Scaffolding Coordinates to Promote Vision-Language Coordination in Large Multi-Modal Models

Anonymous ACL submission

Abstract

State-of-the-art Large Multi-Modal Models (LMMs) have demonstrated exceptional capabilities in vision-language tasks. Despite their advanced functionalities, the performances of LMMs are still limited in challenging scenarios that require complex reasoning with multiple levels of visual information. Existing prompting techniques for LMMs focus on either improving textual reasoning or leveraging tools for image preprocessing, lacking a simple and general visual prompting scheme to promote vision-language coordination in LMMs. In this work, we propose SCAFFOLD prompting 014 that scaffolds coordinates to promote visionlanguage coordination. Specifically, SCAF-016 FOLD overlays a dot matrix within the image as visual information anchors and leverages multi-dimensional coordinates as textual positional references. Extensive experiments on a wide range of challenging vision-language tasks demonstrate the superiority of SCAFFOLD over GPT-4V with the textual CoT prompting.

1 Introduction

004

017

034

040

Large Multi-Modal Models (LMMs) like GPT-4V (Achiam et al., 2023) and Gemini (Team et al., 2023) have demonstrated impressive zero-shot capabilities in processing diverse visual-language tasks. Leveraging the advanced reasoning ability of the language model component, early attempts have been made to deploy LMMs in realistic scenarios, such as autonomous driving (Wen et al., 2023) and anomaly detection (Cao et al., 2023).

However, current LMMs display limited performance when conducting complex reasoning over multiple levels of visual information (Yang et al., 2023b; Wu et al., 2023a; Wu and Xie, 2023). For example, in a spatial reasoning task (Liu et al., 2023a), an LMM needs to verify or generate the statement by elucidating the relationship between different sources of visual information, and aligning its internal workings with textual expressions.



Figure 1: The overall framework of SCAFFOLD. SCAF-FOLD overlays a dot matrix onto the input image with Cartesian coordinates labeled aside the dots. The coordinates are also briefed in the textual prompt, which steers the LMM to leverage the dots on the image as a scaffold and promotes vision-language coordination.

Challenges for LMMs arise in orchestrating precise visual perception with accurate language understanding and generation.

To enhance vision-language coordination, prior efforts for LMMs can be divided into two categories: instruction tuning and prompting. Instruction tuning uses high-quality image-text pairs of either general purposes (Dai et al., 2023; Liu et al., 2023c,b; Xu et al., 2023) or specialized domains (Zhang et al., 2023; Chen et al., 2023a, 2024) to facilitate additional training of LMMs for improved performance. Nevertheless, instructiontuning consumes heavy loads of computing resources and thus sacrifices flexibility in methodology. Prompting for LMMs, on the contrary, steers the improved functioning of LMMs in a nonparametric manner. While techniques for language model prompting like Chain-of-Thought (Wei et al., 2022) apply to LMMs as well, the elicited intermediate reasoning steps mainly take place with the condition of the textual prompts (Wu et al., 2023b). As textual prompting techniques are actively being exploited, few endeavors have been made into visual prompting, which steers the precise visual perception of LMMs for vision-language coordination.



Question: Choose the correct image for the caption. Caption: a big cat is next to a small dog. Options: (A) image I(left)(B) image 2(right)

GPT-4V + SCAFFOLD:

Textual Prompt for SCAFFOLD: Two images are provided, each overlaid with a grid of dots arranged in a matrix with dimensions 6 by 6. Each dot on this grid is assigned a unique set of three-dimensional coordinates labeled as (t, x, y). ... 1. When you mention any key objects in the image, first output their nearest coordinates then identify them. 2. You can use the coordinates to determine the spatial relationships of the objects. ...

Output: In the first image (1), at coordinates (1,4,2) and (1,5,4), there is a small dog wearing a red and white checkered outfit and a large orange and white cat sitting next to each other on the grass. This matches the caption "a big cat is next to a small dog.". In the second image (2), at coordinates (2,3,3) and (2,5,5), there is a large black dog and a small kitten facing each other on the grass. Therefore, the correct image for the caption is: A.

GPT-4V Output: The correct image is: (B) image 2(right). This image shows a large black dog next to a much smaller kitten on a grassy surface. The size difference between the two animals is evident, with the dog being the "big cat" and the kitten being the "small dog" in a metaphorical sense, as the caption likely plays on the term "big cat" which usually refers to large wild felines. The first image (A) shows an actual cat and a small dog, which does not fit the caption.

Figure 2: An example of SCAFFOLD in double-image vision-language tasks, with correct textual reasoning illustrated in green and incorrect textual reasoning illustrated in red. Note that only the original images and questions are visible to **GPT-4V**; **SCAFFOLD** adds the coordinates on images and the corresponding textual prompt guidance.

The challenge of visual prompting for LMMs 067 lies in the mismatch of semantic granularity be-068 tween visual and textual information. While each word is explicitly separated in a textual sentence, different identities in an image are not isolated with clear boundaries. Recent works on visual prompting include leveraging tools to narrow the semantic granularity gap between visual and textual inputs. 074 Yang et al. (2023a) leverage advanced image segmentation models (Kirillov et al., 2023) to construct object segmentation overlays on the input image. Mitra et al. (2023) treat the LMM itself as a scene graph extractor to generate visual information in the textual format. However, tool usage inevitably results in additional resource burdens and potentially erroneous information. An alternative avenue of recent efforts is visual search (Wu and Xie, 2023; Nasiriany et al., 2024), where the so-084 lution to a complex visual task is cast as an iterative search process that accounts for multiple aspects 086 of the image. Nevertheless, the iterative queries of LMMs throughout the visual search process entail considerable expenses, limiting the practical value.

Therefore, it remains elusive whether a simple and general visual prompting scheme exists to promote vision-language coordination in LMMs.

091

093

097

098

100

101

102

103

104

106

108

109

110

111

112

In this work, we present SCAFFOLD, a simple and versatile visual prompting scheme to promote the coordination between vision and language in LMMs. SCAFFOLD overlays a dot matrix onto the input image, and each dot is labeled with its multi-dimensional Cartesian coordinate. The dot matrix on the image forms the scaffold that indicates relative visual positions for LMMs. The overlaid coordinates are also included in the textual prompt, which explicitly strengthens the connection between visual and textual information for LMMs. The LMMs are thus steered to leverage the coordinates to solve different vision-language tasks. In this way, SCAFFOLD provides a scaffold to promote vision-language coordination in LMMs. Extensive experiments on spatial reasoning, compositional reasoning, fine-grained grounding, and hallucination benchmarks demonstrate the superiority of SCAFFOLD over GPT-4V with the textual CoT prompting. We also show that the performance of

210

211

163

164

SCAFFOLD can be further enhanced with region 113 cropping, which reveals the promising future of 114 active perception enabled by SCAFFOLD. 115

Related Work 2

116

117

122

124

127

Large Multi-Modal Models (LMMs). State-ofthe-art LMMs like GPT-4V (Achiam et al., 2023) 118 and Gemini (Team et al., 2023) have excelled in 119 general vision-language tasks (Wu et al., 2023a; Yang et al., 2023b; Fu et al., 2023b). The integra-121 tion of visual capabilities in LMMs with advanced language proficiency and instruction-following 123 skills pave the way for versatile visual interactive agents, both in digital (He et al., 2024; Zheng et al., 125 2024) and embodied environments (Wake et al., 126 2023; Chen et al., 2023b).

