
Enhancing Visual Reasoning with Autonomous Imagination in Multimodal Large Language Models

Jingming Liu* Yumeng Li* Boyuan Xiao Yichang Jian Ziang Qin
Tianjia Shao Yao-Xiang Ding† Kun Zhou

State Key Laboratory of CAD&CG, Zhejiang University

jml754457@gmail.com, yumeng.li@zju.edu.cn

boyuan.xiao@zju.edu.cn, mtdickens1998@gmail.com, qinziang19937@gmail.com

tjshao@zju.edu.cn, dingyx.gm@gmail.com, kunzhou@acm.org

Abstract

There have been recent efforts to extend the Chain-of-Thought (CoT) paradigm to Multimodal Large Language Models (MLLMs) by finding visual clues in the input scene, advancing the visual reasoning ability of MLLMs. However, current approaches are specially designed for the tasks where clue finding plays a major role in the whole reasoning process, leading to the difficulty in handling complex visual scenes where clue finding does not actually simplify the whole reasoning task. To deal with this challenge, we propose a new visual reasoning paradigm enabling MLLMs to autonomously modify the input scene to new ones based on its reasoning status, such that CoT is reformulated as conducting simple closed-loop decision-making and reasoning steps under a sequence of *imagined* visual scenes, leading to natural and general CoT construction. To implement this paradigm, we introduce a novel plug-and-play *imagination space*, where MLLMs conduct visual modifications through operations like focus, ignore, and transform based on their native reasoning ability without specific training. We validate our approach through a benchmark spanning dense counting, simple jigsaw puzzle solving, and object placement, challenging the reasoning ability beyond clue finding. The results verify that while existing techniques fall short, our approach enables MLLMs to effectively reason step by step through autonomous imagination.

Project page: <https://future-item.github.io/autoimagine-site>.

1 Introduction

Recently, Multimodal Large Language Models (MLLMs) have achieved impressive advancements in the domain of visual understanding [38, 31, 3, 52, 75, 7, 36, 2]. Despite these strides, MLLMs still face challenges when tackling complex visual reasoning tasks, often producing erroneous outputs and hallucinations. Recent studies have promoted Chain-of-Thought (CoT) reasoning [55] in MLLMs based on the paradigm of visual clue finding, effectively advancing their reasoning abilities [22, 58, 5, 41, 45, 73]. The central idea lies in constructing CoT steps as letting MLLMs find visual clues by adding visual markers, such as bounding boxes, in the given input visual scene, highlighting the key regions to perceive and understand. In the tasks where clue finding significantly simplifies the whole reasoning task, these approaches achieve promising results.

*Equal contribution.

†Corresponding author.

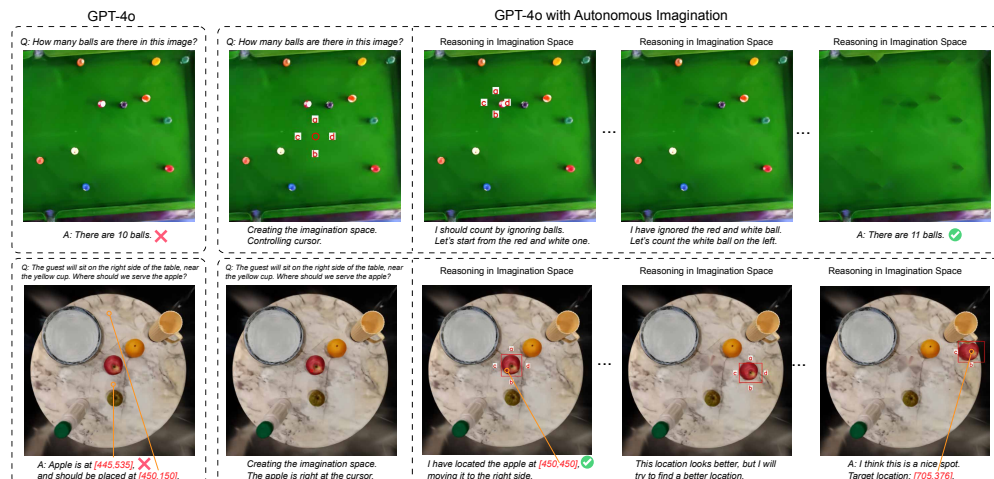


Figure 1: Our autonomous imagination method empowers advanced MLLMs to engage in iterative imaginative reasoning, enabling them to address previously unsolvable tasks without additional training or fine-tuning.

However, there exist three fundamental obstacles for these methods in handling complex visual scenes, where designing CoT based on clue finding does not actually simplify the reasoning task. 1) Adding visual markers could also lead to an increase in visual noise in the input scene. Such visual noise could lead to more serious hallucinations instead of alleviating them. We show an example in Fig. 5c from the dense counting task to illustrate this phenomenon. 2) The process of clue finding does not necessarily lead to the decrease of hardness for single-step reasoning. For example, consider the task of jigsaw puzzle solving shown in Fig. 6, which is to identify the position of a target piece by considering its relationship to all other pieces. Note that identifying one key piece as a trustworthy visual clue also requires considering the complex surrounding relationships. In the worst case, this single step has a similar complexity to the whole task. 3) The reasoning task can involve a significant workload beyond clue finding. As shown in the object placement task in Fig. 4, the reasoning problem challenges the global understanding of object positions and relationships instead of identifying individual visual clues.

On the other hand, inspired by human reasoning, a general way for breaking complex visual reasoning into simpler steps is to gradually modify the input scene, and perform visual reasoning throughout this modification process. For example, for jigsaw puzzle solving, if the CoT process is defined as successively creating the virtual scenes of putting the target piece into one of the candidate positions and justifying its correctness, the whole task is successfully broken into simpler steps. Built on this insight, we introduce a novel visual reasoning paradigm in which MLLMs autonomously modify visual content based on the input scene to create a sequence of *imagined* scenes, and deal with one-step reasoning tasks based on them, all with its native reasoning ability. Given that existing MLLMs are not trained for such imaginative operations, implementing this paradigm presents a substantial challenge.

Our solution is to introduce plug-and-play *imagination spaces*, in which MLLMs autonomously choose to execute modular operations such as *focus*, *ignore*, and *transform*. These operations allow the model to systematically restructure the initially unorganized visual scene, generating a sequence of modified scenes that support multiple rounds of CoT reasoning. This iterative reconfiguration process transforms complex visual tasks into digestible sub-steps aligned with the model’s single-step reasoning capability. Through a closed-loop cycle of imagination and single-step reasoning, MLLMs progressively refine their understanding and derive accurate conclusions—all without requiring additional training or fine-tuning themselves.

To better challenge MLLMs’ visual reasoning ability, we introduce the *Intuitive Visual Reasoning Benchmark (IVRB)*, which focuses on three typical tasks: dense counting, simple jigsaw puzzle solving, and object placement. These tasks represent challenges that are natural and intuitive for humans. More importantly, they all require reasoning beyond clue finding as discussed previously,

outside of what current benchmarks address. Through both quantitative and qualitative experiments, we demonstrate that even though existing techniques fall short, our approach significantly advances the ability of MLLMs to handle these challenging tasks, expanding the boundaries of the visual reasoning ability of MLLMs through autonomous imagination.

2 Related Work

2.1 Multimodal Large Language Models

MLLMs have evolved in visual reasoning tasks, initially leveraging domain-specific expert models like HuggingGPT [47], MM-REACT [61], and VisualChatGPT [56]. The focus later shifted toward training LLMs with an adaptor for the other modal, such as LLaVA [31], BLIP-2 [24, 23], and MoVA [76]. Many recent MLLMs now exhibit native visual understanding through training on text-image pairs, with fine-tuning on question-answering tasks [31, 3, 52, 75, 7, 36, 2], including the state-of-the-art closed-source MLLM GPT-4o [38].

2.2 Visual Reasoning in MLLMs

Numerous studies explore reasoning paradigms in natural language utilizing their in-context learning ability [4], showing that Large Language Models (LLMs) improve through reasoning-based outputs [55, 21]. Subsequent works further enhance reasoning through in-context learning by improving the in-context sample selection [44, 35, 69, 10, 53, 27]. Recently, notable advancements in OpenAI’s o1 [39] have demonstrated that through scaling LLMs for inference-time reasoning before answering, LLMs can be greatly enhanced in their abilities.

