

# GSR-BENCH: A Benchmark for Grounded Spatial Reasoning Evaluation via Multimodal LLMs

Anonymous ACL submission

## Abstract

The ability to understand and reason about spatial relationships between objects in images is an important component of visual reasoning. This skill rests on the ability to recognize and localize objects of interest and determine their spatial relation. Early vision and language models (VLMs) have been shown to struggle to recognize spatial relations. We extend the previously released What’sUp dataset (Kamath et al., 2023) and propose a novel comprehensive evaluation for spatial relationship understanding that highlights the strengths and weaknesses of 27 different models. In addition to the VLMs evaluated in What’sUp, our extensive evaluation encompasses 3 classes of Multimodal LLMs (MLLMs) that vary in their parameter sizes (ranging from 7B to 110B), training/instruction-tuning methods, and visual resolution to benchmark their performances and scrutinize the scaling laws in this task.

## 1 Introduction

Earlier efforts for benchmarking vision and language models (VLMs) were developed for cross-modal and/or dual-encoder, end-to-end models, like LXMERT (Tan and Bansal, 2019), CLIP (Radford et al., 2021), BLIP (Li et al., 2022), with the focus on downstream tasks performances such as VQA (Antol et al., 2015), GQA (Hudson and Manning, 2019), referring expressions (Kazemzadeh et al., 2014), image-text matching or image/text retrieval. While spatial relations are often part of VQA datasets, the evaluation of spatial reasoning is often conflated with grounding referring expressions or objects and their attributes<sup>1</sup>. To isolate these issues, authors in (Kamath et al., 2023) introduced a new benchmark that focuses on spatial relationship understanding only. Using image-text matching evaluation methodology, they showed

<sup>1</sup>VQA example question may be: “Is there a woman to the left of the person that is wearing a wetsuit?”

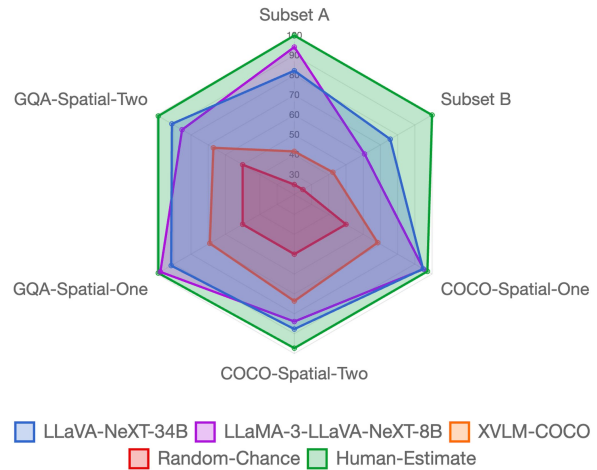


Figure 1: LLAMA-3-LLAVA-NEXT-8B achieves the overall accuracy of 86.1%, compared to 60.4% by XVLM-COCO, in What’sUp benchmark, reaching the best trade-off between accuracy and parameters size, since it performs only 1.1% lower than LLAVA-NEXT-34B, which has  $\times 4.25$  number of parameters.

that contrastive models such as CLIP, BLIP, and their follow-up variants struggle to understand spatial relations with the best accuracy around 61%.

Recent advances in generative large multi-modal models have shown remarkable visual knowledge and reasoning capabilities. We revisit the spatial relationship understanding in the context of MLLMs and extend the existing What’sUp benchmark (Kamath et al., 2023) to include bounding box annotations and depth information. Compositional spatial relationship understanding requires successful recognition of objects and determining their locations. Furthermore, the knowledge of scene depth helps to disambiguate certain relationships (e.g., “in front of” or “behind”). The availability of this information can support a grounded understanding of spatial relations and will contribute to the fine-grained evaluation of large generative MLLMs, which lag behind their earlier counterparts. A few exceptions are multi-task multi-modal