GPT-4V Evaluation. As a leading LMM, GPT-128 4V (Achiam et al., 2023) has significantly expanded 129 the boundaries of LMM capabilities, motivating re-130 searchers to systematically explore its strengths 131 and weaknesses (Wu et al., 2023a; Yang et al., 132 2023b). Despite its proficiency, researchers have 133 proposed challenging benchmarks that reveal large 134 performance gap between GPT-4V and humans, in-135 cluding MMVP (Tong et al., 2024), MMMU (Yue 136 et al., 2023), Mementos (Wang et al., 2024), V* Bench (Wu and Xie, 2023), Contextual (Wad-138 hawan et al., 2024), etc. Extensive evaluations 139 indicate that plenty of room exists for state-of-the-140 141 art LMMs to improve their certain visual-language capabilities. Additionally, previous works (Yan 142 et al., 2023; Liu et al., 2023d; Cao et al., 2023; 143 Wen et al., 2023) have shown a promising future of potential applications of GPT-4V in downstream 145 146 domains like medicine science.

Multi-Modal Prompting Methods. Prompting 147 methods focus on unlocking model potentials by 148 carefully constructing model inputs. Chain-of-149 Thought (Wei et al., 2022; Kojima et al., 2022, 150 CoT) and its variants (Yao et al., 2023; Besta et al., 151 2023) have successfully elicited reasoning capa-152 bilities in language models. However, in multi-153 modal contexts such as compositional reasoning, 154 the original CoT is less effective (Mitra et al., 2023). 155 Consequently, numerous multi-modal prompting 156 methods have been developed for specific visual ca-157 pabilities. For instance, Compositional CoT (Mitra 158 et al., 2023) for compositional reasoning, Spatial 159 CoT (Chen et al., 2024) for spatial understanding, 160 Set-of-Marks prompting (Yang et al., 2023a) for 161 visual grounding. However, these methods tend 162

to be tailored for specific capabilities, calling for simple and general visual prompting schemes.

3 Methodology

In this section, we introduce SCAFFOLD prompting for vision-language coordination in LMMs.

3.1 **Visual Perspective of SCAFFOLD: Dot Matrices and Coordinates**

Visually, we enhance each input image with a uniformly distributed rectangular dot matrix, where each dot is labeled with multi-dimensional coordinates. These dots serve as visual positional anchors, while their coordinates are utilized as textual references in textual responses.

Visual Anchor Implementation. We select rectangular dot matrices as our visual anchor due to their simplicity, flexibility in textual description, and potential adaptability to image sequences. Unlike grids, which divide images into separate regions and may disrupt continuous visual content, dot matrices offer a less intrusive overlay. Additionally, to preserve original information, we deliver both the original image and the coordinates-overlaid image as inputs in single-image tasks.

Coordinates Implementation. For our approach, we use multi-dimensional Cartesian coordinates due to its simplicity and clarity. For a single image with an overlaid dot matrix of size $h \times w$, we assign two-dimensional coordinates (x, y) to each dot, representing its relative visual position. Here, the xcoordinate ascends from 1 to h within each column, while the y-coordinate ascends from 1 to w within each row. For image sequences, we extend these coordinates to three-dimensional (t, x, y). The tcoordinate remains constant within each image but increases sequentially across the sequence, allowing for differentiation between images and enhancing temporal perception.

In comparison, we also consider other coordinate options and identify their limitations. Absolute pixel coordinates, for example, consume excessive space and are complicated to perceive and apply accurately. Furthermore, one-dimensional Cartesian coordinates and alphabetic coordinates fall short in providing adequate positional information.

Other Factors. 1. Matrix Size. The matrix should be visually clear and provide ample space for multidimensional coordinates. For simplicity, we employ a 6×6 matrix for general vision-language tasks. 2. Matrix Density. Without prior visual

knowledge, we choose rectangular dot matrices 212 with a uniform density for general vision-language 213 tasks, providing LMMs equal assistance when rea-214 soning across different regions. 3. Matrix Color. 215 The coordinates are designed to be recognizable by 216 LMMs using their OCR capabilities. Consequently, 217 we color each dot in either black or white according 218 to its contrast against the background. 219

3.2 Textual Perspective of SCAFFOLD: Description and Guidelines

221

223

227

228

229

240

241

242

243

245

246

247

248

249

251

257

To complement the coordinates-overlaid visual inputs, we prepend textual guidance to task instructions for LMMs. This includes a brief description of the dot matrices and coordinates, accompanied by several general guidelines for their effective use, as detailed in Appendix A.1.1. The characteristics of these descriptions and guidelines are as follows: (1) Conciseness: The textual guidance is deliberately brief and clear, ensuring easy comprehension. (2) Generality: Designed to be universally applicable, these guidelines are not specific to any particular scenario, making them suitable for a wide range of vision-language tasks. (3) Extensibility: The guidelines are semantically independent, allowing for the seamless addition of more tailored instructions based on different scenarios. (4) Compositionality: The prepended texts can be easily combined with other prompting methods, such as zero-shot or compositional CoT (Kojima et al., 2022; Mitra et al., 2023).

4 Experiments

To demonstrate the effectiveness of SCAFFOLD, we conduct extensive experiments on top of GPT-4V on a range of challenging vision-language tasks, including *Spatial Reasoning*, *Compositional Reasoning*, *Fine-Grained Visual Understanding* and *Hallucination*. Specifically, we perform systematic evaluations on 11 benchmarks, with detailed information presented in Appendix A.1.2. We set the temperature of GPT-4V as zero in our experiments.

4.1 Benchmarks

This subsection briefly introduces the benchmarks used for evaluation. Due to the limited budget, for some datasets, we sample a subset for experiments.
Spatial Reasoning evaluates LMM capability to infer spatial relationships between objects. Selected benchmarks are as follows. 1. MME (Position split) (Fu et al., 2023a) is a subset of the MME

comprehensive evaluation suite for LMMs to infer object positions. 2. *Visual Spatial Reasoning* (*VSR*) (Liu et al., 2023a) challenges LMMs with 66 types of spatial relations to verify spatial propositions. 3. *EgoThink (Spatial split)* (Cheng et al., 2023) tests the spatial reasoning ability of LMMs from a first-person perspective.

