
Text as Images: Can Multimodal Large Language Models Follow Printed Instructions in Pixels?

Xiujun Li^{∗*} Yujie Lu^{∗*} Zhe Gan[◇] Jianfeng Gao[♡]
William Yang Wang^{*} Yejin Choi[∧]
[∧]University of Washington ^{*}University of California, Santa Barbara
[◇]Apple [♡]Microsoft Research
^{*}Equal contribution

Abstract

Recent multimodal large language models (MLLMs) have shown promising instruction following capabilities on vision-language tasks. In this work, we introduce VISUAL MODALITY INSTRUCTION (VIM)¹, and investigate how well multimodal models can understand textual instructions provided in pixels, despite not being explicitly trained on such data during pretraining or fine-tuning. We adapt VIM to eight benchmarks, including OKVQA, MM-Vet, MathVista, MMMU, and probe diverse MLLMs in both the text-modality instruction (TEM) setting and VIM setting. Notably, we observe a significant performance disparity between the original TEM and VIM settings for open-source MLLMs, indicating that open-source MLLMs face greater challenges when text instruction is presented solely in image form. To address this issue, we train v-MLLM², a generalizable model that is capable to conduct robust instruction following in both text-modality and visual-modality instructions.

1 Introduction

Interleaved image-text data has been increasingly prevalent, ranging from web pages with images and tables, to user interfaces with instructions and forms, in which different modalities interact and blend together. For instance, to perform online shopping, an agent needs to understand the images, instructions and forms. Comprehensive understanding of this multi-modal data demands a range of skills, including recognizing text, understanding images, and also figuring out their interactions.

The current MLLMs are built on top of the pretrained LLMs, and *visual* instruction tuning follows the recipe from its LLM counterparts, specifically, the instruction data is synthesized by the LLMs (mostly from GPT-4 or GPT-4V) in the text format. As illustrated in the left part of Figure 1, the instruction and image are expressed in two modalities, we denote this kind of *visual* instruction data as Text-Modality Instruction (TEM). Under this setting, a pure LLM, for example, Llama 2 or Vicuna in Figure 1 can still make a plausible or correct prediction, even without accessing the image input. All the current benchmarks Fu et al. (2023); Yu et al. (2023); Liu et al. (2023c); Li et al. (2023); Bitton et al. (2023) follow the same format.

This raises a question - *how proficiently these MLLMs can follow instructions if we embed the text instruction into visual format?* As shown in the right part of Figure 1, we name it as VISUAL MODALITY INSTRUCTION, where the image and instruction are in the visual modality. Under the VIM setting, LLMs are not applicable, and LLaVA-1.5 simply repeats the question for the image, may not understand the visual-modality instruction.

¹VIM is short for VISual Modality instruction, code: https://github.com/VIM-Bench/VIM_T00L

²v-MLLM is short for VIM-MLLM, model: <https://huggingface.co/VIM-Bench>

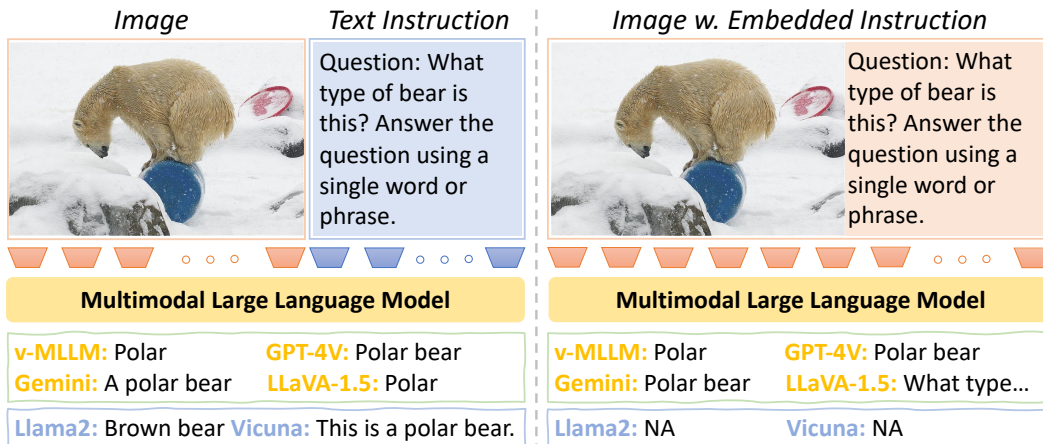


Figure 1: Evaluation paradigm comparison for MLLMs. (a) Left is TEM setting, where Image + Text instruction as two separate modalities are fed into MLLMs for inference; an LLM model (for example, Vicuna) can also make correct prediction, even without accessing to the image. (b) Right: VIM **only** takes the image modality with the text instruction embedded in the image, no additional text prompt is required, LLMs are not applicable. The above example is from OKVQA (question #209725). Note: Image modality input, Text modality input.

