

# VISUAL AGENTS AS FAST AND SLOW THINKERS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Achieving human-level intelligence requires refining cognitive distinctions between *System 1* and *System 2* thinking. While contemporary AI, driven by large language models, demonstrates human-like traits, it falls short of genuine cognition. Transitioning from structured benchmarks to real-world scenarios presents challenges for visual agents, often leading to inaccurate and overly confident responses. To address the challenge, we introduce **FAST**, which incorporates the **F**ast and **S**low Thinking mechanism into visual agents. FAST employs a switch adapter to dynamically select between *System 1/2* modes, tailoring the problem-solving approach to different task complexity. It tackles uncertain and unseen objects by adjusting model confidence and integrating new contextual data. With this novel design, we advocate a *flexible system, hierarchical reasoning* capabilities, and a *transparent decision-making* pipeline, all of which contribute to its ability to emulate human-like cognitive processes in visual intelligence. Empirical results demonstrate that FAST outperforms various well-known baselines, achieving 80.8% accuracy over *VQA<sup>v2</sup>* for visual question answering and 48.7% *GIoU* score over ReasonSeg for reasoning segmentation, demonstrate FAST’s superior performance. Extensive testing validates the efficacy and robustness of FAST’s core components, showcasing its potential to advance the development of cognitive visual agents in AI systems. The code is available at this [link](#).

## 1 INTRODUCTION

In the field of artificial intelligence, *System 2* delineates a cognitive mode distinguished by deliberate, analytical, and consciously reasoned processes (Wei et al., 2022; Wang et al., 2023d; Zelikman et al., 2022; Zhou et al., 2023; Hua & Zhang, 2022). This mode is juxtaposed to *System 1*, which embodies intuitive, automatic, and unconscious cognition. Achieving human-level intelligence in AI systems necessitates the deliberate cultivation and refinement of these cognitive distinctions. This process is crucial for the development of advanced reasoning and decision-making capabilities (Zhang et al., 2023d; Hao et al., 2023).

The emergence of foundation models marks a significant turning point, where Large Language Models (LLMs) based agents have made remarkable strides in many areas, showcasing human-like intelligence across diverse tasks (Brown et al., 2020; Kojima et al., 2022; Ge et al., 2023). However, this achievement is primarily attributed to some features of foundation models: overparameterization and the availability of vast, general-purpose datasets (Kaplan et al., 2020; OpenAI, 2024). It is imperative to note that while these models exhibit human-like traits (*e.g.*, inductive and deductive reasoning (Huang & Chang, 2023b; Dasgupta et al., 2022; Jin et al., 2024)), these characteristics do not equate to the processes of *System 1/2* thinking (Nye et al., 2021; Yao et al., 2023b) and are far less intelligent than human thinking.

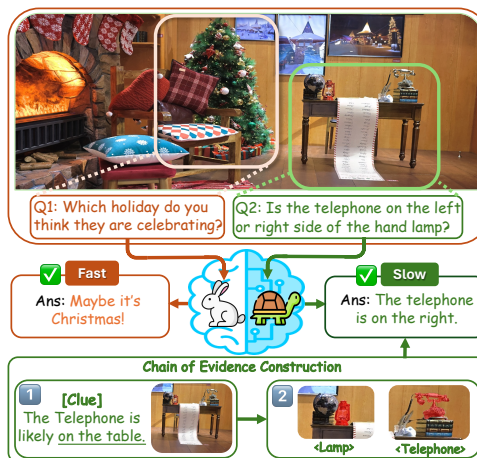


Figure 1. **Working Pipeline.** FAST represents a solution rooted in system switching, demonstrating pronounced capabilities in *hierarchical reasoning* and *ad-hoc explainability*.

In practice, visual agents often encounter challenges when moving from controlled, structured benchmarks to complex, real-world environments (Wu & Xie, 2024; Ge et al., 2023). This problematic circumstance will result in spurious reasoning pathways, akin to hallucinations, where they struggle to acknowledge their limitations or uncertainties (Gunjal et al., 2024; Chen et al., 2023c). Such an issue arises from the absence of explicit modeling of the fast and slow cognitive processes, reminiscent of human *System 1* and *System 2* thinking (Yao et al., 2023b; Kahneman, 2011). Consequently, when faced with intricate inquiries, the Multimodal Large Language Model (MLLM) frequently offers overly confident yet inaccurate responses (Wu & Xie, 2024; Chen et al., 2024b; Tong et al., 2024a). Addressing this problem entails reassessing MLLM algorithms to incorporate insights from the interplay between fast and slow thinking (*System 1* and *System 2*) observed in human cognition. Our design philosophy guides us to incorporate human qualities into our work.

In this study, we introduce the **Fast and Slow Thinking (FAST)** mechanism into visual agents. More concretely, we design a switch adapter to determine whether the encountered problems are best addressed using which thinking mode. Simple tasks require only fast thinking (*System 1*) for a straightforward problem-solving pipeline, while complex tasks necessitate the slow, deliberate processing of *System 2* (see Fig. 1). Specifically, *System 2* is triggered when we encounter visual challenges that have: ① *Uncertainty*: When the model has low confidence in directly identifying the object to which the complex query is referring. For example, the query asks “the appliance for storing and cooling food” instead of “refrigerator;” and ② *Invisibility*: When dealing with minuscule-sized objects that evade detection by standard visual encoders, where normal visual agents cannot tell what it is. This switch adapter is achieved by designing negative contextual data to re-adjust the model’s confidence and ignite world knowledge (as detailed analysis in §3.2). Subsequently, a proposal adapter is engaged to outline regions that are related to the questions. This allows visual agents to leverage the newly acquired data, thereby facilitating a more detailed and precise response. Further, if the inquiry necessitates detailed insights into particular instances, a seg adapter provides segmentation masks, offering additional contextual information for deeper analysis (as detailed analysis in §3.3).

**FAST** enjoys a few attractive qualities. ① **Flexible system**: Building on a foundation that explicitly models *System 1/2 thinking*, our proposed method adeptly handles complex visual tasks, demonstrating competitive performance in a streamlined pipeline (see §2.1). **FAST’s** core epistemology combines an intuitive mechanism for straightforward cases with deliberate analytics for more intricate scenarios, thereby enhancing the development of a human-like visual agent. ② **Hierarchical reasoning**: **FAST** perceives visual tasks with a top-down granularity, encompassing image-level cues, box-level candidates, and pixel-level targets (see Fig. 2). This progressive approach facilitates a sensible understanding of visual content, starting from global concepts, progressing through region-specific candidate assessment, and culminating in precise target identification. Each stage involves developing concrete ideas and establishing a coherent “*chain of evidence*” to support the final inference. ③ **Transparent pipeline**: **FAST’s** decision-making process embodies a neuro-symbolic essence in **System 2 mode**, yielding intermediate step outputs as interpretable symbols (*e.g.*, bounding boxes or masks), facilitating direct visual inspection by humans. This inherent reasoning mechanism enables *ad-hoc explainability* of the model’s behavior (see Fig. 3), distinguishing **FAST** from prior approaches (Liu et al., 2024a) that lack precise explication of their operational mechanisms.

We conducted a series of experiments to validate the efficacy of our proposed method. In §3.1.1, we apply **FAST** to visual question answering and multimodal benchmarks. **FAST** demonstrates significantly improved performance over baselines such as LLaVA-v1.5 (Liu et al., 2024a), achieving performance gains on benchmarks like TextVQA (Singh et al., 2019) with a 2.5% increase in accuracy and a total score improvement of 6.7 on MME (Fu et al., 2024). In §3.1.2, we explore the versatility of our approach through its application to tasks such as referring and reasoning segmentation, with performance gains including an increase of 4.1% *CIoU* with LLaVA-v1.5, and improvements of 3.2% *CIoU* and 2.7% *GIoU* on the ReasonSeg dataset over LISA-7B (Lai et al., 2024). The robustness and effectiveness of the core components of our **FAST** framework are further substantiated through a series of ablation studies, as elaborated in §3.3.

## 2 METHODS

**Notation.** The integration of components in visual agents  $\mathcal{F}$  (based on the Large Language Model) typically involves a visual encoder, denoted as  $\mathcal{E}_V$ , a nature language encoder, represented by  $\mathcal{E}_L$ ,

and a Language Language Model such as Vicuna (Chiang et al., 2023). Initially, the visual agent is presented with an image  $\mathcal{I}$  and an accompanying textual prompt  $\mathcal{Q}$ , which could be a question or instruction. Then the visual agent combines these multimodal tokens into a united space. Finally, the visual agent outputs a textual response  $R$  given the textual and image input. The generation process can be expressed as Eq. 1:

$$R = \mathcal{F} [\mathcal{E}_V(I), \mathcal{E}_L(Q)] \quad (1)$$

**Definition 1** (*System 1 and System 2*) System 1 and 2 are two different systems of thinking proposed by Nobel Laureate Daniel Kahneman in his book *Thinking, Fast and Slow* (Kahneman, 2011).

*System 1 (Fast Thinking): Unconscious, automated thinking processes, fast, intuitive, effortless responsible for automatic responses and basic cognitive operations in daily activities, vulnerable to heuristic biases and errors, e.g., recognizing familiar faces, and knowing the location of objects.*

*System 2 (Slow Thinking): Conscious, energetic thinking processes, slow, effortful, logical, and analytical, responsible for complex calculations, reasoning, and decision-making, can monitor and control System 1 processes, e.g. filling out a tax form, finding the position of a word in a sentence.*

## 2.1 FAST

We present FAST (see Fig. 2), a novel framework designed to efficiently handle both simple and complex visual queries. FAST features a dynamic system switch mechanism that enables rapid responses to straightforward questions (*System 1*) and accommodates deliberate reasoning for intricate scenarios (*System 2*). During slow thinking, the system uses contextual clues to identify a relevant region, facilitated by a proposal adapter. The adapter generates a bounding box around the target object, and if needed, a pixel-level mask adapter refines the proposal for further details. Finally, we summarize the gathered information from the whole system to provide a comprehensive answer.

**System Switch.** Current works on visual agents mostly rely on visual question-answering data, which gives direct answers (*System 1*) after inquiry as Eq. 1. However, attempting to answer questions directly in this way can compromise the reliability of the responses. Agents tend to hallucinate over questions that require more deliberate reasoning and visual details. To reduce hallucination and make the model reliable, we utilize a system switch trigger to tell when to require more visual information. Specifically, for a question  $\mathcal{Q}$  and an image  $\mathcal{I}$ , we define a MLLM with switch adapter  $\mathcal{S}$  and formulate the fast  $\mathcal{F}_{fast}$  and slow thinking process  $\mathcal{F}_{slow}$ . When the query is easy, the frame does not need the switch adapter  $\mathcal{S}_{adapter}$  and only output result  $\mathcal{R}$  by  $\mathcal{F}_{fast}$  as Eq. 2.

$$\mathcal{R} = \mathcal{F}_{fast} [\mathcal{E}_V(I), \mathcal{E}_L(Q)] \quad (2)$$

**Remark 2.1** (*Switching-friendly dataset*) A Negative Data for Target Objects Reasoning Dataset  $\mathcal{D}$  of 100,000 (image, question, answer) triples was constructed to facilitate the identification of target regions or objects required to answer a question. The dataset constructs questions about the absence or details of certain objects, deliberately made too small to be perceived by the visual encoder. [Section A for more details.](#)

**Remark 2.2** (*Switch Adapter*) A light-weight adapter that is fine-tuned with both positive fast-thinking data and negative data ([Remark. 2.1](#)) to acquire system switching capability. When the adapter encounters harder questions, the switch mechanism will be triggered for later slow thinking.