Compositional Reasoning requires LMMs to identify object attributes and their interrelations. Selected benchmarks are as follows. *1. Winoground (Thrush et al., 2022)* is a challenging benchmark that necessitates compositional reasoning of LMMs to match images with captions, reformulated as binary choice questions for our evaluation. *2. WHOOPS! VQA* (Bitton-Guetta et al., 2023) involves compositional reasoning over commonsense-defying images. *3. CLEVR* (Johnson et al., 2017) is designed for assessing compositional reasoning in program-generated scenes.

Fine-Grained Visual Understanding requires LMMs to perform visual search and precisely perceive fine-grained visual details. Selected benchmarks are as follows. *1. V* Bench* (Wu and Xie, 2023) requires LMMs to identify and reason with fine-grained visual details in high-resolution images. *2. Spotting Differences*¹ is our newly-collected dataset challenging LMMs to find and pinpoint differences between two similar images, with further details in Appendix A.2.3.

Hallucination measures the tendency of LMMs to generate hallucinatory or illusory perceptions. Selected benchmarks are as follows. *1. POPE (Adversarial Subset) (Li et al., 2023)* assesses object hallucination by querying the existence of specific objects. *2. HallusionBench (Guan et al., 2023)* consists of meticulously crafted images to measure hallucination and visual illusion in LMMs. *3. Mementos (Wang et al., 2024)* evaluates LMM to conduct precise reasoning over image sequences and measures their performances in terms of object and behavior hallucinations.

4.2 Baselines

This section presents the prompting methods utilized as baselines in our experiments, including **1. Naive Prompting** utilizes original images and user instructions as inputs for LMMs, establishing a straightforward baseline without any prompt optimization. **2. CoT** (Wei et al., 2022) guides LMMs to perform step-by-step reasoning before outputting

¹https://www.crazygames.com/game/find-the-difference

Crucial Capability	Dataset Size		Metric	Naive	СоТ	SCAFFOLD (Ours)
Spatial Reasoning	MME (Position) VSR EgoThink (Spatial)		Accuracy Accuracy LLM as Judge	51.7 67.8 66.0	51.7 70.4 74.0	75.0 (+23.3) 74.4 (+6.6) 76.0 (+10.0)
Compositional Reasoning	positional asoning Winoground 100 Group Score 17. WHOOPS! VQA 200 BEM Score 58. CLEVR 200 LLM as Judge 43.		17.0 58.6 43.5	33.0 57.6 43.0	33.0 (+16.0) 62.7 (+4.1) 48.0 (+4.5)	
Fine-Grained Visual Understanding	V* Bench Spotting Differences	238 100	Accuracy Accuracy	27.2 13.0	30.8 14.0	44.6 (+17.4) 19.0 (+6.0)
Hallucination	POPE (Adversarial) HallusionBench (Hard) Mementos	100 504 100	Accuracy Accuracy LLM as Judge	79.0 45.6 33.3	80.0 48.8 33.5	86.0 (+7.0) 53.0 (+7.4) 36.1 (+2.8)
Overall	All	1852	Average	45.7	48.8	55.3 (+9.6)

Table 1: Results of SCAFFOLD on 11 challenging vision-language benchmarks, with the highest score **bold**.

the final answer. The prompt text "*Let's think step by step*." is prepended to task descriptions.

4.3 Results and Analyses

310

312

313

317

As presented in Table 1, the results demonstrate that SCAFFOLD significantly enhances the visual capabilities of LMMs, surpassing CoT (Wei et al., 2022) in most evaluated benchmarks. With naive prompting and CoT prompting averaging 45.7 and 48.8 respectively, SCAFFOLD successfully obtains an overall improvement of 9.6.

319Spatial Reasoning: SCAFFOLD notably enhances320the spatial reasoning capabilities of LLMs across321the three benchmarks, with an average improve-322ment of 13.3. Fig. 3 illustrates how SCAFFOLD323enabled GPT-4V adeptly identifies crucial objects324and records accurate positional information using325two-dimensional coordinates, leading to the correct326assessment of spatial relations through the numeri-327cal analysis of their x-coordinates.

Compositional Reasoning: With SCAFFOLD, GPT-4V demonstrates improved abilities in compositional reasoning with an average improvement of 8.2, showing enhanced perception of key visual elements and smoother reasoning processes. As Fig. 9 in Appendix A.3 shows, SCAFFOLD associates crucial objects with their textual positions, assisting GPT-4V to accurately identify and localize significant visual details.

Fine-Grained Visual Understanding: On benchmarks such as V* Bench (Wu and Xie, 2023) and Spotting Differences, SCAFFOLD has markedly improved the capabilities of LMMs in conducting visual search and localizing fine-grained visual details, with an average improvement of 11.7. Taking Fig. 4 as an example, with the support of coordinates, GPT-4V effectively pinpoints and documents



Question: Determine whether the following statement is true or false: The bench is at the left side of the potted plant.

Ground-Truth: True. **GPT-4V + CoT:** The statement is false. In the image provided, the bench is on the right side of the potted plant, not the left. Therefore, the statement is false. **GPT-4V + SCAFFOLD:** The bench is located near coordinates (3,2) through (6,4), and the potted plant is near coordinates (3,6) through (6,6). Given the arrangement of the coordinates, the bench is to the left of the potted plant. true.

Figure 3: An example from the Visual Spatial Reasoning (VSR) dataset. SCAFFOLD facilitates better visionlanguage coordination than other methods.

the precise locations of target visual elements, contributing to its accurate perception of the target attribute. Additionally, we notice that without coordinates, GPT-4V is more easily to give up and apologize for its search failure.

Hallucination: With an average improvement of 5.7, utilizing coordinates as a scaffold enables GPT-4V to recognize objects within a scene and further accurately describe their positions, guiding its textual reasoning to focus on precise visual informa-

353



Figure 4: An example from the V* Bench dataset.

tion. With an example shown in Fig. 11 in Appendix A.3, GPT-4V with coordinates is capable of precisely capturing visual details and preventing hallucinating non-existent objects, promoting accurate visual grounding.

5 Ablation Study

355

361

364

368

371

373

We conduct extensive ablation studies on key factors such as matrix size and coordinate color to validate and further explore SCAFFOLD.

5.1 Experimental Setup

Due to limited GPT-4V access quota, we each sample 50 questions from Visual Spatial Reasoning (Liu et al., 2023a), Winoground (Thrush et al., 2022), and POPE (Adversarial Subset) (Li et al., 2023), creating an ablation subset of 150 samples. Overall accuracy per question is adopted as the metric and GPT-4V is used for our experiments. Additionally, for stable results, we run each experiment twice and report average accuracy.

5.2 Effect of Matrix Size

The matrix size h and w may influence the precision of textual reference and the granularity of visual information. Consequently, we incorporate



Figure 5: Impact of matrix sizes. Better accuracies are illustrated in darker green.

