoji understanding problems are prevalent in our dataset. We further discuss the specific properties and challenges in Section 3.3 and Appendix C.

196 3.2 Benchmark Construction

Raw Data Collection As shown in the Figure 2, we collect raw data through two automatic generation methods: *Retrieve from the Internet*. There are a large number of 199 databases for guessing the corresponding words and idioms based on emojis on the internet, 200 which can be obtained freely without commercial usage. We retrieve the relevant web pages 201 of the game, web databases, and the video to get the original emoji images and the corre-202 sponding text answers. Generate Emoji Based on Text. The quality of internet retrieval is 203 not high, due to 1) recurring emoji-text pairs, 2) a relatively higher number of four-word 204 idioms compared with a few multi-word idioms, and 3) a relatively low number of English 205 idioms. We select the texts of common English words, English idioms, and Chinese idioms, 206 and generate the corresponding emoji sequences by the text-to-emoji translation.

Data Annotation and Filtering Automatic filtering. The machine is utilized to auto matically perform data cleaning, including deletion of duplicate values and ethical checking.

- Deletion of duplicate values and missing values, etc. Notably, the machine automatically removes duplicate combinations of the same emoji inputs but retains combinations of the same text result with different emoji inputs.
- Perform ethical checking. Some of the emoji may contain expressions of violence, pornography, or other safety violations, and we utilize GPT-40 to check all the images and remove those that are not ethically safe.

Human Filtering. In this phase, we engage human experts to refine the semantics of emoji text pairs, focusing on the following aspects:

- **Removal of Non-standard Idioms:** Phrases like "蓝天白云 (blue sky and white clouds)" lack historical context and cultural significance, leading to their exclusion.
- Elimination of Low Consistency Pairs: Annotators assess the alignment between emoji and text. Pairs with weak correlations are discarded to make the translation process overly difficult and reducing the image's indicative meaning.
- Exclusion of Unclear Images: Images that are unclear to recognize, such as those from low-resolution video screenshots, are scored on their clarity. Images that score poorly are removed to ensure legibility.
- Mitigation of Repetitive Mappings: Frequent mappings, such as "养" to "结 (jie)"-"节 (jie)", can introduce data bias. To address this, we employ diverse emoji databases, and manually adjust or remove repetitive mappings beyond ten times.
 - Filtering of Unethical Content: We rigorously filter for emoji-text pairs linked to violence, discrimination, or other inappropriate themes. A wide range of emoji including multiple skin tones and gender categories is utilized to promote expression.

To eliminate subjectivity in manual filtering, we provide annotators with detailed guidelines as shown in Appendix A. And additional information about our human filtering can be found in Appendix A.

Table 1: The statistics and image-text pair examples of our Emoji2Idiom in four tasks.

Task	Image-Text Pairs	Emoji Examples	Text Examples
Chinese idioms (Four Characters)	1,876	😟 😟 🗙 😀	闷闷不乐
Chinese idioms (Multi-characters)	334	X ? 3 7 2 1 0 1	不问三七二十一
English Word	842	🛨 🐟	starfish
English Idiom	783	<u></u> = 💰	Health is wealth.

3.3 Data Statistic Analysis and Other Features

We give the statistics of our proposed Emoji2Idiom in Table 1. In the Chinese idioms task, we collect 1,876 and 334 emoji-text pairs of four-character idioms and multi-character idioms, respectively. Among them, since the combinations of four-character idioms are naturally much larger than those of multi-character idioms in dictionaries, such a difference in the distribution is similarly reflected in our dataset. For English words, we set the tasks of emoji to word and idiom, with 842 and 783 sets of image-text pairs, respectively. There are some additional details about Emoji2Idiom in the Appendix B. In our Emoji2Idiom, we observe several interesting linguistic phenomena, and present some examples, with additional details provided in Appendix C. The linguistic phenomena raise great challenges of Emoji2Idiom, and also encourage the exploration of vision-language capabilities of MLLM.

 Word Split In the English word, it is common for multiple emoji to represent one word.
For instance, the word "Panda" can be split into "Pan-" and "-da," where "Pan-" corresponds to Q. Beyond understanding the meaning of individual emojis, the MLLM must also remove unnecessary letters and combine the parts to infer a completely new word.

Harmonic Characters Since it is sometimes difficult to find directly related emoji to represent, harmonic characters with similar pronunciations are often chosen to replace them. For example, "To be loaded", "To" harmonizes with "Two" 2, and "be" harmonizes with "bee" . In the Chinese idiom "难舍难离", "舍" harmonizes with "蛇 (snake)" of emoji
6, "离" harmonizes with "梨 (pear)" of emoji 🍐. The understanding of these harmonics usually requires the model to synthesize the relevant context of the emoji, to reason out the correct expression of the harmonized words.

Table 2: Evaluation results on Chinese idiom task. The Word, Chr-2 and Chr-1 denote the 271 accuracy of guessing the whole word, two or more words, and one or more words correctly. 272

		Idion	with Fe	our wo	rds		Idiom with Multi-words						
	v	Vord-lev	el	Cha	racter-	level	Word-level Character-						
Model	Word	Chr-2	Chr-1	Pre.	Rec.	F-1	Word	Chr-2	Chr-1	Pre.	Rec.	F-1	
Deepseek-VL	0.4	2.3	25.6	6.4	6.4	6.4	1.1	4.9	29.3	8.7	10	9.3	
Qwen-VL	0.5	4.7	30.2	8.7	8.7	8.7	2.4	9.8	31.7	9.1	16.9	11.8	
LLaVa-1.5	0.6	3.8	32.2	10.5	10.5	10.5	2.8	7.9	29.9	9.0	17.3	11.8	
CogAgent	0.6	4.4	34.7	11.6	11.6	11.6	3.6	8.2	30.4	8.7	14.5	10.9	
InternVL-2	0.8	6.3	37.8	9.1	9.1	9.1	3.4	8.3	29.4	8.9	15.6	11.3	
Claude-3.5	1.3	6.7	23.3	8.0	8.0	8.0	1.4	2.9	7.1	6.0	9.7	7.4	
GPT-4V	0.7	1.3	22.1	5.8	5.8	5.8	1.1	6.8	28.4	3.7	9.1	5.3	
GPT-40	3.3	8.7	27.5	10.7	10.7	10.7	9.1	13.6	27.3	7.5	18.1	10.6	

280 281 282

285 286

287

291

270

284

Abstract visual Emoji Understanding. In addition to referring to the direct meanings of the emoji, it is often necessary to deeply infer the semantics of the emoji. For example, in 🗑 🧡 🝃 💪 "同心叶力 (pull together with the same goal) ", 💪 is an arm, but it does not mean "arm" in idioms. Instead, it is a very strong arm, which corresponds to "力 (power)".

288 Cross-cultural Issue Discussion Emoji, as a simple and universally recognized symbol, 289 is widely used across many countries, especially on global social platforms. While cultural 290 nuances are inevitable, emojis generally facilitate cross-cultural understanding. Our dataset tries to minimize ambiguities and emotional shifts caused by complex linguistic contexts, focusing on a sequence instead of a single emoji. In addition, Our dataset is constructed 292 with careful consideration of emoji diversity, covering categories such as smiley faces and emotions, humans and bodies, animals and nature, food and drinks, travel and places, activities, objects, symbols, and flags.

295 296 297

298

299

301

302

294

3.4 Evaluation Metrics

Our dataset computes the precision, recall, F-1, and BLEU value of the results with the ground truth results on the sentence level, word level, and character level to evaluate the MLLM's ability to understand emoji images. We further propose the Chr-2 and Chr-1 to measure in fine-grained evaluation, which denotes the accuracy of guessing two or more words, and one or more words correctly. The details about the evaluation metrics are provided in the Appendix D.

307

308

EXPERIMENT RESULTS 4

4.1 BASELINES

309 We select commercial Claude-3.5-sonnet-20241022, gpt-4-vision-preview and GPT-4o-310 20240513 to evaluate the emoji2idiom benchmark. For a richer evaluation, we select a 311 series of open-source MLLMs for testing. These include: 1) Qwen-VL-7B Bai et al. (2023), 312 DeepSeek-VL-7B Lu et al. (2024a), which have good Chinese language support; 2) LLaVa-313 1.5-7B Li et al. (2023b), CogAgent-18B Hong et al. (2023), InternVL-2-8B which have good 314 visual comprehension capabilities. We provide a detailed description of the baselines, 315 their implementation details, and the prompt template in Appendix E.

316 317

318

4.2 MLLM EVALUATION RESULTS

319 Emoji to Chinese Idiom We evaluate four-character and multi-character idioms shown 320 in Table 2. We observe that all the MLLMs perform poorly on these two tasks. The latest 321 model, GPT-40, achieves accuracy scores of 3.3 and 5.0 at the word level for both tasks. The accuracy at the Chr-1 is significantly higher than at the word level, indicating that MLLMs 322 are equipped with the basic translations of text corresponding to individual emojis, but 323 have limited capability to further infer the corresponding linguistic meanings based on the

Table 3: Evaluation on English word and idiom task. B-1 and B-2 denote the BLEU-1 and BLEU-2 respectively.