Note that a slow thinking process is not always activated. The system switch adapter as [Remark. 2.2](#)  $\mathcal{S}_{adapter}$  will determine whether the question for the particular image is sufficient to give a direct answer. If so, the fast mode  $\mathcal{F}_{fast}$  will give a quick and direct response as [Equation 2](#). If there is any missing information about the question that current agent cannot solve, the switch adapter will be activated and find the pattern to elicit all the possible missing objects  $\mathcal{O}_{missing}$  related to question and context clues  $\mathcal{C}_{clue}$  which is the possible location of the missing objects as [Equation 3](#).

$$\mathcal{O}_{missing}, \mathcal{C}_{clue} = \mathcal{S}_{adapter} [\mathcal{E}_V(I), \mathcal{E}_L(Q)] \quad (3)$$

Specifically, we use negative data that contain missing objects  $\mathcal{O}_{missing}$  and context clues  $\mathcal{C}_{clue}$  for training the system switch adapter for triggering the slow mode as [Remark. 2.2](#). The slow mode

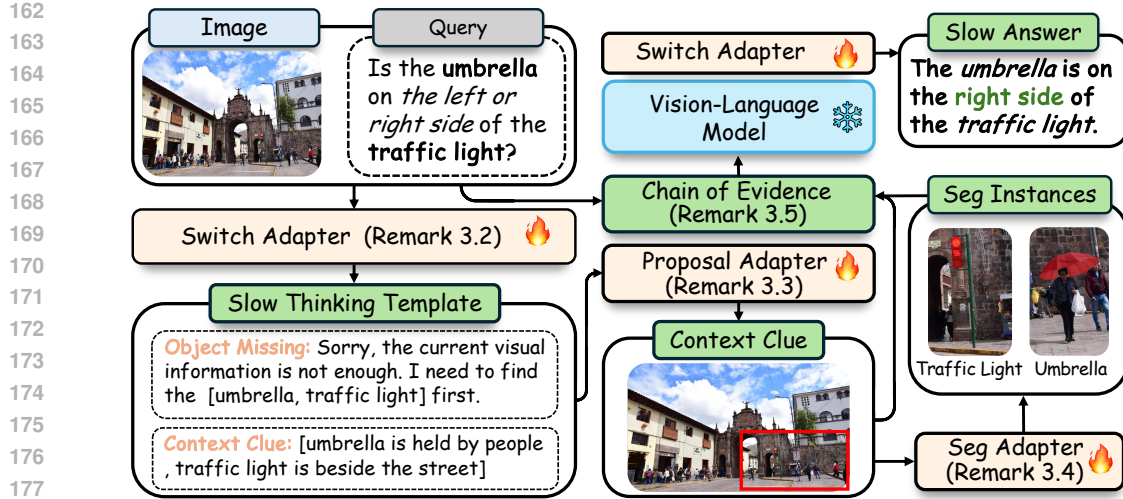


Figure 2. **Slow Thinking Mode of FAST.** Our *slow thinking mode* comprises three core modules: *Switch Adapter*, which selectively activates a slow and analytical thinking mode when encountering complex visual queries, supplementing with extensive world knowledge to provide missing objects and contextual clues; *Proposal Adapter*, which identifies and emphasizes regions of interest within the visual inputs; *Seg Adapter*, which delivers precise pixel-level segmentation, enhancing depth of the visual analysis. The outputs from each module are integrated into a *chain of evidence* (see Fig. 3), providing a methodical and accurate response. FAST represents a neural-symbolic approach that combines the strengths of symbolic reasoning, ensuring that our system is effective and interpretable.

should deal with question-image pairs that are 1) uncertain in pinpointing the specific object in question, and 2) too small to perceive for the standard visual encoder. So we utilize triplet data in dataset Remark. 2.1 (image, question, answer) as where the question requires objects that are not in the image or too small to be perceived by the visual encoder. The threshold is set to be  $20 \times 20$ . We require the model to tell that certain objects are missing instead of a direct answer, and we utilize the world knowledge to list all the objects and also the context information for later deliberate reasoning.

## 2.2 PRELIMINARY

**Hierarchical Reasoning.** We use a top-down scheme to reason over multi-scale granularity images effectively in order to reason and take advantage of world knowledge progressively. Similar to humans would look for some context clue to find specific objects relating to questions and zoom in if they think the answer lies in a particular region, we model this process with system switch adapters as Eq. 2.2 to focus on the context clue  $\mathcal{C}_{clue}$  generated from the switch adapter as Eq. 3.

We denote the MLLM as a proposal adapter  $\mathcal{P}_{adapter}$  (visual agent). In *System 2*, **FaST** uses many visual agents to accomplish hierarchical reasoning. The frame tries to narrow down the search space by using the question  $Q$  and the previously obtained clue  $\mathcal{C}_{clue}$  to let the proposal adapter output a region *Region* that aligns with the question and the context clue as Eq. 4.

$$Region = \mathcal{P}_{adapter} [\mathcal{E}_V(I), \mathcal{E}_L(Q), \mathcal{C}_{clue}] \quad (4)$$

After getting the region, the visual agents  $\mathcal{P}_{adapter}$  will be asked to focus on a more specific target with a bounding box [*Bboxes*] complemented by the context clue  $\mathcal{C}_{clue}$  and region *Region* get from Eq. 4. This process can reveal the step-by-step reasoning and be modeled as Eq. 5.

$$[Bboxes] = \mathcal{P}_{adapter} [\mathcal{E}_L(Q), Region, \mathcal{C}_{clue}] \quad (5)$$

**Remark 2.3 (Proposal Adapter)** A lightweight adapter that is fine-tuned with proposal data to acquire the capability of finding the corresponding region given the context clue or object name.

**Remark 2.4 (Pixel-level mask decoder)** The Pixel-level mask decoder is the decoder of segment anything (SAM (Kirillov et al., 2023)). The pixel-level mask decoder is fine-tuned to produce target masks based on the hidden embeddings.

When we have a more specific target proposal (bounding box  $[Bboxes]$ ), FAST will apply a fine-grained pixel-level mask decoder  $\mathcal{P}_{seg}$  as Eq. 6 to output the specific mask part  $[Mask]$  of the target proposal  $[Bboxes]$  to focus on as Eq. 6. We name this whole process from *Region* to  $[Mask]$  *chain of evidence* as Remark. 2.5 similar to thinking more and more deeply by humans.

$$[Mask] = \mathcal{P}_{seg} [\mathcal{E}_L(Q), [Bboxes], \mathcal{O}_{missing}] \tag{6}$$

**Remark 2.5 (Chain of Evidence)** *Chain of evidence is like the chain of thought in a Large Language Model. But we define it as a deeper and deeper step of thinking based on correct evidence in our frame FAST. The completion of the chain of evidence needs many visual agents to work together.*

After getting the target proposal (bounding box  $[Bboxes]$ ) from context clue  $\mathcal{C}_{clue}$  with proposal adapter and specific mask part  $[Mask]$  by missing objects with seg adapter, a *chain of evidence* is constructed as Remark 2.5 and Fig. 3. Our FAST framework then summarizes all this information ( $\mathcal{I}$  and  $\mathcal{Q}$ ) and the *chain of evidence* with switch adapter to give the final correct reasoning answer  $Ans$  as Eq. 7

$$Ans = \mathcal{F}_{Slow} [\mathcal{E}_L(Q), \mathcal{E}_V(I), [Bboxes/Mask]] \tag{7}$$

The decision-making process in FAST is distinguished by its neuro-symbolic nature, which generates intermediate outputs as easily interpretable symbols, including region-of-interest (RoI) driven boxes and object-driven masks. This capability allows humans to perform direct visual inspections, thereby augmenting the transparency of the model’s operations. Moreover, the intrinsic reasoning mechanism of FAST enhances the ad-hoc explainability of its behavior, see Fig. 2.

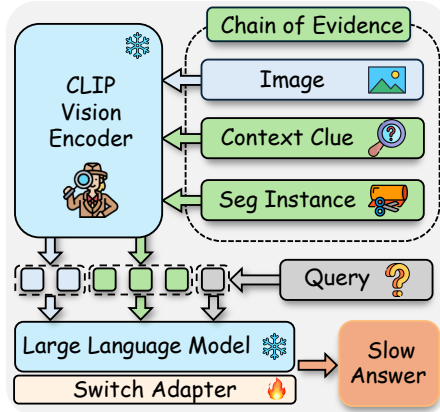


Figure 3. **Chain of Evidence.** FAST represents a solution rooted in switching, demonstrating pronounced capabilities in *hierarchical reasoning* and *ad-hoc explainability*.

### 2.3 IMPLEMENTATION DETAILS

The framework of FAST(as Fig. 2)’s implementation details are shown in this section below.

- *Visual Agents.* We choose the architecture and configuration of LLaVA-v1.5 (Liu et al., 2024a) as our visual agent. The most important component in a visual agent is the visual encoder  $\mathcal{E}_V(I)$ : A CLIP-ViT-L-336px model (Radford et al., 2021) is used, where input images are resized or padded to  $336 * 336$  pixels, learning to associate visual features with corresponding textual descriptions. An MLP projection with channels of [256, 4096, 4096] is used for connecting image representations into the word embedding space.

- *Mask Decoder.* The mask decoder  $\mathcal{P}_{seg}$  architecture is identical to SAM. Besides, it is fully fine-tuned with a collection of semantic segmentation (Caesar et al., 2018; Zhou et al., 2017; Ramanathan et al., 2023; He et al., 2022; Chen et al., 2014) and referring segmentation (Mao et al., 2016; Kazemzadeh et al., 2014) datasets to efficiently map the  $\langle seg \rangle$ -token representations to a mask if the FAST need to segment.

- *Chain of Evidence.* When we apply the *chain of evidence* as Remark 2.5 in the LLM to get the answer as the final step like Eq. 7. The whole sequence of the *chain of evidence* is too long to load in the  $\mathcal{F}_{Slow}$ . So FAST needs a visual sampler based on cross-attention that is trained to decrease the number of image tokens to a suitable length (from 256 to 32), apart from MLP projection.

## 3 EXPERIMENT

We utilize eight popular benchmarks to evaluate our framework FAST comprehensively, categorized into general visual question answering (VQA) datasets and multimodal benchmarks. The VQA benchmarks include VQA-v2 (Goyal et al., 2017), GQA (Hudson & Manning, 2019), ScienceQA (Lu et al., 2022), and TextVQA (Singh et al., 2019) which focus on optical character recognition. For multimodal benchmarks evaluation, we use the hallucination benchmark POPE (Li et al.,

2023c), along with comprehensive benchmarks such as MME (Fu et al., 2024), MM-Vet (Yu et al., 2024), and SEED (Li et al., 2024). We compare our model with the baseline LLaVA-v1.5 (Liu et al., 2023a), and other multimodal large language models. To thoroughly assess our model’s understanding of pixel-level instances, we evaluate its performance on referring segmentation and grounding benchmarks, including refCOCO (Kazemzadeh et al., 2014), refCOCO+ (Kazemzadeh et al., 2014), and refCOCOg (Caesar et al., 2018). Further, to examine the model’s reasoning capabilities on FAST framework, we consider the Reasoning Segmentation benchmark (Lai et al., 2024).