Motivated by this, we introduce a new setting, called VISUAL MODALITY INSTRUCTION (short for VIM), evaluating the capability of MLLMs for visual-modality instruction following. We adapt VIM to various benchmarks Marino et al. (2019); Fu et al. (2023); Yu et al. (2023); Lu et al. (2023); Yue et al. (2023), and compose a new benchmark - VIM-Bench. As highlighted in 2, there exists a performance disparity between the TEM and VIM settings for all open-source MLLMs, all of them are not robust enough at visual-modality instruction following. To summarize, our main contributions are:

- We present VISUAL MODALITY INSTRUCTION, a challenging setting to probe the capability of Multimodal Large Language Models for visual-modality instruction following.
- We adapt the VIM to various benchmarks, and reveal a significant disparity for open-source MLLMs between their text-modality instruction setting and VIM setting.
- We train a v-MLLM, which demonstrates robust visual instruction following capabilities.

2 Method

Instruction following, is viewed as one key capability of high-performing MLLMs. In this section, we first present VIM, to examine the instruction following capability of MLLMs, specifically the visual-modality instruction following. Then, we introduce v-MLLM, enhancing the MLLMs with visual-modality instruction following.

2.1 VIM

2.1.1 Visual-Modality Instruction

As illustrated in the left part of Figure 1, the current evaluation norm of MLLMs takes two modalities as input: image and text (as instruction). The existing MLLMs are built on top of the LLMs, benefiting from its strong text understanding capability. For the current MLLM evaluation paradigm, instruction is presented in the text modality, which can utilize the strong language priors from the LLMs for understanding. As shown in Table 1, even a pure LLM model (GPT-4, Llama 2 or Vicuna) can get some success without accessing to the images. Interestingly, on most of eight tasks, Llama 2 shows better numbers over GPT-4 (gpt-4-1106-preview). We manually check some response, and find