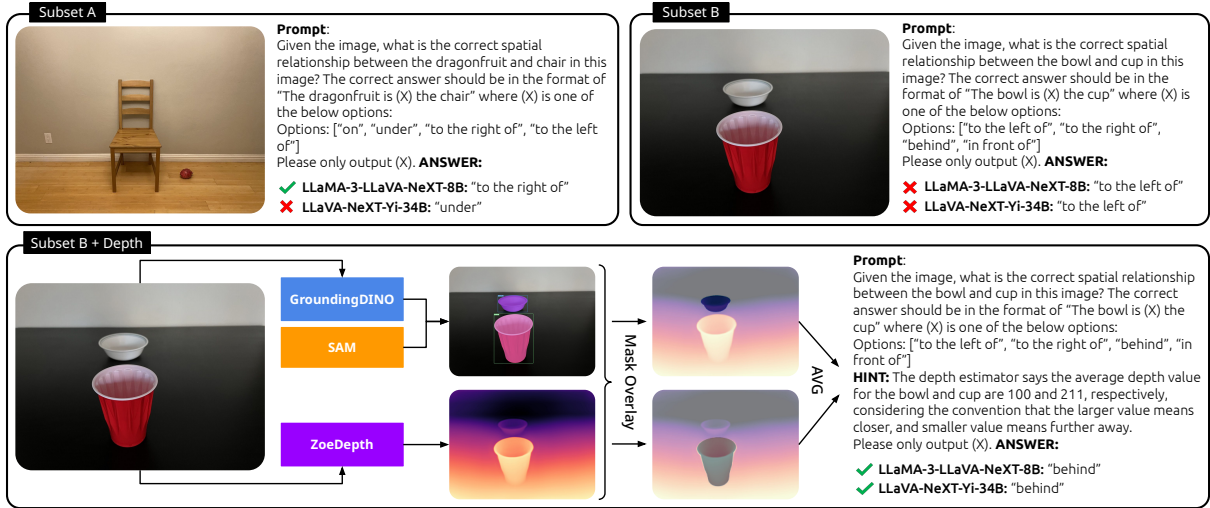


Figure 2: Our pipeline overview for spatial relationship understanding prompting, shown in the top two figures, and our depth-augmented prompting, shown in the bottom figure.

059 benchmarks like MMBench (Liu et al., 2023d) and  
 060 its related benchmarks that focus on evaluating sev-  
 061 eral MLLMs for both visual recognition tasks and  
 062 description generation. Given the simple structure  
 063 of spatial clauses, we can study separately the abil-  
 064 ity of the model to ground the subject and object  
 065 in the clause, and the effect and means of incorpor-  
 066 ating the depth information. The contributions of  
 067 this work can be summarized as follows:

- 068 • Extended What’sUp spatial relationship  
 069 dataset with depth, masks, and bounding box  
 070 annotations.
- 071 • Design of different prompting strategies  
 072 through structured prompting for the evalu-  
 073 ation of grounding and spatial reasoning.
- 074 • Comprehensive evaluation and comparison  
 075 of 18 VLMs and 9 MLLMs, with vari-  
 076 ous sizes, resolutions, pre-training/instruction-  
 077 tuning, and prompting strategies.

## 078 2 GSR Benchmark

079 We extend carefully curated What’sUp dataset (Ka-  
 080 math et al., 2023) that is comprised of Subset A  
 081 containing pairs of objects in unambiguous spatial  
 082 relations, being "on", "under", "left of" or "right  
 083 of" a table, chair, or armchair, and Subset B con-  
 084 taining an object "in front of", "behind", "left" or  
 085 "right" of another object on a tabletop, and subsets  
 086 of COCO-Spatial and GQA-Spatial with either one  
 087 or two objects occurring, accompanied by spatial  
 088 clauses like "on top of", "on the bottom of", "right

089 of", or "left of". To study the grounding in this  
 090 context, we annotate the dataset with bounding box  
 091 coordinates and segmentation masks for all the ob-  
 092 jects mentioned in the captions and the depth maps  
 093 for the images. We leverage GroundingDINO (Liu  
 094 et al., 2023c) as an open-vocabulary object detector,  
 095 Segment Anything (SAM) (Kirillov et al., 2023)  
 096 for the object mask segmentation, and ZoeDepth  
 097 (Bhat et al., 2023) for monocular depth estimation.  
 098 In the next section, we explain in detail how these  
 099 additional annotations enable a more rigorous and  
 100 grounded evaluation of spatial reasoning and its  
 101 components<sup>2</sup>.

## 102 3 GSR-BENCH Experiments

103 Grounded spatial reasoning evaluation is typically  
 104 done using image-text matching, binary VQA, or  
 105 multiple-choice VQA. Further evaluations include  
 106 subject and/or object grounding and localization;  
 107 and exploring the effect of using depth information.  
 108 In addition to 18 VLMs that have been evaluated  
 109 in (Kamath et al., 2023), we focus on the probing  
 110 of open-source generative MLLMs like LLaVA and  
 111 InternVL<sup>3</sup> using structured generation methodolo-  
 112 gies of Multiple choice (MC) and Template-based  
 113 generation (TG). In MC prompting, captions for  
 114 each image are represented as A, B, C, and D op-  
 115 tions for Subset A and Subset B, while A and B  
 116 options for COCO-Spatial and GQA-Spatial Sub-  
 117 sets. Then, the model is prompted to choose the  
 118 correct letter as the final answer. In TG prompting,

<sup>2</sup>All the code and data will be publicly available.

<sup>3</sup>InternVL is the leading model in MMBench.