With advancements in reasoning within text modalities, substantial efforts are now focused on enabling reasoning in visual tasks to enhance performance. Existing MLLMs often face limitations in direct visual perception and are prone to generating answers with hallucinations [70]. Initial approaches employed attention mechanisms to enhance question-answering capabilities [74], and subsequent solutions include strengthening reasoning in visual modality by training [29, 51] or employing auxiliary knowledge [37]. Additionally, researchers have proposed effective visual prompting methods to refine the focus of MLLMs [16, 28, 12, 68, 67, 37, 71, 65, 63]. See [57] for the recent comprehensive survey.

Attempts have been made to incorporate visual reasoning by generating visual promptings during the reasoning steps in CoT. Based on pioneer studies [70, 71, 66, 40], recent approaches build CoT by enabling MLLMs to conduct visual clue finding either through direct model training [22, 58, 5, 41, 45] or by utilizing plug-and-play visual processing models [73]. This paradigm strengthens reasoning by filtering extraneous information and emphasizing essential visual elements, allowing models to leverage visual prompts effectively without human intervention. However, they remain difficult to handle challenging reasoning tasks beyond clue finding as discussed in Sec. 1.

In closed-world action planning where the visual scene is pre-processed and structured, various imagination techniques have been explored, such as using video generation models to simulate control processes as references [1, 9] or employing image generation to visualize target goals [72]. Methods that improve an LLM’s understanding of the current state, either through textual descriptions or visualizations, have shown to enhance decision-making in closed-world sequential action planning as well [32, 59]. However, these approaches are limited to closed-world settings, where imagination or visualization is primarily used to align the closed-world state space with MLLMs. Although these techniques have shown promising results in closed environments, they are designed specifically for constrained settings and cannot be directly extended to open-world contexts, where the visual scene is primitive and unstructured. In open worlds, the state space is undefined, and the solution space is vast. Leveraging generative models that interpret text prompts and operate in open-world scenarios could be a potential approach. However, such an advanced world model is yet to be developed.

The closest studies to ours are from the field of robotics, where recent methods utilize the ability of MLLMs to perform object manipulation [18, 8]. These methods also create a virtual space and use MLLMs as evaluators to judge whether the object’s final state matches the instruction, in order to generate a final state where the object should be moved to. However, current approaches adopt random sampling to conduct exhaustive searches in the virtual environment, which is not feasible

when the possible state space is large. We implement this method under visual reasoning tasks in the experiments to verify this argument.

2.3 Visual Reasoning Benchmarks for MLLMs

Following the rapid development of MLLMs, various benchmarks are proposed to evaluate their visual reasoning abilities from different aspects. See [25] for a recent comprehensive survey. In this section, we focus on discussing two categories of benchmarks, *multi-step reasoning and object hallucination*, which are closest to the major purpose of our work: extending the boundary of the CoT reasoning ability for state-of-the-art MLLMs.

Multi-step reasoning benchmarks. CoT reasoning involves breaking complicated tasks into multiple simpler steps. There have been several benchmarks proposed for evaluating such kind of multi-step reasoning ability. Some of them focus on tasks of math [34], logic [60], and video [54], where our current method is not applicable. On the other hand, the CoM benchmark [41] evaluates the ability to conduct chain of manipulations for fetching detailed evidence from inputs, and the Visual CoT benchmark [45] evaluates the visual clue finding ability. Both these two benchmarks focus on evaluating relatively small-scale MLLMs, requiring specific training to conduct CoT reasoning. Note that our target is to strengthen state-of-the-art large-scale MLLMs, such as GPT-4o, with their native reasoning ability, which is already strong enough to handle many general reasoning problems. Our IVRB benchmark is proposed for this purpose. We also note that our method may not be suitable for small-scale MLLMs, which do not have strong enough native reasoning ability to properly utilize the autonomous imagination operations.

Object hallucination benchmarks. Current MLLMs can suffer from hallucination by perceiving or generating non-existing objects in the input scene. A number of object hallucination benchmarks are proposed for analyzing this phenomenon [64, 30, 17, 50, 11, 48, 43, 49, 13, 33, 26, 6]. Even though such hallucination seems not to be directly related to CoT reasoning, we note that our paradigm can serve as a promising way to address this challenge, especially for the scenes where multiple objects exist. By the closed-loop iterative process of autonomous imagination, MLLMs can deal with the objects one by one as in our counting benchmark instead of handling them altogether, effectively reducing the possibility of hallucination. We conduct additional experiments on the recently proposed multi-object hallucination benchmark ROPE [6] to validate this argument. See Sec. A for details.

3 Method

3.1 Reasoning Paradigm

We propose a new paradigm, *Autonomous Imagination*, to realize CoT visual reasoning of MLLMs. We consider the following visual reasoning task where an *unstructured* input 2D or 3D visual scene is given, meaning that no semantic pre-processing, such as segmentation or grounding, is conducted initially. Denote the input scene as o . Given the input text prompt token r_0 , the target of reasoning is to obtain the prediction probability $P(y|o, r_0)$, where y is the final output text token. We assume MLLMs take both images and text prompts as inputs. For 2D scenes, o can be processed directly. For 3D scenes, we assume to utilize the images rendered from the proposed 3D imagination space, which will be introduced later in Sec. 3.2.

In general, for utilizing CoT, the reasoning problem of predicting $P(y|o, r_0)$ can be transformed into intermediate steps of predicting $P(z_t|o, z_{t-1})$, $t \in 1, 2, \dots, T$. Each $z_t = \{r_t, c_t\}$ denotes the augmented textual prompts r_t and visual information c_t for deepening scene understanding and pushing forward reasoning. After the final reasoning step at T , the final output is set as $y = r_T$. The target transforms into predicting the following joint probability:

$$P(z_{1:T}|o, r_0) = \prod_{t=1}^T P(z_t|o, z_{0:t-1}), \quad (1)$$

where $z_{t_1:t_2} = \{z_{t_1}, z_{t_1+1}, \dots, z_{t_2}\}$ and $z_0 = r_0$. What is essential in CoT design is to make each step simple enough to match the reasoning capability of MLLMs.

In our approach, we propose a new paradigm of CoT design. We introduce the *imagination space*, a virtual visual environment that can render *imagined scenes* \hat{o} , which are also images for MLLMs to

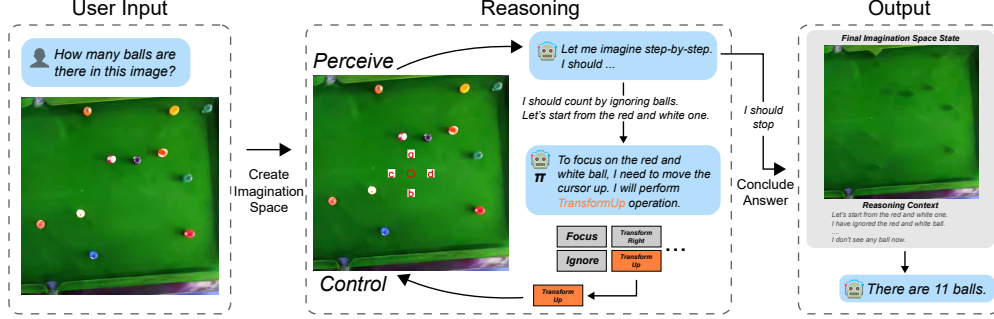


Figure 2: An overview of our *autonomous imagination* method: The imagination space begins with an unstructured input scene and undergoes an iterative reasoning process. In each cycle, MLLMs first perceive the current state of the imagination space, select an operation to apply, and then reassess the updated imagination space. Upon completing this reasoning sequence, MLLMs generate an answer based on the cumulative context of the process and the final state of the imagination space.

perceive and elicit MLLMs to generate new augmented textual prompt r . Furthermore, we introduce a set of *operators* for modifying the imagined scenes into new ones. Denote by $v_t = \{\hat{o}_t, r_{0:t}\}$ and a_t the operator at step t . Our method transforms Eq. 1 into

$$P(\hat{o}_{1:T}, r_{1:T}, a_{1:T} | o, r_0) = \prod_{t=1}^T P(v_t, a_t | v_{t-1}), \quad (2)$$

where $v_0 = \{o, r_0\}$ and finally $y = r_T$. The one-step reasoning task $P(v_t, a_t | v_{t-1})$ is factorized into

$$P(v_t, a_t | v_{t-1}) = \pi(a_t | v_{t-1}) \phi(\hat{o}_t | \hat{o}_{t-1}, a_t) \omega(r_t | \hat{o}_t, r_{1:t-1}). \quad (3)$$

Eq. 3 shows that one-step reasoning is factorized into a decision function $\pi(a_t | v_{t-1})$, a scene modification function $\phi(\hat{o}_t | \hat{o}_{t-1}, a_t)$, and a reasoning function $\omega(r_t | \hat{o}_t, r_{1:t-1})$. Given the current imagined scene \hat{o}_{t-1} and augmented text prompt token r_{t-1} , the decision function chooses a specific scene modification operator a_t . The scene modification function then utilizes a_t to update \hat{o}_{t-1} into \hat{o}_t . The reasoning function finally updates r_{t-1} into r_t based on \hat{o}_t . The decision and reasoning functions rely purely on the native reasoning ability of MLLMs. The scene modification function is implemented inside the imagination space.