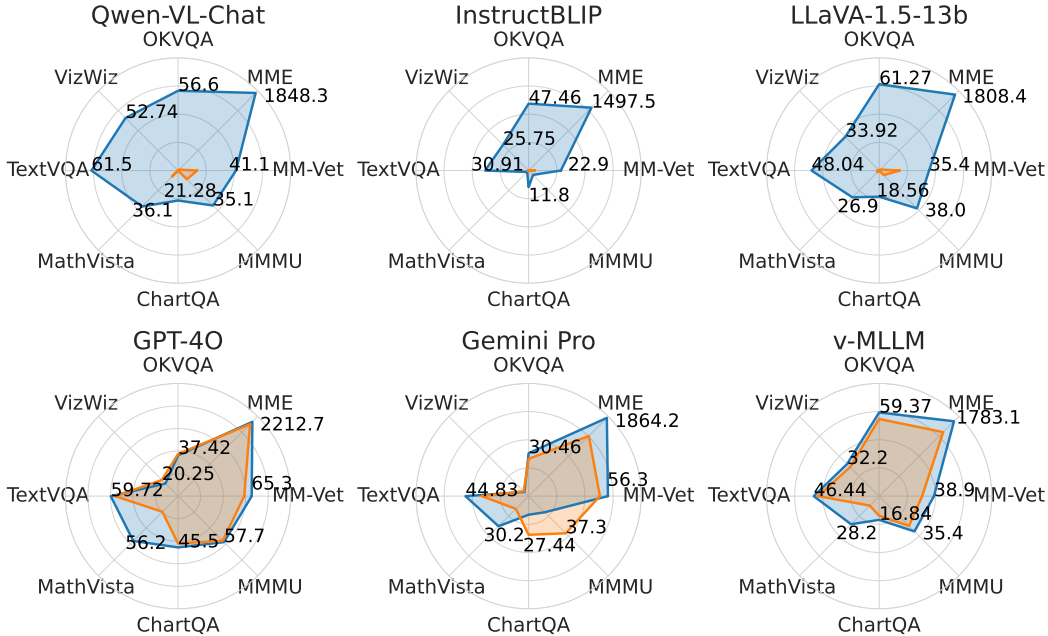


Figure 2: Performance comparison of six selected representative MLLMs for visual instruction following between text-modality instruction (TEM —) and our introduced VISUAL MODALITY INSTRUCTION (VIM —) settings on eight benchmarks. There exists a disparity between TEM and VIM settings for all open-source MLLMs (the first row); GPT-4O, Gemini Pro and our v-MLLM are robust to instruction modality.

that the output from GPT-4 are more reasonable than Llama 2. This might rise several interesting issues, we leave a discussion in Section C.

VIM challenges the MLLMs by rendering the textual instruction into the visual pixel space (image), this enhancement demands not just textual but also strong visual comprehension for instruction understanding. It asks for the strong visual interpretation capability to recognize and follow the embedded instruction in the image.

2.2 v-MLLM

2.2.1 VIM Corpus

One key ingredient of high-performing MLLMs is high-quality instruction tuning data. There are two categories of visual instruction tuning data, one is the synthetic data by LLMs (i.e. GPT-4), like LLaVA Liu et al. (2023b); the other one is the synthetic data generated by GPT-4V, like LVIS-Instruct4V Wang et al. (2023) and ShareGPT-4V Chen et al. (2023). Here, we use the LVIS-Instruct4V-LLaVA-Instruct-mix880k Wang et al. (2023) as our origin instruction tuning corpus D_R , and convert it into the VIM format. We only consider the first turn for the multiple turn conversation data. In total, we get 846k VIM training data D_V after filtering the unavailable image links.

2.2.2 VIM Training

v-MLLM adopts a similar architecture with LLaVA-1.5 Liu et al. (2023a) and LVIS-Instruct4V Wang et al. (2023). The model follows an autoregressive training approach, focusing on optimizing the sequential prediction of the answer words y_1, y_2, \dots, y_n by minimizing the loss function

$$L = \sum_{i=1}^n \text{loss}(\text{LM}(y_{<i}, T, V), y_i)$$

where $y_{<i}$ signifies all tokens preceding the i -th token, T and V represent the textual (e.g., text-modality instruction or prompt) and visual (e.g., visual-modality instruction and image context) tokens. Here the textual token sequence T is optional in the VIM training.

Table 1: Main quantitative results over each benchmark under TEM and VIM settings. : LLM models, : proprietary models, : the proposed models. *We use a more strict evaluation protocol to remove randomness when mapping from open-ended responses to multiple choices.