MODEL	NUM PARAMS	SUBSET A SUB-OBJ	SUBSET B SUB-OBJ	COCO-SPATIAL		GQA-SPATIAL		TOTAL AVERAGE
				ONE-OBJ	TWO-OBJ	ONE-OBJ	TWO-OBJ	
CLIP ViT-B/32 (Radford et al., 2021)	151M	30.3	31.6	43.7	51.1	46.5	47.4	41.8
CLIP ViT-L/14	428M	26.5	25.7	49.2	49.8	46.1	48.5	41.0
NegCLIP (Yuksekgonul et al., 2022)	–	32.5	36.3	47.4	46.4	45.3	46.7	42.4
RoBERTaCLIP (Kamath et al., 2023)	–	25.2	25.0	46.3	53.6	50.8	48.8	41.6
CoCa (Yu et al., 2022)	2.1B	29.4	29.4	48.1	45.2	45.0	49.1	41.0
XVLM 4M (Zeng et al., 2021)	216M	40.0	23.0	58.4	65.0	62.8	54.6	50.6
XVLM 16M	216M	50.7	33.1	65.4	64.5	63.2	53.3	55.0
BLIP 14M (Li et al., 2022)	583M	38.8	38.2	54.2	53.9	49.1	50.5	47.5
BLIP 129M	583M	30.3	30.4	44.8	53.9	50.5	47.4	42.9
BLIP2-ITM (Li et al., 2023)	188M	44.9	30.4	48.3	57.7	46.0	53.6	46.8
BLIP2-ITC	188M	35.9	22.1	55.6	51.8	52.6	49.5	44.6
FLAVA (Singh et al., 2022)	–	33.7	27.2	50.3	55.0	52.2	51.2	44.9
CoCa-Caption	2.1B	25.5	22.8	45.9	51.4	48.5	50.5	40.8
XVLM-Flickr30K	216M	45.1	<u>43.4</u>	63.1	67.3	64.7	58.1	56.9
XVLM-COCO	216M	41.7	42.4	<u>68.4</u>	<u>73.6</u>	<u>69.1</u>	<u>67.0</u>	<u>60.4</u>
BLIP-Flickr30K	583M	29.6	38.0	50.0	58.4	50.3	47.4	45.6
BLIP-COCO	583M	35.7	29.9	46.4	56.4	50.3	52.6	45.2
BLIP-VQA	583M	<u>57.8</u>	37.7	63.6	60.5	63.8	52.9	56.0
LLaVA-1.5-VICUNA	7B	25.0	31.9	90.4	66.6	91.2	62.9	61.3
LLaVA-1.5-VICUNA	13B	58.5	28.2	92.5	78.9	93.1	82.8	72.3
LLaVA-NEXT-MISTRAL	7B	37.4	22.0	81.1	60.4	89.4	57.0	57.9
LLaVA-NEXT-VICUNA	7B	38.6	26.2	95.5	71.8	97.6	79.0	68.1
LLaVA-NEXT-VICUNA	13B	75.0	20.1	<b>95.6</b>	78.6	97.6	84.9	75.3
LLaMA-3-LLaVA-NEXT	8B	<b>94.2</b>	60.8	95.1	83.9	<b>97.8</b>	85.2	<u>86.1</u>
LLaVA-NEXT-YI	34B	82.3	<b>75.7</b>	94.8	<b>87.7</b>	91.5	91.1	<b>87.2</b>
LLaVA-NEXT-QWEN1.5	110B	<u>93.9</u>	54.2	90.6	<u>84.1</u>	96.2	<b>94.2</b>	85.4
INTERN-VL-CHAT-1.5	26B	92.2	<u>61.8</u>	95.1	82.3	<b>97.8</b>	82.8	85.3
Random Chance	–	25.0	25.0	50.0	50.0	50.0	50.0	41.7

Table 1: Template-based generation (TG) results using CircularEval. The first two sections come from What’sUp (Kamath et al., 2023) results. The rest shows our LLaVA 1.5, 1.6, and InternVL-1.5 prompting results. Our best-performing is shown in **bold**, 2nd-best with underline, and What’sUp best-performing with *italic underline*.

as shown in Figure 3, we append the correct format of the entire caption to the prompt, in which the spatial clause acts as the placeholder for the correct spatial relation option. In this way, we are able to leverage LLMs’ open-ended generation capability, handle the models’ verbosity by enforcing the correct answer structure, and overcome the biases observed in MC prompting simultaneously (See Figure 2).

Sample Prompt
Given the image, what is the correct spatial relationship between the <b>subject</b> and <b>object</b> in this image? The correct answer should be in the format of “The <b>subject</b> is (X) the <b>object</b> .”, where (X) is one of the below options: <b>Options:</b> [“on”, “under”, “to the right of”, “to the left of”] Please only output (X), without any other output. <b>ANSWER:</b>

Figure 3: TG sample prompt structure.

We ran each prompt with 4 different permutations so as to vary the position of the answer among the choices in MC and the list of options in TG prompting. An instance is considered correct if all four options are predicted correctly, known as *CircularEval*, introduced in MMBench (Liu et al.,