Comparing Eq. 1 and Eq. 2, our paradigm transforms CoT reasoning into a closed-loop decision-making and reasoning process. Furthermore, Eq. 3 indicates that the reasoning step at step t is only dependent on \hat{o}_{t-1}, \hat{o}_t but independent of all previous scenes. In particular, the reasoning focuses on the current imagined scenes instead of the input scene o . We show that in many complex reasoning tasks, this paradigm effectively alleviates the visual noise issue and leads to effective single-step task decomposition.

3.2 Imagination Space

To handle both 2D and 3D visual inputs, we developed separate 2D and 3D imagination spaces, where 2D space uses image as representation and 3D space is based on the 3D Gaussian splatting model [19]. The imagination space is designed to support visual rendering and supports a minimal set of operators: focus, ignore, and transform. Focus isolates relevant content for further manipulation, ignore enables MLLMs to disregard extraneous information, and transform allows the repositioning of desired content. This compact set of operations enables MLLMs to reason and solve practical challenges by themselves, as demonstrated in our experiments. MLLMs select operators through natural language output, which are then applied to the imagination space. The updated space is rendered and returned to MLLMs for the subsequent reasoning step.

3.3 Focus Operator

The focus operator allows MLLMs to automatically isolate and label the target content from a scene, creating distinct elements for further transformation. Once MLLMs identify an object of interest

Algorithm 1: 3D conditional segmentation

Input : SAM2 mask, 3D scene, thresholds ϵ_1, ϵ_2
Output : Set of marked Gaussians $\mathcal{G}_{\text{marked}}$

```
1 Initialize vote counts:  $\text{votes}[g_i] \leftarrow 0$  for each Gaussian  $g_i$  in the scene
2 foreach pixel  $p$  in masked_pixels do
3   ray  $\leftarrow$  cast_ray( $p$ );  $T \leftarrow 1$ 
4   foreach Gaussian  $g_i$  intersected by ray do
5      $\alpha_i \leftarrow$  get_alpha( $g_i$ , ray);  $C_i \leftarrow T \cdot \alpha_i$ 
6     if  $C_i > \epsilon_1$  then
7       |  $\text{votes}[g_i] \leftarrow \text{votes}[g_i] + 1$ 
8     end
9      $T \leftarrow T \cdot (1 - \alpha_i)$ 
10    if  $T < \epsilon_2$  then break
11  end
12 end
13  $\text{max\_vote} \leftarrow$  maximum of  $\text{votes}[g_i]$  over all  $g_i$ 
14  $\mathcal{G}_{\text{marked}} \leftarrow \emptyset$ 
15 foreach Gaussian  $g_i$  do
16   if  $\text{votes}[g_i] \geq 0.1 \times \text{max\_vote}$  then
17     |  $\mathcal{G}_{\text{marked}} \leftarrow \mathcal{G}_{\text{marked}} \cup \{g_i\}$ 
18   end
19 end
```

using a virtual cursor (as described in Sec. 3.6), the focus operator segments this content for focused manipulation.

In 2D space, this cursor directly conditions the Segment Anything Model (SAM) [20] to perform segmentation. MLLMs verify segmented output to ensure alignment with their intended focus. However, in 3D Gaussian space, segmentation on demand poses unique challenges: Existing methods [62, 46] are designed for unconditional segmentation, where all objects are segmented without the ability to specify conditions. However, our use case requires conditional segmentation, focusing on a specific object selected by an input condition, making these methods unsuitable.

To enable the focus operator in 3D Gaussian-based representations, we introduce a method to selectively segment an object from an unstructured 3D Gaussian scene, which is illustrated in Alg. 1. Given the MLLM’s initial selection block, we first generate a virtual camera trajectory that orbits around the targeted object, creating a video sequence. This sequence is fed into the SAM2 model [42], which applies the selection cursor condition across all rendered frames. For each frame, rays are cast from pixels within the segmentation mask, and path tracing records the radiance contributions of intersecting Gaussians. This process is expressed as volumetric integration illustrated in Line 4-11. Gaussians that receive a contribution higher than a threshold will receive one vote.

We select 3D Gaussian with votes exceeding 10% of the highest vote, forming a shell around the target object’s radiance field. This shell, along with the internal radiance Gaussians, is segmented collectively to define the object’s radiance structure. In both 2D and 3D spaces, the segmented object is placed in a separate layer, while the remaining content is retained in a base layer.

3.4 Ignore Operator

The ignore operation removes a focused object from the imagination space to prevent interference with the reasoning process. Since API access to state-of-the-art MLLMs doesn’t support attention masks, removing the object leaves a hole in the base layer, potentially causing hallucinations. To address this, we simply inpaint the hole created by the object’s removal.

For both 2D and 3D imagination spaces, we project the object mask onto the rendered image as the region to inpaint. We directly use the implementation provided by OpenCV [15] by simply merging the nearby pixels. The inpainted image is then used for further reasoning by the MLLMs, effectively avoiding it from paying attention to the inpainted region as the content has been removed.

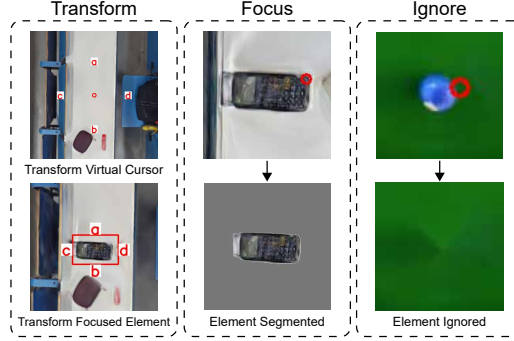


Figure 3: Illustrations of operations within our imagination space: transformations can be applied to focused elements (including the virtual cursor), focus operations allow segmentation of cursor-selected elements, and ignore operations make cursor-selected elements visually invisible.

3.5 Transform Operator

Object transformation can occur in four directions: up, down, left, and right, each represented by an alphabet character (a, b, c, d) to avoid interference from semantically meaningful words. In each step of the iterative imagination process, MLLMs select a direction for movement. Since MLLMs struggle with determining precise distances, we standardize movement units in screen space coordinates and gradually reduce step size throughout the process, where MLLMs only control the direction of movement.

3.6 Reasoning Process

Specifically, the iterative reasoning process described in Sec. 3.1 is executed as follows: After initializing the imagination space, a virtual cursor is positioned at the center and automatically considered as focused. MLLMs can perform transform operations on the cursor to reposition it within the space. Following each operation, MLLMs conduct scene modification and receive an updated image of the newly rendered imagination space.

When MLLMs identify that the cursor has selected a visual element of interest, a focus operation is performed based on the cursor’s location. If MLLMs determine that this operation has correctly segmented the intended element, the focus is shifted to that object, enabling MLLMs to execute either a transformation or ignore operation on it. Once MLLMs have either repositioned or disregarded the object, the focus returns to the virtual cursor. This process iterates continuously as part of the reasoning workflow until all necessary reasoning steps are completed, culminating in a final textual response.

4 Benchmark and Evaluation Protocol

While previous works have demonstrated MLLMs’ capabilities in various visual tasks, a significant gap remains in their ability to perform practical, intuitive reasoning tasks that are simple for humans. Rather than benchmarking existing abilities, our focus is on expanding MLLMs’ capabilities, enabling them to tackle tasks that were previously unsolvable. To this end, we propose the Intuitive Visual Reasoning Benchmark (IVRB), which evaluates three tasks that require practical visual reasoning: dense object counting, simple jigsaw puzzle solving, and object placement.

Dense Object Counting: When faced with numerous densely packed objects, humans often rely on multi-step reasoning to reach an accurate count, as a single glance may not suffice. Existing MLLMs, however, struggle with dense object counting, as their multi-step reasoning is not yet effective for this task. We include this task to assess whether the model can leverage reasoning to offset its limited direct perception, mirroring human strategies.

For evaluation, we directly compare the model’s predicted count to the ground truth. In addition to reporting the **success rate**, we calculate the **mean** and **variance** of counting errors to provide insight

into each model’s accuracy and consistency. We constructed in total 122 images with paired ground truth for our benchmark.