### 3.1 MAIN RESULTS

#### 3.1.1 EXPERIMENTS ON VQA AND MULTIMODAL BENCHMARKS

**Training.** In developing the Switch Adapter, we employed the LLaVA-v1.5 (Liu et al., 2024a) framework, conforming strictly to its established training protocols. We incorporated negative samples from  $V^*$  (Wu & Xie, 2024) with contextual cues to enhance system switching capability to amplify multimodal inferential and world knowledge. This augmented dataset was combined with LLaVA-v1.5’s supervised dataset and trained for one epoch. For the Proposal Adapter, we augmented the LLaVA-v1.5 dataset with region-specific bounding boxes based on contextual cues and queries, then fine-tuned for one epoch to optimize proposal generation. The Segmentation Adapter utilized the LISA (Lai et al., 2024) architecture integrated with the LLaVA-v1.5, employing SAM as the mask decoder. The adapter was fine-tuned using the same datasets as Lisa, including semantic segmentation, referring segmentation, and reasoning segmentation. This fine-tuning process involved 10,000 steps to improve the model’s segmentation capabilities. Throughout developing the Switch Adapter, Proposal Adapter, and Segmentation Adapter, we employed the LoRA (Low-Rank Adaptation) technique (Hu et al., 2022). By leveraging LoRA, we introduce minimal additional parameters while preserving the original multimodal large language model’s architecture and efficiency. All experiments used 8 NVIDIA TESLA A100-80GB GPUs.

**Metric.** In model evaluation across diverse datasets, various performance metrics are utilized.

*Accuracy.* The primary evaluation metric utilized in the  $VQA^{v2}$ , GQA, TextVQA, ScienceQA, and SEED benchmarks is accuracy. Accuracy is a performance measure that quantifies the exact match percentage between predicted and acceptable ground truth answers, indicating a model’s precision.

*F1 Score.* The POPE dataset uses the F1 Score to balance precision and recall, providing a comprehensive assessment by harmonizing the trade-off between positive prediction accuracy and recall.

*Total Score.* The MME evaluation metrics include accuracy (based on individual questions) and accuracy+ (considering both questions per image), reflecting a stricter and more comprehensive model understanding. Random accuracies for these metrics are 50% and 25%, respectively. Perception scores, calculated as the sum of these metrics across subtasks, total 2000 for perception.

*GPT-Evaluation.* In the MM-Vet dataset, performance is evaluated by GPT-4 through a comparative analysis of predicted and ground truth answers, generating a score to quantify alignment.

**Results.** As depicted in Table 1, FAST demonstrates superior performance across multiple VQA datasets and multimodal benchmarks when compared to established methods. To ensure fairness in comparison, all methods in Table 1 share the same visual encoder: basic CLIP (Radford et al., 2021). Remarkably, FAST consistently surpasses the LLaVA-v1.5 model, achieving significant improvements in performance across all evaluated datasets. Specifically, in VQA datasets, our model outperforms LLaVA-v1.5 by 2.3% in  $VQA^{v2}$ , 1.8% in GQA, and 2.5% in  $VQA^T$ . Additionally, FAST excels in multimodal benchmarks, with notable increases of 6.7 in the MME score, 1.5 in the SEED score, and 0.5 in the MM-Vet score, highlighting its versatility and effectiveness in handling a broad range of domains. These results underscore

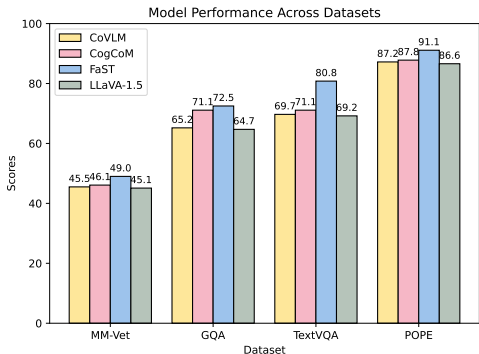


Figure 4. **The Comparison with CoVLM and CogCoM.** These models use the more powerful vision encoder.

324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377

Method	LLM	VQA Datasets				Multimodal Benchmarks			
		$VQA^{v2}$	GQA	$VQA^T$	$SQA^I$	POPE	MME	SEED	MM-Vet
BLIP-2 <sub>[ICML23]</sub>	Vicuna-13B	65.0	32.3	42.5	61.0	85.3	1293.8	46.4	22.4
InstructBLIP <sub>[NeurIPS24]</sub>	Vicuna-13B	-	49.5	50.7	63.1	78.9	1212.8	53.4	25.6
Qwen-VL-Chat <sub>[arXiv23]</sub>	Qwen-7B	78.2	57.5	61.5	68.2	-	1487.5	58.2	-
mPLUG-Owl2 <sub>[CVPR24]</sub>	LLaMA-7B	79.4	56.1	58.2	68.7	-	1450.2	<b>61.6</b>	<b>36.2</b>
Monkey <sub>[CVPR24]</sub>	Qwen-7B	80.3	60.7	-	<b>69.4</b>	67.6	-	-	-
LLaVA-v1.5 <sub>[CVPR24]</sub>	Vicuna-7B	78.5	62.0	58.2	66.8	85.9	<u>1510.7</u>	58.6	30.5
Chain of Spot <sub>[arXiv24]</sub>	Vicuna-7B	<u>80.7</u>	<u>63.7</u>	<u>60.9</u>	68.2	<u>86.4</u>	1501.1	59.7	30.8
V* <sub>[CVPR24]</sub>	Vicuna-7B	-	-	-	-	82.4	1128.9	41.7	27.7
Visual CoT <sub>[arXiv24]</sub>	Vicuna-7B	-	63.1	<b>77.5</b>	-	-	-	-	-
FAST (Ours)	Vicuna-7B	<b>80.8</b>	<b>63.8</b>	60.7	<u>68.9</u>	<b>86.4</b>	<b>1517.4</b>	<u>60.1</u>	<u>31.0</u>
$\Delta$ (vs LLaVA-v1.5)	Vicuna-7B	+2.3	+1.8	+2.5	+2.1	+0.4	+6.7	+1.5	+0.5

Table 1. **Main results on eight VQA and multimodal benchmarks.** Our FAST consistently outperforms the baseline LLaVA1.5 model across all evaluated benchmarks, denoted with line  $\Delta$ .

complex visual and textual tasks. Moreover, Fig. 4 showcases a direct comparison between FAST, CoVLM (Wang et al., 2023b), and CogCoM (Qi et al., 2024), both of which employ the more powerful EVA2-CLIP-E (Sun et al., 2023) model as their visual encoder. As expected, these models exhibit stronger performance due to their enhanced encoder. To align with this, we replaced our original vision encoder with EVA2-CLIP-E, which resulted in further improved performance, ensuring a more rigorous and fair comparison with state-of-the-art methods. This two-tiered comparison—first with basic CLIP and then with the more advanced EVA2-CLIP-E—provides a balanced and comprehensive evaluation of FAST against leading approaches, reinforcing its effectiveness in diverse and challenging tasks.

### 3.1.2 EXPERIMENTS ON REFERRING AND REASONING SEGMENTATION

**Training.** The training settings for the Switch Adapter and Proposal Adapter remain consistent with those previously described as §3.1.1. During the training phase of the Segmentation Adapter, certain specific datasets are intentionally omitted to uphold an unbiased evaluation of referring and reasoning segmentation datasets. This strategic exclusion is a crucial measure implemented to prevent any potential data leakage, thereby ensuring the integrity and reliability of the evaluation results.

**Metric.** Following prior research on segmentation (Kazemzadeh et al., 2014; Mao et al., 2016), two evaluation metrics are employed: Generalized Intersection over Union ( $GIoU$ ) and complete Intersection over Union ( $CIoU$ ).

$CIoU$ . The  $CIoU$  is calculated based on the cumulative intersection over the cumulative union across all images in the dataset. This approach can introduce a significant bias towards larger objects or images with more objects, as they contribute more to the cumulative union area.

$GIoU$ . The  $GIoU$  is computed as the average per image  $IoU$ , where the  $IoU$  is calculated for each image, and then the average is taken across all images in the dataset. This metric provides a balanced assessment by treating all images equally, regardless of their size or the number of objects.

**Results.** Table 2 illustrates the performance of FAST compared to recent visual agents like LISA on referring and reasoning segmentation benchmarks. FAST notably outperforms LISA-7B on the refCOCO+ and refCOCOg benchmarks by 2.0% and 0.6%  $CIoU$ , respectively. For the more complex reasoning segmentation task, FAST shows even stronger results, with a 3.2%  $CIoU$  gain and a

Method	Referring Segmentation			Reasoning Segmentation	
	refCOCO	refCOCO+	refCOCOg	ReasoSeg	
	$CIoU$	$CIoU$	$CIoU$	$CIoU$	$GIoU$
LAVT <sub>[CVPR22]</sub>	72.7	62.1	61.2	-	-
OVSeg <sub>[CVPR23]</sub>	-	-	-	28.5	18.6
GRES <sub>[CVPR23]</sub>	<u>73.8</u>	<b>66.0</b>	65.0	22.4	19.9
X-Decoder <sub>[CVPR23]</sub>	-	-	64.6	22.6	17.9
SEEM <sub>[NeurIPS24]</sub>	-	-	65.7	25.5	21.2
LISA-7B <sub>[CVPR24]</sub>	<b>74.1</b>	62.4	<u>66.4</u>	<u>44.4</u>	<u>46.0</u>
LLaVA w Seg Adapter	70.8	57.5	64.0	43.0	41.0
FAST (Ours)	73.3	<u>64.4</u>	<b>67.0</b>	<b>47.6</b>	<b>48.7</b>

Table 2. **Main results on referring and reasoning segmentation benchmarks.** Our FAST exhibits competitive results in referring segmentation tasks like refCOCOg+ while showcasing superior performance in reasoning segmentation, particularly when evaluated against LISA-7B.

2.7% *GIoU* improvement over LISA. The results highlight FAST’s superior performance and its robustness in handling both straightforward and complex visual reasoning segmentation benchmarks.

### 3.2 ANALYSIS OF SYSTEM SWITCHING ADAPTER

Our study investigates the efficacy of the switch adapter mechanism in balancing accuracy and computational efficiency. As depicted in Fig.5, our analysis illustrates the system’s adeptness in discerning between the *System 1* and *System 2* cognitive modes triggered by query complexity. For queries requiring *System 2* reasoning, the adapter dynamically combines *System 1* reasoning for simpler subcomponents with *System 2* reasoning for more complex aspects. Consequently, the reported accuracy rates under *System 2* mode (52.2% for MME and 56.8% for GQA) reflect a combination of reasoning outcomes, emphasizing the adapter’s ability to differentiate query complexities and optimize task performance accordingly. This highlights the importance of maintaining *System 1* reasoning for prompt and confident responses while effectively utilizing *System 2* reasoning for complex problem-solving.

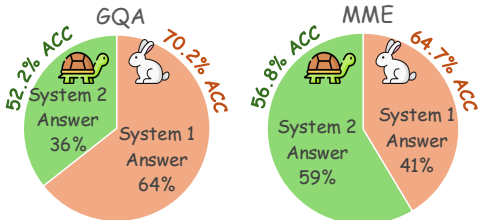


Figure 5. **System 1 Mode Analysis.** We investigate the system switching ratio, along with fast thinking performance on easy or hard queries defined by the switch adapter.

Table 3 compares runtime across system configurations. *System 1 Only*, using a switch adapter, operates efficiently with one-time inference, while *System 2 Only*, which constructs a *chain of evidence* for every query, is significantly more resource-intensive. In contrast, FAST balances efficiency and performance, running 31% faster than *System 2 Only* on MME and 50% faster on GQA, with comparable results. This highlights FAST’s ability to optimize cognitive task processing while conserving computational resources.

Method	MME		GQA	
	Runtime	Result	Runtime	Result
<i>System 1 Only</i>	734ms	1508.7	737ms	61.9
<i>System 2 Only</i>	2938ms	1518.6	2937ms	64.0
<b>OURS</b>	2023ms	1517.4	1475ms	63.8

Table 3. **Runtime Analysis and Comparison** on only *System 1* (fast), our FAST and only *System 2* (slow).