2023d). As opposed to the CircularEval, there exists VanillaEval, which only asks the model to choose the correct answer from a list of options once and has been shown to be prone to bias in recent studies. We first ran our experiments using MC prompting and observed a significant degree of bias among the models when the position of the answer varied among the choices of A, B, C, or D. This bias and sensitivity turned out to be even more detrimental in smaller models, while larger models like LLaVA-NEXT-YI-34B and LLaMA-3-LLaVA-NEXT-8B showed significantly higher robustness (See Figure 4 in the Appendix for details). This phenomenon also corroborates the findings of multiple recent studies in LLMs (Zheng et al., 2023; Pezeshkpour and Hruschka, 2023; Wang et al., 2023; Xue et al., 2024; Wang et al., 2024). According to this observation, we opted for TG prompting, accompanied by the CircularEval methodology, inspired by Gemini 1.5 Pro (Reid et al., 2024). See Table 1 for the TG prompting results, where rows in section 1 and 2 come from the What’sUp benchmark (Kamath et al., 2023), section 3 refers to LLaVA-1.5 models (Liu et al., 2023b), section 4 to the LLaVA-NeXT models (Li et al., 2024; Liu et al., 2024), and section 5 to the

MODEL	SUBSET A		SUBSET B		COCO-SPATIAL			GQA-SPATIAL			AVG
	SUB	OBJ	SUB	OBJ	ONE-OBJ	SUB	OBJ	ONE-OBJ	SUB	OBJ	G-SCORE
LLAVA-1.5-VICUNA-7B	9.7	79.4	51.5	25.7	47.4	49.8	48.0	31.9	55.0	47.8	44.62
LLAVA-1.5-VICUNA-13B	13.8	86.1	77.4	32.3	60.9	61.8	61.0	42.1	72.2	59.8	56.74
LLAVA-NEXT-VICUNA-7B	14.1	99.0	95.8	66.7	81.9	84.5	77.7	45.5	60.1	56.0	68.13
LLAVA-NEXT-MISTRAL-7B	13.1	82.3	93.9	60.0	<u>87.1</u>	86.8	<u>85.7</u>	69.2	85.9	81.8	74.58
LLAVA-NEXT-VICUNA-13B	15.3	84.0	95.3	67.6	<u>87.1</u>	<b>90.2</b>	83.9	69.6	85.9	80.4	75.93
LLAMA-3-LLAVA-NEXT-8B	19.2	99.3	96.6	73.8	85.7	87.5	83.2	69.0	84.5	80.4	77.92
LLAVA-NEXT-YI-34B	<u>21.1</u>	<b>100.0</b>	<u>97.8</u>	<u>78.9</u>	83.7	85.7	81.4	<u>70.0</u>	<b>88.0</b>	<u>83.5</u>	<u>79.01</u>
LLAVA-NEXT-QWEN1.5-110B	<b>29.4</b>	98.5	<b>98.8</b>	<b>80.1</b>	<b>88.7</b>	<u>88.2</u>	<b>86.4</b>	<b>74.9</b>	<u>86.9</u>	<b>84.2</b>	<b>81.61</b>
GroundingDINO [avg( $\rho$ )]	58.8	92.0	78.1	70.1	62.3	62.8	59.3	59.4	70.4	65.2	67.84
GroundingDINO [ $\Sigma(\rho \geq 0.5)/t$ ]	68.9	100.0	90.0	88.7	71.0	73.6	66.4	59.1	76.3	71.1	76.51

Table 2: Grounding/Localization results. AVG G-SCORE refers to the mean accuracy of  $\text{IoU} \geq 0.5$ . The bottom two rows refer to the GroundingDINO mean confidence scores ( $\rho$ ), and mean accuracy of  $\rho \geq 0.5$ , respectively.

InternVL-1.5 results (Chen et al., 2024).

**Grounding/Localization Evaluation.** This experiment aims to measure the MLLMs grounding ability of the objects mentioned in the captions. Recent studies like (Rajabi and Kosecka, 2023) on Visual Spatial Reasoning (VSR) benchmark (Liu et al., 2023a) has demonstrated that there exist multiple cases where the VLM correctly predicts the binary ITM label of 1 using the holistic representations of the image and caption, while the model fails to localize the subject and object correctly. Our experiments aim to quantify these type of behaviors in MLLMs. We prompt MLLMs to extract the normalized bounding box coordinates for the caption’s objects as "Give me the bounding box coordinates for the {object}" and compute the IoU between the model’s output and the GroundingDINO output for each object, assigning the binary accuracy of 1 if  $\text{IoU} \geq 0.5$ , otherwise 0. See Table 2 for the results.

MODEL	W/O DEPTH	WITH DEPTH
INTERNVL-CHAT-1.5-26B	26.5	40.7
LLAMA-3-LLAVA-NEXT-8B	53.4	60.3
LLAVA-NEXT-YI-34B	64.7	81.9

Table 3: DAP results for *behind* & *in front of* cases.

**Depth-Augmented Prompting (DAP).** The experiments in Table 1 revealed that Subset B is the lowest-performing, with many instances requiring reasoning about "*behind*" and "*in front of*" spatial clauses. We propose to incorporate the depth values of subject and object into the prompt, as a hint to the model, utilizing our augmented benchmark annotations, depicted in Figure 2. We show that this minimal change improves the accuracy of top-3 performing models in these instances of Subset B, reported by CircularEval in Table 3.

## 4 Discussion

According to Table 1 and 2, there is a positive correlation, even stronger in grounding, between **scaling the LLM size & visual resolution**, and the **overall accuracy** in both tasks. Conversely, there exist multiple exceptions, which are inevitable to concretely justify due to various intervening factors, such as (1) differences in training/fine-tuning & architectures and (2) release date and further instruction-tuning of the LLMs, like LLAMA-3-8B, which has the most-recent knowledge cut-off.