Solving Simple Jigsaw Puzzles: Jigsaw puzzles are classic tests of visual perception and reasoning, commonly used to assess intelligence. In this task, we evaluate MLLMs’ ability to solve simple jigsaw puzzles, aligning their problem-solving performance with that of humans to assess visual reasoning capabilities. In the Jigsaw Puzzle Solving task, MLLMs are tasked with identifying and placing missing pieces in their correct locations.

For evaluation, we use a digital jigsaw puzzle game with a magnetic mechanism, where pieces automatically snap into place when positioned close to their correct locations. Success is measured by the **completion rate**, defined as the percentage of pieces placed in their correct locations. For evaluation, we constructed 11 cases where four pieces are missing, and 11 cases where six pieces are missing. The jigsaw puzzle’s size ranges from 3×5 to 5×8 .

Object Placement: Previous methods have shown progress in enabling MLLMs to understand and describe static scenes, such as identifying object locations (e.g., "Where is the cup?"). However, for practical use, MLLMs must also interpret dynamic instructions that convey intent, such as "Where should the cup be placed?" In the Object Placement task, MLLMs are required to identify both the current and target locations of objects based on abstract instructions (e.g., "Prepare two cups for guests in the living room"). Given the input scene of 3D Gaussians, MLLMs must determine the original locations of objects and their intended locations based on a provided prompt.

For evaluation, we first measure the **locating success rate**, which reflects the model’s ability to accurately identify the initial location of the correct object. If the model fails to do this, the placement task is automatically marked as a failure. We then measure the **placement success rate** by checking whether the predicted final location falls into the marked ground-truth region. For both success rates, when multiple correct solutions exist, we provide multiple ground truth regions and selecting any valid region is considered correct. We captured a total of 17 scenes, including 271 user prompts and paired ground truth for evaluation.

5 Experiments

5.1 Baselines

We evaluate our approach against state-of-the-art MLLM, namely **GPT-4o** [38]. We provide a clearly defined 2D coordinate system to GPT-4o to ensure output coordinates are generated without any ambiguity. We also adopt the recently proposed representative clue finding method **VCoT** [45] using their open-released pre-trained model. We developed a baseline **GPT-4o Sampling** inspired by the sample-then-evaluate imagination paradigm used in robotics [18, 8]. Since their original methods are designed for robotic manipulation, we re-implement them under our paradigm. Note that this sampling strategy is not designed for general visual reasoning, and can only be applied to some of the tasks in our evaluation. Furthermore, to demonstrate the necessity of altering the visual scene beyond simply finding visual clues, we crafted a robust reasoning baseline named **cursor-only** that utilizes our imagination space but restricts operations solely to transform operations of the virtual cursor. Note that this baseline serves as an ablation of our method, which has the same closed-loop control design except that scene modification is disabled and only virtual cursor moving is enabled.

5.2 Dense Counting

As shown in Tab. 1, our model significantly outperforms all baselines, achieving a higher success rate and lower errors when mistakes occur. Note that the sampling method is not directly applicable to the counting task due to intractability so we omit its comparisons. VCoT only supports drawing one visual marker for clue finding, hence is unsuitable for more complex reasoning tasks involving multiple objects, such as counting. We also compare with our cursor-only baseline in counting, observing that it performs similarly to GPT-4o. We further experimented with reasoning by drawing multiple visual markers, implemented by allowing cursor-only baseline to draw visual markers when performing the focus operation. Although technically able to draw markers on every ball to perform counting accurately, the model quickly enters a negative feedback loop: hallucinations lead it to draw more markers, which introduce additional noise and exacerbate the hallucinations. This results in

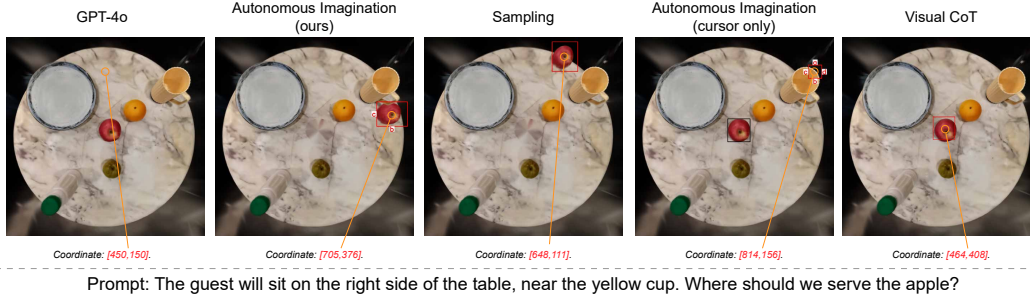


Figure 4: Qualitative comparison of apple placement based on user prompts. Coordinate locations are highlighted on the image with orange lines and circles for improved visualization. Our method aligns more closely with the user’s requirements, placing the apple further to the right side of the table.

Settings		GPT-4o	VCoT	Ours (cursor-only)	GPT-4o Sampling	Ours
Dense Counting	Success Rate	39.8%	15.5%	39.0%/0%	-	85.3%
	Mean Error	0.82	2.90	1.02/ <i>intractable</i>	-	0.17
	Variance	1.49	15.57	2.54/ <i>intractable</i>	-	0.22
Simple Jigsaw Puzzle	4-Piece Missing	29.5%	9.1%	27.3%	43.2%	68.2%
	6-Piece Missing	9.1%	3.3%	24.2%	30.3%	51.5%
Object Placement	Locating	10.9%	10.4%	69.4%*	-	69.4%
	Placement	3.6%	1.5%	27.8%	17.3%	37.3%

Table 1: Quantitative comparison results. See Sec. 4 for details of the metrics. We consider two variants of cursor-only in dense counting, leading to two sets of results. Furthermore, GPT-4o Sampling is intractable in some tasks as illustrated. See Sec. 5.2 and Sec. 5.4 for details. *In object placement, cursor-only functions the same as our method in locating, so their results are identical.

a success rate of 0% and intractable mean error and variance, as the model cannot stop counting, further highlighting the importance of modifying the visual scene.

We further conduct analysis to increase the counting difficulty by adding more objects with partial data. Results are shown in Fig. 5a with an additional baseline **GPT-4o-text-cot**, which performs counting via pure text-based CoT. This reveals perceptual limitations in GPT-4o that cannot be overcome by simply increasing the textual CoT steps, while our method remains unaffected. Fig. 5b shows the percentage of cases solved as reasoning steps increase. For reference, we include the one-step results from GPT-4o depicted as a flat line. Our method performs consistently better as reasoning steps grow, whereas GPT-4o-text-cot fails to do this and may even introduce noise, resulting in poorer performance than one-step GPT-4o.

5.3 Simple Jigsaw Puzzle Solving

As shown in Tab. 1 and Fig. 6, our method consistently outperforms all baselines. It achieves a higher locate accuracy, which contributes to a higher success rate with the same number of attempts. Since this visual reasoning task requires the model to choose the correct target location rather than merely locating visible elements, VCoT is ineffective at reasoning. The sampling method performs reasonably well in this task since we set a large sampling budget to ensure the correct answers can be sampled as the candidate answers. However, despite its extensive reasoning budget, it still does not match the performance of our method. This highlights the effectiveness and robustness of establishing an autonomous perception-control loop. We also demonstrate that the cursor-only baseline performs similarly to primitive GPT-4o under the four-piece missing scenario. However, it performs significantly better under the more challenging six-piece missing scenario, further confirming the effectiveness of the closed-loop control paradigm.

5.4 Object Placement

As shown in Tab. 1, our method consistently outperforms baseline methods in both object locating and placement performance. It is important to note that our cursor-only ablation baseline is identical

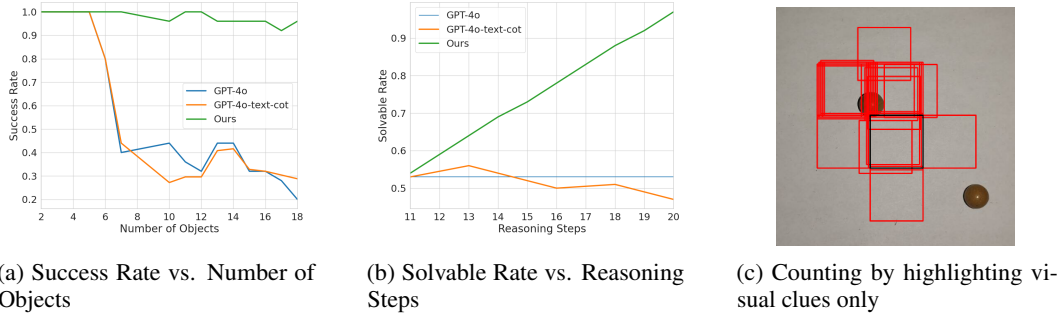


Figure 5: We demonstrate in (a)(b) that as counting task difficulty increases linearly, scaling inference-time textual reasoning (implemented as GPT-4o-text-cot, See Sec. 5.2) fails—and even performs worse than vanilla GPT-4o—as complexity exceeds perception limits. In contrast, our methods remain unaffected by these limits, achieving correct counting even as difficulty rises. Additionally, highlighting visual clues can introduce noise, causing MLLMs to enter hallucination loops; for instance, in (c), the model incorrectly concluded there were 196 balls when only two were present. Qualitative comparisons of counting are provided in the appendix.

