



## ***II-Bench: An Image Implication Understanding Benchmark for Multimodal Large Language Models***

Ziqiang Liu, Feiteng Fang, Xi Feng, Xinrun Du, Chenhao Zhang,  
Zekun Wang, Yuelin Bai, Qixuan Zhao, Liyang Fan, Chengguang Gan,  
Hongquan Lin, Jiaming Li, Yuansheng Ni, Haihong Wu, Yaswanth Narsupalli,  
Zhigang Zheng, Chengming Li, Xiping Hu, Ruifeng Xu, Xiaojun Chen, Min  
Yang, Jiaheng Liu, Ruibo Liu, Wenhao Huang, Ge Zhang, Shiwen Ni.

Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences  
University of Chinese Academy of Sciences  
01.ai

NeurIPS 2024 | Presenter: Ziqiang Liu | November, 2024

<https://huggingface.co/datasets/m-a-p/II-Bench>

# Contents

---



中国科学院大学  
University of Chinese Academy of Sciences

- 01.** Motivation
- 02.** II-Bench
- 03.** Experiments
- 04.** Analysis



中国科学院大学  
University of Chinese Academy of Sciences

Can MLLMs Understand the **Deep Implication** Behind Images?





- Numerous challenging and comprehensive benchmarks have been proposed to more accurately assess the capabilities of MLLMs.
- There is a dearth of exploration of the **higher-order perceptual capabilities** of MLLMs.
- We propose the **Image Implication understanding Benchmark, II-Bench**, which aims to evaluate the model's higher-order perception of images.
- We believe that II-Bench will inspire the community to develop the next generation of MLLMs, advancing the journey towards expert artificial general intelligence (AGI).



中国科学院大学

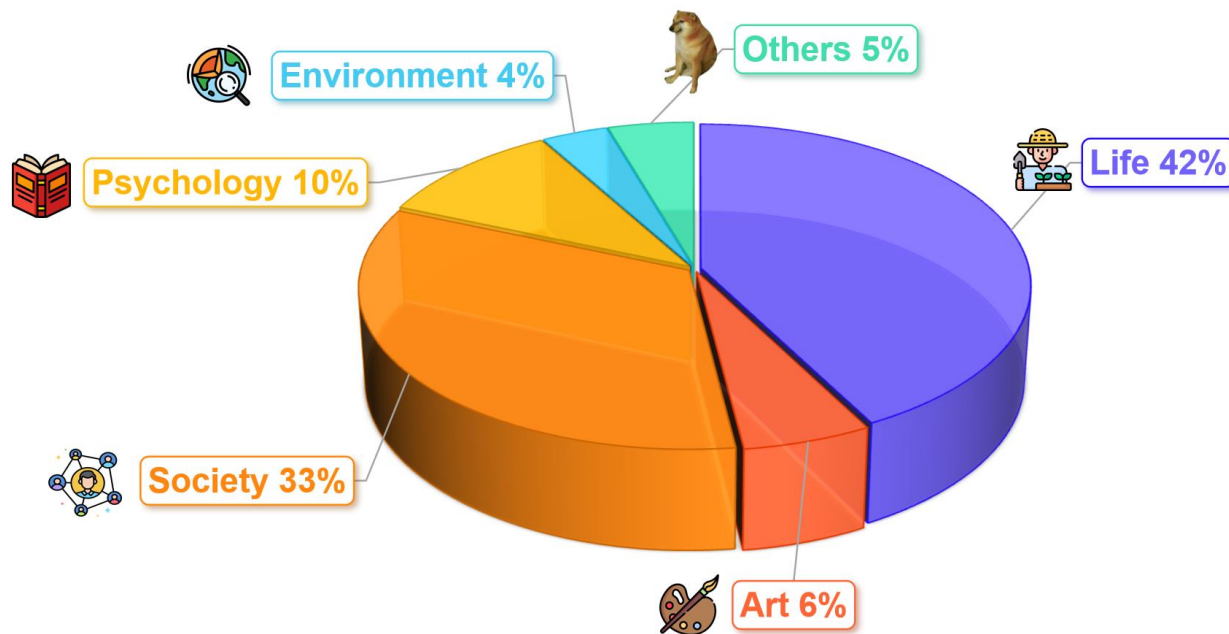
University of Chinese Academy of Sciences

How to construct II-Bench?