### 3.3 ABLATION STUDY

Algorithm Component	GQA	POPE	MME
BASELINE	62.1	85.7	1509.2
+ Proposal Adapter	63.2	86.0	1516.5
+ Seg Adpater	62.8	85.8	1514.4
<b>OURS (both)</b>	63.8	86.2	1517.4

Table 4. **Key Component Analysis**

Output Component	MME	refCOCOg
BASELINE*	1511.8	66.0
+ Missing Objects	1513.4	66.8
+ Context Clue	1516.6	66.4
<b>OURS (both)</b>	1517.4	67.0

Table 5. **Switch Adapter Output Analysis**

**Key Component Analysis.** We undertake a detailed investigation into the core elements of our novel framework, FAST, with particular emphasis on the proposal adapter for contextual region localization and the seg adapter for pixel-level mask segmentation. To establish a comparative baseline, we design a model configuration that excludes both the proposal and seg adapters, instead relying solely on a switch adapter to provide missing objects and context clues. This baseline model serves as the foundation for evaluating the impact of the individual and combined components of the framework. As demonstrated in Table 4, the introduction of the proposal adapter, the seg adapter, or both, results in progressive and substantial improvements in performance across various evaluation metrics. For instance, accuracy on the *VQA<sup>v2</sup>* dataset improves from 62.1% to 63.8%, showcasing the considerable value these components add. This underscores the pivotal roles of the proposal and seg adapters in enhancing the model’s overall capability, further affirming their importance within the FAST framework.

Further, we evaluate the switch adapter’s role in incorporating missing objects and context clues using a variant BASELINE\*, which omits these features. Table 5 shows that adding missing objects or context clues improves metrics like MME and refCOCOg, with the best performance achieved



when both are included. These results confirm the importance of all components in optimizing FAST’s effectiveness.

### 3.4 QUALITATIVE COMPARISONS OF FAST

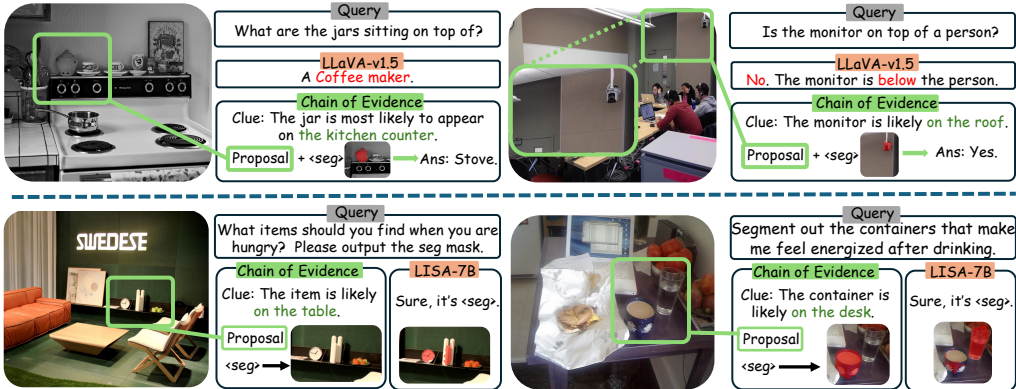


Figure 6. **Qualitative Comparisons of FAST.** The top row shows the VQA results on FAST compared to LLaVA-v1.5. The bottom row presents the segmentation results compared to LISA-7B.

In Fig. 6, we present qualitative comparisons that highlight the enhancements introduced by FAST. The top row, above the dotted line, shows results from the VQA task, comparing FAST with LLaVA-v1.5. LLaVA-v1.5 often fails to focus on key areas within the image, leading to incorrect or incomplete responses. In contrast, FAST builds a *chain of evidence* by identifying key objects and elements (e.g., detecting a woman on the street or a monitor on the roof) and then applying object-level pixel masks to accurately determine the focus areas. This enables FAST to provide more precise and deliberate answers. The bottom row, below the dotted line, shows segmentation results, comparing FAST with LISA-7B. LISA-7B struggles with segmenting smaller objects or those requiring more complex reasoning, often causing confusion. In contrast, FAST excels at isolating relevant objects with greater accuracy and granularity, particularly with smaller or less obvious items. This demonstrates FAST’s superior performance in both VQA and segmentation tasks, showcasing its ability to handle a wide range of visual and reasoning challenges more effectively than its counterparts.

## 4 RELATED WORKS

**LLM as Visual Agents.** With the capabilities that LLMs have demonstrated in language understanding and generation (Ouyang et al., 2022; OpenAI, 2024; Zheng et al., 2023; Touvron et al., 2023a;b; Wang et al., 2024a), the research community has progressed to explore how LLMs can be enhanced with vision input for multimodal tasks as visual agents (Alayrac et al., 2022; Driess et al., 2023; Li et al., 2023a; Ge et al., 2023; Dai et al., 2024; Liu et al., 2023b; 2024a). There are two paradigms for LLM-based visual agents: end-to-end based and tool-using visual agents. Following the principle of instruction tuning, end-to-end visual agents are trained with a curated visual instruction tuning dataset to digest features from multi-modality, unlocking the capability to answer visual questions (Huang et al., 2023; Luo et al., 2023; Zhu et al., 2023a; Bai et al., 2023; Zhang et al., 2023b;c; Chen et al., 2024a; Ye et al., 2023; Singh et al., 2019). For other visual tasks (e.g., Segmentation, Detection, etc), end-to-end trained tailored agents can further perform downstream tasks (Pi et al., 2023; Peng et al., 2024; Lai et al., 2024; Chen et al., 2023b; Wang et al., 2023c; Dai et al., 2024; Wang et al., 2024b; 2023b; Jiang et al., 2023; Chen et al., 2023a). Recent research has focused on leveraging improved vision encoders and fostering more detailed visual understanding, yielding promising results (Fan et al., 2024; Xu et al., 2024a; Shi et al., 2024). While these approaches can be implemented with direct instruction tuning data, they represent a ‘*System 1*’ type of training. This type of training primarily relies on the dataset’s quality and tends to provide direct answers that are prone to hallucinations, a consequence inherent to the nature of *System 1* instruction tuning data. For the second paradigm, tool-using models are built on top of a frozen LLM with access to pretrained visual perception tools (Surís et al., 2023; Shen et al., 2023; Lu et al., 2023a). In this scenario, the LLM first selects visual tools and then decides by thoroughly analyzing the fine-grained information extracted by visual tools (Lu et al., 2023a; You et al., 2023; Wu et al., 2024).

486 While external visual tools enhance the interpretability of the reasoning process, their complexity  
487 can introduce inaccuracies. Moreover, the abundance of information generated during reasoning  
488 may overshadow key details relevant to the query, resulting in incorrect answers.

489 Our research introduces a novel and adaptable framework designed to enhance response accuracy  
490 by adopting distinct slow thinking cognitive modes. Unlike traditional end-to-end visual agents,  
491 our framework, FAST, systematically assesses information sufficiency, thereby mitigating the risk  
492 of overconfidence. When *System 2* (slow, analytical thinking) is activated, FAST employs multiple  
493 experts to construct a coherent *chain of evidence*. This approach ensures the generation of accurate  
494 and interpretable responses, significantly advancing the reliability and transparency of visual agents.

495 **System 2 in AI.** Recently, LLMs have been engineered to produce text that mimics the step-by-step  
496 reasoning process characteristic of human cognition, akin to the analytical and deliberate thought  
497 processes associated with what is termed as *System 2* in the human cognition process (Qiao et al.,  
498 2023; Huang & Chang, 2023a; Wang et al., 2023a; Shaikh et al., 2023; Shao et al., 2024). The system-  
499 atic approach to problem-solving is a hallmark across various domains, including mathematical  
500 word problems (Kojima et al., 2022; Wang et al., 2023d; Lightman et al., 2023; Cobbe et al., 2021;  
501 Liu et al., 2023c; Zhu et al., 2023b; Lu et al., 2023b), logical reasoning (Yao et al., 2023d;a; Besta  
502 et al., 2024; Wen et al., 2023; Lei et al., 2023; Cheng et al., 2024; Jin et al., 2024), and multi-modal  
503 reasoning (Chen et al., 2024d; You et al., 2023; Wu & Xie, 2024). In Explainable AI, this system-  
504 atic method is emulated by the model as it generates a text-based elucidation of its reasoning and  
505 decision-making process through step by step reasoning process (e.g., chain of thought) (Han et al.,  
506 2024; Zhao et al., 2024; Jacovi & Goldberg, 2020; Hua & Zhang, 2022). However, it is crucial to  
507 recognize that LLMs, while powerful, are not exempt from encountering challenges when facing  
508 complex problems. One such challenge is the issue of hallucination (Zhou et al., 2024; Cui et al.,  
509 2023; Li et al., 2023b; Zhang et al., 2023a; Chen et al., 2024c; Guan et al., 2024a), which can distort  
510 the model’s reasoning process and lead to inaccuracies in the explanations provided. Initially,  
511 LLM reasoning is seen as only a linear chain of thoughts, where each step in the reasoning process  
512 is clearly articulated. As models evolve, they adopt more complex structures like hierarchical  
513 trees (Geng et al., 2023; Yao et al., 2023b) and intricate graphs (Besta et al., 2024), which enable  
514 them to handle much more complex problems but also restrict their general applicability because  
515 of increased topological complexity (Yao et al., 2023b;a; Besta et al., 2024; Lei et al., 2023; Wen  
516 et al., 2023). Moreover, these complex structures can lead to errors propagating through the model’s  
517 reasoning, causing a cascade of mistakes (Xu et al., 2024b). To counter this, incorporating feedback  
518 from intermediate reasoning steps and employing iterative refinement, which is similar to human reflection,  
519 could help mitigate errors (Chu et al., 2023; Tong et al., 2024b; Guan et al., 2024b; Madaan  
520 et al., 2023; Yuan et al., 2024; Wu & Xie, 2024). In unsupervised scenarios, such feedback is vital  
521 for enhancing the reasoning capabilities of LLMs and reducing errors (Yao et al., 2023c).

522 Our key contribution is the introduction of the *chain of evidence* within multimodal reasoning frame-  
523 works. This methodology enriches each reasoning step with accurate, image-based cascading in-  
524 formation, effectively mirroring human visual and *System 2* cognitive processes. Our approach  
525 enhances accuracy and significantly improves interpretability and generalization capabilities.

## 526 5 CONCLUSION

527  
528  
529  
530  
531 In this study, we introduced FAST, a framework that combines *System 1* (which is fast and intu-  
532 itive) and *System 2* (which is slow and deliberate) thinking to improve visual agents’ reasoning and  
533 decision-making. FAST adapts to queries of varying complexity with a flexible system switch, deliver-  
534 ing quick responses for simple tasks and using hierarchical reasoning for more complex scenarios.  
535 The FAST leverages neuro-symbolic decision-making transparent pipeline delivering interpretable  
536 intermediate outputs that enable explainability. Our results show significant improvements across  
537 benchmarks, demonstrating the effectiveness of FAST’s *chain of evidence* in reducing hallucina-  
538 tions and improving interpretability. Furthermore, ablation studies highlight the critical importance  
539 of contextual clues, symbolic reasoning, and pixel-level adapters in refining visual reasoning and  
understanding, marking a step forward in creating more reliable and accurate AI cognition.

## 6 ETHICAL SAFEGUARDS

In our paper introducing a novel framework FAST, we implement rigorous ethical measures to prevent potential misuse and promote responsible application. These measures are delineated in comprehensive protocols accompanying the final release of models and datasets. Our protocols encompass stringent usage guidelines, access controls, incorporation of safety filters, and monitoring systems. These concerted efforts reflect our steadfast dedication to upholding the utmost ethical standards in scientific exploration. Our objective is to protect the rights and privacy of all stakeholders involved, thereby fostering a culture of responsible and ethical research within our community.

## 7 REPRODUCIBILITY

Our FAST framework is implemented in PyTorch (Paszke et al., 2019). All the experiments are conducted on eight NVIDIA A100-80GB GPUs. Our full implementation shall be publicly released upon paper acceptance to guarantee reproducibility. The codes are available at the anonymous link <https://anonymous.4open.science/r/Sys2-LLaVA-8B0F/> for the review process.