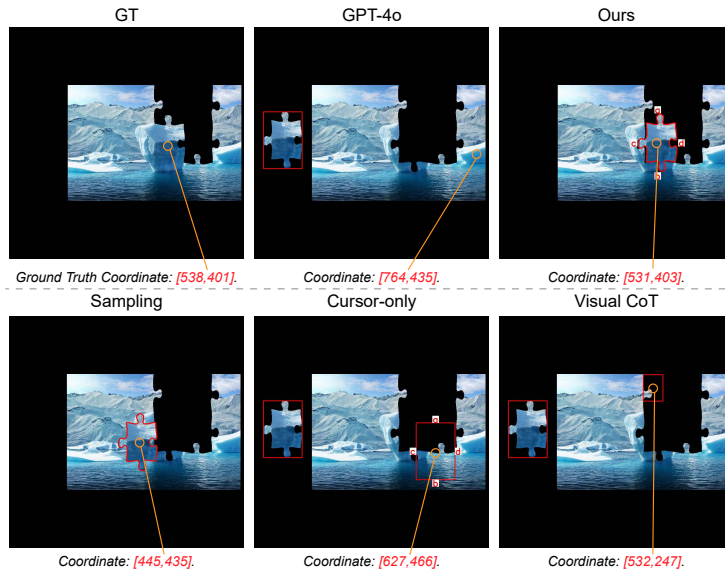


Figure 6: Qualitative comparison on simple jigsaw puzzle solving, we use a black background to make jigsaw pieces more visible to MLLMs. We illustrate the final visual state achieved by different methods after completing their reasoning processes and producing a solution. For clarity, the actual location of each coordinate in the image is highlighted with an orange line and circle.

to our full method when locating an object. Therefore, we assign them with the same locating result. Additionally, the sampling method is intractable for locating, we directly utilize the locating results of our method for placement.

Methods lacking advanced visual reasoning capabilities perform poorly on the placement task, as the successful placement requires inferencing about the correct spatial position for that element beyond visual recognition. The cursor-only baseline, which moves visual marker instead of scene modification, does not perform as effectively as our complete method. This difference underscores the advantage of scene manipulation in simplifying the reasoning process by revealing additional visual information. Although the sampling baseline receives substantial visual information, it still under-performs relative to our method, and even falls short of the cursor-only baseline. Conversely, in simpler jigsaw puzzle solving, the sampling method achieves a relatively high success rate compared to the cursor-only baseline. This contrast highlights the importance of a structured reasoning pathway, particularly in complex, open-world scenarios where the abundant visual information could

overwhelm the MLLM, preventing it from identifying the correct answer despite its presence in the sample. By following the reasoning process of closed-loop control, MLLMs can progressively approach the correct answer without requiring an exhaustive number of samples, resulting not only in greater efficiency but also in improved accuracy.

6 Limitations and Future Work

While our method enables MLLMs to perform visual reasoning tasks by operating within a structured framework autonomously, its expressiveness is constrained by the defined set of operators. Future work could focus on extending the operator set to address a broader range of tasks as MLLMs' capability to handle complex operations improves. Furthermore, if future MLLMs natively possess robust and precise image editing capabilities, the reasoning process could become more efficient and streamlined, potentially removing the need for a specially designed imagination space.

7 Conclusion

Current MLLMs still struggle with visual reasoning tasks beyond clue finding that humans perform intuitively. We propose a new reasoning paradigm to conduct closed-loop decision-making and reasoning with imagination, thus creating an interactive imagination space where MLLMs can refine reasoning. Our method employs plug-and-play imagination space and operator design, enabling MLLMs to modify visual content autonomously, enabling general and natural CoT construction. We propose the IVRB benchmark across three challenging visual reasoning tasks, simple jigsaw puzzle solving, dense counting, and object placement, demonstrating the limitations of existing methods for reasoning beyond clue finding. Our approach, however, enables MLLMs to handle these tasks with autonomous imagination from their native reasoning abilities, pushing the boundaries of visual reasoning capabilities for MLLMs.

8 Acknowledgement

This work is supported by National Natural Science Foundation of China (U23A20311,62206245).

References

- [1] Anurag Ajay, Seungwook Han, Yilun Du, Shaung Li, Abhi Gupta, Tommi Jaakkola, Josh Tenenbaum, Leslie Kaelbling, Akash Srivastava, and Pulkit Agrawal. Compositional foundation models for hierarchical planning. *arXiv preprint arXiv:2309.08587*, 2023.
- [2] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022.
- [3] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. 2023.
- [4] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [5] Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multimodal llm's referential dialogue magic. *arXiv preprint arXiv:2306.15195*, 2023.
- [6] Xuweiyi Chen, Ziqiao Ma, Xuejun Zhang, Sihan Xu, Shengyi Qian, Jianing Yang, David Fouhey, and Joyce Chai. Multi-object hallucination in vision language models. In *Advances in Neural Information Processing Systems 37*, 2024.
- [7] Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. Pali: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2209.06794*, 2022.

- [8] Yufei Ding, Haoran Geng, Chaoyi Xu, Xiaomeng Fang, Jiazhao Zhang, Qiyu Dai, Songlin Wei, Zhizheng Zhang, and He Wang. Open6DOR: Benchmarking open-instruction 6-dof object rearrangement and a VLM-based approach. In *ICRA 2024 Workshop on 3D Visual Representations for Robot Manipulation*, 2024.
- [9] Yilun Du, Mengjiao Yang, Bo Dai, Hanjun Dai, Ofir Nachum, Joshua B Tenenbaum, Dale Schuurmans, and Pieter Abbeel. Learning universal policies via text-guided video generation. *arXiv e-prints*, pages arXiv–2302, 2023.
- [10] Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and Tushar Khot. Complexity-based prompting for multi-step reasoning. In *The Eleventh International Conference on Learning Representations*, 2022.
- [11] Anisha Gunjal, Jihan Yin, and Erhan Bas. Detecting and preventing hallucinations in large vision language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 18135–18143, 2024.
- [12] Wenyi Hong, Weihang Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxiao Dong, Ming Ding, et al. Cogagent: A visual language model for gui agents. *arXiv preprint arXiv:2312.08914*, 2023.
- [13] Hongyu Hu, Jiyuan Zhang, Minyi Zhao, and Zhenbang Sun. Ciem: Contrastive instruction evaluation method for better instruction tuning. *arXiv preprint arXiv:2309.02301*, 2023.
- [14] Yuzhou Huang, Liangbin Xie, Xintao Wang, Ziyang Yuan, Xiaodong Cun, Yixiao Ge, Jiantao Zhou, Chao Dong, Rui Huang, Ruimao Zhang, et al. Smartedit: Exploring complex instruction-based image editing with multimodal large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8362–8371, 2024.
- [15] Itseez. Open source computer vision library. <https://github.com/itseez/opencv>, 2015.
- [16] Songtao Jiang, Yan Zhang, Chenyi Zhou, Yeying Jin, Yang Feng, Jian Wu, and Zuozhu Liu. Joint visual and text prompting for improved object-centric perception with multimodal large language models. *arXiv preprint arXiv:2404.04514*, 2024.
- [17] Liqiang Jing, Ruosen Li, Yunmo Chen, Mengzhao Jia, and Xinya Du. Faithscore: Evaluating hallucinations in large vision-language models. *arXiv preprint arXiv:2311.01477*, 2023.
- [18] Ivan Kapelyukh, Yifei Ren, Ignacio Alzugaray, and Edward Johns. Dream2Real: Zero-shot 3D object rearrangement with vision-language models. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2024.
- [19] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), July 2023.
- [20] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023.
- [21] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213, 2022.
- [22] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. *arXiv preprint arXiv:2308.00692*, 2023.
- [23] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023.
- [24] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022.
- [25] Lin Li, Guikun Chen, Hanrong Shi, Jun Xiao, and Long Chen. A survey on multimodal benchmarks: In the era of large ai models. *arXiv preprint arXiv:2409.18142*, 2024.
- [26] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023.
- [27] Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. On the advance of making language models better reasoners. *arXiv preprint arXiv:2206.02336*, 2022.