Models	LLM	Res.	MM-Vet	MME	OKVQA	VizWiz	TextVQA	MathVista*	ChartQA	MMMU*
<i>TEM Setting</i>										
GPT-4	-	-	9.8	74.6	8.37	2.76	3.36	18.7	4.12	28.8
Llama 2	Llama2-7b	-	11.1	1609.5	16.21	5.67	7.18	23.2	0	6.2
Vicuna	Vicuna-7b	-	11.7	1120.6	4.5	1.88	1.88	18.1	0	2.0
GPT-4V	-	-	67.7	1926.6	22.28	17.59	43.14	46.1	28.00	42.9
GPT-4O	-	-	65.3	2212.7	36.20	15.79	59.72	56.2	45.50	57.7
Gemini Pro	-	-	56.3	1864.2	30.46	4.17	44.83	30.2	13.20	16.2
Qwen-VL-Chat	-	-	41.1	1848.3	56.6	52.74	61.5	36.1	21.28	35.1
InstructBLIP	FlanT5 _{XXL}	224	22.9	1497.5	47.46	25.75	30.91	1.4	11.80	4.40
LLaVA-1.5	Vicuna-7b	336	30.5	1851.5	58.41	32.08	45.36	25.1	18.08	36.1
LLaVA-1.5	Vicuna-13b	336	35.4	1808.4	61.27	33.92	48.04	26.9	18.56	38.0
LLaVA-1.6	Vicuna-7b	-	44	1828.6	58.6	34.29	63.61	31.6	22.76	29.3
LLaVA-1.6	Vicuna-13b	-	49.2	1880.5	62.01	35.81	65.79	34.1	26.68	31.2
v-MLLM	Vicuna-7b	336	29.9	1771.1	56.09	30.48	43.38	25.7	16.72	34.0
v-MLLM	Vicuna-13b	336	38.9	1783.1	59.37	32.20	46.44	28.2	16.84	35.4
<i>VIM Setting</i>										
GPT-4V	-	-	63.5	1713.1	28.32	22.18	42.50	12.8	27.44	37.3
GPT-4O	-	-	58.7	2144.3	37.42	20.25	55.88	19.7	42.00	56.0
Gemini Pro	-	-	50.6	1434.6	26.43	4.93	33.24	11.7	15.44	21.9
Qwen-VL-Chat	-	-	13.5	21.2	0.01	0.15	0.27	6.1	0	8.8
InstructBLIP	FlanT5 _{XXL}	224	4.40	0	0.07	0	0.04	0.6	0	0
LLaVA-1.5	Vicuna-7b	336	11.0	2.9	0	0	0	0.9	0	1.4
LLaVA-1.5	Vicuna-13b	336	14.6	24.4	0.38	0	1.51	1.8	0	4.6
LLaVA-1.6	Vicuna-7b	-	20.7	0	0	0	0	6.5	0	8.8
LLaVA-1.6	Vicuna-13b	-	34.2	8.33	0	0	0.01	7.5	0	9.1
v-MLLM	Vicuna-7b	336	25.9	1474.6	52.10	26.40	38.96	7.2	12.24	22.0
v-MLLM	Vicuna-13b	336	30.5	1525.1	54.76	29.15	43.40	9.5	13.96	29.9

3 Experiments

We first build our VIM-Bench based on eight existing representative benchmarks, then compare the v-MLLM with six representative MLLMs under two settings (TEM and VIM) across all the tasks.

3.1 VIM-Bench


Benchmarks To assess the generalization capability of MLLMs, we adapt VIM to eight representative benchmarks, including MME Fu et al. (2023), MM-Vet Yu et al. (2023), OKVQA Marino et al. (2019), VizWiz Bigham et al. (2010), TextVQA Singh et al. (2019), MathVista Lu et al. (2023), ChartQA Masry et al. (2022), and MMMU Yue et al. (2023). The details of source datasets, data processing pipeline, and evaluations can be found in Appendix A.

Data Reformatting Given the above mentioned benchmarks, we try to do minimal changes (i.e., keeping the image resolution) for evaluation. This process involves reformatting text instruction into visual-modality instruction by moving the text instruction into the image modality. In reformatting, we retain the original task’s goal while maintaining the original images with text instructions rendering at the bottom of the image, see the example in Figure 1. These repurposed benchmarks are integrated into our VIM-Bench. Theoretically, VIM can be applied to any existing benchmarks, even for pure NLP tasks. We choose eight representative MLLM benchmarks, although our selections are not exhaustive, they provide a broad basis for MLLM evaluation.

3.2 Main Results

Table 1 summarizes the overall results for two settings. 1). In the original TEM setting, the backbone LLM models can get decent performance on these benchmarks, even without access to the image modality. Interestingly, on six out of eight tasks, Llama 2 is much better than GPT-4, we will briefly

Table 2: MLLMs’ instruction recognition response to the question #3575865 in OKVQA.

<p>Image w. Embedded Instruction</p>  <p>Question: What toy is this? Answer the question using a single word or phrase.</p>	<p>Text Prompt: What is the text in the image?</p> <p style="text-align: center;">Recognized Instructions</p> <p>LLaVA-1.5-13B: The image shows a man sitting in a pew with a teddy bear on his back. The teddy bear is wearing a backpack, and the man appears to be looking at it. The scene takes place in a church, with several other people present in the background.</p> <p>GPT-4V: The text in the image says: "Question: What toy is this? Answer the question using a single word or phrase."</p>
---	---

discuss this issue in Section C. 2). For all open-source MLLMs, there is a significant performance disparity between the TEM setting and VIM setting. 3). GPT-4V and Gemini Pro are robust to the instruction modality, while, open-source MLLMs struggle in the VIM setting, achieve significantly low scores. 4). Our proposed v-MLLM shows robust instruction following capabilities in two settings across all the tasks, especially in the VIM setting, significant gain over open-source MLLMs.