All Experiments are conducted on eight NVIDIA A100-80GB SXM GPUs<sup>1</sup>. Reproducing the fine-tuning process would require approximately 15 A100 GPU days.

## REFERENCES

- 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. In *NeurIPS*, 2022.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-VL: A Frontier Large Vision-Language Model with Versatile Abilities. *arXiv preprint arXiv:2308.12966*, 2023.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. Graph of Thoughts: Solving Elaborate Problems with Large Language Models. In *AAAI*, 2024.
- 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. In *NeurIPS*, 2020.
- Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. COCO-Stuff: Thing and Stuff Classes in Context. In *CVPR*, 2018.
- Hila Chefer, Shir Gur, and Lior Wolf. Generic Attention-Model Explainability for Interpreting Bi-Modal and Encoder-Decoder Transformers. In *ICCV*, 2021.
- Delong Chen, Jianfeng Liu, Wenliang Dai, and Baoyuan Wang. Visual Instruction Tuning with Polite Flamingo. In *AAAI*, 2024a.
- Jiaxing Chen, Yuxuan Liu, Dehu Li, Xiang An, Ziyong Feng, Yongle Zhao, and Yin Xie. Plug-and-Play Grounding of Reasoning in Multimodal Large Language Models. *arXiv preprint arXiv:2403.19322*, 2024b.
- Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. MiniGPT-V2: Large Language Model as a Unified Interface for Vision-Language Multi-Task Learning. *arXiv preprint arXiv:2310.09478*, 2023a.
- 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*, 2023b.

<sup>1</sup><https://www.nvidia.com/en-sg/data-center/a100/>

- 594 Xiang Chen, Chenxi Wang, Yida Xue, Ningyu Zhang, Xiaoyan Yang, Qiang Li, Yue Shen, Jinjie  
595 Gu, and Huajun Chen. Unified Hallucination Detection for Multimodal Large Language Models.  
596 *arXiv preprint arXiv:2402.03190*, 2024c.
- 597 Xianjie Chen, Roozbeh Mottaghi, Xiaobai Liu, Sanja Fidler, Raquel Urtasun, and Alan Yuille. De-  
598 tect What You Can: Detecting and Representing Objects Using Holistic Models and Body Parts.  
599 In *CVPR*, 2014.
- 600 Yuyan Chen, Qiang Fu, Yichen Yuan, Zhihao Wen, Ge Fan, Dayiheng Liu, Dongmei Zhang, Zhixu  
601 Li, and Yanghua Xiao. Hallucination Detection: Robustly Discerning Reliable Answers in Large  
602 Language Models. In *CIKM*, 2023c.
- 603 Zhenfang Chen, Rui Sun, Wenjun Liu, Yining Hong, and Chuang Gan. GENOME: Generative  
604 Neuro-Symbolic Visual Reasoning by Growing and Reusing Modules. In *ICLR*, 2024d.
- 605 Pengyu Cheng, Tianhao Hu, Han Xu, Zhisong Zhang, Yong Dai, Lei Han, and Nan Du. Self-Playing  
606 Adversarial Language Game Enhances LLM Reasoning. *arXiv preprint arXiv:2404.10642*, 2024.
- 607 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng,  
608 Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna:  
609 An Open-Source Chatbot Impressing GPT-4 with 90%\* ChatGPT Quality, March 2023. URL  
610 <https://lmsys.org/blog/2023-03-30-vicuna/>.
- 611 Zheng Chu, Jingchang Chen, Qianglong Chen, Weijiang Yu, Tao He, Haotian Wang, Weihua Peng,  
612 Ming Liu, Bing Qin, and Ting Liu. A Survey of Chain of Thought Reasoning: Advances, Frontiers  
613 and Future. *arXiv preprint arXiv:2309.15402*, 2023.
- 614 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
615 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training Verifiers to  
616 Solve Math Word Problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 617 Chenhang Cui, Yiyang Zhou, Xinyu Yang, Shirley Wu, Linjun Zhang, James Zou, and Huaxiu Yao.  
618 Holistic Analysis of Hallucination in GPT-4V (ision): Bias and Interference Challenges. *arXiv*  
619 *preprint arXiv:2311.03287*, 2023.
- 620 Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang,  
621 Boyang Li, Pascale N Fung, and Steven Hoi. InstructBLIP: Towards General-Purpose Vision-  
622 Language Models with Instruction Tuning. In *NeurIPS*, 2024.
- 623 Ishita Dasgupta, Andrew K Lampinen, Stephanie CY Chan, Antonia Creswell, Dharshan Kumaran,  
624 James L McClelland, and Felix Hill. Language Models Show Human-Like Content Effects on  
625 Reasoning. *arXiv preprint arXiv:2207.07051*, 2022.
- 626 Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter,  
627 Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. PaLM-E: An Embodied Mul-  
628 timodal Language Model. *arXiv preprint arXiv:2303.03378*, 2023.
- 629 Xiaoran Fan, Tao Ji, Changhao Jiang, Shuo Li, Senjie Jin, Sirui Song, Junke Wang, Boyang Hong,  
630 Lu Chen, Guodong Zheng, et al. MouSi: Poly-Visual-Expert Vision-Language Models. *arXiv*  
631 *preprint arXiv:2401.17221*, 2024.
- 632 Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu  
633 Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. MME: A Comprehensive Evaluation  
634 Benchmark for Multimodal Large Language Models. *arXiv preprint arXiv:2306.13394*, 2024.
- 635 Yingqiang Ge, Wenyue Hua, Kai Mei, Jianchao Ji, Juntao Tan, Shuyuan Xu, Zelong Li, and  
636 Yongfeng Zhang. OpenAGI: When LLM Meets Domain Experts. In *NeurIPS*, 2023.
- 637 Shijie Geng, Jianbo Yuan, Yu Tian, Yuxiao Chen, and Yongfeng Zhang. Hiclip: Contrastive  
638 language-image pretraining with hierarchy-aware attention. In *ICLR*, 2023.
- 639 Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V  
640 in VQA Matter: Elevating the Role of Image Understanding in Visual Question Answering. In  
641 *CVPR*, 2017.

- 648 Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang  
649 Chen, Furong Huang, Yaser Yacoob, et al. Hallusionbench: An Advanced Diagnostic Suite  
650 for Entangled Language Hallucination & Visual Illusion in Large Vision-Language Models. In  
651 *CVPR*, 2024a.
- 652 Xinyan Guan, Yanjiang Liu, Hongyu Lin, Yaojie Lu, Ben He, Xianpei Han, and Le Sun. Mitigating  
653 Large Language Model Hallucinations via Autonomous Knowledge Graph-Based Retrofitting. In  
654 *AAAI*, 2024b.
- 656 Anisha Gunjal, Jihan Yin, and Erhan Bas. Detecting and Preventing Hallucinations in Large Vision  
657 Language Models. In *AAAI*, 2024.
- 658 Cheng Han, James C Liang, Qifan Wang, Majid Rabbani, Sohail Dianat, Raghuvver Rao, Ying Nian  
659 Wu, and Dongfang Liu. Image Translation as Diffusion Visual Programmers. In *ICLR*, 2024.
- 661 Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu.  
662 Reasoning with Language Model is Planning with World Model. In *EMNLP*, 2023.
- 663 Ju He, Shuo Yang, Shaokang Yang, Adam Kortylewski, Xiaoding Yuan, Jie-Neng Chen, Shuai Liu,  
664 Cheng Yang, Qihang Yu, and Alan Yuille. PartImageNet: A Large, High-Quality Dataset of Parts.  
665 In *ECCV*, 2022.
- 667 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,  
668 and Weizhu Chen. LoRA: Low-Rank Adaptation of Large Language Models. In *ICLR*, 2022.
- 669 Wen Yue Hua and Yongfeng Zhang. System 1+ System 2= Better World: Neural-Symbolic Chain of  
670 Logic Reasoning. In *EMNLP*, 2022.
- 672 Jie Huang and Kevin Chen-Chuan Chang. Towards Reasoning in Large Language Models: A Sur-  
673 vey. In *ACL*, 2023a.
- 674 Jie Huang and Kevin Chen-Chuan Chang. Towards Reasoning in Large Language Models: A Sur-  
675 vey. In *ACL*, 2023b.
- 677 Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv,  
678 Lei Cui, Owais Khan Mohammed, Barun Patra, et al. Language is Not All You Need: Aligning  
679 Perception with Language Models. In *NeurIPS*, 2023.
- 680 Drew A Hudson and Christopher D Manning. GQA: A New Dataset for Real-World Visual Reason-  
681 ing and Compositional Question Answering. In *CVPR*, 2019.
- 682 Alon Jacovi and Yoav Goldberg. Towards Faithfully Interpretable NLP Systems: How Should We  
683 Define and Evaluate Faithfulness? In *ACL*, 2020.
- 684 Dongsheng Jiang, Yuchen Liu, Songlin Liu, Xiaopeng Zhang, Jin Li, Hongkai Xiong, and Qi Tian.  
685 From CLIP to DINO: Visual Encoders Shout in Multi-Modal Large Language Models. *arXiv*  
686 *preprint arXiv:2310.08825*, 2023.
- 687 Mingyu Jin, Qinkai Yu, Dong Shu, Haiyan Zhao, Wen Yue Hua, Yanda Meng, Yongfeng Zhang, and  
688 Mengnan Du. The Impact of Reasoning Step Length on Large Language Models. In *ACL*, 2024.
- 689 Daniel Kahneman. *Thinking, Fast and Slow*. Farrar, Straus and Giroux, 2011.
- 690 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child,  
691 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling Laws for Neural Language  
692 Models. *arXiv preprint arXiv:2001.08361*, 2020.
- 693 Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. ReferItGame: Referring to  
694 Objects in Photographs of Natural Scenes. In *EMNLP*, 2014.
- 695 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete  
696 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment Anything. In *ICCV*,  
697 2023.

- 702 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large  
703 Language Models are Zero-Shot Reasoners. In *NeurIPS*, 2022.  
704
- 705 Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Rea-  
706 soning Segmentation via Large Language Model. In *CVPR*, 2024.
- 707 Bin Lei, Chunhua Liao, Caiwen Ding, et al. Boosting Logical Reasoning in Large Language Models  
708 Through a New Framework: The Graph of Thought. *arXiv preprint arXiv:2308.08614*, 2023.  
709
- 710 Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-Bench: Bench-  
711 marking Multimodal LLMs with Generative Comprehension. In *CVPR*, 2024.
- 712 Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. BLIP: Bootstrapping Language-Image  
713 Pre-Training for Unified Vision-Language Understanding and Generation. In *ICML*, 2022.  
714
- 715 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: Bootstrapping Language-Image  
716 Pre-training with Frozen Image Encoders and Large Language Models. In *ICML*, 2023a.
- 717 Junyi Li, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. HALUEVAL: A Large-  
718 Scale Hallucination Evaluation Benchmark for Large Language Models. In *EMNLP*, 2023b.  
719
- 720 Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating  
721 Object Hallucination in Large Vision-Language Models. In *ACL*, 2023c.
- 722 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan  
723 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s Verify Step by Step. *arXiv preprint*  
724 *arXiv:2305.20050*, 2023.
- 725 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual Instruction Tuning. In *NeurIPS*,  
726 2023a.  
727
- 728 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual Instruction Tuning. In *NeurIPS*,  
729 2023b.
- 730 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved Baselines with Visual Instruc-  
731 tion Tuning. In *CVPR*, 2024a.  
732
- 733 Yixin Liu, Avi Singh, C Daniel Freeman, John D Co-Reyes, and Peter J Liu. Improving Large  
734 Language Model Fine-Tuning for Solving Math Problems. *arXiv preprint arXiv:2310.10047*,  
735 2023c.
- 736 Zuyan Liu, Yuhao Dong, Yongming Rao, Jie Zhou, and Jiwen Lu. Chain-of-Spot: Interactive Rea-  
737 soning Improves Large Vision-Language Models. *arXiv preprint arXiv:2403.12966*, 2024b.  
738
- 739 Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization. In *ICLR*, 2019.
- 740 Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord,  
741 Peter Clark, and Ashwin Kalyan. Learn to Explain: Multimodal Reasoning via Thought Chains  
742 for Science Question Answering. In *NeurIPS*, 2022.  
743
- 744 Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu,  
745 and Jianfeng Gao. Chameleon: Plug-and-Play Compositional Reasoning with Large Language  
746 Models. In *NeurIPS*, 2023a.
- 747 Pan Lu, Liang Qiu, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, Tanmay Rajpurohit, Peter  
748 Clark, and Ashwin Kalyan. Dynamic Prompt Learning via Policy Gradient for Semi-Structured  
749 Mathematical Reasoning. In *ICML*, 2023b.
- 750 Gen Luo, Yiyi Zhou, Tianhe Ren, Shengxin Chen, Xiaoshuai Sun, and Rongrong Ji. Cheap and  
751 Quick: Efficient Vision-Language Instruction Tuning for Large Language Models. In *NeurIPS*,  
752 2023.  
753
- 754 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri  
755 Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-Refine: Iterative Refinement  
with Self-Feedback. In *NeurIPS*, 2023.