- [28] Weifeng Lin, Xinyu Wei, Ruichuan An, Peng Gao, Bocheng Zou, Yulin Luo, Siyuan Huang, Shanghang Zhang, and Hongsheng Li. Draw-and-understand: Leveraging visual prompts to enable mllms to comprehend what you want. *arXiv preprint arXiv:2403.20271*, 2024.
- [29] Ziyi Lin, Chris Liu, Renrui Zhang, Peng Gao, Longtian Qiu, Han Xiao, Han Qiu, Chen Lin, Wenqi Shao, Keqin Chen, et al. Sphinx: The joint mixing of weights, tasks, and visual embeddings for multi-modal large language models. *arXiv preprint arXiv:2311.07575*, 2023.
- [30] Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Mitigating hallucination in large multi-modal models via robust instruction tuning. In *The Twelfth International Conference on Learning Representations*, 2023.
- [31] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023.
- [32] Ruibo Liu, Jason Wei, Shixiang Shane Gu, Te-Yen Wu, Soroush Vosoughi, Claire Cui, Denny Zhou, and Andrew M. Dai. Mind’s eye: Grounded language model reasoning through simulation, 2022.
- [33] Holy Lovenia, Wenliang Dai, Samuel Cahyawijaya, Ziwei Ji, and Pascale Fung. Negative object presence evaluation (nope) to measure object hallucination in vision-language models. *arXiv preprint arXiv:2310.05338*, 2023.
- [34] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023.
- [35] Pan Lu, Liang Qiu, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, Tanmay Rajpurohit, Peter Clark, and Ashwin Kalyan. Dynamic prompt learning via policy gradient for semi-structured mathematical reasoning. *arXiv preprint arXiv:2209.14610*, 2022.
- [36] Gen Luo, Yiyi Zhou, Tianhe Ren, Shengxin Chen, Xiaoshuai Sun, and Rongrong Ji. Cheap and quick: Efficient vision-language instruction tuning for large language models. *Advances in Neural Information Processing Systems*, 36, 2023.
- [37] Chancharik Mitra, Brandon Huang, Trevor Darrell, and Roei Herzig. Compositional chain-of-thought prompting for large multimodal models. *arXiv preprint arXiv:2311.17076*, 2023.
- [38] OpenAI. Hello GPT-4o, 2024.
- [39] OpenAI. Learning to reason with llms, 2024.
- [40] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*, 2023.
- [41] Ji Qi, Ming Ding, Weihang Wang, Yushi Bai, Qingsong Lv, Wenyi Hong, Bin Xu, Lei Hou, Juanzi Li, Yuxiao Dong, et al. Cogcom: Train large vision-language models diving into details through chain of manipulations. *arXiv preprint arXiv:2402.04236*, 2024.
- [42] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024.
- [43] Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. Object hallucination in image captioning. *arXiv preprint arXiv:1809.02156*, 2018.
- [44] Ohad Rubin, Jonathan Herzig, and Jonathan Berant. Learning to retrieve prompts for in-context learning. *arXiv preprint arXiv:2112.08633*, 2021.
- [45] Hao Shao, Shengju Qian, Han Xiao, Guanglu Song, Zhuofan Zong, Letian Wang, Yu Liu, and Hongsheng Li. Visual cot: Unleashing chain-of-thought reasoning in multi-modal language models. *Advances in Neural Information Processing Systems*, 37, 2024.
- [46] Qihong Shen, Xingyi Yang, and Xinchao Wang. Flashsplat: 2d to 3d gaussian splatting segmentation solved optimally. In *European Conference on Computer Vision*, pages 456–472. Springer, 2024.
- [47] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving ai tasks with chatgpt and its friends in hugging face. *Advances in Neural Information Processing Systems*, 36, 2023.
- [48] Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. Aligning large multimodal models with factually augmented rlhf. *arXiv preprint arXiv:2309.14525*, 2023.
- [49] Junyang Wang, Yuhang Wang, Guohai Xu, Jing Zhang, Yukai Gu, Haitao Jia, Ming Yan, Ji Zhang, and Jitao Sang. An llm-free multi-dimensional benchmark for mllms hallucination evaluation. *arXiv preprint arXiv:2311.07397*, 2023.

- [50] Junyang Wang, Yiyang Zhou, Guohai Xu, Pengcheng Shi, Chenlin Zhao, Haiyang Xu, Qinghao Ye, Ming Yan, Ji Zhang, Jihua Zhu, et al. Evaluation and analysis of hallucination in large vision-language models. *arXiv preprint arXiv:2308.15126*, 2023.
- [51] Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, et al. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. *Advances in Neural Information Processing Systems*, 36, 2023.
- [52] Weihai Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. Cogvlm: Visual expert for pretrained language models. *arXiv preprint arXiv:2311.03079*, 2023.
- [53] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- [54] Yan Wang, Yawen Zeng, Jingsheng Zheng, Xiaofen Xing, Jin Xu, and Xiangmin Xu. Videocot: A video chain-of-thought dataset with active annotation tool. *arXiv preprint arXiv:2407.05355*, 2024.
- [55] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- [56] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*, 2023.
- [57] Junda Wu, Zhehao Zhang, Yu Xia, Xintong Li, Zhaoyang Xia, Aaron Chang, Tong Yu, Sungchul Kim, Ryan A Rossi, Ruiyi Zhang, et al. Visual prompting in multimodal large language models: A survey. *arXiv preprint arXiv:2409.15310*, 2024.
- [58] Penghao Wu and Saining Xie. V*: Guided visual search as a core mechanism in multimodal llms. *arXiv preprint arXiv:2312.14135*, 17, 2023.
- [59] Wenshan Wu, Shaoguang Mao, Yadong Zhang, Yan Xia, Li Dong, Lei Cui, and Furu Wei. Visualization-of-thought elicits spatial reasoning in large language models. *arXiv preprint arXiv:2404.03622*, 2024.
- [60] Yijia Xiao, Edward Sun, Tianyu Liu, and Wei Wang. Logicvista: Multimodal llm logical reasoning benchmark in visual contexts. *arXiv preprint arXiv:2407.04973*, 2024.
- [61] Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. Mm-react: Prompting chatgpt for multimodal reasoning and action. *arXiv preprint arXiv:2303.11381*, 2023.
- [62] Mingqiao Ye, Martin Danelljan, Fisher Yu, and Lei Ke. Gaussian grouping: Segment and edit anything in 3d scenes. In *European Conference on Computer Vision*, pages 162–179. Springer, 2024.
- [63] Runpeng Yu, Weihao Yu, and Xinchao Wang. Api: Attention prompting on image for large vision-language models, 2024.
- [64] Bohan Zhai, Shijia Yang, Xiangchen Zhao, Chenfeng Xu, Sheng Shen, Dongdi Zhao, Kurt Keutzer, Manling Li, Tan Yan, and Xiangjun Fan. Halle-switch: Rethinking and controlling object existence hallucinations in large vision language models for detailed caption. *arXiv preprint arXiv:2310.01779*, 2023.
- [65] Daoan Zhang, Junming Yang, Hanjia Lyu, Zijian Jin, Yuan Yao, Mingkai Chen, and Jiebo Luo. Cocot: Contrastive chain-of-thought prompting for large multimodal models with multiple image inputs. *arXiv preprint arXiv:2401.02582*, 2024.
- [66] Shilong Zhang, Peize Sun, Shoufa Chen, Min Xiao, Wenqi Shao, Wenwei Zhang, Kai Chen, and Ping Luo. Gpt4roi: Instruction tuning large language model on region-of-interest. *arXiv preprint arXiv:2307.03601*, 2023.
- [67] Yuechen Zhang, Shengju Qian, Bohao Peng, Shu Liu, and Jiaya Jia. Prompt highlighter: Interactive control for multi-modal llms. *arXiv preprint arXiv:2312.04302*, 2023.
- [68] Yuanhan Zhang, Kaiyang Zhou, and Ziwei Liu. What makes good examples for visual in-context learning? *Advances in Neural Information Processing Systems*, 2023.
- [69] Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493*, 2022.
- [70] Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. Multimodal chain-of-thought reasoning in language models. *arXiv preprint arXiv:2302.00923*, 2023.