Grounding small objects, which refers to the SUB column in Subset A, seems challenging for all, and worst in smaller models, according to Table 2. We also observed a plateau in Table 1, especially in QWEN-1.5-110B, which is the largest ever released open-source MLLM at the moment. This could be a sign of saturation where the reasoning capability flattens out, although scaling still improves grounding, shown in Table 2.

## 5 Conclusions

In this work, we introduce a new benchmark for grounded spatial reasoning by enriching the What’sUp dataset with additional supervision for a more fine-grained assessment of MLLM’s spatial understanding. We also propose a new compositional evaluation methodology for (1) a stricter assessment of spatial relationship understanding through CircularEval, and (2) measuring the model’s grounding capability using the labels we generate through our cost-effective auto-annotation pipeline. Our evaluations reveal the superiority of LLaVA MLLMs over the best-performing VLMs evaluated in What’sUp, like XVLM, by a significant margin of  $\sim +26.8\%$ . Future works may investigate the remaining gap between the top open-source MLLMs and human-level accuracy.

## 228 Limitations

229 **Small-scale Dataset:** Our split sizes remain the  
230 same as the What’sUp dataset in which Subset A  
231 has 412, Subset B has 408, COCO-Spatial-One has  
232 2247, COCO-Spatial-Two has 440, GQA-Spatial-  
233 One has 1160, and GQA-Spatial-Two has 291 in-  
234 stances. Although this benchmark includes 4,958  
235 image instances in total, each instance covering one  
236 or two objects, with various domain shifts in each 6  
237 split, it is smaller than already existing benchmarks  
238 related to spatial reasoning, like Visual Genome  
239 (Krishna et al., 2017), GQA (Hudson and Manning,  
240 2019), VSR (Liu et al., 2022), SpatialSense (Yang  
241 et al., 2019), MMBench (Liu et al., 2023d), etc.  
242 The reason is that this work aims to provide a care-  
243 fully curated benchmark for spatial relationship  
244 understanding evaluation in a controlled setting  
245 to abstract away intervening factors that make the  
246 evaluations noisy.

247 **Limited Spatial Prepositions:** Following the  
248 What’sUp dataset, our benchmark is also confined  
249 to the primitive spatial clauses of *on*, *under*, *behind*,  
250 *in front of*, *to the left of*, *to the right of*, *below* and  
251 *above*, when having two objects involved in the  
252 caption, and, *on the top*, *on the bottom*, *on the left*  
253 and *on the right* when having only one object in  
254 the caption, like in COCO-Spatial-One and GQA-  
255 Spatial-One.

256 **Lack of Robustness in MC Prompting:** In ad-  
257 dition to the similar findings of MC noisiness in  
258 LLMs that we discussed earlier, we hypothesize  
259 that the higher degree of variance in multiple-  
260 choice results in the last two subsets (COCO and  
261 GQA), which is more significant in the smaller  
262 models, could be due to the language domain dis-  
263 tribution shift. Most of the LLMs and MLLMs  
264 are being trained and evaluated with 4 options in  
265 the multiple-choice settings. Conversely, in the  
266 last two subsets, we have two captions per image,  
267 which means we only provide options A and B to  
268 the model in the prompt instead of ABCD without  
269 any fine-tuning for this task or this specific type of  
270 prompting.

271 **Intern-VL-1.5 Poor Grounding Observation:**  
272 An unexpected, significant noisiness in the output  
273 of grounding/localization prompting of InternVL-  
274 1.5 model prevented us from analyzing and report-  
275 ing the results for this model, which requires fur-  
276 ther investigation since a similar behavior has been

observed through our interaction with the InternVL-  
1.5 demo, as well.

**Depth Augmentation Nuances:** The issue we  
noticed in the DAP experiment was the distraction  
the depth hint can cause in cases where multiple  
correct relationships hold in the image. For in-  
stance, object A can be *to the left of* object B, and  
also *in front of* object B, at the same time. So, in  
these ambiguous cases, incorporating depth could  
make the model’s decision biased towards the *in  
front of* preposition, while the ground-truth might  
be *to the left of* in this case. Therefore, we believe  
that trying both prompts, with and w/o depth hint,  
would be helpful for disambiguation in such cases.

**No Human Annotation:** Due to the resource  
constraints, our extended benchmark relies on the  
pseudo-labels we generate using state-of-the-art,  
off-the-shelf models like GroundingDINO, SAM,  
and ZoeDepth. Future works could incorporate hu-  
man inspection and labeling for further robustness  
in annotations.