- 756 Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L Yuille, and Kevin Murphy.  
757 Generation and Comprehension of Unambiguous Object Descriptions. In *CVPR*, 2016.  
758
- 759 Maxwell Nye, Michael Tessler, Josh Tenenbaum, and Brenden M Lake. Improving Coherence  
760 and Consistency in Neural Sequence Models with Dual-System, Neuro-Symbolic Reasoning. In  
761 *NeurIPS*, 2021.
- 762 OpenAI. GPT-4 Technical Report. *arXiv preprint arXiv:2303.08774*, 2024.  
763
- 764 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
765 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training Language Models to Follow  
766 Instructions with Human Feedback. In *NeurIPS*, 2022.
- 767 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor  
768 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. PyTorch: An Imperative Style,  
769 High-Performance Deep Learning Library. In *NeurIPS*, 2019.  
770
- 771 Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei.  
772 Kosmos-2: Grounding Multimodal Large Language Models to the World. In *ICLR*, 2024.
- 773 Renjie Pi, Jiahui Gao, Shizhe Diao, Rui Pan, Hanze Dong, Jipeng Zhang, Lewei Yao, Jianhua Han,  
774 Hang Xu, and Lingpeng Kong Tong Zhang. DetGPT: Detect What You Need via Reasoning. In  
775 *EMNLP*, 2023.  
776
- 777 Ji Qi, Ming Ding, Weihang Wang, Yushi Bai, Qingsong Lv, Wenyi Hong, Bin Xu, Lei Hou, Juanzi  
778 Li, Yuxiao Dong, and Jie Tang. Cogcom: Train large vision-language models diving into details  
779 through chain of manipulations. *arXiv preprint arXiv:2402.04236*, 2024.
- 780 Shuofei Qiao, Yixin Ou, Ningyu Zhang, Xiang Chen, Yunzhi Yao, Shumin Deng, Chuanqi Tan, Fei  
781 Huang, and Huajun Chen. Reasoning with Language Model Prompting: A Survey. In *ACL*, 2023.  
782
- 783 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
784 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning Transferable Visual  
785 Models from Natural Language Supervision. In *ICML*, 2021.
- 786 Vignesh Ramanathan, Anmol Kalia, Vladan Petrovic, Yi Wen, Baixue Zheng, Baishan Guo, Rui  
787 Wang, Aaron Marquez, Rama Kovvuri, Abhishek Kadian, et al. PACO: Parts and Attributes of  
788 Common Objects. In *CVPR*, 2023.
- 789 Zhongwei Ren, Zhicheng Huang, Yunchao Wei, Yao Zhao, Dongmei Fu, Jiashi Feng, and Xiaojie  
790 Jin. Pixellm: Pixel reasoning with large multimodal model. In *CVPR*, 2024.  
791
- 792 Anna Rohrbach, Marcus Rohrbach, Ronghang Hu, Trevor Darrell, and Bernt Schiele. Grounding of  
793 Textual Phrases in Images by Reconstruction. In *ECCV*, 2016.
- 794 Omar Shaikh, Hongxin Zhang, William Held, Michael Bernstein, and Diyi Yang. On Second  
795 Thought, Let’s Not Think Step by Step! Bias and Toxicity in Zero-Shot Reasoning. In *ACL*,  
796 2023.  
797
- 798 Hao Shao, Shengju Qian, Han Xiao, Guanglu Song, Zhuofan Zong, Letian Wang, Yu Liu, and Hong-  
799 sheng Li. Visual cot: Advancing multi-modal language models with a comprehensive dataset and  
800 benchmark for chain-of-thought reasoning. *arXiv preprint arXiv:2403.16999*, 2024.
- 801 Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugging-  
802 GPT: Solving AI Tasks with ChatGPT and Its Friends in Hugging Face. In *NeurIPS*, 2023.  
803
- 804 Baifeng Shi, Ziyang Wu, Maolin Mao, Xin Wang, and Trevor Darrell. When Do We Not Need  
805 Larger Vision Models? *arXiv preprint arXiv:2403.13043*, 2024.
- 806 Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh,  
807 and Marcus Rohrbach. Towards VQA Models that Can Read. In *CVPR*, 2019.  
808
- 809 Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. Eva-clip: Improved training  
techniques for clip at scale. *arXiv preprint arXiv:2303.15389*, 2023.

- 810 Dídac Surís, Sachit Menon, and Carl Vondrick. ViperGPT: Visual Inference via Python Execution  
811 for Reasoning. In *ICCV*, 2023.
- 812
- 813 Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes Wide  
814 Shut? Exploring the Visual Shortcomings of Multimodal LLMs. In *CVPR*, 2024a.
- 815
- 816 Yongqi Tong, Dawei Li, Sizhe Wang, Yujia Wang, Fei Teng, and Jingbo Shang. Can LLMs Learn  
817 from Previous Mistakes? Investigating LLMs’ Errors to Boost for Reasoning. *arXiv preprint*  
818 *arXiv:2403.20046*, 2024b.
- 819 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
820 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
821 Efficient Foundation Language Models. *arXiv preprint arXiv:2302.13971*, 2023a.
- 822
- 823 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
824 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open Founda-  
825 tion and Fine-Tuned Chat Models. *arXiv preprint arXiv:2307.09288*, 2023b.
- 826
- 827 Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen, You Wu, Luke Zettlemoyer, and Huan Sun.  
828 Towards Understanding Chain-of-Thought Prompting: An Empirical Study of What Matters. In  
829 *ACL*, 2023a.
- 830
- 831 Taowen Wang, Yiyang Liu, James Chenhao Liang, Yiming Cui, Yuning Mao, Shaoliang Nie, Jiahao  
832 Liu, Fuli Feng, Zenglin Xu, Cheng Han, et al. M<sup>2</sup>PT: Multimodal Prompt Tuning for Zero-shot  
833 Instruction Learning. In *EMNLP*, 2024a.
- 834
- 835 Weihang Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang,  
836 Lei Zhao, Xixuan Song, et al. CogVLM: Visual Expert for Pretrained Language Models. *arXiv*  
837 *preprint arXiv:2311.03079*, 2023b.
- 838
- 839 Weiyun Wang, Min Shi, Qingyun Li, Wenhao Wang, Zhenhang Huang, Linjie Xing, Zhe Chen, Hao  
840 Li, Xizhou Zhu, Zhiguo Cao, et al. The All-Seeing Project: Towards Panoptic Visual Recognition  
841 and Understanding of the Open World. In *ICLR*, 2024b.
- 842
- 843 Wenhao Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong  
844 Lu, Jie Zhou, Yu Qiao, et al. VisionLLM: Large Language Model is Also an Open-Ended Decoder  
845 for Vision-Centric Tasks. In *NeurIPS*, 2023c.
- 846
- 847 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowd-  
848 hery, and Denny Zhou. Self-Consistency Improves Chain of Thought Reasoning in Language  
849 Models. In *ICLR*, 2023d.
- 850
- 851 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny  
852 Zhou, et al. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In  
853 *NeurIPS*, 2022.
- 854
- 855 Yilin Wen, Zifeng Wang, and Jimeng Sun. MindMap: Knowledge Graph Prompting Sparks Graph  
856 of Thoughts in Large Language Models. *arXiv preprint arXiv:2308.09729*, 2023.
- 857
- 858 Penghao Wu and Saining Xie. Guided Visual Search as a Core Mechanism in Multimodal LLMs.  
859 In *CVPR*, 2024.
- 860
- 861 Yixuan Wu, Yizhou Wang, Shixiang Tang, Wenhao Wu, Tong He, Wanli Ouyang, Jian Wu, and  
862 Philip Torr. DetToolChain: A New Prompting Paradigm to Unleash Detection Ability of MLLM.  
863 *arXiv preprint arXiv:2403.12488*, 2024.
- 864
- 865 Ruyi Xu, Yuan Yao, Zonghao Guo, Junbo Cui, Zanlin Ni, Chunjiang Ge, Tat-Seng Chua, Zhiyuan  
866 Liu, Maosong Sun, and Gao Huang. Llava-UHD: An LMM Perceiving Any Aspect Ratio and  
867 High-Resolution Images. *arXiv preprint arXiv:2403.11703*, 2024a.
- 868
- 869 Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli. Hallucination is Inevitable: An Innate Limitation  
870 of Large Language Models. *arXiv preprint arXiv:2401.11817*, 2024b.