		English Word						English Idiom						
	W	Word-level Character-level					Sentence-level Word-level						zel	
Model	Pre.	Rec.	F-1	Pre.	Rec.	F-1	Pre.	Rec.	F-1	Pre.	Rec.	F-1	B-1	B-2
Deepseek-VL	23.2	26.3	24.7	46.2	47.5	46.8	11.9	11.9	11.9	15.1	14.6	14.8	15.1	11
Qwen-VL	28.6	29.1	28.8	51.2	50.4	50.8	12.1	12.1	12.1	17.1	12.5	14.4	17.1	11.3
LLaVa-1.5	30.1	30.1	30.1	54.6	55.7	55.1	14.4	14.4	14.4	19.7	21.3	20.5	19.7	16.4
CogAgent	29.8	29.8	29.8	52.8	51.9	52.3	13.2	13.2	13.2	18.3	19.5	18.9	18.3	15.2
InternVL-2	31.1	31.1	31.1	56.6	57.2	56.9	15.3	15.3	15.3	19.3	22.1	20.6	18.4	16.1
Claude-3.5	42.3	42.3	42.3	63.9	73.8	68.5	29.8	29.8	29.8	48.0	42.7	45.2	42.3	39.7
GPT-4V	38.5	38.5	38.5	60.3	69.2	64.4	26.4	26.4	26.4	41.1	43.1	42.1	39.4	37.5
GPT-40	55.8	55.8	55.8	68.5	77.5	72.7	35.2	35.2	35.2	46.8	47.3	47.0	45.0	41.6

Table 4: Evaluation of the semantic similarity scores of Chinese task, where 1 is categorized as dissimilar and 5 is categorized as perfect similarity.

	Chines	e Idio	m witl	n Four	word	s	Chinese	Idion	ı with	Mult	i-wor	\mathbf{ds}
	Average		Distr	ibutic	n(%)		Average		Distr	ibutio	n(%)	
Model	Semantics	1	2	3	4	5	Semantics	1	2	3	4	5
InternVL-2	1.41	66.9	26.3	6.1	0	0.7	1.47	61.4	32.9	4.5	0	1.1
GPT-40	1.66	56.7	29.1	9.5	0.6	4.1	1.76	59.1	25.0	5.7	1.2	9.0
		Eng	lish W	/ord				Engl	ish Idi	iom		
InternVL-2	2.75	46.2	11.5	19.2	17.2	40.4	2.75	60.3	19.0	11.1	4.8	4.9
GPT-40	3.55	27.5	7.8	2.0	7.5	55.2	2.99	28.5	22.0	7.7	5.5	36.3

relevant emoji context, especially for the harmonization reasoning. Thus, our Emoji2Idiom is challenging to MLLMs due to a huge number of harmonization word mapping with emojis.

Emoji to English Word and English Idiom MLLM's overall accuracy is higher compared to the two Chinese idiom tasks. In Table 3, GPT-40 achieves impressive F-1 values of 55.8 and 35.2 at the word and sentence levels, in emoji-to-English word and English idiom respectively. This is likely because the model has encountered more similar English texts during training, making it more adept at reasoning about English words. However, MLLMs always suffer from hallucination problems. When they catch a linguistic meaning of a single emoji, they quickly focus on the word or idiom related to this emoji from the inner knowledge they have, and ignore the relevant context of the emojis. Based on our Emoji2Idiom, the community can explore the hallucination problem and improve the inference ability.

Evaluation of the semantic similarity of the response We further experiment the semantic similarity between the responses and the ground truth. We input the model output answers and ground truth into LLM and let LLM score the semantic similarity from 1 to 5. As shown in the Table 4.2, we observe that the average scores of the model on the English task are significantly higher than those on the Chinese task. In addition, when we carefully observe the distribution, we find that 1)for the Chinese task, most of the scores are concentrated in 1 and 2, which indicates that the MLLM can almost barely guess; 2)while for the English task, most of the scores are concentrated in 1 and 5, which indicates that the MLLM can either predict the answer correctly, or get irrelevant answers.

4.3 Further Exploration on MLLM Learning

Exploration with In-context Learning In addition to evaluating the direct inference abilities of MLLMs, we further explore their performance using in-context learning. We select the open-source Qwen-VL and the closed-source GPT-40, evaluating each task with 3, 5, and 7 context examples, as shown in Table 5 and Table 6. MLLMs improve across various tasks with the addition of contextual examples, indicating the high quality of our Emoji2Idiom that the randomly chosen examples can improve the performance a lot. However, in the Chinese task, performance decreases when using too many samples (7 in-context examples). This decline indicates that MLLMs learn incorrect mappings in this complex

Table 5: Exploration on in-context learning in Chinese idiom tasks. The Word, Chr-2 and Chr-1 denote the accuracy of the whole word, two or more words, and one or more words.

	Idiom with Four words						Idiom with Multi-words							
	v	Word-level Character-level				Word-level Character-lev								
Model	Word	Chr-2	Chr-1	Pre.	Rec.	F-1	Word	Chr-2	Chr-1	Pre.	Rec.	F-1		
Qwen-VL	0.5	4.7	30.2	8.7	8.7	8.7	2.4	9.8	31.7	9.1	16.9	11.8		
+3 in-context example	0.5	5.1	31.3	9.3	9.3	9.3	2.2	10.1	28.6	10.4	13.1	11.6		
+5 in-context example	0.6	5.3	31.6	9.4	9.4	9.4	3.3	12.3	32.1	11.7	16.9	13.8		
+7 in-context example	0.5	4.9	32.1	9.4	9.4	9.4	2.8	10.7	31.4	11.4	15.4	13.1		
GPT-40	3.3	8.7	27.5	10.7	10.7	10.7	9.1	13.6	27.3	7.5	18.1	10.6		
+3 in-context example	2.6	11.3	33.9	12.6	12.6	12.6	9.5	23.8	36.9	17.0	21.9	19.1		
+5 in-context example	3.5	12.2	35.7	13.7	13.7	13.7	13.1	27.4	42.9	20.7	29.1	24.2		
+7 in-context example	3.5	8.7	31.3	12.0	12.0	12.0	10.7	19.0	34.5	16.2	23.1	19.0		

Table 6: Exploration on in-context learning in English tasks.

	English Words English Idiom											
	W	ord-lev	vel	Cha	racter-	level	W	ord-lev	vel	Cha	racter-	level
Model	Pre.	Rec.	F-1	Pre.	Rec.	F-1	Pre.	Rec.	F-1	Pre.	Rec.	F-1
Qwen-VL	28.6	29.1	28.8	51.2	50.4	50.8	12.1	12.1	12.1	17.1	12.5	14.4
+3 in-context example	29.3	29.3	29.3	53.6	52.1	52.8	12.0	12.0	12.0	17.6	17.9	17.7
+5 in-context example	30.6	30.6	30.6	55.9	54.2	55.0	13.0	13.0	13.0	18.9	18.4	18.6
+7 in-context example	32.5	32.5	32.5	57.8	55.7	56.7	15.2	15.2	15.2	19.7	20.2	19.9
GPT-40	55.8	55.8	55.8	68.5	77.5	72.7	35.2	35.2	35.2	46.8	47.3	47.0
+3 in-context example	57.6	57.6	57.6	72.3	75.0	73.6	36.3	36.3	36.3	47.6	47.3	47.4
+5 in-context example	54.5	54.5	54.5	77.5	79.0	78.2	37.4	37.4	37.4	48.2	50.0	49.1
+7 in-context example	60.6	60.6	60.6	79.4	73.9	76.5	38.5	38.5	38.5	49.5	50.5	50.0

task and suffer from hallucination issues. We further provide some insights of fine-tuneing and reasoning approaches on Emoji2Idiom in Appendix. I.

Exploration with Chain-of-Thought We further investigate the enhancement of CoT inference. This task prompts the MLLM to think step by step, with the detailed prompt in the Appendix E.3.2. In Figure 3, we evaluate GPT-40 and qwen-vl. Our findings indicate that the CoT design enables the MLLM to produce better answers without additional training, improving accuracy at both character and word levels while improving semantically and visually aligned responses. This demonstrates the method's effectiveness and the high quality of our data. Furthermore, the CoT framework mitigates hallucination issues in GPT-40 while avoiding significant semantic bias. By mimicking human reasoning processes, the CoT design offers insights into MLLM errors, guiding future research.

Exploration on Input Length Effects We discuss that how the length of emoji se-quences might impact model performance. The length of chinese four-character idioms and English words exhibit short, with average lengths of 4.11 and 4.23, respectively, and lead to minimal impact from image size variations. However, Chinese multi-character idioms and English idioms have longer sequences (averaging 7.48 and 5.32), resulting in more elongated images with higher variance in length. The resizing methods employed by different MLLMs can distort longer images, degrading performance. To address this, further work can propose an additional preprocessing step.

4.4 HUMAN EVALUATION

Human Performance on Chinese Idiom Tasks Due to the limited performance of MLLMs on the Chinese idiom task, we invite human experts to participate and assess the task's difficulty, thereby determining the upper limit of machine performance on this benchmark. Humans are tested by the same evaluation metrics, and task complexity is rated on a scale from one (very easy) to five (very difficult), with evaluation details provided in Appendix G. The results in Figure 4 show that MLLM still has significant room for improvement, and our Emoji2Idiom presents significant challenges.







Figure 4: Human performance on Chinese idiom task. To better show the task complexity, we map the score to the 1-100 interval.

Human Evaluation on MLLMs We conduct a human assessment of the answers generated by MLLMs. This evaluation considers the normality, semantic similarity to the ground truth, emotional similarity, visual similarity to the original emojis, and the fluency of the 459 generated text, with details provided in Appendix. G. As shown in Figure 5, MLLMs perform well in generating standardized idiom expressions. However, lower scores in semantic similarity and visual similarity suggest that emoji comprehension and idiom reasoning remain challenging areas for MLLMs. 463

4.5CASE STUDY

466 We provide a case study based on our experimental results, which contains four types of 467 prevalent challenges, and propose some potential training methods in Appendix I. 468

469 Harmonization Problem MLLMs often fail on harmonization problems. As shown in 470 Figure 6, 🐍 - "蛇 (snake, sound like "she")" homonym to "舍 (leave, sound like "she")" 471 but the MLLM fails to recognize. In the English idiom, 🐋 - "whale" homonym to "well". 472 Our Emoji2Idiom includes many harmonic character phenomena. Current MLLMs are not 473 yet capable of effectively capturing emoji context and reasoning with harmonic words, and 474 struggle with our challenging Emoji2Idiom.

