

Can MLLMs Understand the Deep Implication Behind Chinese Images?



Chenhao Zhang\*, Xi Feng\*, Yuelin Bai\*, Xinrun Du\*  
Jinchang Hou, Kaixin Deng, Guangzeng Han, Qinrui Li, Bingli Wang, Jiaheng Liu  
Xingwei Qu, Yifei Zhang, Qixuan Zhao, Yiming Liang, Ziqiang Liu, Feiteng Fang  
Min Yang, Wenhao Huang, Chenghua Lin, Ge Zhang†, Shiwen Ni†



*Huazhong University of Science and Technology,  
Shenzhen Institute of Advanced Technology, CAS  
01.ai, M-A-P, Hokudai, UCSB, University of Manchester, etc.*

Introduction

As the capabilities of Multimodal Large Language Models (MLLMs) improve, the need for higher-order evaluation of them is increasing. However, there is a lack of work evaluating MLLM for **higher-order perception** and understanding of **Chinese visual content**. To address this, we introduce the **CII-Bench**, which aims to assess MLLMs' such capabilities for Chinese images. To ensure the authenticity of the Chinese context, images in CII-Bench are sourced from the Chinese Internet and manually reviewed, with corresponding answers also manually crafted. Additionally, CII-Bench incorporates images that represent Chinese traditional culture, such as famous Chinese traditional paintings, which can deeply reflect the model's understanding of Chinese traditional culture. Through experiments on multiple MLLMs using CII-Bench, significant findings emerged. There is a large gap between MLLMs and humans in performance. The highest MLLM accuracy is 64.4%, while the human average is 78.2% and the peak is 81.0%. MLLMs perform poorly on traditional culture images, indicating limitations in understanding high-level semantics and lacking a deep knowledge base of Chinese traditional culture. Moreover, most models have higher accuracy when image emotion hints are added to the prompts. We believe CII-Bench will help MLLMs better understand Chinese semantics and specific images, and move forward the development of expert artificial general intelligence (AGI). Our project is publicly available at <https://cii-bench.github.io>.

Overview

We introduce the **Chinese Image Implication Understanding Benchmark CII-Bench**, a new benchmark measuring the **higher-order perceptual, reasoning** and **comprehension** abilities of MLLMs when presented with **complex Chinese implication images**. These images, including abstract artworks, comics, memes, posters and paintings, possess visual implications that require an understanding of visual details and reasoning ability. CII-Bench reveals whether current MLLMs, leveraging their inherent comprehension abilities, can accurately decode the metaphors embedded within the complex and abstract information presented in these images.