- 864 Fanglong Yao, Changyuan Tian, Jintao Liu, Zequn Zhang, Qing Liu, Li Jin, Shuchao Li, Xiaoyu Li,  
865 and Xian Sun. Thinking Like an Expert: Multimodal Hypergraph-of-Thought (HOT) Reasoning  
866 to Boost Foundation Modals. *arXiv preprint arXiv:2308.06207*, 2023a.
- 867
- 868 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik  
869 Narasimhan. Tree of Thoughts: Deliberate Problem Solving with Large Language Models. In  
870 *NeurIPS*, 2023b.
- 871 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.  
872 React: Synergizing reasoning and acting in language models. In *ICLR*, 2023c.
- 873
- 874 Yao Yao, Zuchao Li, and Hai Zhao. Beyond Chain-of-Thought, Effective Graph-of-Thought Rea-  
875 soning in Large Language Models. *arXiv preprint arXiv:2305.16582*, 2023d.
- 876 Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen  
877 Hu, Pengcheng Shi, Yaya Shi, et al. mPLUG-Owl: Modularization Empowers Large Language  
878 Models with Multimodality. *arXiv preprint arXiv:2304.14178*, 2023.
- 879
- 880 Haoxuan You, Rui Sun, Zhecan Wang, Long Chen, Gengyu Wang, Hammad A Ayyubi, Kai-Wei  
881 Chang, and Shih-Fu Chang. Idealgpt: Iteratively Decomposing Vision and Language Reasoning  
882 via Large Language Models. In *EMNLP*, 2023.
- 883 Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang,  
884 and Lijuan Wang. MM-Vet: Evaluating Large Multimodal Models for Integrated Capabilities. In  
885 *ICML*, 2024.
- 886
- 887 Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Sainbayar Sukhbaatar, Jing Xu, and Jason  
888 Weston. Self-Rewarding Language Models. *arXiv preprint arXiv:2401.10020*, 2024.
- 889 Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D. Goodman. Star: Self-Taught Reasoner Boot-  
890 strapping Reasoning with Reasoning. In *NeurIPS*, 2022.
- 891
- 892 Muru Zhang, Ofir Press, William Merrill, Alisa Liu, and Noah A Smith. How Language Model  
893 Hallucinations Can Snowball. In *arXiv preprint arXiv:2305.13534*, 2023a.
- 894 Shilong Zhang, Peize Sun, Shoufa Chen, Min Xiao, Wenqi Shao, Wenwei Zhang, Kai Chen, and  
895 Ping Luo. GPT4ROI: Instruction Tuning Large Language Model on Region-of-Interest. *arXiv*  
896 *preprint arXiv:2307.03601*, 2023b.
- 897
- 898 Yanzhe Zhang, Ruiyi Zhang, Jiuxiang Gu, Yufan Zhou, Nedim Lipka, Diyi Yang, and Tong  
899 Sun. LLAVAR: Enhanced Visual Instruction Tuning for Text-Rich Image Understanding. *arXiv*  
900 *preprint arXiv:2306.17107*, 2023c.
- 901 Yifan Zhang, Yang Yuan, and Andrew Chi-Chih Yao. Meta Prompting for AI Systems. *arXiv*  
902 *preprint arXiv:2311.11482*, 2023d.
- 903
- 904 Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic Chain of Thought Prompting  
905 in Large Language Models. In *ICLR*, 2023e.
- 906 Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang,  
907 Dawei Yin, and Mengnan Du. Explainability for Large Language Models: A Survey. *ACM*  
908 *Transactions on Intelligent Systems and Technology*, 15(2):1–38, 2024.
- 909
- 910 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
911 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging LLM-as-a-Judge with MT-Bench and  
912 Chatbot Arena. In *NeurIPS*, 2023.
- 913 Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene  
914 Parsing through ADE20K Dataset. In *CVPR*, 2017.
- 915
- 916 Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuur-  
917 mans, Claire Cui, Olivier Bousquet, Quoc Le, et al. Least-to-Most Prompting Enables Complex  
Reasoning in Large Language Models. In *ICLR*, 2023.

918 Yiyang Zhou, Chenhang Cui, Jaehong Yoon, Linjun Zhang, Zhun Deng, Chelsea Finn, Mohit  
919 Bansal, and Huaxiu Yao. Analyzing and Mitigating Object Hallucination in Large Vision-  
920 Language Models. In *ICLR*, 2024.

921  
922 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. MiniGPT-4: Enhanc-  
923 ing Vision-Language Understanding with Advanced Large Language Models. *arXiv preprint*  
924 *arXiv:2304.10592*, 2023a.

925 Xinyu Zhu, Junjie Wang, Lin Zhang, Yuxiang Zhang, Yongfeng Huang, Ruyi Gan, Jiaying Zhang,  
926 and Yujiu Yang. Solving Math Word Problems via Cooperative Reasoning Induced Language  
927 Models. In *ACL*, 2023b.

928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

## SUMMARY OF THE APPENDIX

This appendix contains additional details for the ICLR 2025 submission, titled “*Visual Agents as Fast and Slow Thinkers*”. The appendix is organized as follows:

- §A provides **Implementation Details and Pseudo Code**.
- §B reports more **Results** for **Different Thinking Modes**.
- §C reports more **Quantitative Results** for **Visual Question Answering**.
- §D shows more **Quantitative Results** for **Segmentation**.
- §E analyzes **Failure Case**.
- §F examines the **Limitation and Future Work** of our research.
- §G discusses the **Social Impact** of our research.
- §H offers **Ethical Guard** or our dataset.
- §I claims **Reproducibility** of our approach.
- §J supplies **Data License** for the methods we used for comparison.

## A IMPLEMENTATION DETAILS AND PSEUDO-CODE OF FAST

**Visual Resampler.** The resampler (Alayrac et al., 2022) compresses high-dimensional visual features into a fixed-size latent space using a cross-attention mechanism. It begins with a set of learnable latent embeddings, which query the vision encoder’s output features through scaled dot-product attention. Each latent embedding attends selectively to the most relevant visual tokens, guided by attention weights computed via the query-key interaction. The process iterates across multiple layers of cross-attention, followed by feedforward transformations, refining the latent representations at each step. This approach ensures efficient dimensionality reduction while retaining critical information, producing a compact set of visual tokens for downstream tasks.

**Hyper-parameters.** We follow established methodologies and utilize LLaVA-v1.5 (Liu et al., 2024a) as the foundational visual agent. The image resolution is preprocessed to  $336 \times 336$  pixels to accommodate the clip-vit-large-patch14-336 vision encoder (Radford et al., 2021). The AdamW optimizer (Loshchilov & Hutter, 2019) is employed with the DeepSpeed ZeRO 2<sup>2</sup> configuration for fine-tuning the switch, proposal, and segmentation adapters with LoRA (Hu et al., 2022). For the LoRA configuration, we set the rank to 128 and alpha to 256, consistent with the settings of LLaVA-v1.5. Additionally, we adjust the learning rate of the vision encoder projection layer to  $2e-5$  to achieve better alignment. An MLP projection with channels of [256, 4096, 4096] is used to connect image representations into the word embedding space for the projection layer. An additional resampler projection layer is used to reduce the number of image tokens.

**Training Data for Switch Adapter.** Consistent with the pretraining stage of LLaVA-v1.5, we initially pretrain Vicuna-v1.5 as a base frozen large language model and for the MLP projection layer and sampler layer of the CLIP vision encoder using a 558K subset of the LAION-CC-SBU dataset<sup>3</sup> with BLIP (Li et al., 2022) captions. During the fine-tuning stage, we integrate the negative dataset acquired from  $V^*$  (Wu & Xie, 2024) and PixellM (Ren et al., 2024) and with the original LLaVA-v1.5 instruction tuning 665k data<sup>4</sup> for LoRA based finetuning.

Specifically, The dataset for fine-tuning the switch adapter was carefully constructed to emphasize scenarios requiring precise object recognition and complex reasoning. For the GQA subset of the 167k VQA data, we specifically targeted questions where the annotated objects mentioned in the query were critical for deriving the correct answer. Initially, the InstructBLIP model was used to evaluate GQA questions with annotated objects. Only questions that the model could correctly answer were retained. To ensure the importance of these annotated objects, we applied the LaMa image inpainting model to erase the mentioned objects from the corresponding images. The modified images were re-evaluated using InstructBLIP, and only questions that the model failed to answer

<sup>2</sup><https://github.com/microsoft/DeepSpeed>

<sup>3</sup><https://huggingface.co/datasets/liuhaotian/LLaVA-Pretrain>

<sup>4</sup><https://huggingface.co/datasets/liuhaotian/LLaVA-Instruct-150K>

1026 after object removal were included. This process ensured that the curated subset focused exclusively  
1027 on questions where the annotated objects were essential, forming a robust component of the VQA  
1028 data.

1029 For the VAW object attribution dataset, both open-ended and binary questions about object attributes  
1030 were synthesized. Open-ended questions were formulated around attributes such as “color,” “mate-  
1031 rial,” and “pose,” while binary questions incorporated additional attributes like “state” and “optical  
1032 property.” Answer formats adhered to predefined structures to ensure consistency. The same ob-  
1033 ject removal and re-evaluation strategy as used in the GQA subset was applied, filtering the data to  
1034 include only questions where the absence of objects rendered the query unanswerable.

1035 From the LLaVA-80K instruction tuning data, noun phrases were extracted from the text of questions  
1036 or instructions and matched with object category names defined by COCO, augmented with common  
1037 synonyms such as “man” and “woman” for the “person” category. Images were retained only if the  
1038 identified categories had annotated instances with bounding boxes. These annotated instances, along  
1039 with their bounding box coordinates, were used as target objects during training.

1040 In addition, we incorporated a data generation approach inspired by LISA, utilizing GPT-4 and  
1041 GPT-4V to expand and diversify the dataset. Initially, LLAVA was used for image captioning, and  
1042 GPT-4 generated questions about multiple regions in the image. While this approach utilized pre-  
1043 existing mask annotations to reduce costs, its diversity was limited to the scope of the captions. To  
1044 address these limitations, we refined the pipeline with GPT-4V, leveraging its advanced capabilities  
1045 in visual understanding. Image captions, object names, and bounding box coordinates were input  
1046 into GPT-4V, which, using dynamically crafted prompts, autonomously selected instances and gen-  
1047 erated nuanced question-answer pairs tailored to the image content. This refinement significantly  
1048 improved the diversity and contextual relevance of the data. An illustrative example of such prompts  
1049 is provided below:

1050  
1051  
1052 **Prompt:** Imagine you need to query a machine agent about an  
1053 image. The image has a height of 720 pixels and a width of  
1054 1280 pixels. You are given several entities described by a  
1055 list, each identifying an object in the image along with its  
1056 location. The class names and corresponding coordinates are  
1057 as follows:

- 1058 • Dog at [350.12, 450.45, 480.89, 600.67];
- 1059 • Ball at [200.33, 300.22, 250.78, 350.56];
- 1060 • Grass at [0.0, 500.0, 1280.0, 720.0];

1061 Coordinates are represented as (top-left x, top-left  
1062 y, bottom-right x, bottom-right y). The question must  
1063 incorporate at least two of these objects and require  
1064 reasoning about the relationships or interactions between  
1065 them. Additional requirements for the generated question are  
1066 as follows:

- 1067 1.The answer to the question must explicitly reference each  
1068 included object or its equivalent and avoid implying the  
1069 presence of any other objects not listed.
- 1070 2.The question must be precise, meaningful, and avoid being  
1071 overly general.
- 1072 3.The question should frame a single cohesive activity  
1073 or relationship rather than merely combining independent  
1074 sub-queries.
- 1075 4.When answering, the class names should be rephrased to  
1076 indicate their position, role, or interaction in the image.

1077  
1078 This multi-faceted dataset construction process ensured the generation of diverse and challenging  
1079 samples, providing a robust foundation for fine-tuning the switch adapter on complex reasoning  
tasks.

**Training Data for Proposal Adapter** To determine the corresponding region for a query, we use LRP++ (Chefer et al., 2021) for data construction, similar to Chain of Spot (Liu et al., 2024b). Our initial prompt is as follows:

```
<Image>
To answer the question: [Q],
where is the region of interest in the image based on [C]?

Ans.str[w0, w1, h0, h1]
```

The question  $Q$  and the context clue  $C$  are formatted to get the answer in terms of a bounding box. In this format,  $w_0, w_1$  represent the left and right boundaries, respectively, while  $h_0, h_1$  denote the upper and lower boundaries. To identify the correct region, we sampled one question per image from the LLaVA instruction tuning data, consisting of a total of 665k data for proposal finetuning.

**Training Data for Seg Adapter.** Adopting an approach similar to LISA Lai et al. (2024), the training data for our model comprises three distinct segments: a semantic segmentation dataset, a referring segmentation dataset, and a reasoning segmentation dataset. We deliberately exclude visual question-answering datasets to enhance the model’s segmentation performance. The semantic segmentation segment includes the ADE20K (Zhou et al., 2017), COCO-Stuff (Caesar et al., 2018), and LVIS-PACO (Ramanathan et al., 2023) part segmentation datasets. The referring segmentation datasets encompass refCOCO (Kazemzadeh et al. (2014), refCOCO+ (Kazemzadeh et al., 2014), refCOCOg (Caesar et al., 2018), and refCLEF (Rohrbach et al., 2016). The reasoning segmentation dataset includes ReasonSeg (Lai et al., 2024). It is important to note that the referring segmentation and reasoning segmentation datasets are carefully excluded during the evaluation of the segmentation benchmarks to prevent any potential data leakage.