## References

- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Mar-  
garet Mitchell, Dhruv Batra, C. Lawrence Zitnick,  
and Devi Parikh. 2015. VQA: Visual Question An-  
swering. In *International Conference on Computer  
Vision (ICCV)*.
- Shariq Farooq Bhat, Reiner Birkel, Diana Wofk, Peter  
Wonka, and Matthias Müller. 2023. Zoedepth: Zero-  
shot transfer by combining relative and metric depth.  
*arXiv preprint arXiv:2302.12288*.
- Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye,  
Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi  
Hu, Jiapeng Luo, Zheng Ma, et al. 2024. How far  
are we to gpt-4v? closing the gap to commercial  
multimodal models with open-source suites. *arXiv  
preprint arXiv:2404.16821*.
- Drew A Hudson and Christopher D Manning. 2019.  
Gqa: A new dataset for real-world visual reasoning  
and compositional question answering. In *Proceeed-  
ings of the IEEE/CVF conference on computer vision  
and pattern recognition*, pages 6700–6709.
- Amita Kamath, Jack Hessel, and Kai-Wei Chang. 2023.  
What’s “up” with vision-language models? investi-  
gating their struggle with spatial reasoning. In *Pro-  
ceedings of the 2023 Conference on Empirical Meth-  
ods in Natural Language Processing*, pages 9161–  
9175, Singapore. Association for Computational Lin-  
guistics.
- Sahar Kazemzadeh, Vicente Ordonez, Mark Matten,  
and Tamara Berg. 2014. Referitgame: Referring to

328	objects in photographs of natural scenes. In <i>Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)</i> , pages 787–798.	
329		
330		
331		
332	Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. 2023. Segment anything. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 4015–4026.	
333		
334		
335		
336		
337		
338	Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. <i>International journal of computer vision</i> , 123:32–73.	
339		
340		
341		
342		
343		
344	Bo Li, Kaichen Zhang, Hao Zhang, Dong Guo, Renrui Zhang, Feng Li, Yuanhan Zhang, Ziwei Liu, and Chunyuan Li. 2024. <a href="#">Llava-next: Stronger llms supercharge multimodal capabilities in the wild.</a>	
345		
346		
347		
348	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. <i>arXiv preprint arXiv:2301.12597</i> .	
349		
350		
351		
352	Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In <i>International Conference on Machine Learning</i> , pages 12888–12900. PMLR.	
353		
354		
355		
356		
357	Fangyu Liu, Guy Emerson, and Nigel Collier. 2022. <a href="#">Visual spatial reasoning.</a> <i>arXiv preprint arXiv:2205.00363</i> .	
358		
359		
360	Fangyu Liu, Guy Emerson, and Nigel Collier. 2023a. Visual spatial reasoning. <i>Transactions of the Association for Computational Linguistics</i> , 11:635–651.	
361		
362		
363	Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024. <a href="#">Llava-next: Improved reasoning, ocr, and world knowledge.</a>	
364		
365		
366	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. Visual instruction tuning. In <i>NeurIPS</i> .	
367		
368	Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. 2023c. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. <i>arXiv preprint arXiv:2303.05499</i> .	
369		
370		
371		
372		
373	Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. 2023d. Mmbench: Is your multi-modal model an all-around player? <i>arXiv preprint arXiv:2307.06281</i> .	
374		
375		
376		
377		
378	Pouya Pezeshkpour and Estevam Hruschka. 2023. Large language models sensitivity to the order of options in multiple-choice questions. <i>arXiv preprint arXiv:2308.11483</i> .	
379		
380		
381		
	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In <i>International Conference on Machine Learning</i> , pages 8748–8763. PMLR.	382 383 384 385 386 387 388
	Navid Rajabi and Jana Kosecka. 2023. Towards grounded visual spatial reasoning in multimodal vision language models. <i>arXiv preprint arXiv:2308.09778</i> .	389 390 391 392
	Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. <i>arXiv preprint arXiv:2403.05530</i> .	393 394 395 396 397 398
	Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. 2022. Flava: A foundational language and vision alignment model. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 15638–15650.	399 400 401 402 403 404
	Hao Tan and Mohit Bansal. 2019. Lxmert: Learning cross-modality encoder representations from transformers. <i>arXiv preprint arXiv:1908.07490</i> .	405 406 407
	Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. 2023. Large language models are not fair evaluators. <i>arXiv preprint arXiv:2305.17926</i> .	408 409 410 411
	Xinpeng Wang, Chengzhi Hu, Bolei Ma, Paul Röttger, and Barbara Plank. 2024. Look at the text: Instruction-tuned language models are more robust multiple choice selectors than you think. <i>arXiv preprint arXiv:2404.08382</i> .	412 413 414 415 416
	Mengge Xue, Zhenyu Hu, Meng Zhao, Liqun Liu, Kuo Liao, Shuang Li, Honglin Han, and Chengguo Yin. 2024. <a href="#">Strengthened symbol binding makes large language models reliable multiple-choice selectors.</a>	417 418 419 420
	Kaiyu Yang, Olga Russakovsky, and Jia Deng. 2019. Spatialsense: An adversarially crowdsourced benchmark for spatial relation recognition. In <i>International Conference on Computer Vision (ICCV)</i> .	421 422 423 424
	Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. 2022. Coca: Contrastive captioners are image-text foundation models. <i>arXiv preprint arXiv:2205.01917</i> .	425 426 427 428
	Mert Yuksekogunul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. 2022. When and why vision-language models behave like bag-of-words models, and what to do about it? <i>arXiv preprint arXiv:2210.01936</i> .	429 430 431 432 433
	Yan Zeng, Xinsong Zhang, and Hang Li. 2021. Multi-grained vision language pre-training: Aligning texts with visual concepts. <i>arXiv preprint arXiv:2111.08276</i> .	434 435 436 437

438 Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and  
439 Minlie Huang. 2023. Large language models are not  
440 robust multiple choice selectors. In *The Twelfth Inter-*  
441 *national Conference on Learning Representations*.