生活 Life	艺术 Art	社会 Society
<p><b>Question:</b> 这张图片有什么隐喻?</p> <p><b>Option:</b></p> <p>(A)有的人不充分考虑自身的情况就盲目的去追求某些事物。</p> <p>(B)羽毛熄灭蜡烛是一种天马行空的想法，讽刺了不切实际的胡乱尝试</p> <p>(C)每个人都有追求时尚的权利。</p> <p>(D)羽毛熄灭蜡烛是一种天马行空的想法，虽然失败了，但这种创新值得赞扬。</p> <p>(E)图中人物多次用羽毛熄灭蜡烛，赞扬了坚持不懈、百折不挠的精神。</p> <p>(F)用羽毛熄灭蜡烛这种不合理的行为，讽刺了有的人做事不考虑周到，盲目的尝试。</p>	<p><b>Question:</b> 远处的小岛暗示了什么?</p> <p><b>Option:</b></p> <p>(A)远处的小岛被描绘为精神寄托的象征，代表了人们在现实世界中寻找精神慰藉和寄托的地方。</p> <p>(B)远处的小岛与远处的棕榈树共同营造出一种与自然和谐共处的氛围，暗示着人与自然之间的和谐关系。</p> <p>(C)远处的小岛象征着希望和<b>目标，虽然距离遥远，但依旧可以到达。</b></p> <p>(D)远处的小岛上有着特定的文化景观或历史遗迹，象征着特定的文化背景或历史时期，提醒人们关注和尊重历史与文化的重要性。</p> <p>(E)远处的小岛作为远方的地标，象征着未知的领域或新的探索方向，鼓励人们勇敢地去探索未知。</p> <p>(F)远处的小岛象征着个人内心深处的平静之地，是人们在面对外界压力和挑战时寻求内心平静和恢复的地方。</p>	<p><b>Question:</b> 这张图片有什么隐喻?</p> <p><b>Option:</b></p> <p>(A)坚持不懈是一种重要的美德。</p> <p>(B)父母的行为习惯决定了孩子的未来。</p> <p>(C)教育的失败是因为家长没有起到足够的监督作用。</p> <p>(D)钢琴的学习应该从小做起并坚持下来，这才能够走向成功。</p> <p>(E)<b>有的家长把教育失败的原因归咎于孩子，却忽略了自身的原因。</b></p> <p>(F)如果父母不以身作则成为榜样，那么将来孩子的教育一定失败。</p>
<p><b>Image Type:</b> 多格漫画(Multi-panel Comic)</p> <p><b>Rhetoric:</b> 隐喻</p> <p><b>Emotion:</b> 消极</p> <p><b>Difficulty Level:</b> 简单</p>	<p><b>Image Type:</b> 插画(Illustration)</p> <p><b>Rhetoric:</b> 象征</p> <p><b>Emotion:</b> 积极</p> <p><b>Difficulty Level:</b> 简单</p>	<p><b>Image Type:</b> 单格漫画(Single-panel Comic)</p> <p><b>Rhetoric:</b> 对比</p> <p><b>Emotion:</b> 消极</p> <p><b>Difficulty Level:</b> 简单</p>
中华传统文化 Chinese Traditional Culture	环境 Environment	政治 Politics
<p><b>Question:</b> 这张图片有什么隐喻?</p> <p><b>Option:</b></p> <p>(A)萧瑟的冬景暗示了人物对于春天到来、万物复苏的渴望。</p> <p>(B)<b>孤身赏雪景暗示了图片中人物淡然、超脱世俗的心境。</b></p> <p>(C)独自一人欣赏雪景暗示了人物内心的孤独和知己难求的悲伤。</p> <p>(D)抬头的动作暗示了人物的思考。</p> <p>(E)独自一人暗示了人物对于亲人和家乡的怀念。</p> <p>(F)萧瑟的冬景暗示了人物内心的悲伤。</p>	<p><b>Question:</b> 这张图片有什么隐喻?</p> <p><b>Option:</b></p> <p>(A)象征着自然界的生物受到人类活动的严重影响，甚至面临灭绝的威胁。</p> <p>(B)这张图片表现了工业技术的飞速发展，暗示着未来生活将更加便利和富裕。</p> <p>(C)这张图片旨在宣传新型环保技术的应用，表现工业与自然和谐共处的美好愿景。</p> <p>(D)暗示了人们有能力通过改变行为模式、采用新技术、实施环保政策等方式，来减轻对自然环境的破坏，实现可持续发展和生态平衡的可能。</p> <p>(E)表达了人类对自然界的彻底征服，通过技术改变地表环境。<b>(F)表达了对环境污染和生态破坏的深刻忧虑，它提醒观者在追求工业发展的同时，不应忽视对自然环境的保护和珍惜。</b></p>	<p><b>Question:</b> 这张图片有什么隐喻?</p> <p><b>Option:</b></p> <p>(A)个体在面对群体或更高权威时，所面临的道德困境和选择。</p> <p>(B)天使和士兵形象之间的冲突暗示了信仰与现实之间的张力，以及个体在面对残酷现实时，如何坚持自己的信仰。</p> <p>(C)图片象征了人类对宗教信仰的追求，表达了对精神世界的渴望。</p> <p>(D)<b>图片可能讽刺了那些以战争干预其他国家或地区的行为，表达了对和平的渴望与对战争后果的担忧。</b></p> <p>(E)个人的命运既受到外力的影响，也取决于个人的选择。</p> <p>(F)即使在平时时期，战争的威胁也可能随时存在；而即使在战争中，人们也可能怀抱着对和平的渴望。</p>
<p><b>Image Type:</b> 绘画(Painting)</p> <p><b>Rhetoric:</b> 隐喻</p> <p><b>Emotion:</b> 积极</p> <p><b>Difficulty Level:</b> 困难</p>	<p><b>Image Type:</b> 海报(Poster)</p> <p><b>Rhetoric:</b> 象征</p> <p><b>Emotion:</b> 消极</p> <p><b>Difficulty Level:</b> 中等</p>	<p><b>Image Type:</b> 插画(Illustration)</p> <p><b>Rhetoric:</b> 隐喻、对比</p> <p><b>Emotion:</b> 消极</p> <p><b>Difficulty Level:</b> 困难</p>