- **Data Collection:** We collect 20,150 raw images from various renowned illustration websites, ensuring a sufficiently extensive raw dataset.
- **Data Filtration:** image deduplication -> text-to-image ratio control -> visual inspection
- **Data Annotation:** The annotators mark the images with their difficulty, image type, domain, and corresponding rhetoric first. An explanation of contained visual implications is then drafted for each image, Finally, the annotators devise 1-3 fine-grained questions per image, each with only one correct answer and five distractor options related to the implication nuances.
- **Data Quality Assurance:** Each question and option undergoes multiple rounds of meticulous manual annotation to ensure the distractors are sufficiently challenging and not easily distinguishable from the correct option and ensure consistency across different annotators.



- **1222** images
- **1434** questions
- **6** domains
- **6** categories
- **3** sentiments
- **3** difficulties
- **9** rhetorics





Model	Size	ViT	Projection Module	LLM
CogVLM2-Llama3-Chat [58]	19.5B	EVA2-CLIP-E	MLP	Llama-3-8B + Visual Expert
MiniCPM-Llama3-2.5 [23]	8.5B	SigLip-400M	Perceiver Resampler	Llama3-8B
InternVL-Chat-1.5 [8]	25.5B	InternViT-6B	MLP	InternLM2-20B
InternLM-XComposer2-VL [13]	7B	OpenAI ViT-Large	PLoRA	InternLM-2
DeepSeek-VL-Chat-7B [40]	7.3B	SAM-B + SigLIP-L	MLP	DeepSeek-LLM-7B
InstructBLIP-T5 [11]	4.0B/12.3B	ViT-g/14	MLP	FLAN T5 XL/XXL
BLIP-2 FLAN-T5 [33]	4.1B/12.1B	ViT-g/14	MLP	FLAN T5 XL/XXL
mPLUGw-OWL2 [62]	8.2B	ViT-L/14	Visual Abstractor	Llama-2-7B
Qwen-VL-Chat [3]	9.6B	ViT-bigG	VL Adapter	Qwen-7B
Yi-VL-34B-Chat [63]	7.1B/35.4B	CLIP ViT-H/14	MLP	Yi-34B-Chat
LLaVA-1.6-34B [35]	34.8B	ViT-L/14	MLP	Nous-Hermes-2-Yi-34B
Mantis-8B-siglip-llama3 [26]	8.5B	SigLIP	MLP	Llama-3-8B
Idefics2-8B [28]	8.4B	SigLIP	MLP	Mistral-7B

- Zero(Few)-Shot Prompting: 0, 1, 2, 3 shot(s)
- Chain of Thought Prompting
- Domain: give the image's domain (e.g. life, environment) in the prompt
- Emotion: give the image's emotion(e.g. positive, negative) in the prompt
- Rhetoric: give the rhetorical devices (e.g. metaphor, personification)