4 Conclusion

In this work, we review the existing MLLMs from a visual perspective, and present VIM, a challenging setting to assess the visual instruction following ability of Multimodal Large Language Models. We adapt VIM to eight benchmarks, leading to VIM-Bench. Through in-depth probing under zero-shot setting, we observed a common issue for the existing open-source MLLMs: all fall short in the VIM setting, in most cases performing not as good as those in the original TEM setting. Furthermore, we train v-MLLM, which demonstrates robust instruction following capabilities under text and visual modality instruction settings on all the tasks.

References

- Jeffrey P Bigham, Chandrika Jayant, Hanjie Ji, Greg Little, Andrew Miller, Robert C Miller, Robin Miller, Aubrey Tatarowicz, Brandyn White, Samuel White, et al. Vizwiz: nearly real-time answers to visual questions. In *Proceedings of the 23rd annual ACM symposium on User interface software and technology*, pp. 333–342, 2010.
- Yonatan Bitton, Hritik Bansal, Jack Hessel, Rulin Shao, Wanrong Zhu, Anas Awadalla, Josh Gardner, Rohan Taori, and Ludwig Schimdt. Visit-bench: A benchmark for vision-language instruction following inspired by real-world use. *arXiv preprint arXiv:2308.06595*, 2023.
- Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*, 2023.
- Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Zhenyu Qiu, Wei Lin, Jinrui Yang, Xiawu Zheng, et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023.
- Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Benchmarking multimodal llms with generative comprehension. *arXiv preprint arXiv:2307.16125*, 2023.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023a.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023b.

- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? *arXiv preprint arXiv:2307.06281*, 2023c.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023.
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*, pp. 3195–3204, 2019.
- Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*, 2022.
- Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 8317–8326, 2019.
- Junke Wang, Lingchen Meng, Zejia Weng, Bo He, Zuxuan Wu, and Yu-Gang Jiang. To see is to believe: Prompting gpt-4v for better visual instruction tuning. *arXiv preprint arXiv:2311.07574*, 2023.
- Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*, 2023.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. *arXiv preprint arXiv:2311.16502*, 2023.

Part I

Appendix

A Benchmark Details

A.1 Source Datasets

We consider eight representative datasets, MM-Vet Yu et al. (2023), MME Fu et al. (2023), OKVQA, VizWiz, TextVQA, MathVista, ChartQA, and MMMU. We also provide probing analysis on VQA_{v2} *test-dev* split, RefCOCO *testA* split, RefCOCO+ *testA* split, and RefCOCOg *test* split to better illustrate how the data is formatted.

In Figure 3, we showcase an example from the source dataset VQA_{v2}, which is with instruction probing setting “Text”. We also consider other two probing settings, “MIX” that have an additional text prompt to guide the visual embedded instruction following, and our proposed “VIM” that only allow the image with embedded instruction as input. We also showcases dataset examples sampled from source datasets RefCOCO, MME, and MM-Vet. All the VIM test samples does not include any additional instruction input in text modality (noted as “NA”).

A.2 Dataset Processing Pipeline

For each source dataset, we start by building up zero shot by embedding instructions into the input image to concatenate as a new image which contains instructions in image modality. In this way, we obtain a new image with embedded instructions for each image-question pair.

A.3 Evaluation Details

Metrics We follow the evaluation pipeline for each benchmark. We use parsing and accuracy for MathVista and MMMU with a more strict protocol. For example, when the model is outputting an empty or random string for a multiple-choice question, in MathVista, the original evaluation protocol will use Levenshtein distance to map to a most similar prediction option, and in MMMU, a random choice from the candidate list will be applied. This will introduce noise and randomness for the evaluation, may not correctly reflect the model performance. In our strict evaluation protocol, we eliminate this random match strategy.

For OKVQA and TextVQA, we follow the leaderboard evaluation to use an evaluation metric that is robust to inter-human variability: $\text{Acc}(\text{ans}) = \min \left\{ \frac{\#\text{humans that said ans}}{3}, 1 \right\}$. For ChartQA, we use relaxed accuracy on human and augmented split.