## 442 **A Appendix**

443 The appendix is organized as follows:

- 444 • Figure 4 demonstrates the biases of multiple-  
445 choice (MC) prompting.
- 446 • Figures 5 - 11 depict the distributions of ob-  
447 jects occurring in the captions.
- 448 • Figure 12 shows sample failures in grounding  
449 small objects in Subset A.

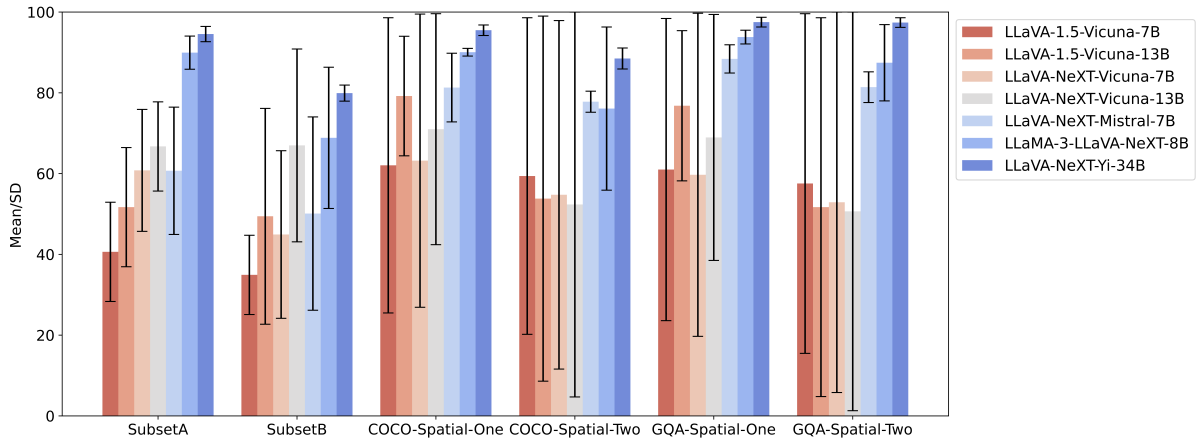


Figure 4: Sensitivity of the models to different permutations of choice order, in the multiple-choice (MC) experiment, which is more significant in the smaller models, and when having two choices of A and B instead of regular 4-choice of A, B, C, and D. LLAVA-NEXT-YI-34B demonstrates an excellent robustness against this issue.

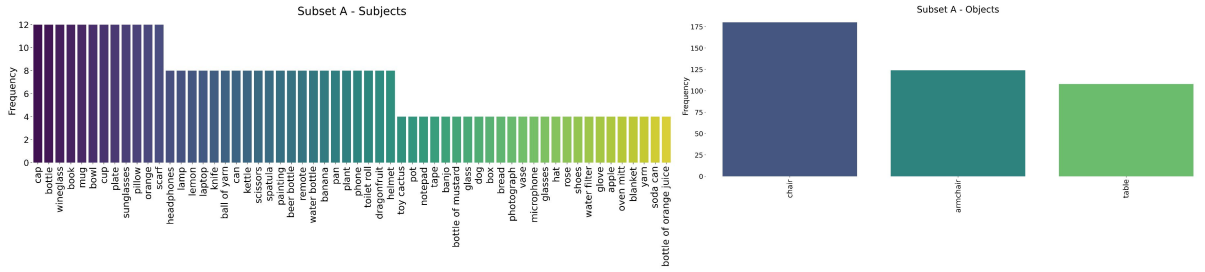


Figure 5: Subset A - subjects and objects distributions.

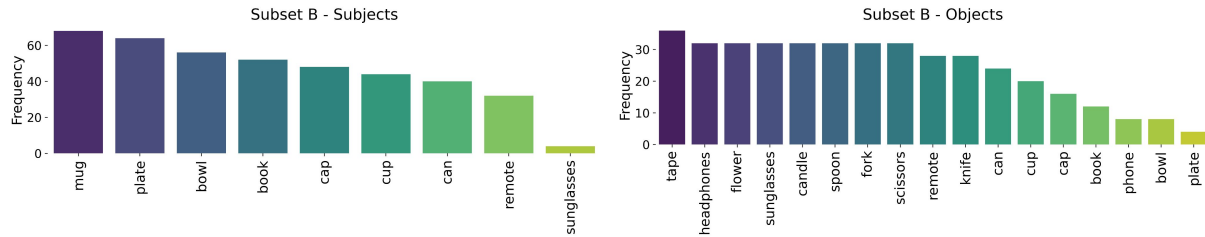


Figure 6: Subset B - subjects and objects distributions.