# Experiments

## Main Results

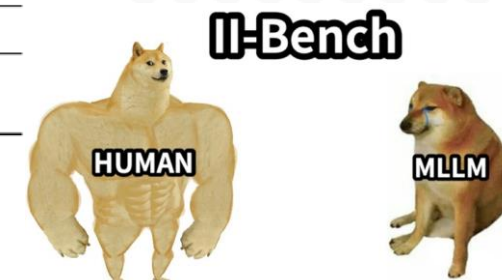


中国科学院大学

University of Chinese Academy of Sciences

	Overall (1,399)	Life (585)	Art (85)	Society (461)	Psy. (152)	Env. (51)	Others (65)	Positive (196)	Neutral (789)	Negative (414)
<i>Open-source Models</i>										
InstructBLIP-T5-XL	47.3	45.6	48.2	48.8	44.7	52.9	50.8	46.9	48.3	45.4
BLIP-2 FLAN-T5-XL	52.8	53.0	58.8	52.5	42.8	64.7	58.5	56.1	52.9	51.0
mPLUGw-OWL2	53.2	54.0	56.5	50.5	52.0	60.8	56.9	55.6	52.6	53.1
Qwen-VL-Chat	53.4	53.2	49.4	52.1	50.0	60.8	72.3	56.1	52.6	53.6
InstructBLIP-T5-XXL	56.7	56.2	58.8	58.6	45.4	64.7	64.6	63.3	56.1	54.6
Mantis-8B-siglip-Llama3	57.5	56.8	61.2	57.5	53.9	64.7	61.5	59.2	58.0	55.6
BLIP-2 FLAN-T5-XXL	57.8	57.1	63.5	57.0	53.3	66.7	66.2	67.9	57.2	54.3
DeepSeek-VL-Chat-7B	60.3	59.0	58.8	58.4	61.8	68.6	76.9	65.8	60.1	58.0
Yi-VL-6B-Chat	61.3	60.9	63.5	60.7	56.6	66.7	72.3	61.7	61.7	60.1
InternLM-XComposer2-VL	62.1	61.7	62.4	62.3	58.6	70.6	66.2	65.8	63.0	58.7
InternVL-Chat-1.5	66.3	63.6	65.9	68.5	65.8	64.7	76.9	73.5	65.4	64.5
Idefics2-8B	67.7	67.2	<b>74.1</b>	67.7	62.5	74.5	70.8	68.9	67.0	68.4
Yi-VL-34B-Chat	67.9	67.5	70.6	67.7	63.8	70.6	76.9	74.0	68.2	64.5
MiniCPM-Llama3-2.5	69.4	68.4	<u>71.8</u>	69.4	64.5	<b>80.4</b>	78.5	<u>75.0</u>	69.3	66.9
CogVLM2-Llama3-Chat	70.3	<u>68.9</u>	<u>68.2</u>	<u>70.9</u>	<u>67.8</u>	72.5	<b>86.2</b>	69.9	<u>71.1</u>	<u>69.1</u>
LLaVA-1.6-34B	<b>73.8</b>	<b>73.8</b>	<u>71.8</u>	<b>73.3</b>	<b>71.1</b>	<u>78.4</u>	<u>81.5</u>	<b>79.1</b>	<b>72.9</b>	<b>72.9</b>
<i>Closed-source Models</i>										
GPT-4V	65.9	65.0	69.4	65.3	59.9	<u>76.5</u>	80.0	69.4	66.0	64.0
GPT-4o	72.6	72.5	72.9	73.3	<u>68.4</u>	<u>76.5</u>	75.4	<u>78.6</u>	<u>71.2</u>	72.5
Gemini-1.5 Pro	73.9	<u>73.7</u>	<b>74.1</b>	<u>74.4</u>	63.2	<b>80.4</b>	<u>83.1</u>	<b>80.1</b>	70.8	<b>75.4</b>
Qwen-VL-MAX	<b>74.8</b>	<b>74.7</b>	<u>71.8</u>	<b>74.6</b>	<b>73.0</b>	<u>76.5</u>	<b>84.6</b>	<b>80.1</b>	<b>74.5</b>	<u>72.9</u>
<i>Humans</i>										
Human_avg	90.3	90.0	88.2	91.4	86.6	96.1	92.3	84.7	89.1	92.2
Human_best	<b>98.2</b>	<b>97.9</b>	<b>98.8</b>	<b>98.3</b>	<b>97.4</b>	<b>100.0</b>	<b>100.0</b>	<b>98.0</b>	<b>98.0</b>	<b>98.8</b>

There is still a huge gap between humans and MLLMs.