475

446

447

448

449

450

451

452

453 454

455

456 457

458

460

461

462

464

465

476 Hallucination Problem In Figure 6, the model recognizes \gg and immediately outputs 477 "horsing around", without considering other emojis. Another example shows GPT-4v and 478 GPT-40 recognizing the number 3 and associating it with the idioms "朝三暮四 (change 479 one's mind often)" and "颠三倒四 (disorderly)", both containing the number 3, without 480 considering the surrounding emojis. That is due to the hallucination problem. MLLMs 481 often think narrowly, focusing only on words or idioms directly related to a single emoji. Our Emoji2Idiom is concerned about this issue, and look forward to further exploration of 482 the poor performance of MLLMs that we have discovered. 483

- 484
- Multi-emoji to One Character Mapping. Emoji2Idiom presents a huge challenge on 485 this mapping issue, where MLLMs fail to perform a multi-to-one or one-to-multi mapping.



Figure 5: Human evaluation on GPT-4v and GPT-4o. The Std., Sem Sim., Emj Sim., Emo Sim., and Flu. denote the normality, semantic similarity, emotional similarity, visual similarity to the original emojis, and fluency.



Figure 6: Four typical problems the GPT-4v and GPT-4o suffer in our Emoji2Idiom.

For example, **1 0 0** are four emojis, but the model does not successfully combine them into one character " \uparrow (one thousand)".

Abstract Visual Image Understanding of the Emoji Symbol MLLMs struggle to align emoji semantics with intricate meanings when it comes to deep comprehension. For example, in "receive a kickback", the model simply captures is, the meaning of "box", and interprets it as "out of the box", but does not combine the package attributes of "receiving something" with the hint of money to generate the correct answer. Our Emoji2Idiom highly focuses on this deeper understanding, evaluating and exploring the capabilities of MLLMs.

5 Conclusion

We propose the Emoji2Idiom benchmark, containing emoji to Chinese idioms, English
words, and English idioms. It provides a way to measure the ability of MLLM to understand complex emoji symbol sequences on images. We design a measurement framework
containing harmonic characters, abstract visual understanding, and many-to-one mapping
problems, to validate the ability of the MLLM to synthesize the understanding of emoji
contexts with emoji-to-text coupled reasoning and generation. We evaluate advanced opensource and closed-source MLLMs with our dataset, analyze the results, and highlight future
research directions with case studies.

540 **Reproducibility Statement** 6 541

In order to ensure that other researchers can better reproduce our work in us, we put a lot of effort into reproducibility. We describe in detail our data collection and data building process in Section 3 and Appendix A, and provide full experimental details in the Section 4.1 and Appendix E,G, including the parameter details of the model we used with the prompt template. All data and source code can be found on the Github link Emoji2Idiom. We promise to continue to maintain our Github repository, discuss this research with other researchers, and contribute to the entire multimodal large language model community.

548 549 550

551

553

555

556

557

558 559

560

564

565

566

542

543

544

545

546

547

7 ETHICS STATEMENT

552 We introduce a novel benchmark, Emoji2Idiom, incorporating a thorough description of data collection, annotation, and filtration processes. We emphasize that the dataset's creation adheres strictly to ethical guidelines, with vigilant measures against any breach or 554 impropriety. Great care has been taken to uphold ethical standards in the dataset, employing anonymization, desensitization, and data cleaning. The text samples pose no risk to public welfare. Hence, the innovative research directions and tasks proposed are ethically robust and harmless to society.

- References
- 561 Jonathan E Abel. Not everyone s: Or, the question of emoji as 'universal' expression. In 562 Emoticons, Kaomoji, and Emoji, pp. 25–43. Routledge, 2019. 563
 - Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. arXiv preprint arXiv:2308.12966, 2023.
- Rizhao Cai, Zirui Song, Dayan Guan, Zhenhao Chen, Xing Luo, Chenyu Yi, and Alex Kot. 567 Benchlmm: Benchmarking cross-style visual capability of large multimodal models. arXiv 568 preprint arXiv:2312.02896, 2023. 569
- 570 Dongping Chen, Ruoxi Chen, Shilin Zhang, Yinuo Liu, Yaochen Wang, Huichi Zhou, Qihui 571 Zhang, Pan Zhou, Yao Wan, and Lichao Sun. Mllm-as-a-judge: Assessing multimodal 572 llm-as-a-judge with vision-language benchmark. arXiv preprint arXiv:2402.04788, 2024.
- 573 Yuxiao Chen, Jianbo Yuan, Quanzeng You, and Jiebo Luo. Twitter sentiment analysis via 574 bi-sense emoji embedding and attention-based lstm. In Proceedings of the 26th ACM 575 international conference on Multimedia, pp. 117–125, 2018. 576
- Zhenpeng Chen, Sheng Shen, Ziniu Hu, Xuan Lu, Qiaozhu Mei, and Xuanzhe Liu. Emoji-577 powered representation learning for cross-lingual sentiment classification. In The world 578 wide web conference, pp. 251–262, 2019. 579
- 580 Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, Yang Zhou, Kaizhao Liang, Jintai Chen, 581 Juanwu Lu, Zichong Yang, Kuei-Da Liao, et al. A survey on multimodal large language models for autonomous driving. In Proceedings of the IEEE/CVF Winter Conference on 582 Applications of Computer Vision, pp. 958–979, 2024. 583
- 584 Chenhang Cui, Yiyang Zhou, Xinyu Yang, Shirley Wu, Linjun Zhang, James Zou, and 585 Huaxiu Yao. Holistic analysis of hallucination in gpt-4v (ision): Bias and interference 586 challenges. arXiv preprint arXiv:2311.03287, 2023.
- Marcel Danesi. Emotional wellbeing and the semiotic translation of emojis. In Exploring 588 the Translatability of Emotions: Cross-Cultural and Transdisciplinary Encounters, pp. 589 323-344. Springer, 2022. 590
- 591 Ben Eisner, Tim Rocktäschel, Isabelle Augenstein, Matko Bosnjak, and Sebastian Riedel. emoji2vec: Learning emoji representations from their description. In Proceedings of The 592 Fourth International Workshop on Natural Language Processing for Social Media. Association for Computational Linguistics, 2016.

- Brigitte Fischer and Cornelia Herbert. Emoji as affective symbols: affective judgments of
 emoji, emoticons, and human faces varying in emotional content. Frontiers in psychology,
 12:645173, 2021.
- Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. Mme: A comprehensive evaluation benchmark for multimodal large language models. arXiv preprint arXiv:2306.13394, 2023a.
- Chaoyou Fu, Renrui Zhang, Haojia Lin, Zihan Wang, Timin Gao, Yongdong Luo, Yubo
 Huang, Zhengye Zhang, Longtian Qiu, Gaoxiang Ye, et al. A challenger to gpt-4v? early
 explorations of gemini in visual expertise. arXiv preprint arXiv:2312.12436, 2023b.
- Will Gibson, Pingping Huang, and Qianyun Yu. Emoji and communicative action: The semiotics, sequence and gestural actions of 'face covering hand'. Discourse, Context & Media, 26:91–99, 2018.
- Wenyi Hong, Weihan Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang,
 Zihan Wang, Yuxuan Zhang, Juanzi Li, Bin Xu, Yuxiao Dong, Ming Ding, and Jie Tang.
 Cogagent: A visual language model for GUI agents. CoRR, abs/2312.08914, 2023. doi:
 10.48550/ARXIV.2312.08914. URL https://doi.org/10.48550/arXiv.2312.08914.
- Lars Henning Klein, Roland Aydin, and Robert West. Emojinize: Enriching any text with emoji translations. arXiv preprint arXiv:2403.03857, 2024.
- Jing Yu Koh, Daniel Fried, and Russ R Salakhutdinov. Generating images with multimodal
 language models. Advances in Neural Information Processing Systems, 36, 2024.
- Vanessa Leonardi. Communication challenges and transformations in the digital era: emoji language and emoji translation. Language and Semiotic Studies, 8(3):22-44, 2022.
- Bohao Li, Yuying Ge, Yixiao Ge, Guangzhi Wang, Rui Wang, Ruimao Zhang, and Ying
 Shan. Seed-bench-2: Benchmarking multimodal large language models. arXiv preprint
 arXiv:2311.17092, 2023a.
- Bohao Li, Yuying Ge, Yi Chen, Yixiao Ge, Ruimao Zhang, and Ying Shan. Seed-bench-2-plus: Benchmarking multimodal large language models with text-rich visual comprehension. arXiv preprint arXiv:2404.16790, 2024a.
- Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench:
 Benchmarking multimodal llms with generative comprehension. In *CVPR*, 2024b.

- 630 Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, 631 Tristan Naumann, Hoifung Poon, and Jianfeng Gao. Llava-med: Training a large language-and-vision assistant for biomedicine in one day. In Alice Oh, Tristan Nau-632 mann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), Advances 633 in Neural Information Processing Systems 36: Annual Conference on Neural Infor-634 mation Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 635 10 - 16, 2023, 2023b. URL http://papers.nips.cc/paper_files/paper/2023/hash/ 636 5abcdf8ecdcacba028c6662789194572-Abstract-Datasets_and_Benchmarks.html. 637
- Yingshu Li, Yunyi Liu, Zhanyu Wang, Xinyu Liang, Lingqiao Liu, Lei Wang, Leyang Cui,
 Zhaopeng Tu, Longyue Wang, and Luping Zhou. A comprehensive study of gpt-4v's
 multimodal capabilities in medical imaging. *medRxiv*, pp. 2023–11, 2023c.
- Chuchu Liu, Fan Fang, Xu Lin, Tie Cai, Xu Tan, Jianguo Liu, and Xin Lu. Improving sentiment analysis accuracy with emoji embedding. *Journal of Safety Science and Resilience*, 2(4):246–252, 2021.
- Fuxiao Liu, Tianrui Guan, Zongxia Li, Lichang Chen, Yaser Yacoob, Dinesh Manocha, and
 Tianyi Zhou. Hallusionbench: You see what you think? or you think what you see? an
 image-context reasoning benchmark challenging for gpt-4v (ision), llava-1.5, and other
 multi-modality models. arXiv preprint arXiv:2310.14566, 2023.

- Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Yaofeng Sun, et al. Deepseek-vl: towards real-world visionlanguage understanding. arXiv preprint arXiv:2403.05525, 2024a.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao
 Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating math
 reasoning in visual contexts with gpt-4v, bard, and other large multimodal models. In *ICLR*, 2024b.
- Johanna Monti, Federico Sangati, Francesca Chiusaroli, Benjamin Martin, Mansour Sina,
 et al. Emojitalianobot and emojiworldbot-new online tools and digital environments for
 translation into emoji. In Proceedings of Third Italian Conference on Computational
 Linguistics (CLiC-it 2016), 2016.
- Moeen Mostafavi and Michael D Porter. How emoji and word embedding helps to unveil emotional transitions during online messaging. In 2021 IEEE International Systems Conference (SysCon), pp. 1–8. IEEE, 2021.
- Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and
 Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. arXiv
 preprint arXiv:2306.14824, 2023.
- Jonathan Roberts, Timo Lüddecke, Rehan Sheikh, Kai Han, and Samuel Albanie. Charting new territories: Exploring the geographic and geospatial capabilities of multimodal llms. arXiv preprint arXiv:2311.14656, 2023.
- Mustafa Shukor, Alexandre Rame, Corentin Dancette, and Matthieu Cord. Beyond task
 performance: Evaluating and reducing the flaws of large multimodal models with incontext learning. arXiv preprint arXiv:2310.00647, 2023.
- 674
 675
 675
 676
 676
 676
 677
 676
 677
 676
 677
 678
 679
 679
 679
 670
 670
 670
 671
 672
 674
 672
 674
 673
 674
 674
 674
 674
 674
 675
 676
 677
 676
 677
 676
 677
 676
 677
 678
 678
 679
 679
 679
 670
 670
 670
 670
 671
 672
 672
 673
 674
 674
 674
 674
 675
 676
 677
 676
 677
 676
 677
 678
 678
 678
 678
 678
 679
 679
 679
 670
 670
 670
 670
 670
 670
 671
 672
 672
 674
 675
 675
 676
 677
 676
 677
 678
 678
 678
 678
 678
 678
 679
 679
 679
 670
 670
 670
 671
 672
 672
 674
 674
 675
 675
 676
 677
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
- Yiqi Wang, Wentao Chen, Xiaotian Han, Xudong Lin, Haiteng Zhao, Yongfei Liu, Bohan Zhai, Jianbo Yuan, Quanzeng You, and Hongxia Yang. Exploring the reasoning abilities of multimodal large language models (mllms): A comprehensive survey on emerging trends in multimodal reasoning. arXiv preprint arXiv:2401.06805, 2024.
- Jiayang Wu, Wensheng Gan, Zefeng Chen, Shicheng Wan, and S Yu Philip. Multimodal large language models: A survey. In 2023 IEEE International Conference on Big Data (BigData), pp. 2247–2256. IEEE, 2023.

686

687 688

689

690

691

- Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on multimodal large language models. *arXiv preprint arXiv:2306.13549*, 2023.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of CVPR*, 2024.
- ⁶⁹³ Duzhen Zhang, Yahan Yu, Chenxing Li, Jiahua Dong, Dan Su, Chenhui Chu, and Dong
 ⁶⁹⁴ Yu. Mm-llms: Recent advances in multimodal large language models. arXiv preprint arXiv:2401.13601, 2024a.
- Ge Zhang, Xinrun Du, Bei Chen, Yiming Liang, Tongxu Luo, Tianyu Zheng, Kang Zhu,
 Yuyang Cheng, Chunpu Xu, Shuyue Guo, et al. Cmmmu: A chinese massive multidiscipline multimodal understanding benchmark. arXiv preprint arXiv:2401.11944, 2024b.
- Jiawei Zhang, Tianyu Pang, Chao Du, Yi Ren, Bo Li, and Min Lin. Benchmarking large multimodal models against common corruptions. In NAACL, 2024c.

702 A Additional Details of Data Filtering

A.1 AUTOMATIC FILTERING

708

709

710

711

724

725

726

727

728

730

731

732 733

734

735

736

737

738

739

In this phase, we mainly utilize machines and large language models to filter large-scaledata, which includes the following steps in total:

- 1. Deletion of default values. We utilize a machine to automatically remove incomplete emoji-idiom pairs, including those with missing corresponding emoji images and those with missing corresponding idioms. It is guaranteed that each emoji image corresponds to the standard idiom answer one by one.
- 712 Control points to the standard latent and any of the by one.
 713 2. Removing Duplicate Values. We utilize the machine to automatically remove duplicate emoji-idiom pairs. Here, we only need to remove the emoji-idiom pairs corresponding to identical emoji sequences while retaining the pairs with the same idiom text but corresponding to different emoji representations, which helps to enhance the diversity of the dataset. Note that we will first filter the pairs corresponding to the same idiom text by the machine with additional labels, and make a manual decision on whether to perform the deletion in the next stage of manual filtering.
- 719
 3. Image Quality Check. We utilize LLM (specifically GPT-40 is used), to perform image quality checking, which entails marking and removing: images that are too blurry and those that do not meet the ethical norms (images that contain elements of violence, abusive language, discrimination, etc.) along with their corresponding idioms.
 - 4. Text Ethics Checking. We utilize LLM (specifically GPT-40) to perform text ethics checking, which involves tagging and deleting idiom with elements of violence, discrimination, abuse, etc. For example, "红颜祸水" is a sexist idiom, and we will delete its corresponding emoji-idiom pair.
- 729 A.2 HUMAN FILTERING

In this phase, we invited human experts in Chinese and English languages to perform manual data filtering, which included the following steps in total:

- 1. Duplicate value checking: for the automatic filtering phase, the machine flags a portion of emoji-idiom pairs where the text is the same but the corresponding images are not the same. the human expert needs to further check whether the emoji expressions here are really different. For the pairs with identical emoji images, the human expert will delete them.
 - 2. Image quality check: Human experts further check whether the emoji images are unclear and illegible, and remove the illegible images.
- 740 3. Idiom standardization check: Human experts need to check whether the idiom text 741 expression is standardized, including the format of the idiom, whether it has a 742 specific linguistic meaning, and whether it is in line with common human usage, 743 etc., to ensure that our dataset meets the real-world usability. For example, for the 744 idiom "blue sky and white clouds", although it is a four-word idiom that conforms to 745 the norms of human usage, it does not have a specific allusion, mythological story, 746 traditional story background, or special semantic meaning, and does not belong to the standard idioms. For example, although "流水高山" is a four-letter word 747 with a specific historical background, people more often use the expression "high 748 mountains and flowing water". Therefore, "流水高山" is not an expression that 749 conforms to human language usage and will be deleted. 750
- 4. Emoji and Idiom Relevance Check: Since in emoji to idiom expression, many times the representation of harmonic characters will be utilized, which will increase the difficulty of emoji comprehension and the difficulty of generating the final idioms. Human experts will evaluate the relevance of emoji to idioms:
 - If too many or too complex harmonic characters are used with the emoji representation, at this time the task will be too difficult for not only MLLM but

756	also humans to understand. At this point, the human expert will consider the
757	relevance of this emoji sequence to the idiom to be too low and delete the
758	emoji-idiom pair.
759	• It is noteworthy that we conducted an evaluation of human ability on this
760	benchmark in Sec. 4.4 and found that humans achieved an average score of
761	66.5 on the word-level accuracy of Chinese idioms. This score demonstrates
762	both that our dataset is challenging and that the task is accomplishable, and
763	that there is still much room for improvement in the performance of the current
764	MLLM on this task.
765	5. Repeated harmonic word mapping check: due to the limited expression of emoji,
766	when using emoji to replace textual expressions, harmonic words are often used to
767	find the corresponding emoji for expression. emoji2idiom also has a large number
768	of harmonic words. However, we found that if the mapping of the same emoji
769	corresponding to a certain harmonic word occurs too many times, it may cause data
770	bias to LLM in subsequent training, i.e., when LLM sees this emoji it automatically
771	thinks of this harmonic word that occurs multiple times. To mitigate the bias caused
772	by this narmonic word mapping, we performed.
773	• Count the repeated emoji-character harmonic word mappings, and when there
774	are more than ten occurrences, we manually replace the expression of the emoji
775	(find other harmonic word counterparts to replace the original repeated emoji), or just delete the redundant emoji idiom pair
776	In addition, we also considered this issue during the original data collection
777	• In addition, we also considered this issue during the original data conection.
778	can reduce this duplicate mapping. We also take different generation methods
779	when manually constructing text-to-emoji data, which also helps to increase
780	the diversity of harmonic word mappings.
781	6 Safety and Ethics Check: Based on the automatic detection, the human experts
782	further conducted a safety and ethics check of the emoii images and idiom text
783	checking whether there are any issues such as violent gore, abusive language, sexism.
784	racial discrimination, stereotyping, and so on, in the data.
785	
786	To eliminate subjectivity in manual filtering, we provide annotators with detailed guidelines
787	as shown in Figure 7 and 8, including scoring criteria for each item (1-5 points) covering
788	idiomatic normality, graphic consistency, image legibility, repetition mapping, and ethical

794

805

789

790

B Additional Details of Data statistics

emoji and text levels to guide judgments.