Coloring Strategy	WNG	VSR	POPE	Overall
None (Baseline)	70.0	64.0	73.0	69.0
White	78.0	72.0	78.0	76.0
Black	79.0	73.0	82.0	78.0
Complementary	77.0	71.0	81.0	76.3
Binary (Ours)	81.0	73.0	84.0	79.3

Table 2: Results of different coloring strategies in ablation experiments, where *WNG* denotes the Winoground (Thrush et al., 2022) dataset and *POPE* denotes the POPE (Li et al., 2023) Adversarial subset.

matrices of difference sizes varying from 3×3 to 7×7 and measure their performances.

378

379

380

381

383

384

386

388

390

391

392

393

394

396

397

398

400

Fig. 5 depicts the performance variations with different matrix sizes, suggesting 6×6 as the optimal size for our ablation dataset. Additionally, the sizes in the upper right section tend to perform better than those in the lower left section. It may be due to the sampled images usually possessing equal or larger widths than heights, suggesting matrix sizes may ideally align with image sizes.

Additionally, the 6×6 size did not yield the best results across all three subsets, hinting at potential improvements by customizing matrix sizes for specific tasks. It may be beneficial to automatically and dynamically adjust the matrix size and we leave this open problem to future research.

5.3 Effect of Matrix Color

In terms of matrix color, we design different coloring strategies and compare their performances. As illustrated in Fig. 6, uniform coloring strategies adopt the same color for various scenes, occasionally blending into the surroundings. Complementary colors introduce large amounts of colors and



Figure 6: Examples of different color configurations of dot matrices and coordinates.

Coordinates	WNG	VSR	POPE	Overall
None (Baseline)	70.0	64.0	73.0	69.0
Alphabet	80.0	72.0	79.0	77.0
Pixel	78.0	75.0	77.0	76.7
One-Dimensional	72.0	71.0	81.0	74.7
Cartesian (Ours)	81.0	73.0	84.0	79.3

Table 3: Results of different coordinates designs in ablation experiments, where *WNG* denotes the Winoground (Thrush et al., 2022) dataset and *POPE* denotes the POPE (Li et al., 2023) Adversarial subset.

may mislead model attention. Consequently, for simplicity and visibility, we choose the most contrasting color from black and white at each dot location. To assess the efficacy of our approach, we compared it against baseline coloring strategies, including *uniform black*, *uniform white*, and *complementary coloring*. As shown in Table 2, our *binary coloring* strategy slightly surpasses the alternatives in performance.

5.4 Effect of Coordinate Format

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

Coordinates, as textual references for dots, are vital for aligning visual inputs with textual outputs. To assess the effectiveness of our implementation, we experiment with various coordinate formats, including alphabetic, one-dimensional numerical, and pixel absolute coordinates. The examples of these formats are exhibited in Fig. 12 of Appendix A.3.

The results, detailed in Table 3, reveal that our approach surpasses alternative coordinate designs in performance. Additionally, all the coordinates designs perform better than the baseline without coordinates, indicating the flexibility of coordinates design and the adaptability to different scenarios.

5.5 Effect of Dot Perturbations

To assess the resilience of SCAFFOLD, we introduce Gaussian noise to the dots, slightly adjusting their positions without significantly changing their relative placements, as illustrated in Fig. 7. We model the original dot positions as (X, Y), with l_h



Figure 7: The results of perturbed and standard coordinates (left) and a perturbed coordinatesoverlaid example (right), where *WNG* denotes the Winoground (Thrush et al., 2022) dataset and *POPE* denotes the POPE (Li et al., 2023) Adversarial subset.

and l_w representing the distances between neighboring dots along the x and y axes, respectively. The perturbed coordinate $(X_{\text{new}}, Y_{\text{new}})$ reads:

$$\begin{bmatrix} \mathbf{X}_{\text{new}} \\ \mathbf{Y}_{\text{new}} \end{bmatrix} = \begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} + \begin{bmatrix} \mathcal{N}\left(0, \left(\frac{1}{4} \cdot l_{h}\right)^{2}\right) \\ \mathcal{N}\left(0, \left(\frac{1}{4} \cdot l_{w}\right)^{2}\right) \end{bmatrix}$$
(1)

· -

According to the findings depicted in Fig. 7, the coordinates perturbed with noise not only retain the enhancements provided by standard coordinates but also, in the VSR subset, exceed their performance. This indicates the substantial robustness against perturbation of the coordinates and suggests the potential for further optimizing their placement.

6 Integration with Other Methods

This section describes integration experiments of SCAFFOLD combined with active perception and Chain-of-Thought (Wei et al., 2022).

6.1 SCAFFOLD + Active Perception

7

In complex visual environments, humans would proactively engage with their surroundings to enhance scene understanding, like zooming in or changing perspectives. Similarly, we recognize that LMMs should possess such capabilities in realistic 433

434

435

436

430

- 437 438
- 439 440
- 441
- 442
- 443 444
- 445 446

447

448

449



Figure 8: The procedure of combined SCAFFOLD and active perception techniques on V* Bench.

Method	$ $ N.F.R. \downarrow (%)	S.R. ↑ (%)
Naive	72.2	21.7
CoT	71.3	21.7
SCAFFOLD	26.2	31.3
Scaffold + A.P.	14.8	45.2

Table 4: Results of SCAFFOLD + active perception on V* Bench (Wu and Xie, 2023) direct_attributes subset, where A.P. denotes active perception, N.F.R. denotes Not Found Rate and S.R. denotes Success Rate.

scenarios and propose that SCAFFOLD can function as a scaffold for effective active perception.

To validate this, we integrate SCAFFOLD with active perception in the experiments on the direct_attributes subset of V* Bench (Wu and Xie, 2023), which requires LMMs to perceive finegrained details in high-resolution images. This challenge encompasses both the localization of target objects and the identification of their attributes under resolution constraints. Consequently, we adopt two metrics to measure LMM performance, including *Not Found Rate* representing the percentage of invalid responses, and *Success Rate* representing the percentage of correct responses.

As depicted in Fig. 8, our combined method unfolds in two phases: initial visual search to locate the target details, followed by cropping the image around the pinpointed coordinates to closely examine and identify the target attributes.

The results, presented in Table 4, reveal a performance enhancement of 14.1% compared with SCAFFOLD alone, underscoring the utility of coordinates in facilitating active perception. Furthermore, the results reveal two notable performance leaps. The initial improvement (CoT \rightarrow SCAF-FOLD) is attributed to the use of coordinates, sig-

Dataset	Naive	СоТ	SCAFFOLD	SCAFFOLD + CoT
Wino.	17.0	33.0	33.0	41.0
V*	27.2	30.8	44.6	47.9

Table 5: Results of SCAFFOLD combined with Chainof-Thought (Wei et al., 2022) on V* Bench (Wu and Xie, 2023) and Winoground (Thrush et al., 2022).

nificantly reducing the *Not Found Rate* and thereby aiding in the visual search process. The subsequent gain (SCAFFOLD \rightarrow SCAFFOLD + A.P.) results from the combined implementation of active perception, which enables LMMs to accurately discern target attributes within the cropped regions. 477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

502

6.2 SCAFFOLD + Chain-of-Thought

Our prompting method, characterized by its simplicity, can seamlessly integrate with zero-shot CoT by appending *Let's think step by step.* to user instructions. To test its effectiveness, we conduct experiments on Winoground (Thrush et al., 2022) and V* Bench (Wu and Xie, 2023). Results from Table 5 demonstrate that combining our method with CoT enhances LMM performance beyond what either method achieves independently. These findings underscore our method's substantial compatibility and potential for performance improvement when combined with other methods.