- [71] Ge Zheng, Bin Yang, Jiajin Tang, Hong-Yu Zhou, and Sibeil Yang. Ddcot: Duty-distinct chain-of-thought prompting for multimodal reasoning in language models. *Advances in Neural Information Processing Systems*, 36:5168–5191, 2023.
- [72] Enshen Zhou, Yiran Qin, Zhenfei Yin, Yuzhou Huang, Ruimao Zhang, Lu Sheng, Yu Qiao, and Jing Shao. Minedreamer: Learning to follow instructions via chain-of-imagination for simulated-world control. *arXiv preprint arXiv:2403.12037*, 2024.
- [73] Qiji Zhou, Ruochen Zhou, Zike Hu, Panzhong Lu, Siyang Gao, and Yue Zhang. Image-of-thought prompting for visual reasoning refinement in multimodal large language models, 2024.
- [74] Yiyi Zhou, Tianhe Ren, Chaoyang Zhu, Xiaoshuai Sun, Jianzhuang Liu, Xinghao Ding, Mingliang Xu, and Rongrong Ji. Trar: Routing the attention spans in transformer for visual question answering. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2074–2084, 2021.
- [75] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.
- [76] Zhuofan Zong, Bingqi Ma, Dazhong Shen, Guanglu Song, Hao Shao, Dongzhi Jiang, Hongsheng Li, and Yu Liu. Mova: Adapting mixture of vision experts to multimodal context. *arXiv preprint arXiv:2404.13046*, 2024.



Figure 8: We show failure cases in our method caused by hallucination. These include instances where MLLMs mistakenly perceive a fry that does not exist, misinterpret the wrongly focused floor as a dumbbell, and incorrectly believe that a puzzle piece has been placed in the correct spot.

A More Experiments

ROPE Benchmark. We evaluate our method using the recently proposed ROPE benchmark for Multi-Object Hallucination [6], focusing on the most challenging subset where GPT-4o exhibits the poorest performance, which is answering about single object given multi-object image, heterogeneous object types, on unseen data split. Our method uses GPT-4o as the base model with our imagination space and operations equipped. Please refer to the original paper for more details. Given that this benchmark emphasizes the identification of abstract objects that are not suitable for segmentation, we turn the focus operation into drawing a large rectangular region by MLLMs natively, through selecting top-left and bottom-right corners. The benchmark indicates that MLLMs are prone to hallucinations when confronted with multiple objects that introduce additional visual distractions. We demonstrate that using our method, MLLMs can autonomously focus on the important region despite these distractions, leading to more accurate responses due to the elimination of irrelevant visual information.

Models	Acc.	Models	Acc.
LLaVA-34B	30.81%	GroundHOG*	38.13%
IDEFICS	6.50%	Qwen VL-C	15.37%
Yi-VL-34B	0.41%	MiniCPM-V	14.39%
CogVLM-C	13.50%	GPT-4o	53.74%
GLaMM*	52.28%	GPT-4o+Ours	65.00%

Table 2: Comparison of our method with baselines in ROPE benchmark, see Sec. A for details. Except for our method, the results of other baselines are cited from the latest leaderboard in <https://multi-object-hallucination.github.io/> before our arXiv submission time. For each base model family, we take the one with the highest performance. Mechanistically grounded LVLMs (marked with *) take visual prompts by dedicated pointer tokens. Please refer to the original paper [6] for more details.

Imagine with Image Editing Models. We experimented with utilizing image editing models that accept natural language input, such as the state-of-the-art model SmartEdit [14]. Ideally, such models would provide effective guidance by generating the target goal based on natural language descriptions. However, as shown in Fig. 7, we find that the image editing model is not ideal when it comes to following complex instructions. This underscores the significant challenge of “imagining” the target state in open-world problems, suggesting that a world model capable of accurately providing such visual guidance requires further development.

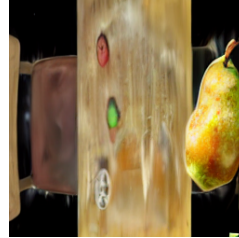
Please also refer to the following figures for additional qualitative experiment results. We present qualitative results on tasks including dense counting, simple jigsaw puzzle solving (please also see the video in project page), and object placement. These images more effectively illustrate the detailed visual reasoning processes and the changes within the imagination space.

B More Experiment Details

B.1 Dense Counting Experiments in Fig. 5a and 5b

The aim of these experiments was to evaluate the impact of difficulty and reasoning steps on performance. To achieve this, we selected a subset of cases from dense counting and supplemented them with simpler cases involving fewer objects. We then re-conducted the experiments on this partial dataset using GPT-4o, alongside a newly introduced textual CoT reasoning baseline, as discussed in the main paper.

In the first experiment (Fig. 5a), we show the success rate of different methods when faced with varying numbers of objects. In the second experiment (Fig. 5b), we test how different methods improve when given different limits on the number of reasoning steps. The result is shown as the solvable rate, which indicates the percentage of cases the method can solve under the current budget. When imposing such a limit, our model is restricted to performing no more than the specified number of steps, with each step defined as a combination of the focus operation and the subsequent operation performed after focusing. For the textual CoT baseline, we explicitly instructed the model on the maximum number of reasoning steps it was allowed to use.



(a) SmartEdit placing pear



(b) SmartEdit placing apple

Figure 7: We show that existing image editing models lack the capability to effectively comprehend and execute complex commands, such as placing the apple/pear at the left/right side on the table.

B.2 Simple Jigsaw Puzzles

The rotations of pieces are not considered in the task. As multiple trials are allowed, we restrict the total number of attempts to 20 and report the finish rate, defined as the proportion of missing pieces successfully placed. The sampling methods used involve comparing pairs of target locations and iteratively selecting the better option until a single target location remains. We guarantee that the correct answer exists within the sampling process. In our method, once a jigsaw piece is moved to the imagined target location, it is considered placed. The GPT-4o baseline selects the target location by outputting coordinates, while the VCoT baseline determines the target location by drawing a block, leveraging its specialized training for such tasks. For baselines that directly output a location, such as GPT-4o and VCoT, we ensure that the jigsaw puzzle piece is provided in the given image. Our method locates and selects a piece independently, as it possesses this capability, while other methods are automatically given a piece.

B.3 Object Placement

Similar to the approach used in solving simple jigsaw puzzles, the methods compared include the following: GPT-4o, which outputs the target location by specifying coordinates; the sampling method, which iteratively compares sampled results until a single option remains; VCoT, which determines placement by drawing blocks; and our method, which outputs the final placement after completing a transformation operation. In alignment with how related works in robotics handle sampling-based methods [18, 8], the sampling baseline restricts its search to a sub-region of the space, filtering out certain incorrect answers. For example, regions occupied by existing objects are excluded from consideration. This principle is also incorporated into the transformation operation in our method. Specifically, a focus operation is first applied to identify the platform where the target object should reside, creating a region mask. The transformation operation then ensures that no movement is suggested if it would lead the object outside the defined mask. This approach enhances efficiency and ensures a fair comparison between the sampling baseline and our method by reducing unnecessary exploration of invalid regions.