For MME, the standard metric (*Score*) proposed in Fu et al. (2023) is the summed up Accuracy (*Acc*) and Accuracy+ (*Acc+*) as: $\text{Score} = \text{sum}(\text{Acc} \times 100\%, \text{Acc}_+ \times 100\%)$, where the former one count each correct answer as correct, while the latter one only considers correct when both “Yes” and “No” questions for each image are answered correctly.

For MM-Vet, we use GPT-4 (“gpt-4-0613” version) to automatically provide the score for each sample. The final Accuracy reported as: $\text{Acc} = \frac{\sum_{i=1}^N s_i}{N} \times 100\%$, where s_i is the score at scale 0 – 1 for sample i .

B Full Analysis

B.1 Robustness Analysis

B.1.1 Prompt for MIX Probing Setting

In *mix* probing setting, the MLLMs can accept an extra text instruction input as guidance. The model performance may vary when given different prompts. We report the results using four relevant but diversified prompts (Prompt #1-#5) in Table 7. To be specific, the detailed prompts we use are: 1)

Table 3: Zero shot evaluation results of Text probing setting on VQA_{v2}, MME, MM-vet. This is the popular setting for evaluating text instruction following capability of MLLMs, where the input image and text instruction are both provided.

Models	LLM	Embedded Instruction	<i>Zero shot</i>		
			VQA _{v2}	MME	MM-Vet
Sub set					
LLaVA-1.5	Vicuna-7b	w.o.	60.75	108	31.3
		w.	57.88	88	27.9
	Vicuna-13b	w.o.	61.00	106	35.2
		w.	58.00	87	32.7

Prompt #1: “Answer the question in the image.”, 2) Prompt #2: “Please answer the question that is written in the image.”, 3) Prompt #3: “Follow the instruction embedded in the image.”, 4) Prompt #4: “Detect the question in the image and directly answer to it.”.

B.1.2 Image Embedded with Instruction

To investigate whether the model performance is robust to the minimal changes introduced by the embedded instruction, we give both the original instruction in the text modality and the image with instruction embedded as the image modality to the model.

In Table 3, we present comparative results of LLaVA-1.5 using Vicuna-7b and Vicuna-13b language backbones. It’s observed that embedding images with instructions leads to a marginal decline in performance. This trend suggests that current MLLMs may not be entirely robust to variations in images. However, this performance degradation is minor and within acceptable limits. This implies that the disparity in performance between the Text and VIM probing settings is not solely attributable to changes in the images, but is largely due to the models’ capacity to follow visual embedded instructions.

B.2 Details of Instruction Recognition

Based on Section ??, we conduct an ablation to verify the instruction recognition capability of these MLLMs. Table 9 showcases some example results for zero shot instruction recognition. GPT-4V can recognize the text instruction in both settings, LLaVA can detect some words of the instructions, but may not perfectly recognize the instructions, especially in the one shot setting. Table 4 shows a failure example of GPT-4V for Instruction Recognition, it recognizes the logo texts on the bus as the text instruction.

B.3 Qualitative Observations

B.3.1 Grounding of GPT-4V


Table 5 shows that GPT-4V’s grounding capability can be unleashed when carefully prompted.

C Limitations

We discuss the limitations of our work as follows: 1). Though v-MLLM exhibits robust instruction following capability in both the TEM and VIM settings, it still has a gap with proprietary models, especially GPT-4V, there is still space to improve to be a generalist model. 2). In this work, the evaluation protocols and metrics of VIM setting follow these from the original TEM setting, it also inherits the shortcoming of these evaluations, we leave these for future work. 3). For proprietary models (GPT-4V and Gemini Pro), they are evolving with in-context learning as more queries are fed to the models, the results from these models may not be stable or potentially reproducible.

Table 4: One “failure” examples to show the instruction recognition of GPT-4V.

Image w. Embedded Instruction



**Question: How many girls are on the advertisement?
Answer the question using a single word or phrase.**

Text Prompt:
What is the text in the image?

GPT-4V:

1. "anwb"
2. "Feel Good Bus"
3. "Verrassend leuk vrijwilligerswerk!"