7 Conclusion

In this work, we propose SCAFFOLD, a simple and general visual prompting method that utilizes scaffolding coordinates to promote vision-language coordination in LMMs. Extensive experiments show that SCAFFOLD successfully elicits LMM capabilities in several challenging vision-language tasks.

476

451

609

554

555

Limitations

503

Here we discuss two limitations of this work.

(1) To automatically adjust dot matrix attributes. 505 In this work, for simplicity and clarity, we adopt 506 matrices of size 6×6 in our implementation. How-507 ever, our ablation study in Section 5 suggests that a one-size-fits-all matrix size can yield good, but not 509 the best results across all datasets. Given the diver-510 sity of visual tasks and the varying granularity of 511 information in different scenes, it stands to reason 512 that tailoring the matrix attributes, such as size and 513 coordinates format, to the specific requirements of 514 each task or even each sample could improve per-515 formance. Addressing the dynamic and automatic 516 adjustment of these attributes to better suit different 517 scenarios remains an area for future exploration. 518

(2) To enhance precision in visual localization. 519 By integrating dot matrices with coordinates onto 520 images, we aimed to facilitate improved visionlanguage coordination by associating key objects 522 with their closest coordinates. However, our ob-523 servations indicate that, particularly in complex or 524 clustered scenes, GPT-4V occasionally struggles to accurately associate textual reasoning with the nearest coordinates. This challenge underscores the need for LMMs to achieve improved visual lo-528 calization and grounding capabilities in intricate 530 environments. With SCAFFOLD, we expect the future of LMMs and visual prompting techniques to 531 be further improved in terms of visual localization. 532

References

534

535

537

539

541

542

543

544

545

546

547

548

549

552

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Michal Podstawski, Hubert Niewiadomski, Piotr Nyczyk, et al. 2023. Graph of thoughts: Solving elaborate problems with large language models. arXiv preprint arXiv:2308.09687.
- Nitzan Bitton-Guetta, Yonatan Bitton, Jack Hessel, Ludwig Schmidt, Yuval Elovici, Gabriel Stanovsky, and Roy Schwartz. 2023. Breaking common sense: Whoops! a vision-and-language benchmark of synthetic and compositional images. In *Proceedings* of the IEEE/CVF International Conference on Computer Vision, pages 2616–2627.
- Yunkang Cao, Xiaohao Xu, Chen Sun, Xiaonan Huang, and Weiming Shen. 2023. Towards generic anomaly

detection and understanding: Large-scale visuallinguistic model (gpt-4v) takes the lead. *arXiv preprint arXiv:2311.02782*.

- Boyuan Chen, Zhuo Xu, Sean Kirmani, Brian Ichter, Danny Driess, Pete Florence, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. 2024. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. *arXiv preprint arXiv:2401.12168*.
- Chi Chen, Ruoyu Qin, Fuwen Luo, Xiaoyue Mi, Peng Li, Maosong Sun, and Yang Liu. 2023a. Position-enhanced visual instruction tuning for multimodal large language models. *arXiv preprint arXiv:2308.13437*.
- Liang Chen, Yichi Zhang, Shuhuai Ren, Haozhe Zhao, Zefan Cai, Yuchi Wang, Peiyi Wang, Tianyu Liu, and Baobao Chang. 2023b. Towards end-to-end embodied decision making via multi-modal large language model: Explorations with gpt4-vision and beyond. *arXiv preprint arXiv:2310.02071*.
- Sijie Cheng, Zhicheng Guo, Jingwen Wu, Kechen Fang, Peng Li, Huaping Liu, and Yang Liu. 2023. Can vision-language models think from a first-person perspective? *arXiv preprint arXiv:2311.15596*.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. 2023. Instructblip: Towards general-purpose vision-language models with instruction tuning. *arXiv preprint arXiv:2305.06500*.
- Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, et al. 2023a. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*.
- Chaoyou Fu, Renrui Zhang, Haojia Lin, Zihan Wang, Timin Gao, Yongdong Luo, Yubo Huang, Zhengye Zhang, Longtian Qiu, Gaoxiang Ye, et al. 2023b. A challenger to gpt-4v? early explorations of gemini in visual expertise. *arXiv preprint arXiv:2312.12436*.
- Tianrui Guan, Fuxiao Liu, Xiyang Wu Ruiqi Xian Zongxia Li, Xiaoyu Liu Xijun Wang, Lichang Chen Furong Huang Yaser Yacoob, and Dinesh Manocha Tianyi Zhou. 2023. Hallusionbench: An advanced diagnostic suite for entangled language hallucination & visual illusion in large vision-language models. *arXiv preprint arXiv:2310.14566*.
- Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu, Yong Dai, Hongming Zhang, Zhenzhong Lan, and Dong Yu. 2024. Webvoyager: Building an end-toend web agent with large multimodal models. *arXiv preprint arXiv:2401.13919*.
- Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. 2017. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE conference*

719

720

- 610 611
- 612 613 614
- 615 616
- 617 618 619
- 621 622 623 624
- 626 627 628
- 631 632 633
- 635 636 637 638 639 640 641

642

- 643 644 645 646 647 648 649
- 651 652
- 6
- 655 656 657

- 659 660 661
- 661 662
- 663 664

- on computer vision and pattern recognition, pages 2901–2910.
- 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. *arXiv preprint arXiv:2304.02643*.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199– 22213.
 - Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*.
 - Fangyu Liu, Guy Emerson, and Nigel Collier. 2023a. Visual spatial reasoning. *Transactions of the Association for Computational Linguistics*, 11:635–651.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023b. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023c. Visual instruction tuning. In *NeurIPS*.
- Zhengliang Liu, Hanqi Jiang, Tianyang Zhong, Zihao Wu, Chong Ma, Yiwei Li, Xiaowei Yu, Yutong Zhang, Yi Pan, Peng Shu, et al. 2023d. Holistic evaluation of gpt-4v for biomedical imaging. *arXiv preprint arXiv*:2312.05256.
- Chancharik Mitra, Brandon Huang, Trevor Darrell, and Roei Herzig. 2023. Compositional chain-of-Thought prompting for large multimodal models. *arXiv preprint arXiv:2311.17076*.
- Soroush Nasiriany, Fei Xia, Wenhao Yu, Ted Xiao, Jacky Liang, Ishita Dasgupta, Annie Xie, Danny Driess, Ayzaan Wahid, Zhuo Xu, Quan Vuong, Tingnan Zhang, Tsang-Wei Edward Lee, Kuang-Huei Lee, Peng Xu, Sean Kirmani, Yuke Zhu, Andy Zeng, Karol Hausman, Nicolas Heess, Chelsea Finn, Sergey Levine, and Brian Ichter. 2024. Pivot: Iterative visual prompting elicits actionable knowledge for vlms. *arXiv preprint arXiv:2402.07872*.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.
- Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Candace Ross. 2022. Winoground: Probing vision and language models for visio-linguistic compositionality. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 5238– 5248.

- Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. 2024. Eyes wide shut? exploring the visual shortcomings of multimodal llms. *arXiv preprint arXiv:2401.06209*.
- Rohan Wadhawan, Hritik Bansal, Kai-Wei Chang, and Nanyun Peng. 2024. Contextual: Evaluating contextsensitive text-rich visual reasoning in large multimodal models. *arXiv preprint arXiv:2401.13311*.
- Naoki Wake, Atsushi Kanehira, Kazuhiro Sasabuchi, Jun Takamatsu, and Katsushi Ikeuchi. 2023. Gpt-4v (ision) for robotics: Multimodal task planning from human demonstration. *arXiv preprint arXiv:2311.12015*.
- Xiyao Wang, Yuhang Zhou, Xiaoyu Liu, Hongjin Lu, Yuancheng Xu, Feihong He, Jaehong Yoon, Taixi Lu, Gedas Bertasius, Mohit Bansal, et al. 2024. Mementos: A comprehensive benchmark for multimodal large language model reasoning over image sequences. *arXiv preprint arXiv:2401.10529*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-Thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Licheng Wen, Xuemeng Yang, Daocheng Fu, Xiaofeng Wang, Pinlong Cai, Xin Li, Tao Ma, Yingxuan Li, Linran Xu, Dengke Shang, et al. 2023. On the road with gpt-4v (ision): Early explorations of visual-language model on autonomous driving. *arXiv* preprint arXiv:2311.05332.
- Penghao Wu and Saining Xie. 2023. V*: Guided visual search as a core mechanism in multimodal llms. *arXiv preprint arXiv:2312.14135*, 17.
- Yang Wu, Shilong Wang, Hao Yang, Tian Zheng, Hongbo Zhang, Yanyan Zhao, and Bing Qin. 2023a. An early evaluation of gpt-4v (ision). *arXiv preprint arXiv:2310.16534*.
- Yifan Wu, Pengchuan Zhang, Wenhan Xiong, Barlas Oguz, James C Gee, and Yixin Nie. 2023b. The role of chain-of-Thought in complex vision-language reasoning task. *arXiv preprint arXiv:2311.09193*.
- Zhiyang Xu, Ying Shen, and Lifu Huang. 2023. Multi-Instruct: Improving multi-modal zero-shot learning via instruction tuning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11445– 11465, Toronto, Canada. Association for Computational Linguistics.
- Zhiling Yan, Kai Zhang, Rong Zhou, Lifang He, Xiang Li, and Lichao Sun. 2023. Multimodal chatgpt for medical applications: an experimental study of gpt-4v. *arXiv preprint arXiv:2310.19061*.
- Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chunyuan Li, and Jianfeng Gao. 2023a. Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v. *arXiv preprint arXiv:2310.11441*.

- 721 722
- 724
- 726
- 727
- 731 732
- 733 734
- 735

- 740 741
- 742 743
- 744
- 745
- 746 747

749 750

751

754

752

- Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. 2023b. The dawn of lmms: Preliminary explorations with gpt-4v (ision). arXiv preprint arXiv:2309.17421, 9(1):1.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. arXiv preprint arXiv:2305.10601.
- 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. 2023. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. arXiv *preprint arXiv:2311.16502.*
- Shilong Zhang, Peize Sun, Shoufa Chen, Min Xiao, Wenqi Shao, Wenwei Zhang, Kai Chen, and Ping Luo. 2023. Gpt4roi: Instruction tuning large language model on region-of-interest. arXiv preprint arXiv:2307.03601.
 - Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. 2024. Gpt-4v (ision) is a generalist web agent, if grounded. arXiv preprint arXiv:2401.01614.
 - Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. arXiv preprint arXiv:2306.05685.

Appendix А

- A.1 Prompts
- This section exhibits the prompts used in SCAF-755 FOLD implementation and our experiments.
 - **SCAFFOLD** implementation A.1.1

758 In SCAFFOLD implementation, we use the following textual guidelines to describe the effective use of coordinates.

Single-Image Setting. In the single-image setting, 761 we label all the dots with two-dimensional coordi-762 nates and deliver both the original image and the coordinates-overlaid image to the model. Consequently, we use the following guidelines.

I will provide you with two images of the same scene. The second image is overlaid with a dot matrix of the shape of HEIGHT * WIDTH to help you with your task, and each dot is labeled with two-dimensional coordinates (x,y). 1. When you mention any key objects in the image, first output their nearest coordinates then identify them. 2. You use the coordinates to determine the spatial relationships of the objects. Within each column, the x-coordinate increases from top to bottom, and within each row, the y-coordinate increases from left to right. 3. You can search and reason region by region with the help of the dots.

4. Finally, conclude your answer in format [[AN-SWER]], such as [[A]], [[B]], [[C]] or [[D]].

Note that the fourth guideline serves as a constraint for specific output formats and may vary among different tasks.

Double-Images Setting. In a double-images setting, we label all the dots with three-dimensional coordinates, with the first coordinate serving to distinguish between two images. Consequently, we use the following guidelines.

Two images are provided, each overlaid with a grid of dots arranged in a matrix with dimensions h by w. Each dot on this grid is assigned a unique set of three-dimensional coordinates labeled as (t, x, y). The first coordinate, 't', serves to distinguish between the two images: '1' is assigned to the first image on the left, and '2' to the second image on the right. The other two coordinates, 'x' and 'y', are used to specify the dot's spatial location within its respective image. This labeling system is designed to assist you in identifying and referring to specific points within each image.

1. When you mention any key objects in the image, first output their nearest coordinates then identify them.

2. You use the coordinates to determine the spatial relationships of the objects. within each column, the x-coordinate increases from top to bottom, and Within each row, the y-coordinate increases from left to right.

3. You can search and reason region by region with the help of the dots.

4. Finally, you must conclude your answer in format [[ANSWER]], such as [[A]] or [[B]].

Note that the fourth guideline serves as a constraint for specific output formats and may vary among different tasks.

Image-Sequence Setting. In the image sequence setting, we label all the dots with three-dimensional coordinates. For simplicity and efficiency, we only deliver the coordinates-overlaid image to the model. Consequently, we use the following guidelines.