795
796 In addition to the numerical statistics, we further do some statistics to better show our
797 Emoji2Idiom.

safety checks. We also provide at least three examples for each item. For ethical safety

checks, we distinguish between subcategories such as violence, abusive language, gender

discrimination, stereotyping, and racial discrimination. We provide examples at both the

Word Frequency and Word Cloud Statistic of Chinese idiom To better present our dataset, we perform word frequency statistics on Chinese idioms and display the word cloud and word rectangle tree graphs, as shown in Figure 9. We first perform word frequency statistics on all characters, filter out the top 1,000 characters, and discard low-frequency words.
From the filtered top 1,000 characters, we conduct lexical analysis and plot word cloud and word rectangle diagrams for adjective and adverbial morphemes, noun morphemes, and verb morphemes, respectively.

Word Frequency and Word Cloud Statistic of English idiom Similarly, we perform
word frequency statistics on English idioms and display word cloud maps with word rectangle
tree diagrams, as shown in Figure 10. We first perform word frequency statistics on all words,
filter out the top 180 words, and discard low-frequency words. From the filtered top 180
words, we create word cloud maps with word rectangle mapping.



The purpose of this work is to screen out emoji-text pairs that do not comply with the rules. For each indicator, there will be a corresponding criterion and examples, which you will need to score the emoji-text pairs, and only the pairs that meet the requirements of each indicator can be retained.

	С	ase							
Image: 👯 즞 🗟 🐔 🧟	50	Normality: 5 Consistency: 4 Jezibility: 4							
Text: 捷雷不)	及掩耳	Emoji Ethical security: 4 Text Ethical security: 5							
	Me	etrics							
Normality historical a	: Whether the text conforms to the illusions and specific cultural backg	idiom's specifications. This includes whether it has rounds, or does not conform to human usage habits							
Options	 Completely non-standard Mostly standard Completely non-standard 	. Mostly non-standard 3. Fairly standard tely standard							
Examples	 "朝三暮四" shows comple "蓝天白云" shows comple "红红火火" shows mostly 	1. "朝三幕四" shows completely standard to the normality. 2. "蓝天白云" shows completely non-standard. 3. "红红火火" shows mostly non-standard.							
Consistent consistence	cy : The consistency of the emoji any y of the example, the easier it is to	d the image is scored, and the higher the get the final translation result							
Options	 Completely inconsistent Mostly consistent Completely 	Mostly inconsistent 3. Fairly consistent tely consistent							
Examples	 2 5 3 ℃ - "两面三刀" 2. ▲ ♥ ♥ ♥ - "出尔反尔" 3. 〒 - "I'm bored." pair is 	pair is mostly consistent. pair is mostly inconsistent. completely inconsistent.							
legibility:	Whether the image is very blurry a	nd illegible is difficult for MLLM to process.							
Options	1. Completely illegible 2. Mos 4. Mostly legible 5. Completel	tly illegible 3. Fairly legible y legible							
Examples	 1.	y legible. gible. Je.							
Duplicate emojis.	emoji-character mapping: Remove	or modify the duplicate emoji-character mapping of							
Examples	 ○○回調 五字琼楼 ▲ ○□ ○○○○○ ◆○○○○ ◆○○○○○○ ◆○○○○○○ ◆○○○○○○○○ ◆○○○○○○○○○○○○ ◆○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○	 ▲ 玉减香消 ○ \$\$ \$ \$ \$ \$ \$ \$ \$ \$ \$ \$ \$ \$ \$ \$ \$ \$ \$							

Figure 7: The first page guidelines for human filtering.

	Guideline of Human Filtering of Data - 2				
The purpose c indicator, the emoji-text pai	of this work is to screen out emoji-text pairs that do not comply with the rules. For eac re will be a corresponding criterion and examples, which you will need to score the rs, and only the pairs that meet the requirements of each indicator can be retained.				
	Case				
Image: 戦争名名 Text: 捷雷オ	Normality: 5 Consistency: 4 legibility: 4 医及掩耳 Emoji Ethical security: 4 Text Ethical security: 5				
	Evaluation Metrics				
Ethical se emoji-te:	ecurity check: Remove emoji-text pairs that are not ethically safe. We rigorously vet tt pairs for issues such as violence, name-calling, and gender bias.				
	Emoji images filtering				
Possible issues	Contains elements of violence, abusiveness, racial discrimination, gender discrimination, and stereotypes.				
Options	 Completely insecure 2. Mostly insecure 3. Fairly secure Mostly secure 5. Completely secure 				
Examples	 Shows completely insecure to the ethics, due to the abusiveness. Shows mostly insecure to the ethics, because it does not conform to social order and good customs. Shows completely insecure to the ethics, due to violence. 				
	Texts filtering				
Possible issues	Contains elements of violence, abusiveness, racial discrimination, gender discrimination, stereotypes, expressions of partiality and passion, and does not conform to social order and good customs.				
Options1. Completely insecure2. Mostly insecure3. Fairly secure4. Mostly secure5. Completely secure					
Examples	 "头发长见识短" shows mostly insecure to the ethics, due to the gender discrimination on the women. "男主外女主内" shows completely insecure to the ethics, due to the stereotypes. "feisty woman" shows completely insecure to the ethics, due to the gender discrimination. 				

Figure 8: The guidelines for human filtering.



972 C Additional Details of Data Attributes and Linguistic 973 PHENOMENON

C.1 Chinese Idiom Task

977 **Harmonization Word** Since it is sometimes difficult to find directly related emoji to 978 represent, harmonic characters with similar pronunciations are often chosen to replace them. 979 For example, Usually, for characters that can't be represented directly by emoji, we will first 980 search for their harmonized characters, then find an emoji that can directly represent the 981 harmonized character, and replace it with this emoji. For example, "捷" does not have a 982 direct emoji, but it harmonizes with 结", which corresponds to "bow" %, and so, we select 983 🎀 chosen to represent the character "捷". There are a large number of harmonic characters 984 in our data. This poses a great challenge to MLLM's understanding and reasoning ability. 985 The reasoning of harmonic words needs the help of related contexts, and in our data scenario, the model is required to analyze the context of emoji in depth instead of understanding 986 individual emoji alone. The understanding of these harmonics usually requires the model 987 to synthesize the relevant context of the emoji, to reason out the correct expression of the 988 harmonized words. 989

990 991

992

975

976

C.2 Abstract visual Emoji Understanding.

The model shows better performance in simply recognizing the shallow meanings of indi-993 vidual emoji, but in Abstract visual in-depth understanding, it is difficult for the model to 994 work with the contextual emoji information to get the real corresponding relevant emoji 995 meanings. For example, vs means match, PK, duel, competition, and so on. In Chinese, 996 "决" represents duel, and then harmonized to "绝" to get the idiom "精才绝艳". In "African 997 elephant", the superficial meaning of the emoji is the earth, but further combined with the 998 specific location of the earth map in the figure and the hint of an elephant, the emoji rep-999 resents the African elephant. Abstract visual understanding in conjunction with its textual 1000 meaning to further reason about the correct answer. In 🗑 🤎 🝃 💪 "同心叶力 (pull together 1001 with the same goal) ", **6** is an arm, but it does not mean "arm" in idioms. Instead, it is a very strong arm, which corresponds to " \mathcal{I} (power)". 1003

1004

1006

1005 C.3 Chinese Idiom Format

Chinese idioms are a special kind of words, which often have specific formats and seman-1007 tic information, so they cannot directly translate the meaning of a single emoji and con-1008 catenate words into sentences. The most common format is four-character idioms, which 1009 often come from ancient Chinese myths, historical stories, classics, etc., consisting of four 1010 Chinese characters, with a Chinese literary style, and often a symmetrical structure. In 1011 addition, multi-character idioms, although far fewer in number than four-character idioms, 1012 are equally important components. Some of them have less than four words (e.g., three-1013 character idioms) and some have more than four words. Generally speaking, whether it 1014 is a four-character idiom or a multi-character idiom, it follows the one-to-one relationship 1015 between emoji and characters, but there are special cases.

1016

1017 Chinese character mapping Usually, idioms follow a one-to-one relationship between 1018 emojis and characters, but there are special cases. First of all, there will be multiple emo-1019 jis corresponding to one character. Often, many numbers will have this correspondence, 1020 especially those with large digits. For instance, " $\begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & \\ 0 & \end{bmatrix}$ " denotes the " \mathcal{T} " 1021 (ten thousand). In addition, there is a mapping relationship between multiple characters in 1022 an emoji. This kind of correspondence is relatively rare, usually in multi-character idioms, 1023 and this one-to-many mapping relationship occurs when two or more characters can form a new word represented by an emoji. The above two mapping relationships require MLLM to 1024 further complete the understanding and reasoning of multiple emoji contexts on the basis 1025 of recognizing the meaning of a single emoji.

1026 C.4 English Word Task

1028 **English Word Split** In the English word task, unlike the regular one word corresponding to one emoji, it is common for multiple emoji to represent one word. The task usually splits 1029 the word, corresponds multiple emoji to different parts, and finally synthesizes them into 1030 one word. For example, "blackBerry" is split into "black-", "ber-", "-ry", and then the 🐻 is utilized to represent "ber-", and finally, the box of black and the letter "E" is added to 1032 get "blackBerry". The word "Panda" can be split into "Pan-" and "-da," where "Pan-" 1033 corresponds to \mathbf{Q} . This kind of word splitting usually does not occur alone but is also 1034 accompanied by the linguistic phenomenon of harmonic words with many-to-one mapping. 1035 For example, in the word "lemon", the word is split into "le-" and "-mon", then "mon-" 1036 is harmonized as "man", and 👨 is chosen to represent the split syllable "mon". Beyond 1037 understanding the meaning of individual emojis, the MLLM must also remove unnecessary 1038 letters and combine the parts to infer a completely new word.

1039 1040 1041

C.5 ENGLISH IDIOM TASK

Harmonization Word Similar to Chinese idioms, there are also a lot of harmonic characters in English idioms. Sometimes difficult to find directly related emoji to represent, harmonic characters with similar pronunciations are often chosen to replace them. For example, "To be loaded", "To" harmonizes with "Two" 2, and "be" harmonizes with "bee"
Most English idioms still keep the simple direct correspondence between emoji and words. What is more challenging for English idioms is their Abstract visual comprehension and word mapping reasoning problem.