MLLM Evaluation In the LLM exploration setting, Table 1 shows that Llama 2 is much better than GPT-4 on six of eight multimodal tasks, without image input. We manually check some results, and find that the responses from GPT-4 are more reasonable than Vicuna and Llama 2. For example, in Table 8, we only provide the text Prompt to the LLMs, the output from Vicuna and Llama 2 is more like text continuation based on the training corpus, while, GPT-4 makes more reasonable response to the Prompt question.

Another potential issue exposed from this exploration experiment is about MLLM evaluation, though GPT-4 makes more reasonable response, its score is low on all six tasks, which may be contrary to the objective of these benchmarks. In-depth analysis and discussion of evaluation protocol and metrics are beyond the scope of this work, which may leave for future work.

D Broader Impacts

Evaluation and Safety Safety is a critical aspect of MLLMs, particularly when these models are used in the real-world applications. The safety of MLLMs is assessed through various benchmarks that test their robustness against unsafe instructions and harmful content. VIM identified a common issue for the existing open-source MLLMs, it may help to improve the robustness of the MLLMs; and also improve the current evaluation benchmarks of MLLMs.

Table 5: Carefully curated prompting to unleash the grounding capability of GPT-4V.

Image w.o. Embedded Instruction



Text Prompt:
How many people in the image? Give them a label and show me their normalized coordinate in the format [x1, y1, x2, y2] rounded to four decimal places.

GPT-4V:
There are 3 people in the image.
Labels and normalized coordinates for the people are:
1. Person A: [0.1081, 0.0432, 0.3784, 0.7973]
2. Person B: [0.3838, 0.0000, 0.7027, 0.8108]
3. Person C: [0.7027, 0.2973, 0.9676, 0.8108]

Table 6: Zero shot evaluation results on MM-Vet with detail scores across six core vision and language capabilities. highlights the numbers for mix instruction setting. The numbers are presented in %. Rec: Recognition, Know: Knowledge, Gen: Language Generation, Spat: Spatial Awareness.


	Rec	OCR	Know	Gen	Spat	Math	Total
LLaVA-1.5-7b	12	6.6	5.6	5.6	6.4	3.8	10.1
	9.7	7.7	5.1	3.1	6.9	3.8	8.5
LLaVA-1.5-13b	15.2	13.6	6.9	10.9	11.9	3.8	14.4
	18.5	15.7	7.3	9.5	15.3	9.6	16.9
InstructBLIP	6.4	1.6	1.8	1.2	2.7	0	4.4
	14.5	7.8	2.6	0.9	9.3	11.5	12.5
GPT-4V	61.4	65.2	51.2	53.7	67.6	59.2	63.5

Ethical Considerations The deployment of MLLMs necessitates careful consideration of ethical implications, including privacy, bias, and the potential misuse of technology. Ensuring that these models are properly developed and used responsibly is crucial to mitigate risks and maximize their positive impact on society.

Table 7: Zero shot evaluation results on MME subset under mix setting. We compare the performance of LLaVA-1.5-7b and LLaVA-1.5-13b across four different prompts.

Task	LLaVA-1.5-7b				LLaVA-1.5-13b			
	Prompt #1	Prompt #2	Prompt #3	Prompt #4	Prompt #1	Prompt #2	Prompt #3	Prompt #4
artwork	5	5	0	0	0	0	0	0
celebrity	5	5	0	0	1	0	0	0
code reasoning	4	3	0	0	0	0	0	0
color	5	5	0	0	1	1	0	0
commonsense reasoning	4	4	0	0	0	0	0	0
count	5	5	0	0	1	0	0	0
existence	6	6	0	0	0	0	0	0
landmark	5	5	0	0	3	0	0	0
numerical calculation	5	5	0	0	0	0	0	0
OCR	5	5	0	0	0	0	0	0
position	4	5	0	0	0	0	0	0
posters	5	5	0	0	1	0	0	0
scene	5	5	0	0	0	0	0	0
text translation	5	5	0	0	0	0	0	0
<i>Correct</i>	68	68	0	0	7	1	0	0

Table 8: LLMs’ responses to the questions in MM-Vet and OKVQA. (*Image* is not feed into the LLMs.)