# Experiments

## Main Results



中国科学院大学

University of Chinese Academy of Sciences

	Overall (1,399)	Life (585)	Art (85)	Society (461)	Psy. (152)	Env. (51)	Others (65)	Positive (196)	Neutral (789)	Negative (414)
<i>Open-source Models</i>										
InstructBLIP-T5-XL	47.3	45.6	48.2	48.8	44.7	52.9	50.8	46.9	48.3	45.4
BLIP-2 FLAN-T5-XL	52.8	53.0	58.8	52.5	42.8	64.7	58.5	56.1	52.9	51.0
mPLUGw-OWL2	53.2	54.0	56.5	50.5	52.0	60.8	56.9	55.6	52.6	53.1
Qwen-VL-Chat	53.4	53.2	49.4	52.1	50.0	60.8	72.3	56.1	52.6	53.6
InstructBLIP-T5-XXL	56.7	56.2	58.8	58.6	45.4	64.7	64.6	63.3	56.1	54.6
Mantis-8B-siglip-Llama3	57.5	56.8	61.2	57.5	53.9	64.7	61.5	59.2	58.0	55.6
BLIP-2 FLAN-T5-XXL	57.8	57.1	63.5	57.0	53.3	66.7	66.2	67.9	57.2	54.3
DeepSeek-VL-Chat-7B	60.3	59.0	58.8	58.4	61.8	68.6	76.9	65.8	60.1	58.0
Yi-VL-6B-Chat	61.3	60.9	63.5	60.7	56.6	66.7	72.3	61.7	61.7	60.1
InternLM-XComposer2-VL	62.1	61.7	62.4	62.3	58.6	70.6	66.2	65.8	63.0	58.7
InternVL-Chat-1.5	66.3	63.6	65.9	68.5	65.8	64.7	76.9	73.5	65.4	64.5
Idefics2-8B	67.7	67.2	<b>74.1</b>	67.7	62.5	74.5	70.8	68.9	67.0	68.4
Yi-VL-34B-Chat	67.9	67.5	70.6	67.7	63.8	70.6	76.9	74.0	68.2	64.5
MiniCPM-Llama3-2.5	69.4	68.4	<u>71.8</u>	69.4	64.5	<b>80.4</b>	78.5	<u>75.0</u>	69.3	66.9
CogVLM2-Llama3-Chat	70.3	<u>68.9</u>	68.2	<u>70.9</u>	<u>67.8</u>	72.5	<b>86.2</b>	69.9	<u>71.1</u>	<u>69.1</u>
LLaVA-1.6-34B	<b>73.8</b>	<b>73.8</b>	<u>71.8</u>	<b>73.3</b>	<b>71.1</b>	<u>78.4</u>	<u>81.5</u>	<b>79.1</b>	<b>72.9</b>	<b>72.9</b>
<i>Closed-source Models</i>										
GPT-4V	65.9	65.0	69.4	65.3	59.9	<u>76.5</u>	80.0	69.4	66.0	64.0
GPT-4o	72.6	72.5	72.9	73.3	<u>68.4</u>	<u>76.5</u>	75.4	<u>78.6</u>	<u>71.2</u>	72.5
Gemini-1.5 Pro	73.9	<u>73.7</u>	<b>74.1</b>	<u>74.4</u>	63.2	<b>80.4</b>	<u>83.1</u>	<b>80.1</b>	70.8	<b>75.4</b>
Qwen-VL-MAX	<b>74.8</b>	<b>74.7</b>	<u>71.8</u>	<b>74.6</b>	<b>73.0</b>	<u>76.5</u>	<b>84.6</b>	<b>80.1</b>	<b>74.5</b>	<u>72.9</u>
<i>Humans</i>										
Human_avg	90.3	90.0	88.2	91.4	86.6	96.1	92.3	84.7	89.1	92.2
Human_best	<b>98.2</b>	<b>97.9</b>	<b>98.8</b>	<b>98.3</b>	<b>97.4</b>	<b>100.0</b>	<b>100.0</b>	<b>98.0</b>	<b>98.0</b>	<b>98.8</b>