1049

Abstract visual Emoji Understanding In English, for emoji that cannot be represented by direct correspondence, the data do not tend to choose harmonic words, but further associate related emoji, putting further demands on the reasoning ability of MLLM.
For example, in "As genuine as a three-dollar bill", "genius" is usually accompanied by intellect and inspiration, and so a shining star + is used to represent the image of sparkling inspiration of such genius. This deeper level of image comprehension requires a greater understanding of the meaning of the image and the text behind it.

1050

English Word Mapping Unlike most one-to-one relationships in Chinese idioms, there 1058 are a large number of non-one-to-one correspondences in English idioms. Due to the large 1059 number of articles, prepositions, conjunctions, and other words in English that are difficult to directly use emojis, such words are usually omitted in the emoji representation of English 1061 idiom, and only the most critical nouns, adjectives, verbs, etc., are retained to express the core meaning. Therefore, the prediction process of English idiom is not a one-to-one 1062 translation mapping, which also poses more challenges to the ability of MLLM. For example, 1063 in "An apple a day keeps the doctor away.", for MLLM, it is necessary to reason out such 1064 common idioms just for the emojis of 🍎 and 🏥. This examines the internal knowledge-1065 mining ability of the large language model and the strong reasoning ability. However, this 1066 kind of reasoning is also very easy to cause the hallucination problem. 1067

- 1068
- 1069 1070

D Additional Details of Evaluation Metrics

Since our primary goal is to propose the emoji-to-idiom task and assess MLLM's ability to understand and reason about the textual semantics corresponding to abstract visual information, our work primarily focuses on task formulation, data construction, and the underlying assessment approach. We believe this task fills a crucial gap in evaluating MLLM's visual capabilities in representing abstract symbols and bridging the visual-verbal divide. Therefore, our current assessment metrics compare predicted answers with standardized answers that have undergone rigorous automated and manual filtering across multiple granularities.

1078 When we calculate the word-level metrics, we need to match the correct answers exactly,1079 and here we also include the consideration of structural information. The accuracy between the output response and the ground truth of the character-level model does not take into

1080 account the structural one-to-one correspondence, but rather divides and acquires the answer 1081 by character, and calculates it at the character level, as long as the character level can be 1082 matched with the ground truth, it can be regarded as a correct character.

1083 1084 1085

D.1 OVERVIEW OF THE DESIGN OF METRICS AND HOW TO USE

1086 1087 1088

> 1099 1100

Word-level (in Chinese idiom and English word) / Sentence-level (in English id-1089 iom): This is the most direct measure of MLLM's ability to fully understand the semantic 1090 information of the symbols in the image. When MLLM's output and the standard answer 1091 can be matched exactly at word-level or sentence-level (including structural matches), i.e., 1092 when MLLM successfully outputs the correct complete idiom, MLLM is considered to have 1093 answered the question correctly. At this level, we computed the associated precision, re-1094 call, and F-1 values. At this point, MLLM possesses both the understanding of individual 1095 emoji, and moreover the corresponding reasoning ability and text generation ability, which is the one that satisfies our initial motivation and truly realizes the ability of unified visual-1096 linguistic understanding. Therefore, this is the most direct indicator of MLLM's ability. 1097 1098

1101 Character-level (in Chinese idiom and English word)/Word-level (in English 1102 idiom): due to the greater challenge of this benchmark, we found that without additional 1103 training, it is more difficult for MLLM to fully answer the correct and complete idiom. In 1104 order to better analyze which part of emoji-to-idiom comprehension is more challenging for 1105 MLLM, we evaluated at character-level/word-level and calculated Precision, recall, and F-1 1106 values. Specifically, for English idiom, we computed BLEU-1 vs. BLEU-2 to better measure MLLM correctness at this level. Since we did not consider structural information in this 1107 segment, the Character/word-level metrics reflect more on MLLM's ability to understand 1108 individual emoji, due to which there are still a large number of emoji that just need to 1109 understand their meanings directly without additional reasoning. Therefore, MLLM's ability 1110 to understand the emoji themselves is reflected when MLLM receives a higher score in this 1111 item. If MLLM's score in the first item slips very significantly compared to the second item, 1112 we can conclude that MLLM possesses basic emoji comprehension skills but lacks further 1113 reasoning skills.

- 1114
- 1115

1116

1117 **Semantic similarity:** After we computed the character/word-level with exploring Chain-1118 of-thought reasoning, we could not help but notice that sometimes MLLM is actually better at understanding individual emoji, predicting one or two characters correctly, but performs 1119 poorly at the full idiomorphic level. but poorer performance on the complete idiom level. 1120 There are even some MLLMs that correctly determine the meaning of each emoji during 1121 the CoT process, but when outputting the idiom, they output an idiom that has similar 1122 semantics but is completely different at the character level, resulting in serious semantic 1123 drift or even hallucination. Therefore, we added an extra step of calculating the metrics for 1124 the semantic similarity of the output response to the standard answer. 1125

1126 • We use LLM (specifically GPT-40) as an expert to score the semantic similarity 1127 of the output response to the standard answer. The specific scoring criteria are 1128 as follows: scoring is done on a scale of 1-5, with 1 being completely dissimilar, 2 1129 being relatively dissimilarity, 3 fairly similar, 4 being relatively similar, and 5 being 1130 completely similar. The specific prompt we use for scoring is: "Please measure the semantic similarity between the given standard answer and the model output on a 1131 scale of 1 to 5, where 1 means completely dissimilar, 2 means relatively dissimilarity, 1132 3 means fairly similar, 4 means relatively similar, and 5 means completely similar. 1133 Output only a numerical score."

The semantic similarity metrics can be complemented with Character-level/word-level metrics, both of which play an important role when the MLLM is unable to fully match the standard answer at the idiomorphic level. The semantic similarity metric focuses more on whether the answers output by the model are semantically similar to the standard answers, and does not focus on the understanding of individual emoji, but rather reflects an overall comprehension of the semantics of the text directly from the images.

1141
1142 We provide a detailed description of the evaluation metrics and the different capabilities of
1143 MLLM they embody in the emoji2idiom, which helps researchers to use our benchmark and
1144 assess the specific capability bottlenecks of MLLM.

1145

1146 D.2 DETAILS OF EVALUATION METRICS FOR DIFFERENT TASKS

1147 D.2.1 CHINESE IDIOM TASK

1149 In the task of Chinese idioms, we evaluate them separately at the word level and at the 1150 character level. At the word level, we first calculate the Word level accuracy, which is the ratio of the number of words that exactly match the ground truth to the total number 1151 of words. In order to further validate the image-to-language comprehension and reasoning 1152 ability of MLLM, we further propose the Chr-1 and Chr-2 indicators at the word level, which 1153 represent the ratio of the number of words with one or more characters correctly and two 1154 or more characters correctly compared to ground truth, to the total number of words. At 1155 the character level, we compare the difference between each character in the predicted word 1156 and each character in the ground truth to calculate the Precision, Recall, and F-1 values. 1157

1158 D.2.2 ENGLISH WORD TASK

In the task of English words, we evaluate them separately at the word level and at the character level. At both levels, we compare the difference between the predicted word/character and each word/character in the ground truth, calculating the Precision, Recall, and F-1 values.

- 1164
- 1165 D.2.3 English Idiom Task

In the task of English idioms, we evaluate them separately at the sentence level and at the word level. At both levels, we compare the difference between the predicted sentence/word and each sentence/word in the ground truth, calculating the Precision, Recall, and F-1 values. In addition, to further measure the similarity of the generated sentences to ground truth, we further calculated BLEU-1 and BLEU-2 values.

- 1171
- 1172 D.3 DISCCUSION OF METRICS

Since our primary goal is to propose the emoji-to-idiom task and assess MLLM's ability to understand and reason about the textual semantics corresponding to abstract visual information, our work primarily focuses on task formulation, data construction, and the underlying assessment approach. We believe this task fills a crucial gap in evaluating MLLM's visual capabilities in representing abstract symbols and bridging the visual-verbal divide. Therefore, our current assessment metrics compare predicted answers with standard-ized answers that have undergone rigorous automated and manual filtering across multiple granularities, and we have not yet explored further metrics in our evaluation.

- 1181 1182
- 1183 D.3.1 DISCUSSION ABOUT GROUND TRUTH

1184 It is worth noting that an emoji has different meanings in different cultures and contexts, 1185 which is one of the key challenges in emoji-to-idiom task. Therefore, instead of focusing on 1186 understanding the direct meaning of a **single** emoji (in fact, the current MLLM can directly 1187 give multiple possible meanings for a single emoji), we provide a specific contextual **context** (a sequence of multiple emojis with a specific semantic meaning) to limit the semantic of single emoji. In addition, the correct answer of emoji sequence needs to meet the meaning of each emoji in the sequence and the structural information of the sequence, which largely avoids the generation of multiple possible answers.

1191 Certainly, in the process of data collection, we did encounter a very small number of scenarios 1192 where other answers were barely acceptable. For example, " $\mathcal{L} = \mathbf{A}$ ", the standard answer 1193 is "Health is wealth", while the other possible answer is "Money is power".But there are 1194 two problems here: 1) the predicted answer does not fully satisfy the structural information 1195 of the sequence, i.e., translating the idiom from left to right.2) The length of this emoji 1196 sequence is very short, which makes the possible prediction results more variable. As the 1197 length of the sequence becomes longer, the less likely it is that other matching answers 1198 will appear. In our data, the average length of the series is 4.11, 4.23, 7.48, 5.32 in Chinese 1199 four-character idioms, English words, Chinese multi-character idioms, and English idioms. Therefore, we believe that it is feasible to provide a standard answer to predict the outcome 1201 for evaluation, and to measure the consistency of emoji and text.