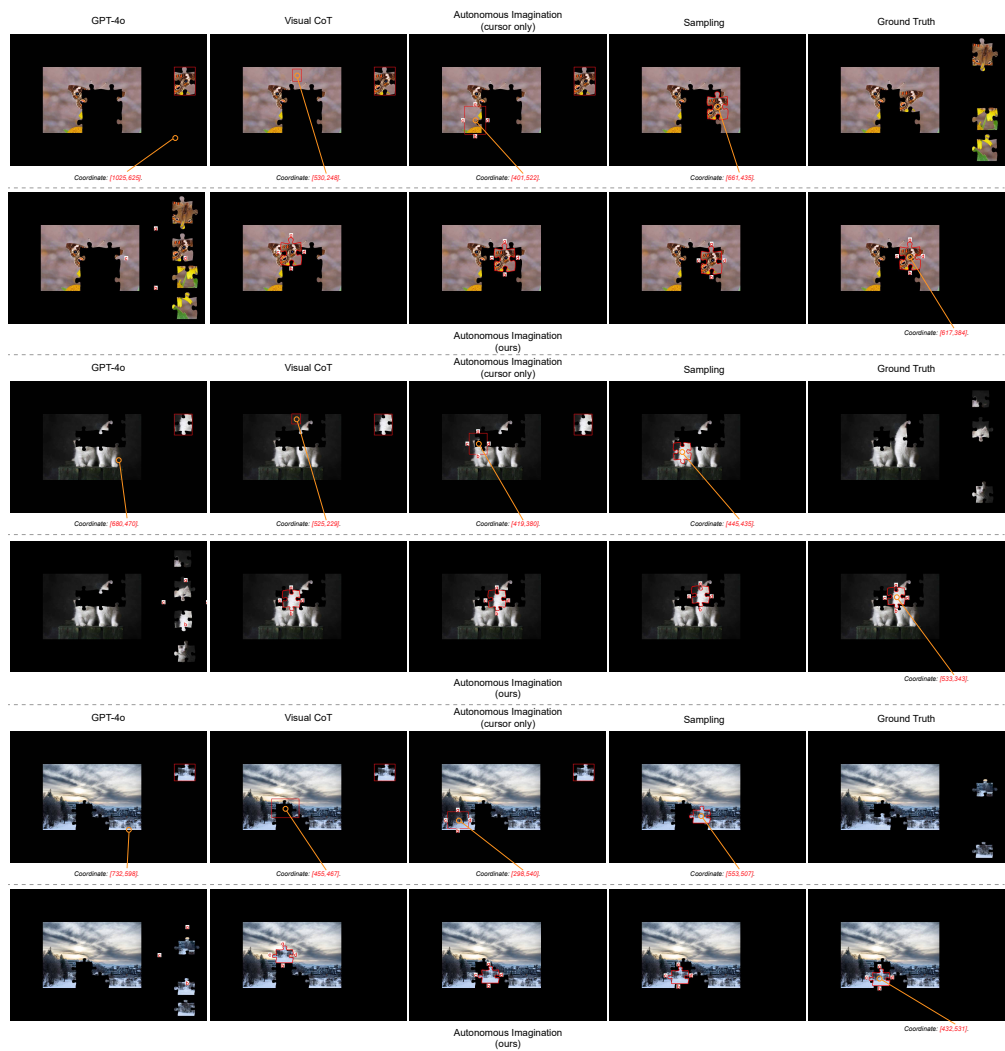


Figure 9: Additional qualitative results of simple puzzle solving. Please zoom in for a clearer view.

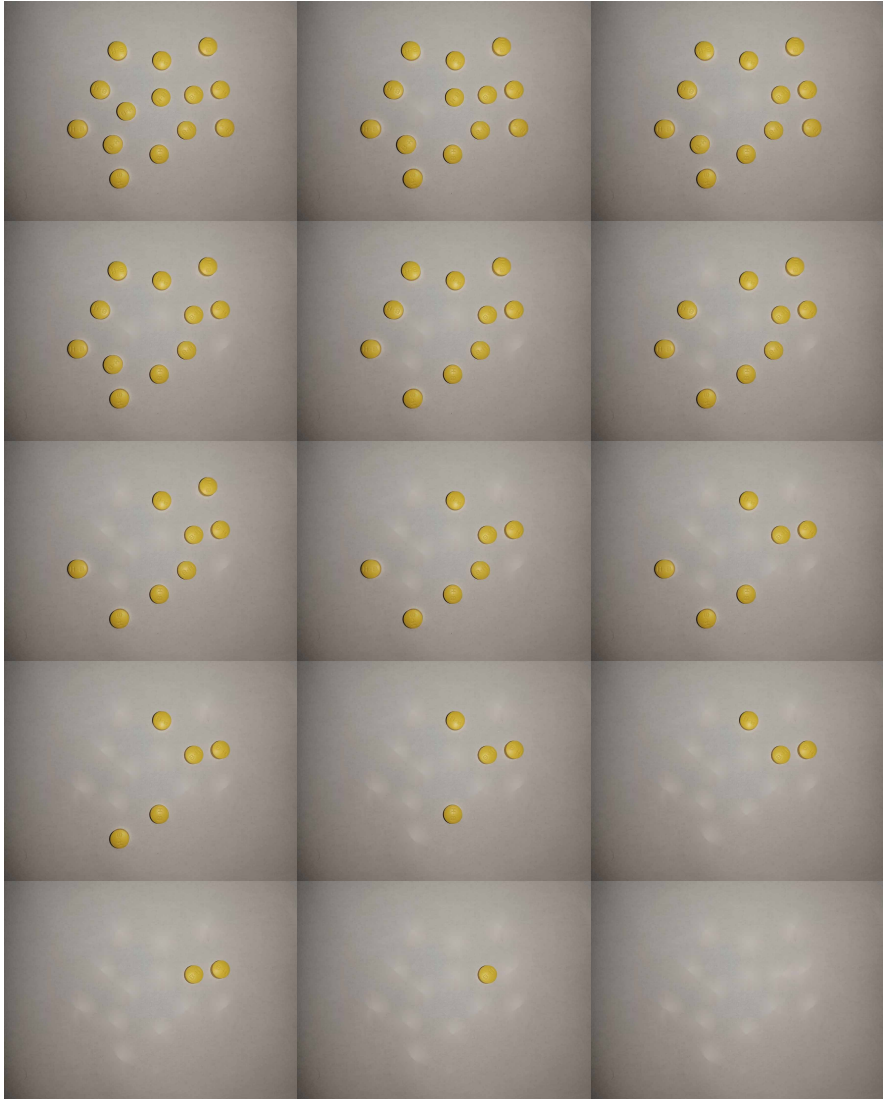


Figure 10: Qualitative demonstration of the counting process autonomously performed by MLLMs.



Figure 11: Qualitative demonstration of the counting process autonomously performed by MLLMs.



Figure 12: Additional qualitative results of object placement are provided. We also illustrate some steps in our method during the search for the target object and its placement. Please zoom in for a clearer view.

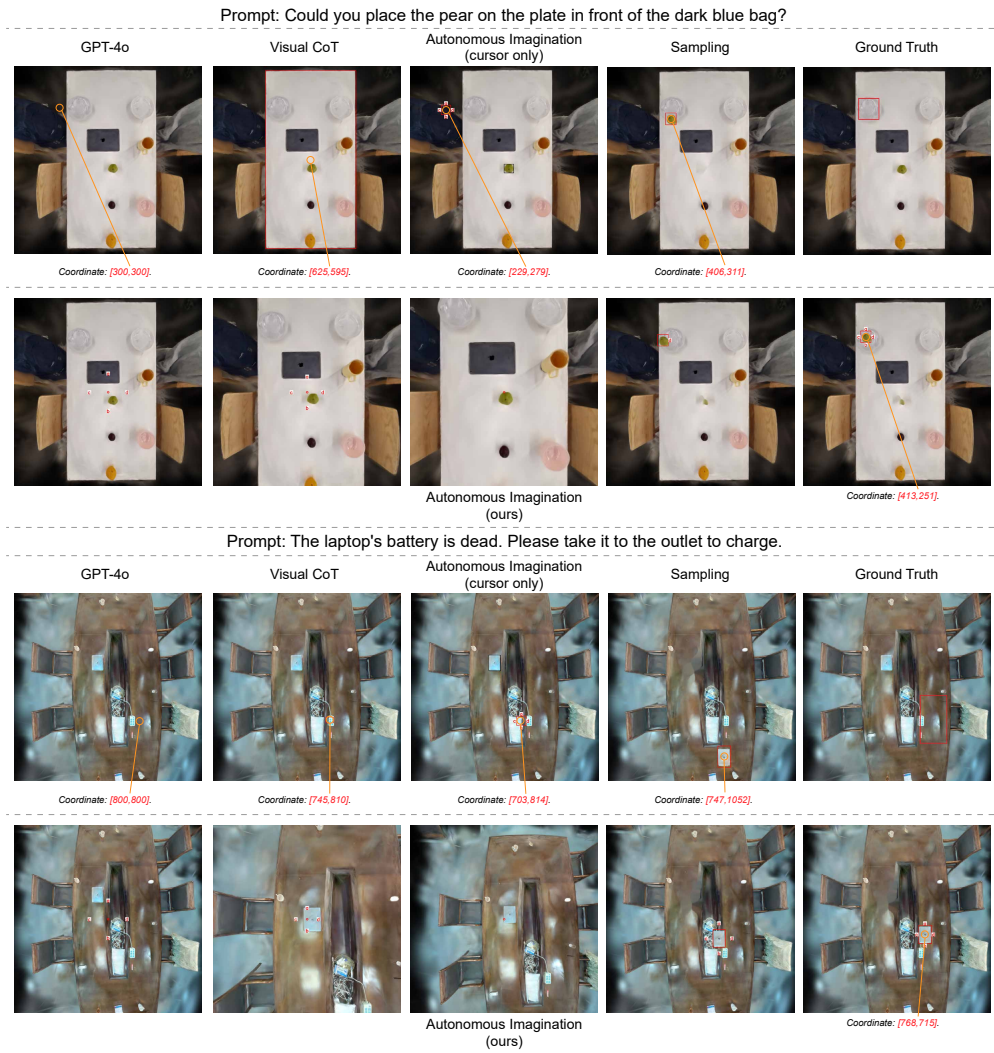


Figure 13: Additional qualitative results of object placement are provided. Please zoom in for a clearer view.