**Small Disparity** between Open-source and Closed-source Models

# Experiments

## Main Results

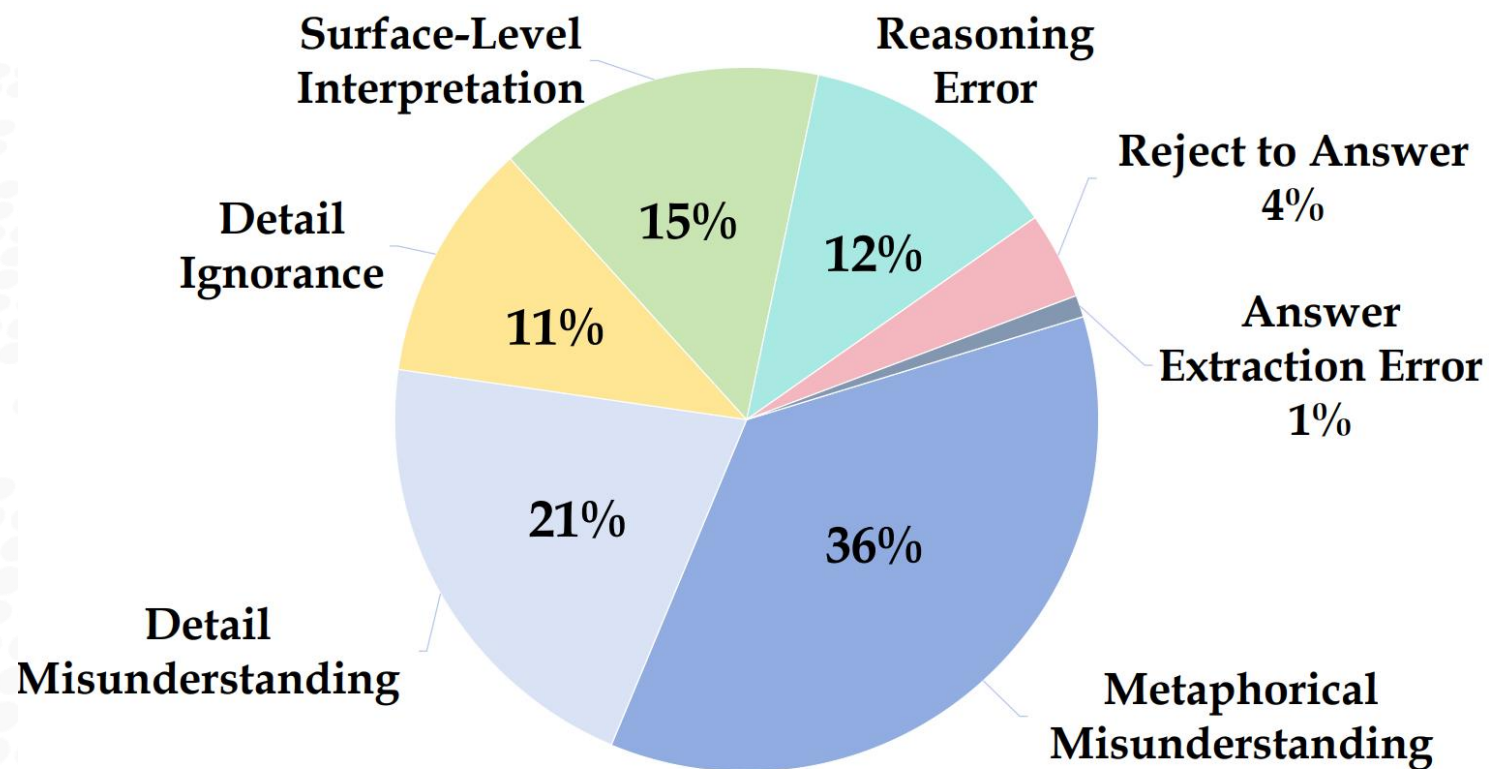


中国科学院大学

University of Chinese Academy of Sciences

	Overall (1,399)	Life (585)	Art (85)	Society (461)	Psy. (152)	Env. (51)	Others (65)	Positive (196)	Neutral (789)	Negative (414)
<i>Open-source Models</i>										
InstructBLIP-T5-XL	47.3	45.6	48.2	48.8	44.7	52.9	50.8	46.9	48.3	45.4
BLIP-2 FLAN-T5-XL	52.8	53.0	58.8	52.5	42.8	64.7	58.5	56.1	52.9	51.0
mPLUGw-OWL2	53.2	54.0	56.5	50.5	52.0	60.8	56.9	55.6	52.6	53.1
Qwen-VL-Chat	53.4	53.2	49.4	52.1	50.0	60.8	72.3	56.1	52.6	53.6
InstructBLIP-T5-XXL	56.7	56.2	58.8	58.6	45.4	64.7	64.6	63.3	56.1	54.6
Mantis-8B-siglip-Llama3	57.5	56.8	61.2	57.5	53.9	64.7	61.5	59.2	58.0	55.6
BLIP-2 FLAN-T5-XXL	57.8	57.1	63.5	57.0	53.3	66.7	66.2	67.9	57.2	54.3
DeepSeek-VL-Chat-7B	60.3	59.0	58.8	58.4	61.8	68.6	76.9	65.8	60.1	58.0
Yi-VL-6B-Chat	61.3	60.9	63.5	60.7	56.6	66.7	72.3	61.7	61.7	60.1
InternLM-XComposer2-VL	62.1	61.7	62.4	62.3	58.6	70.6	66.2	65.8	63.0	58.7
InternVL-Chat-1.5	66.3	63.6	65.9	68.5	65.8	64.7	76.9	73.5	65.4	64.5
Idefics2-8B	67.7	67.2	<b>74.1</b>	67.7	62.5	74.5	70.8	68.9	67.0	68.4
Yi-VL-34B-Chat	67.9	67.5	70.6	67.7	63.8	70.6	76.9	74.0	68.2	64.5
MiniCPM-Llama3-2.5	69.4	68.4	<u>71.8</u>	69.4	64.5	<b>80.4</b>	78.5	<u>75.0</u>	69.3	66.9
CogVLM2-Llama3-Chat	<u>70.3</u>	<u>68.9</u>	68.2	<u>70.9</u>	<u>67.8</u>	72.5	<b>86.2</b>	69.9	<u>71.1</u>	<u>69.1</u>