1202 1203

D.3.2 Discussion about Furthur Metrics

In future work, we plan to develop additional evaluation metrics to better assess MLLM's ability to bridge the multimodal divide between vision and language. Our goals include:

- Adding semantic similarity metrics: Our current metrics primarily quantify the direct correspondence between the standard answer and the generated result, with a relative lack of semantic similarity calculation. Incorporating semantic similarity into our evaluation, particularly in human assessments, will provide a more comprehensive measure of performance.
- 1212 • Measuring the similarity between the emoji's original visual information 1213 and the final prediction: By annotating emojis with a standardized language 1214 base, we can compare results to predictions more effectively. For example, the emoji "禁" might correspond to the textual interpretation "sun, 阳 (read as 'yang')" and 1215 1216 relate to the harmonic word "养 (read as 'yang')" in the final ground truth. While 1217 a predicted result like "日 (sun)" might not match the direct character level, it 1218 captures the initial visual information of the emoji and should be scored 1219 accordingly. This step will help identify specific bottlenecks MLLM faces in this task, whether in visual understanding, harmonic character mapping, or textual 1220 reasoning. 1221
- Including GENERATION metrics: In addition to common generative metrics (e.g., ROUGE, METEOR, diversity, complexity), we will consider task-specific metrics, such as adherence to idiomatic format specifications, like meeting the four-character idiom requirement.
- 1225 1226

1227

1228

1230

- E Additional Details of Baselines and Implementation details
- 1229 E.1 BASELINES
- We select close-source MLLMs, GPT-4V and GPT-40, to evaluate the emoji2idiom benchmark.

GPT-4V Building on the work done for GPT-4, GPT-4 with vision (GPT-4V) enables
users to instruct GPT-4 to analyze image inputs provided by the user.

- 1236
 1237
 1238
 1239
 GPT-40 GPT-40 ("o" for "omni") accepts as input any combination of text, audio, image, and video, which is similar to human response time(opens in a new window) in a conversation. In our work, we choose the GPT-4o-20240513 as our baseline.
- To conduct a richer evaluation, we select a series of open-source MLLMs for testing, including
 Qwen-VL Bai et al. (2023), DeepSeek-VLLu et al. (2024a), LLaVa Li et al. (2023b), and
 CogAgent Hong et al. (2023).

1242
1243
1243
1244
1244
1245
Qwen-VL-7B Qwen-VL (Qwen Large Vision Language Model), proposed by Alibaba Cloud, accepts images, text, and bounding boxes as inputs. It provides Multi-lingual LVLM supporting text recognition and Abstract visual recognition and understanding.

1246
 1247
 1248
 1248
 1249
 DeepSeek-VL-7B DeepSeek-VL is an open-source MLLM designed for real-world vision and language understanding applications, which possesses general multimodal understanding capabilities.

LLaVA-1.5-7B LLaVA is a MLLM that connects a vision encoder and a language model for visual and language understanding, which uses instruction tuning data generated by GPT-4.

 1254
 1255
 CogAgent-18B CogAgent-18B supports image understanding based on CogVLM, which further possesses GUI image Agent capabilities.

1257

1259

1258 E.2 IMPLEMENTATION DETAILS

In our experiments, we explore the inference capabilities of MLLM to accomplish multiple 1260 tasks. In the GPT-4v and GPT-40 tests, we call the official API and use the original 1261 temperature coefficient for the experiment. The time of the GPT4v and GPT4o experiments 1262 in this work has been updated to May 30, 2024. It is important to note that since the closed-1263 source model GPT series will be updated over time, the reproduction of results in future 1264 studies may be affected by the GPT version. In the experiments of the closed-source model, 1265 we use the original official weights for evaluation without additional training. For Qwen-VL, we use the open-source model of Qwen-VL-7B and experiment on a single NVIDIA RTX 1267 3090. For DeepSeek-VL, we experiment with DeepSeek-VL-7B-chat on an NVIDIA RTX 1268 3090. We implement CogAgent-18B on 2 NVIDIA RTX 3090 cards for FP16 inference, and 1269 LLaVA-1.5-7B is also implemented with 2 NVIDIA RTX 3090 cards. For all the evaluations, 1270 we set the temperature as 0.7 and top-k as 0.9. We further provide the computation source and time usage in Table 7. The Emoji2Idiom data and evaluation scripts can be found on 1271 GitHub https://anonymous.4open.science/r/Emoji2Idiom-0CCA. 1272

1273

1277

1274 E.3 PROMPT TEMPLATE

1276 E.3.1 GENERAL PROMPT

For different MLLMs, the templates of the input prompt and message are naturally different due to the different ways the models were originally called. In the MLLM assessment, our prompt design mainly follows the following principles. (1) Keep it as short as possible. Provide effective information in a short prompt to avoid interfering with the understanding of MLLM. (2) Ensure the consistency of the prompts of different MLLMs as much as possible. This ensures that our evaluation results are not affected by the prompt. (3) The design of the different models is designed to give the task concerns more clearly. We show our prompt as shown in Figure 12.

1285 1286

1200 E.3.2 Cot Prompt

¹²⁸⁸ We design the CoT process, inspired by human thinking when seeing the emoji2idiom task.

1289 1290

- 1. Understand each emoji and provide a directly related textual representation.
- 1291
 1292
 1293
 2. Generate possible harmonic words, fine-grained comprehension, and idiom associations.
- 1294 3. Combine multiple emojis to ensure the idioms align or find other possible matches.
 - 4. Finalize the text and check for grammatical errors.



Figure 11: The error bar graphs of different evaluation results of MLLM, which illustrate the Word accuracy of Chinese idiom with four words and Multi-words, Word-level precision of English Word, and Sentence-level precision of English idiom task, respectively.

Table 7: The usage of the computation source and time of MLLMs.

	ÿ				
Model	Hardware	Time Usage	Model	Hardware	Time Usage
GPT-4v GPT-4o Qwen-VL-7b-chat	API API 1 RTX 3090	$\begin{array}{c} 156 \mathrm{min} \\ 149 \mathrm{min} \\ 623 \mathrm{min} \end{array}$	DeepSeek-VL-7b CogAgent-18b LLaVA-1.5-7b	1 RTX 3090 2 RTX 3090 2 RTX 3090	$719 \mathrm{min}$ $503 \mathrm{min}$ $562 \mathrm{min}$

1308

1309

1310 1311 1312

F Additional Details of Automatic Evaluation

To ensure the reliability and robustness of the results, we set up three different random seeds for the experiment in the automatic evaluation of the open-source model and take the average value as the final experimental result. The resulting error bar diagram is shown in Figure 11.

The detailed information of the total amount computed and the type of resources used is shown in Table 7.

1328

1329 G Additional Details of Human Evaluation

1330 1331 G.1 HUMAN EVALUATION GUIDELINE

We invite human experts to conduct human assessments, one for human performance on
 Emoji2Idiom and one for scoring MLLM results. The specific evaluation guideline is shown
 in the Figure 13.

1336 G.2 DETAILED EVALUATION RESULTS

Based on these evaluation guidelines, human experts were able to obtain results from the
evaluation of human performance on the Emoji2Idiom and the evaluation results of MLLMs.
The specific results are shown in the Table 8 and Table 9.

1341

1343

1345

1342 H Additional Details of Case Study

1344 H.1 Typical Case Study

Harmonization Problem There are a large number of harmonic character phenomena in our dataset, which poses a great challenge to the understanding and reasoning of the large language model. The MLLM is also significantly hampered by these harmonic words during emoji understanding. As shown in the Fig. 15 ****** stands for "捷" and 合 stands for "河", and the model does not succeed in recognizing any of these harmonic words. The

	Prompt Template for Evaluation of MLLMs
We provide the d tasks in our Emoji	etails of our prompt designed for different MLLMs for evaluation on different i2Idiom benchmark.
	Task Definition
Task: Emoji-to-Ch Identifier: Chinese Output: Chinese I	inese Idiom / Emoji-to-English Word / Emoji-to-English Idiom e Idiom, English Word, English Idiom diom, English Word, English Idiom
	Prompt Template
Without In-c	ontext learning and additional training, we evaluate the inference ability.
Qwen-VL	'text': 'What is the <identifier> represented by the emojis in this image? Output format: The <output> is', 'image': file_path</output></identifier>
DeepSeek-VL	"content": "'What is the <identifier> represented by the emojis in this image? Output format: The <output> is", "images": ["file_path"]</output></identifier>
LLaVA-1.5	'text': 'What is the <identifier> represented by the emojis in this image? Output format: The <output> is', 'image': file_path</output></identifier>
CogAgent	'text': 'What is the <identifier> represented by the emojis in this image? Output format: The <output> is, 'image': file_path</output></identifier>
InternVL-2	'text': 'What is the <identifier> represented by the emojis in this image? Output format: The <output> is, 'image': file_path</output></identifier>
GPT-4v/GPT- 4o/Claude-3.5- sonnet	<pre>{"type": "text", "text": "What is the <identifier> represented by the emojis in this image? Output format: 'The <output> is'."}, { "type": "image_url", "image_url": { "url": f"data:image/jpeg;base64,{base64_image}"}</output></identifier></pre>
> Without add	itional training, we evaluate the inference ability and In-context learning.
Qwen-VL	<pre>'text': 'What is the <ldentifier> represented by the emojis in this image? Output format: The <output> is' 'image': file_path 'text': 'Here are some <task> examples of the emoji images and the corresponding idioms. Emojis come first, and follow the corresponding <ldentifier> .' 'image': example_image_1 'text': 'The idiom is <ground truth=""> .'</ground></ldentifier></task></output></ldentifier></pre>
GPT-4o	"text": "What is the <ldentifier> represented by the emojis in this image? Output format: 'The <output> is'." "image_url": {"url": f"data:image/jpeg;base64,{base64_image}" "text": "Here are some <task> examples of the emoji images and the corresponding idioms. Emojis come first, and follow the corresponding <ldentifier> ." "image_url": {"url": f"data:image/jpeg;base64,{base64_example_image_1}"} "text": "The idiom is <ground truth="">"</ground></ldentifier></task></output></ldentifier>

Figure 12: Our prompt template is designed for evaluation on MLLMs.