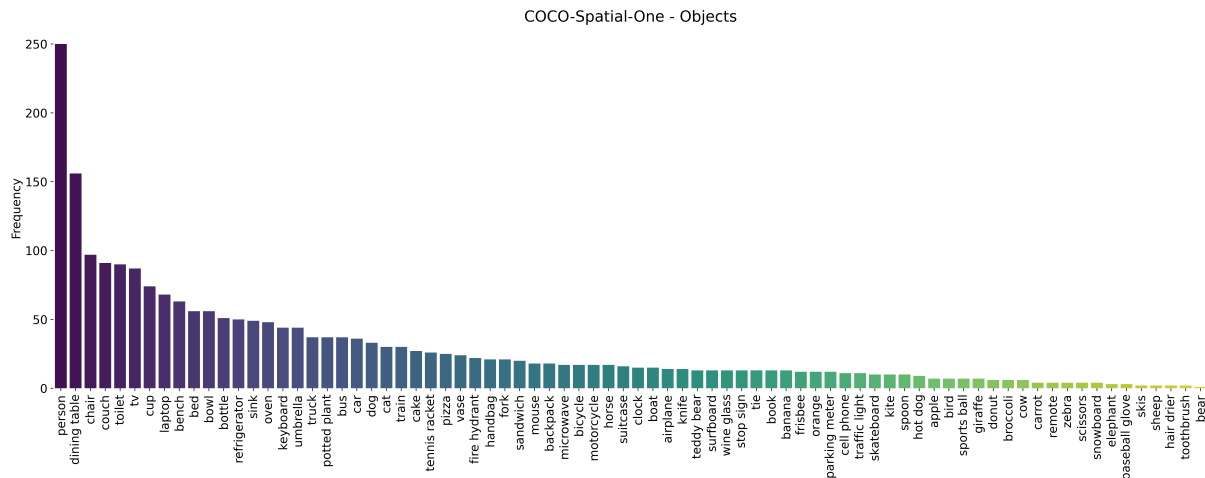
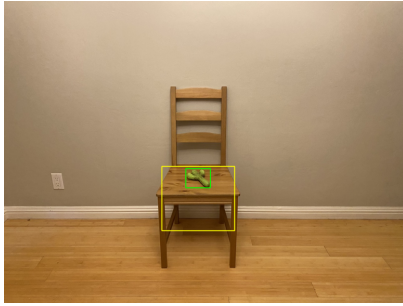


Figure 7: COCO Spatial One - objects distribution.







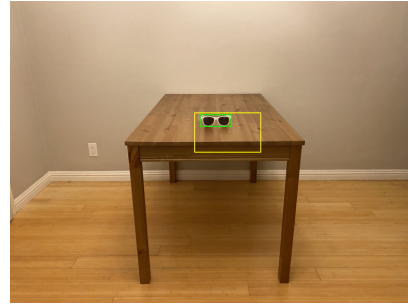
(a) **toy cactus** on chair



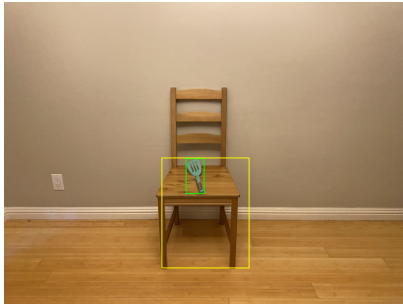
(b) **wineglass** under armchair



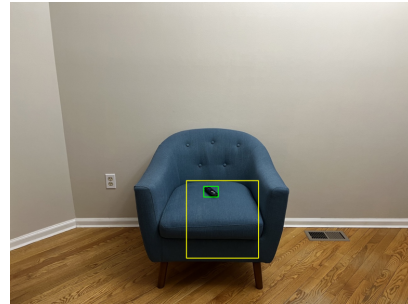
(c) **tape** under armchair



(d) **sunglasses** on table



(e) **spatula** on chair



(f) **remote** on armchair



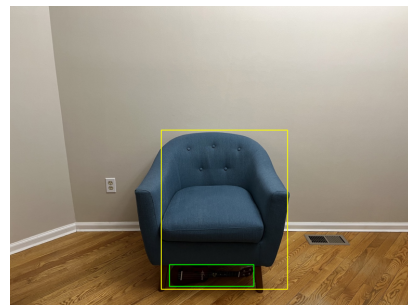
(g) **orange** right of armchair



(h) **ball of yarn** left of table



(i) **wineglass** on table



(j) **banjo** under armchair

Figure 12: Sample failures in small objects grounding (i.e.,  $\text{IoU} < 0.5$ ), which refers to the SUB column results of Subset A in Table 2. The pseudo-ground-truth bounding box, which is the GroundingDINO output, is indicated in **green**, and the output of LLAVA-NEXT-QWEN-1.5-110B, which is the best-performing MLLM in our grounding/localization experiment, is demonstrated in **yellow**.