775 776

766

767

768

769

770

771

773

774

779 780

781

782

200

A sequence of images is provided, each overlaid with a grid of dots arranged in a matrix with dimensions HEIGHT by WIDTH. Each dot on this grid is assigned a unique set of three-dimensional coordinates labeled as (t, x, y). The first coordinate, 't', serves to distinguish between the images, for instance, '1' is assigned to the first image, and '8' to the last image. The other two coordinates, 'x' and 'y', are used to specify the dot's spatial location within its respective image. This labeling system is designed to assist you in identifying and referring to specific points within each image.

1. When you mention any key objects in the image, first output their nearest coordinates then identify them.

2. You use the coordinates to determine the temporal and spatial relationships of the objects. Within the image sequence, the t-coordinate increases as time grows; within each column, the x-coordinate increases from top to bottom; within each row, the y-coordinate increases from left to right.

3. You can search and reason region by region with the help of the dots.

4. you need to keep your descriptions concise and clear.

Note that the fourth guideline serves as a constraint for specific output formats and may vary among different tasks.

A.1.2 Datasets

Spotting Differences dataset is our newly formulated dataset challenging LMMs to precisely localize different details in two similar images. We exhibit the prompts as follows. Spotting Differences 2

[Question] Spot ten differences between the images and answer the dot position closest to every difference in the format: [spot index, x, y]. I am going to tip \$100 for a better answer. [Answer] Let's identify ten differences between the left and right images: 1. The bird in the sky is missing in the right image.

Closest dot: [1, 2, 1]...

A.1.3 Integration Experiments

We list the prompts used in the integration experiments.

SCAFFOLD + Active Perception Firstly, we guide the LMM to visually search the image and localize target objects using the following prompt.

Based on the question: question you should first identify key objects in the question and link them with their nearest coordinate, and finally conclude the coordinates in format [[(x,y)]] in the end(you don't need to answer the question).

Secondly, we crop the image based on the output coordinates and guide the LMM to answer the question based on the cropped images in the second turn of the conversation.

Here are the cropped images from the scene according to your selected coordinates, you can take a closer look and answer the question. If I don't provide cropped images or the target does not exist in the cropped image, please visually search the original image and answer the question. Question: QUESTION Options: OPTIONS

A.2 Benchmarks

This section details the benchmarks used in our experiments to evaluate our method. The benchmarks are divided into four categories and elaborated respectively.

A.2.1 Spatial Reasoning

Spatial reasoning is a crucial capability of LMMs to determine spatial relationships between objects in the image. To evaluate the effectiveness of our method to elicit spatial reasoning capabilities, we select several challenging datasets including MME (Fu et al., 2023a) Position split, Visual Spatial Reasoning (Liu et al., 2023a) dataset, Ego-Think (Cheng et al., 2023) Spatial split. Note that we only select a subset related to spatial reasoning in comprehensive evaluation benchmarks like MME (Fu et al., 2023a) and EgoThink (Cheng et al., 2023a) and EgoThink (Cheng et al., 2023a) because the other subsets don't constitute spatial reasoning and our GPT-4V access quota is limited. The details of the benchmarks are as follows.

MME (Position) (Fu et al., 2023a)

Dataset Introduction. The MME (Fu et al., 2023a) benchmark is a comprehensive evaluation suite designed to measure the perception and recognition capabilities of LMMs in 14 tasks. To evaluate spatial reasoning capabilities, we only select Position split from 14 subtasks. It contains 60 questions about the spatial relationship among the objects in the scene, such as *left, above*, etc..

Metric. All questions are formatted as general interrogative sentences, which can be answered with either "yes" or "no". Therefore, we guide the model

806

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

827

828

829

830

831

832

833

834

835

836

837

838

²https://www.crazygames.com/game/find-the-difference

brackets [[]], and then use accuracy as the metric to evaluate the model's performance. *Example.*



to output the final answer within double square

Visual Spatial Reasoning (VSR) (Liu et al., 2023a)

Dataset Introduction. The VSR (Liu et al., 2023a) dataset is designed to comprehensively evaluate LMM capabilities to perform spatial reasoning in images, containing 66 types of spatial descriptions in language. Each sample provides an image and a corresponding spatial description in terms of two individual objects in the scene. Because of the limited GPT-4V (Achiam et al., 2023) access quota, we randomly sample 200 samples from the dataset for evaluation.

Metric. The task is to determine the correctness of the given spatial description, which can be answered with either "true" or "false". Therefore, we guide the model to output the final answer within double square brackets [[]], and then use accuracy as the metric to evaluate the model's performance. *Example.*



EgoThink (Spatial) (Cheng et al., 2023) *Dataset Introduction.* The Egothink (Cheng et al., 2023) dataset is intended for first-perspective vision-language problem-solving capabilities. To evaluate spatial reasoning capabilities, we only select Spatial split from it. The Spatial split contains 50 questions about the spatial relationship among the objects, particularly requiring LMMs to perceive and reason from the first perspective.

Metric. The task is to answer spatial questions from the first perspective. We adopt the evaluation script in EgoThink official implementation ³, which uses GPT-4 (Achiam et al., 2023) as the judge model to score the answers.

Example.



A.2.2 Compositional Reasoning

Compositional reasoning represents the capability of LMMs to perceive and reason in terms of objects' attributes and their relationships, significant for visual perception. To evaluate the effectiveness of our method to elicit compositional reasoning capabilities, we select several challenging datasets including Winoground (Thrush et al., 2022), WHOOPS! (Bitton-Guetta et al., 2023) and CLEVR (Johnson et al., 2017).

Winoground (Thrush et al., 2022)

Dataset Introduction. Winoground (Thrush et al., 2022) proposes a novel dataset that challenges LMMs to correctly match two images and two captions. The captions use the same words but in a different order, requiring a precise understanding of both images and captions. The challenging dataset is a suitable benchmark for compositional reasoning.

Metric. Given two images and two captions, we compose a sample into four binary-choice questions for the effective evaluation of LMMs, including choosing the correct caption given two images and choosing the correct image given two captions respectively. Finally, we adopt *group score* to measure LMM performance: only when the model an-

³https://github.com/AdaCheng/EgoThink

911

912

913

914

915

916

917

918

919

920

923

swers all four questions correctly is it considered completely correct for this sample. *Example*.



Q1: [[Image1]] [[Image2]] Choose the correct image for the caption. Caption: an old person kisses a young person. Options: (A) image 1(left) (B) image 2(right) Q2: [[Image1]] [[Image2]] Choose the correct image for the caption. Caption: a young person kisses an old person. Options: (A) image 1(left) (B) image 2(right)

Q3: [[Image1]] Choose the correct caption for the image. (A) an old person kisses a young person. (B) a young person kisses an old person.

Q4: [[Image2]] Choose the correct caption for the image. (A) an old person kisses a young person. (B) a young person kisses an old person.

WHOOPS! (Bitton-Guetta et al., 2023)

Dataset Introduction. The WHOOPS! (Bitton-Guetta et al., 2023) is designed to challenge LMMs to perform compositional reasoning in terms of purposefully commonsense-defying images. It remains challenging for LMMs to recognize and interpret these unconventional images. Furthermore, several tasks are posed over the dataset, and we select the visual question-answering task.