LLaVA-1.6-34B	<b>73.8</b>	<b>73.8</b>	<u>71.8</u>	<b>73.3</b>	<b>71.1</b>	<u>78.4</u>	<u>81.5</u>	<b>79.1</b>	<b>72.9</b>	<b>72.9</b>
<i>Closed-source Models</i>										
GPT-4V	65.9	65.0	69.4	65.3	59.9	<u>76.5</u>	80.0	69.4	66.0	64.0
GPT-4o	72.6	72.5	72.9	73.3	<u>68.4</u>	<u>76.5</u>	75.4	<u>78.6</u>	<u>71.2</u>	72.5
Gemini-1.5 Pro	<u>73.9</u>	<u>73.7</u>	<b>74.1</b>	<u>74.4</u>	63.2	<b>80.4</b>	<u>83.1</u>	<b>80.1</b>	70.8	<b>75.4</b>
Qwen-VL-MAX	<b>74.8</b>	<b>74.7</b>	<u>71.8</u>	<b>74.6</b>	<b>73.0</b>	<u>76.5</u>	<b>84.6</b>	<b>80.1</b>	<b>74.5</b>	<u>72.9</u>
<i>Humans</i>										
Human_avg	90.3	90.0	88.2	91.4	86.6	96.1	92.3	84.7	89.1	92.2
Human_best	<b>98.2</b>	<b>97.9</b>	<b>98.8</b>	<b>98.3</b>	<b>97.4</b>	<b>100.0</b>	<b>100.0</b>	<b>98.0</b>	<b>98.0</b>	<b>98.8</b>

Model Performance across **Different Domains and Emotions**



**Metaphorical Misunderstanding** is a most common error that GPT-4V makes when generating responses based on image comprehension.



# Future



中国科学院大学

University of Chinese Academy of Sciences



**AGI**



中国科学院大学

University of Chinese Academy of Sciences

# Thanks for your attention!



**Contact:**

[zq.liu4@siat.ac.cn](mailto:zq.liu4@siat.ac.cn)



**Paper:**

<https://arxiv.org/pdf/2406.05862>



**Code:**

<https://github.com/II-Bench/II-Bench>



**Huggingface:**

<https://huggingface.co/datasets/m-a-p/II-Bench>