Statistics

CII-Bench contains a total of **698** various Chinese images. These images are manually collected and annotated by 30 undergraduate students from various disciplines and institutions, with sources from multiple renowned Chinese illustration websites. Each image is manually designed with one to three multiple-choice questions, each with six options and only one correct answer. The questions cover the metaphors, symbolism, and detailed understanding of the images. The benchmark includes a total of **800** multiple-choice questions, with 765 questions used to construct the test set and 35 questions used to construct the development and validation set for few-shot tasks.

Statistics	Statistics
Total Questions800	Life216 (30.95%)
Total Images698	Art123 (17.62%)
Dev : Validation : Test15 : 20 : 765	Society157 (22.49%)
Easy : Medium : Hard305 : 282 : 111	Environment51 (7.31%)
Average Question Length10.54	Politics21 (3.01%)
Average Option Length28.31	Chinese Traditional Culture130 (18.62%)
Average Explanation Length121.06	Positive220 (31.52%)
Metaphor562	Neutral247 (35.39%)
Exaggerate121	Negative231 (33.09%)
Symbolism236	Illustration178 (25.50%)
Visual Dislocation42	Meme145 (20.77%)
Antithesis13	Poster87 (12.46%)
Analogy19	Multi-panel Comic34 (4.87%)
Personification73	Single-panel Comic143 (20.49%)
Contrast87	Painting119 (17.05%)

Main Results

Model	Overall (800)	Life (216)	Art (123)	Society (157)	Politics (21)	Env. (51)	CTC (130)	Positive (220)	Negative (247)	Neutral (231)
Open-source Models										
Qwen-VL-Chat	34.3	27.9	34.7	32.5	45.8	55.2	36.5	34.0	35.1	33.6
idefics2-8b	36.3	25.0	46.3	38.1	41.7	56.9	32.9	32.8	39.1	36.4
MiniCPM-Llama3-2.5	40.4	36.3	45.6	37.1	50.0	51.7	40.2	43.2	37.0	41.3
CogVLM2-Llama3-Chinese-Chat	43.4	37.1	48.3	42.3	54.2	63.8	40.2	40.3	45.7	43.8
MiniCPM-v2.6	45.0	37.5	47.6	49.5	58.3	55.2	42.3	45.6	44.6	44.9
LLaVA-1.6-34B	46.0	40.8	55.1	42.8	45.8	62.1	43.1	44.4	48.2	45.2
LLaVA-1.6-72B	48.0	43.8	48.3	49.5	<u>70.8</u>	60.3	43.8	41.5	52.5	49.2
Qwen2-VL-7B	49.6	42.5	51.7	54.1	62.5	65.5	44.5	50.2	47.5	51.2
GLM-4V-9b	50.3	46.7	48.3	53.6	54.2	62.1	48.2	51.9	52.9	46.3
InternVL2-Llama3-76B	52.9	50.8	53.7	51.0	58.3	67.2	51.1	<u>54.8</u>	51.8	52.3
InternVL2-8B	53.1	49.2	53.1	55.7	62.5	63.8	50.4	50.6	53.3	55.1
InternVL2-40B	<u>57.9</u>	<u>55.8</u>	<u>55.1</u>	61.9	62.5	<u>70.7</u>	<u>52.6</u>	54.4	<u>58.0</u>	<u>60.8</u>
Qwen2-VL-72B	<b>64.4</b>	<b>61.7</b>	<b>61.2</b>	<b>68.0</b>	<b>79.2</b>	<b>75.9</b>	<b>59.9</b>	<b>62.7</b>	<b>63.8</b>	<b>66.4</b>
Closed-source Models										
GPT-4o	54.1	54.1	55.8	52.1	50.0	63.8	51.8	51.9	56.2	54.1
Claude-3.5-Sonnet	54.1	52.1	<u>61.9</u>	52.6	62.5	46.6	<u>53.3</u>	52.7	56.5	53.0
Qwen-VL-MAX	56.9	53.3	59.2	58.8	62.5	<u>67.2</u>	52.6	53.9	58.3	58.0
Gemini-1.5 Pro	<u>60.1</u>	<b>60.0</b>	<b>63.3</b>	<u>62.4</u>	<b>70.8</b>	62.1	51.1	<u>54.8</u>	<b>65.6</b>	<b>59.4</b>
GLM-4V	<b>60.9</b>	<u>55.0</u>	59.9	<b>66.5</b>	<u>66.7</u>	<b>79.3</b>	<b>55.5</b>	<b>58.5</b>	<u>64.5</u>	<b>59.4</b>
Text-Only Models										
Llama-3-8B-Instruct	21.7	22.2	26.9	18.6	<u>25.0</u>	27.8	<u>20.4</u>	21.2	24.4	19.5
DeepSeek-67B-Chat	<u>27.1</u>	<u>26.6</u>	<u>32.7</u>	<b>30.9</b>	20.0	<u>35.2</u>	18.2	<u>25.7</u>	22.2	<u>33.2</u>
Qwen2-7B-Instruct	<b>32.5</b>	<b>33.2</b>	<b>34.6</b>	<b>30.9</b>	<b>35.0</b>	<b>40.7</b>	<b>28.5</b>	<b>33.6</b>	<b>30.4</b>	<b>33.6</b>
Humans										
Human_avg	78.2	81.0	67.7	82.7	87.7	84.0	65.9	77.9	75.2	81.6
Human_best	<b>81.0</b>	<b>83.2</b>	<b>73.6</b>	<b>87.2</b>	<b>89.5</b>	<b>86.0</b>	<b>66.7</b>	<b>78.2</b>	<b>78.8</b>	<b>83.3</b>