Metric. We adopt the BEM score in WHOOPS! official implementation ⁴ to evaluate LMM performances.

Example.



CLEVR (Johnson et al., 2017)

Dataset Introduction. The CLEVR (Johnson et al., 2017) is designed as a standard evaluation suite for compositional vision-language reasoning with program-rendered images and correctly generated annotations. Containing objects like a metal cube or red sphere, the dataset poses questions in terms of object existence, object attributes, and object relationships. For effective evaluation, we randomly sampled 200 samples in the dataset.

Metric. We guide the LMMs to answer the questions in an open-ended generative manner. Due to the complexity of the question and the answer, we adopt GPT-4 (Achiam et al., 2023) as a judge model to determine the correctness of answers. Our judge prompt is as follows, adapted from MT-Bench (Zheng et al., 2023).

[Instruction] Please act as an impartial judge and evaluate the quality of the response provided by an AI assistant to the user question displayed below. Your evaluation should consider correctness and helpfulness. You will be given a reference answer and the assistant's answer. Begin your evaluation by comparing the assistant's answer with the reference answer. Identify and correct any mistakes. The assistant has access to an image along with questions but you will not be given images. Therefore, please consider only how the answer is close to the reference answer. If the assistant's answer is not exactly the same as or similar to the answer, then he must be wrong. Be as objective as possible. Discourage uninformative answers. Also, equally, treat short and long answers and focus on the correctness of answers. After providing your explanation, you must rate the response with either 0, 0.5, or 1 by strictly following this format: "[[rating]]", for example: "Rating: [[0.5]]". [Question] question [The Start of Reference Answer] ground_truth [The

End of Reference Answer] [The Start of Assistant's Answer] answer [The End of Assistant's Answer]



Example.

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

⁴https://whoops-benchmark.github.io/

957

961

962

963

964

965

968

969

970 971

972

973

974

975

976

977

978

947

A.2.3 Fine-Grained Visual Understanding

Fine-grained Visual Understanding represents the capability of LMMs to precisely capture, perceive, and describe certain visual details in the scene. To evaluate the effectiveness of our method to elicit Fine-grained Visual Understanding, we select several challenging datasets including V* Bench (Wu and Xie, 2023) and Spotting Differences.

V* Bench (Wu and Xie, 2023)

Dataset Introduction. The V* Bench (Wu and Xie. 2023) is designed to challenge LMMs to perform visual search and identify fine-grained visual details in high-resolution images. The challenges lie in the necessity of visual search, while GPT-4V sometimes fails to capture visual details, responding with I'm sorry, I couldn't find XXX in the image. Note that we conduct our evaluation in an openended generative manner, GPT-4V sometimes fails to identify targets and refuses to choose an option. *Metric.* The questions are multiple-choice questions. Due to the generative nature of current LMMs, we pose a question to the target LMM and let it generate an open-ended response. Finally, we guide the model to output the final answer within double square brackets [[]], and then use accuracy as the metric to evaluate the model's performance. Example.



Options: (A) The shop is a yoga studio. (B) The shop is a cafe. (C) The shop is a seven-eleven. (D) The shop is a milk tea shop.

Spotting Differences ⁵

Dataset Introduction. The dataset was designed to challenge the spatial analysis capabilities and object localization abilities of LLMs. Inspired and derived by the "Spot the Difference" web game, each level consists of two images. The goal is to locate the differences between the two images at 10 specific locations. The first 50 levels, increasing in difficulty, were selected as samples for this dataset. The challenge lies in the LLM's ability to first analyze the images and then accurately localize the correct object positions. GPT-4V may sometimes return *I'm sorry, but I cannot assist with this request*.

Metric. We guide GPT-4V to answer the pixel or matrix positions of the difference by dividing the prompt into question text + standardized answer requirement + prompt trick. We use OpenCV's Hough Circle Transform to locate the correct image position of the difference. When the distance between the correct position and the position given by GPT-4V is less than 50 pixels, the difference positioning is considered successful. *Example.*



1. [1, 150, 100] - The bird flying in the sky is missing...

A.2.4 Hallucination

Hallucination in vision-language contexts represents a misalignment between visual inputs and textual outputs in terms of object existence, attributes, and relationships. To evaluate the effectiveness of our method to mitigate hallucination in

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

⁵https://www.crazygames.com/game/find-the-difference

vision-language tasks, we select several challenging datasets including POPE Adversarial subset (Li
et al., 2023), HallusionBench (Hard) (Guan et al.,
2023) and Mementos (Wang et al., 2024). The
details of these benchmarks are as follows.

POPE (Li et al., 2023)

1009

1012

1013

1016

1017

1019

1020

1021

1023

1025

1026

1029

1033

1037

Dataset Introduction. The POPE (Li et al., 2023) benchmark is designed for measuring object hallucination in images. The questions are built based on the object annotations, challenging LMMs to determine object existence in images.

Metric. All questions are formatted as general interrogative sentences, which can be answered with either "yes" or "no". Therefore, we guide the model to output the final answer within double square brackets [[]], and then use accuracy as the metric to evaluate the model's performance.

Example.



HallusionBench (Li et al., 2023)

Dataset Introduction. The HallusionBench (Li et al., 2023) is designed for evaluating language hallucination and visual illusion in LMMs. The images and questions are meticulously crafted by human experts, posing great challenges to current LMMs. Due to the limited GPT-4V quota, we evaluate our method on the *Hard* set.

Metric. All questions are formatted as general interrogative sentences, which can be answered with either "yes" or "no". Therefore, we simplify the evaluation process ⁶ and guide the model to output the final answer within double square brackets [[]], and then use accuracy as the metric to evaluate the model's performance.

Example.



Mementos (Wang et al., 2024)

Dataset Introduction. The Mementos (Wang et al., 2024) is designed for evaluating reasoning capabilities across image sequences. Featuring 4,761 diverse image sequences with varying lengths, the dataset adopts a GPT-4_assisted metric to evaluate the correctness of objects and behaviors in generated descriptions, reflecting both reasoning capabilities and hallucination levels.

Metric. The task is to generate a description of the image sequences. We adopt the evaluation script in Mementos official implementation ⁷, which uses GPT-4 (Achiam et al., 2023) as the judge model to extract the objects and behaviors in the description, then calculates F1 score in terms of objects and behaviors respectively. Finally, we use the average F1 score to measure the performance.

Example.



⁷https://github.com/umd-huang-lab/Mementos

⁶the official evaluation process of HallusionBench involves GPT-4 as judge model

A.3 Complementary Cases



Figure 9: An example from the Winoground dataset.



Figure 10: An example from the Spotting Differences dataset.



Figure 11: An example from the POPE (adversarial) dataset.



Figure 12: Examples of different coordinate formats.