<p><i>Image:</i></p> <p>Solve the following equations:</p> <p>1) $8x + 11 = 4x + 14$</p> <p>2) $7d - 4 = 11d - 9$</p>	<p>MM-Vet #1</p> <p><i>Text Prompt:</i> What is d in the last equation? Answer the question in one word or phrase.</p> <p>GPT-4: I'm sorry, but I need more context or information about the "last equation" you are referring to in order to provide an answer. Could you please provide the equation or more details surrounding it?</p> <p>Vicuna: d is the derivative of the function with respect to x.</p> <p>Llama 2: D is the constant of integration.</p>
<p><i>Image:</i></p> 	<p>OKVQA #2971475</p> <p><i>Text Prompt:</i> What sport can you use this for? Answer the question in one word or phrase.</p> <p>GPT-4: Please provide the context or the object you're referring to in order to answer the question about which sport it can be used for.</p> <p>Vicuna: Basketball.</p> <p>Llama 2: Basketball</p>

Practical Applications MLLMs have a wide range of applications, from enhancing accessibility in technology to improving human-computer interactions. There are many real application scenarios for VIM, like UI interface navigation, agent development etc. This can improve efficiency, accuracy, and safety in manufacturing, logistics, and other sectors where visual instruction is crucial. We hope the VIM can benefit for the MLLM’s evolution.

Table 9: Zero Shot Instruct Recognition: MLLMs’s recognition to the example questions in VQA.

<p><i>Image w. Embedded Instruction</i> #42000 in VQA</p>  <p>Question: What color are the gym shoes? Answer the question using a single word or phrase.</p> <p>Text Prompt: What is the text in the image?</p> <p>GPT-4V: Question: What color are they gym shoes? Answer the question using a single word or phrase.</p> <p>LLaVA-1.5-7B: The text in the image is a question asking, “What color are the gym shoes?”</p> <p>LLaVA-1.5-13B: The text in the image is a question asking about the color of the gym shoes.</p>	<p><i>Image w. Embedded Instruction</i> #757000 in VQA</p>  <p>Question: How many elephants are pictured in this photo? Answer the question using a single word or phrase.</p> <p>Text Prompt: What is the text in the image?</p> <p>GPT-4V: Question: How many elephants are pictured in this photo? Answer the question using a single word or phrase.</p> <p>LLaVA-1.5-7B: The text in the image is a question asking how many elephants are pictured in the photo.</p> <p>LLaVA-1.5-13B: The text in the image is a question asking how many elephants are pictured in the photo.</p>
<p><i>Image w. Embedded Instruction</i> #132132000 in VQA</p>  <p>Question: How many men are shopping? Answer the question using a single word or phrase.</p> <p>Text Prompt: What is the text in the image?</p> <p>GPT-4V: Question: How many men are shopping? Answer the question using a single word or phrase.</p> <p>LLaVA-1.5-7B: The text in the image is a question asking, “How many men are shopping?”</p> <p>LLaVA-1.5-13B: The text in the image is a question asking how many men are shopping.</p>	<p><i>Image w. Embedded Instruction</i> #131612000 in VQA</p>  <p>Question: What is the man doing with the horse? Answer the question using a single word or phrase.</p> <p>Text Prompt: What is the text in the image?</p> <p>GPT-4V: Question: What is the man doing with the horse? Answer the question using a single word or phrase.</p> <p>LLaVA-1.5-7B: The text in the image is a question asking, “What is the man doing with the horse?”</p> <p>LLaVA-1.5-13B: The text in the image is a caption that reads “Question: What is the man doing with the horse? Answer the question using a single word or phrase.”</p>







Setting	Source Dataset	Image Modality Input	Text Modality Input
Text	VQAv2		Question: What sport is the man participating in? Answer the question using a single word or phrase.
MIX	VQAv2	 Question: What sport is the man participating in? Answer the question using a single word or phrase.	Answer the question in the image.
VIM (Zero Shot)	VQAv2	 Question: What sport is the man participating in? Answer the question using a single word or phrase.	NA
	RefCOCO series	 Question: What is the normalized coordinate of "the catcher" in the format [x1, y1, x2, y2] rounded to four decimal places?	NA
	MME	 Question: Is this artwork titled the adoration of the shepherds? Answer the question using a single word or phrase.	NA
	MM-Vet	 Question: What will the girl on the right write on the board? Answer the question using a single word or phrase.	NA

Figure 3: Dataset example comparison of three instruction probing settings: Text, MIX and VIM.