- Gap between Humans and MLLMs

There is a notable performance gap between MLLMs and humans. Models demonstrate the highest accuracy of **64.4%**, while human accuracy average at **78.2%** and best at **81.0%**.

- Disparity between Open-source and Closed-source Models

Closed-source models generally **outperform** open-source models, but the best-performing open-source model surpasses the top closed source model, with a difference of more than **3%**.

- Model Performance across Different Domains and Emotions

Models perform significantly worse in **Chinese traditional culture** compared to other domains, indicating that current models still lack sufficient understanding of Chinese culture. Further analysis shows that GPT-4o can only observe the **surface-level information**, it's difficult to deeply interpret the complex cultural elements contained in Chinese traditional painting.

- Analysis on different prompt skills

Incorporating image **emotion hints** into prompts generally improves model scores, indicating that models struggle with **emotional understanding**, leading to misinterpretation of the implicit meanings in the images.

Illustration	Meme	Poster
Life (19, 10.68%)	Life (138, 95.17%)	Life (8, 9.20%)
Art (79, 44.38%)	Art (0, 0%)	Art (33, 37.93%)
Society (58, 32.58%)	Society (4, 2.76%)	Society (3, 3.45%)
CTC (0, 0%)	CTC (0, 0%)	CTC (7, 8.05%)
Environment (12, 6.74%)	Environment (3, 2.07%)	Environment (36,41.37%)
Politics (10, 5.62%)	Politics (0, 0%)	Politics (0, 0%)
Multi-panel Comic	Single-panel Comic	Painting
Life (25, 73.53%)	Life (27, 18.88%)	Life (0, 0%)
Art (0, 0%)	Art (11, 7.69%)	Art (0, 0%)
Society (8, 23.53%)	Society (88, 61.54%)	Society (0, 0%)
CTC (0, 0%)	CTC (4, 2.80%)	CTC (119, 100%)
Environment (0, 0%)	Environment (2, 1.40%)	Environment (0, 0%)
Politics (1, 2.94%)	Politics (11, 7.69%)	Politics (0, 0%)

Other Information

■ Paper & Code

PDF: <https://arxiv.org/abs/2410.1385>  
Code: <https://github.com/MING-ZCH/CII-Bench>

■ Authors

Chenhao Zhang: [ch\\_zhang@hust.edu.cn](mailto:ch_zhang@hust.edu.cn)  
Shiwen Ni: [sw.ni@.siat.ac.cn](mailto:sw.ni@.siat.ac.cn)