	Guideline of Human Evaluation of MLLM Performance						
This study aim case provides the generated	ns to evaluate the quality of MLLM performance on our benchmark Emoji2Idiom. I you with task type, an emoji image, answer, and ground truth, You need to evalue I answer from the following aspects.						
	Case						
Task: Emoji-to Generated Ar Ground Truth	p-Chinese Idiom Image: Iswer: 闻鸡起舞 : 捷雷不及掩耳 疑 💬 😤 🐔 🔊 🖗						
	Evaluation Metrics						
Normalit formattin	ty: Whether the generated answer conforms to the idiom's specifications, includin ng specifications and semantic specifications						
Options	 Completely non-standard Mostly non-standard Fairly standard Mostly standard Completely standard 						
Examples	 mples 1. "朝三暮四" shows completely standard to the normality. 2. "天鹅绒门盘" shows completely non-standard. 3. "眼大眼小耳朵瞎" shows mostly non-standard. 						
> Semanti	c similarity : Whether the generated answers are semantic similar to the ground tr						
Options	 Completely dissimilar Mostly dissimilar Fairly similar Completely similar 						
Examples	 "眼见为实" shows completely dissimilar to the ground truth "星星点点 "杞人忧天" shows fairly similar to the ground truth "闷闷不乐". 						
> Emotion	al similarity: Whether the generated answers are emotional similar to the ground						
Options	 Completely dissimilar Mostly dissimilar Fairly similar Completely similar 						
Examples	 "狐假虎威" shows mostly similar to the ground truth "阴魂不散". "班门弄斧" shows mostly dissimilar to the ground truth "当务之急". 						
Visual sin	milarity: Whether the generated answers are visually similar to the origin emoji im						
Options	 Completely dissimilar Mostly dissimilar Fairly similar Mostly similar Completely similar 						
Examples	 "Money is power." is mostly similar to the origin emoji image "bright idea." is fairly similar to the emoji image 						
Fluency:	Whether the generated answers are fluency and easy to understand.						
Options	 Completely influent Mostly influent Fairly fluent Mostly fluent Completely fluent 						
Examples	 "Break the ice." is mostly fluent. "日日山如故" is completely influent. 						
> Complex	kity: Whether this task is complex for the MLLM and thus difficult to solve.						
Options	 Completely easy Mostly easy Fairly complex Mostly complex Completely complex 						
Examples	1. "☆☆♀♀ corresponds to 星星点点" is mostly easy to solve.						





Figure 14: Four typical problems the MLLM suffer in English word and idiom tasks.



1567	Table 8: Human performance on Chinese idiom task.								
1568 1569	Human performance	word-level	charact character-2	er-level character-1	Complexity	Time usage			
1570 1571 1572 1573 1574	Human expert-1 Human expert-2 Human expert-3 Human expert-4 Average	$ \begin{array}{c c} 56 \\ 69 \\ 64 \\ 77 \\ 66.5 \end{array} $	$65 \\ 74 \\ 70 \\ 85 \\ 73.5$	$76 \\ 81 \\ 85 \\ 95 \\ 84.25$	$ \begin{array}{c c} 3.5 \\ 3.1 \\ 3.2 \\ 3.4 \\ 3.3 \\ \end{array} $	15s per image 22s per image 28s per image 24s per image 22s per image			
1575									

о. т¹

Table 9: Human evaluation on Chinese idiom task.											
Madal	Expert	Chinese idiom				English idiom					
model		Std.	Sem.	Emj.	Emo.	Flu.	Std.	Sem.	Emj.	Emo.	Flu.
	1	3.7	1.1	1.3	2.4	3.7	4.4	2.3	2.5	2.1	4.3
	2	4.5	1.6	1.8	2.6	3.9	4.2	2.1	2.4	2.2	4.5
GPT-4v	3	3.9	1.1	1.4	2.1	3.6	4.3	2.2	2.5	2.3	4.4
	4	3.9	1.2	2.4	2.2	3.8	4.3	2.2	2.7	2.4	4.5
	Avg.	4.0	1.3	1.7	2.3	3.8	4.3	2.2	2.5	2.3	4.4
	1	4.1	1.4	1.8	2.3	3.8	4.4	2.2	2.6	2.5	4.4
	2	4.8	1.6	2.1	2.3	4.0	4.3	2.2	2.6	2.3	4.7
GPT-40	3	4.1	1.4	1.5	2.4	3.6	4.5	2.3	2.4	2.2	4.2
	4	4.5	1.7	2.5	2.3	3.9	4.2	2.4	2.6	2.5	4.5
	Avg.	4.1	1.4	1.8	2.3	3.8	4.4	2.3	2.6	2.4	4.5

1591 inference of such harmonic words requires the help of relevant contexts, and in our data 1592 scenario, the model is required not to understand individual emoji alone, but to deeply and 1593 comprehensively analyze the context of the emoji. Obviously, under this task requirement, 1594 current multimodal large language models are not well equipped to capture emoji context 1595 with harmonic word reasoning. 1596

1597 Hallucination Problem During the process of recognizing emoji, the model can usually 1598 recognize the corresponding meaning of individual emoji better. At this point, the models 1599 are prone to hallucinations. After recognizing the meaning of a single emoji, they think diffusely about this emoji and only consider words or idioms directly related to the emoji, ignoring the involvement of emoji in other contexts. For example, in Fig. 15 the model recognizes and starts thinking about idioms related to horse and directly outputs "horsing 1603 around" without considering another emoji. Similarly, when MLLMs capture 🔥, they search for the Chinese idiom with the character " χ (fire)". Another example shows that the 1604 1605 GPT-4v and GPT-4o recognize the number 3 and directly associate it with the idiom "朝 1606 三暮四" and "颠三倒四", which contains the number 3, without considering the information of the rest of the emoji around.

- 1608 1609 Multi-to-One or One-to-Multi Character Mapping. For the MLLM, it is customary to perform a one-to-one mapping operation where an emoji corresponds to a Chinese char-1610 acter or English word. In many scenarios, however, it is necessary to perform a multi-to-one 1611 or one-to-multi mapping. For example, in Figure 15the number 1 0 0 0 is composed 1612 of four emojis, but the model does not successfully combine them into one character "f 1613 1614 (one thousand)". And in Figure 14, 🔔 not just indicates a single "bell" or "alarm", but the idiom "ring a bell". This reasoning relies on the capability of knowledge ming and the 1615 reasoning based on the emojis in images and their corresponding linguistic meanings. 1616
 - 1617

1576

1579 1580

1582

1585 1586 1587

1589

Abstract visual Image Understanding of the Emoji Symbol The model shows good 1618 performance in simply recognizing the shallow meanings of individual emoji, but in Abstract 1619 visual understanding, it is difficult to match the emoji information with the context to get the deep corresponding emoji meanings. For example, in Figure 14, the prediction of the idiom "receive a kickback", the model simply captures the emoji idiom, the meaning of "box", and interprets it as "think outside of the box" or "out of the box", but does not combine the package attributes of "receiving something" with the hint of money to generate the correct answer.

I Additional Details of Training and Finetuning for Future WORK

Based on these results and error case studies, we propose potential training methods and frameworks that could significantly improve MLLM performance in visual-linguistic tasks, drawing inspiration from human approaches to joint visual-semantic reasoning:

- Direct fine-tuning: We can incrementally pre-train MLLMs on an emoji-rich corpus to build a basic understanding of emoji. Our initial tests indicate that MLLMs already demonstrate a foundational grasp of emoji, performing well in many cases. Following pre-training, we suggest a 4:1 division of the fine-tuning dataset and test set, with direct fine-tuning on the pre-trained MLLM. This method mirrors human learning, where repeated practice after initial knowledge acquisition leads to mastery in a specific domain.
- 1640 • Incorporating Chain of Thought (CoT) design: When translating emoji to idioms, we can model the process after human reasoning. This CoT design references the 1641 process of human thinking and reasoning, which can assist MLLM to think about 1642 idiom generation in a structured way, and is better able to further analyze where 1643 exactly MLLM goes wrong and provide inspiration for subsequent research work. We hope that such reasoning can be further generalized to more general symbol 1645 understanding, and our emoji2idiom data can also be used as part of general symbol 1646 understanding to evaluate the general symbol understanding capability of the large 1647 language model. 1648
- Adding a symbol mapping set as external knowledge: A single emoji may correspond to multiple characters. By constructing an emoji-to-character mapping set, we can enable MLLM to learn possible alignments. This approach is similar to how humans use external knowledge to accomplish tasks that might be challenging without it.
- Multi-agent invocation: Referring to the CoT process, we can utilize multiple intelligences for tasks like emoji comprehension, harmonic word association, and emoji combination, allowing for integrated task planning, memory iteration, and refined reasoning.

1657 Finally, our work significantly contributes to enhancing the visual comprehension and rea-1658 soning capabilities of MLLMs. Most current unified evaluation metrics focus on MLLM's understanding of natural images, often overlooking abstract visual information and symbolic 1660 representations—areas that receive less attention during training. Additionally, MLLMs struggle with recognizing complex textual information in images, particularly handwritten text or intricate symbols. We believe our emoji2idiom task not only complements existing 1662 evaluations of abstract symbolic representations but also offers a solution for deeper visual 1663 reasoning, thus promoting the development of visual-textual alignment and multimodal uni-1664 fication architecture. 1665

666

1633

1634

1635

1636

1637

1639

1667

1668

1669

1670

1070

1671 1672