**Pseudo-code Implementation.** The pseudo-code of FAST is given in Pseudo-code 1.

## B MORE RESULTS FOR DIFFERENT THINKING MODES

As shown in Table 6, FaST demonstrates strong performance in both System 1 and System 2 reasoning on VQA datasets, outperforming baseline methods in most cases. Notably, FaST achieves the highest accuracy in challenging System 2 tasks across GQA,  $VQA^T$ , and  $SQA^T$ , which require advanced reasoning capabilities. This highlights the effectiveness of the switch adapter mechanism in dynamically allocating tasks based on complexity. While maintaining competitive performance in simpler System 1 tasks, FaST leverages its adaptive architecture to excel in more complex scenarios, as evidenced by its superior System 2 results.

Further, in Table 7, FaST’s robustness extends to reasoning segmentation tasks, where it achieves significant improvements in System 2 performance compared to baseline models such as LLaVA with segmentation and LISA-7B. For example, in the ReasonSeg dataset, FaST records a remarkable 48.2 CIoU in System 2 tasks, significantly outperforming LISA-7B and LLaVA, which achieve 43.3 CIoU and 42.4 CIoU, respectively. This result underscores FaST’s ability to generalize effectively across diverse task families and reasoning paradigms.

Overall, the results validate the universal applicability and robustness of the FaST framework. By effectively utilizing the switch adapter to allocate tasks dynamically, FaST demonstrates a strong capability to balance performance across both simple and complex reasoning tasks, making it a reliable solution for diverse real-world applications.

## C MORE QUALITATIVE RESULTS FOR VISUAL QUESTION ANSWERING

Figure 7 presents additional qualitative results for Visual Question Answering (VQA). Our FAST framework consistently demonstrates remarkable performance across various challenging scenarios. Notably, in the bottom right corner of Figure 7, our FAST leverages extensive world knowledge to identify the keyboard, which subsequently aids in discovering the hidden computer mouse and providing the correct answer. This ability to integrate and utilize contextual information showcases the

**Algorithm 1:** Pseudo-code of FAST in a PyTorch-like style.

```

1134
1135
1136 class FaST:
1137     def __init__(self, switch_llm, proposal_llm, seg_llm):
1138         self.switch_llm = switch_llm
1139         self.proposal_llm = proposal_llm
1140         self.seg_llm = seg_llm
1141
1142     def get_contextual_clues(self, image, question):
1143         # Get missing objects and context clues using switch adapter
1144         return self.switch_llm(image, question)
1145
1146     # Construct Chain of Evidence
1147     def construct_coe(self, image, question, context_clues, missing_objects):
1148         # Step 1: Get region proposals
1149         region = self.proposal_llm(image, question, context_clues)
1150
1151         # Step 2: Get pixel-level mask for the missing objects
1152         mask = self.seg_llm(region, missing_objects)
1153
1154         return (context_clues, region, missing_objects, mask)
1155
1156     # Main Function
1157     def forward(self, image, question):
1158
1159         # Get initial answer
1160         initial_answer = self.switch_llm(image, question)
1161
1162         # Check if slow thinking is needed based on the initial answer
1163         if "sorry, i can not answer" in initial_answer.lower():
1164             # Perform slow thinking
1165             missing_objects, context_clues = initial_answer['obj'], initial_answer
1166             ['clue']
1167             chain_of_evidence = self.construct_coe(image, question, context_clues,
1168             missing_objects)
1169
1170             # Generate the final answer using the constructed chain of evidence
1171             final_answer = self.switch_llm(image, question, chain_of_evidence)
1172         else:
1173             # Perform fast thinking
1174             final_answer = initial_answer
1175
1176         return final_answer
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

```

Method	LLM	$VQA^{v2}$		GQA		$VQA^T$		$SQA^I$	
		Sys 1	Sys 2	Sys 1	Sys 2	Sys 1	Sys 2	Sys 1	Sys 2
BLIP-2	Vicuna-13B	67.3	53.1	37.8	22.4	44.3	39.7	63.4	59.2
LLaVA-v1.5	Vicuna-7B	<u>81.2</u>	68.0	<b>70.3</b>	47.0	<u>61.1</u>	53.7	<u>68.4</u>	65.7
Chain of Spot	Vicuna-7B	<b>82.1</b>	<u>74.5</u>	<u>70.9</u>	<u>50.7</u>	<b>62.1</b>	<u>59.0</u>	<b>68.6</b>	<u>67.8</u>
FAST (Ours)	Vicuna-7B	81.1	<b>75.5</b>	70.2	<b>52.3</b>	61.2	<b>60.2</b>	68.2	<b>70.2</b>

Table 6. **System 1 and System 2 performance on VQA datasets.** FaST demonstrates superior performance in both reasoning modes compared to baselines.

model’s advanced capabilities and highlights its potential for practical applications. The qualitative results further underscore FAST’s robustness and versatility in handling diverse VQA tasks.

### D MORE QUALITATIVE RESULTS FOR SEGMENTATION

Figure 8 showcases further qualitative results for the Segmentation task. Our FAST model excels in various challenging scenarios, accurately locating difficult targets and performing complex reasoning for more demanding queries. For instance, in the bottom right corner, the model successfully identifies an appliance that can be turned on when feeling hot by recognizing relevant contextual clues that suggest the appliance should probably appear on the wall, thereby resulting in the correct answer. This example demonstrates the model’s advanced understanding, adaptability, and precision.

Method	refCOCOg		ReasonSeg	
	Sys 1	Sys 2	Sys 1	Sys 2
LISA-7B	70.2	63.4	46.6	43.3
LLaVA w Seg	68.4	60.2	44.2	42.4
FaST (Ours)	<b>70.8</b>	<b>64.1</b>	46.4	<b>48.2</b>

Table 7. **System 1 and System 2 performance on reasoning segmentation tasks.** FaST achieves strong performance across both tasks, demonstrating its robustness and effectiveness in dynamic task allocation between System 1 and System 2.

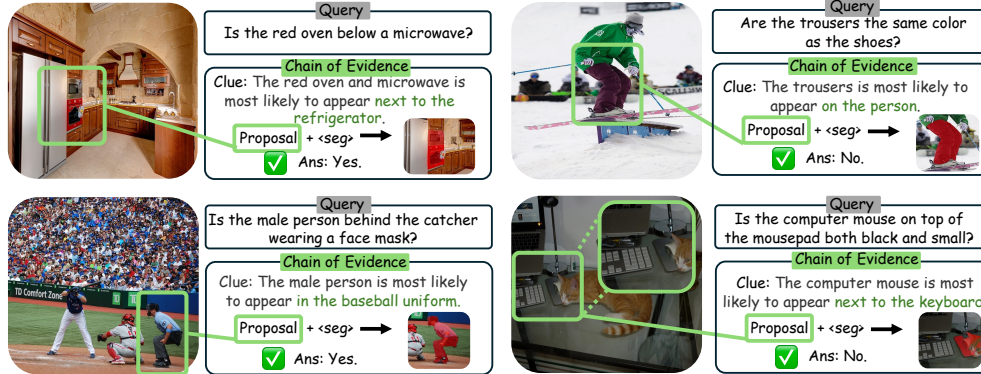


Figure 7. **More qualitative results for Visual Question Answering.**

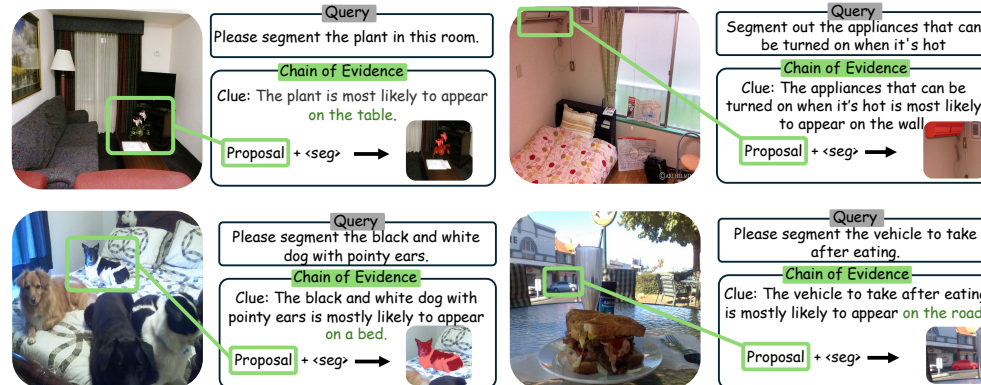


Figure 8. **More qualitative results for Visual Question Answering.**

## E FAILURE CASE

In Figure 9, we present an overview of the most notable failure cases, providing insights into the distinct patterns that lead to suboptimal outputs in our FAST model. These challenges include difficulty in triggering the *System 2* thinking mode, constructing adequate contextual clues, generating appropriate proposals, and providing accurate pixel masks. The model often fails to recognize the need for deliberate reasoning, relying instead on *System 1* thinking, which leads to incorrect responses, as seen in Figure 9a. Inadequate contextual clues generated by the switch adapter impair the model’s focus on the correct region, resulting in vague or incorrect responses, as illustrated in Fig. 9b. The proposal adapter’s inaccurate identification of regions of interest, as shown in Figure 9c, leads to proposals that do not correspond to the query. Additionally, the segmentation adapter struggles with producing precise masks, particularly for small or occluded objects, causing erroneous conclusions, as highlighted in Figure 9d. These failure cases underscore the urgent need for refinement in our FAST framework, emphasizing the importance of significantly enhancing the precision of the system switch adapter, improving contextual clue construction, and optimizing the proposal and segmentation adapters to achieve more reliable and consistently accurate responses in complex visual and textual tasks.


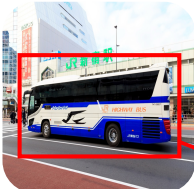
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

	<p>Query Is this a picture of Canelles de Baix (la Vall de Binya)?</p> <p>GT Answer Ans: No.</p> <p>Fast Answer ✗ Ans: Yes.</p>		<p>Query Does the flower look yellow?</p> <p>GT Answer Ans: No.</p> <p>Fast Answer ✗ Ans: Yes, the flower is yellow.</p>
---	---	--	--

(a) The model fails to trigger *System 2* thinking mode.

	<p>Query What is the rope attached to?</p> <p>GT Answer Ans: Coat.</p> <p>Chain of Evidence Clue: The rope is most likely to appear not present in the image.</p> <p>Proposal + &lt;seg&gt; → ✗ Ans: Hat.</p>		<p>Query Is this artwork displayed in St. vitus's cathedral, Prague?</p> <p>GT Answer Ans: No.</p> <p>Chain of Evidence Clue: The St. vitus's cathedral, Prague is most likely to appear near the cliff.</p> <p>Proposal + &lt;seg&gt; → ✗ Ans: Yes.</p>
---	---	---	--

(b) The model fails to construct adequate contextual clues.

	<p>Query What is the color of the dog?</p> <p>GT Answer Ans: White.</p> <p>Chain of Evidence Clue: The dog is most likely to appear next to a person.</p> <p>Proposal + &lt;seg&gt; → ✗ Ans: Black</p>		<p>Query Is the blue luggage on the left or right side of the bus?</p> <p>GT Answer Ans: Right.</p> <p>Chain of Evidence Clue: The blue luggage, bus is most likely to appear inside the bus.</p> <p>Proposal + &lt;seg&gt; → ✗ Ans: Left.</p>
---	--	--	--

(c) The model fails to generate appropriate proposals.

	<p>Query Is the lamp to the right or to the left of the car?</p> <p>GT Answer Ans: Left.</p> <p>Chain of Evidence Clue: The cars and lamp is most likely to appear on the street.</p> <p>Proposal + &lt;seg&gt; → ✗ Ans: Right.</p>		<p>Query Are there any yellow shoes in the image?</p> <p>GT Answer Ans: No.</p> <p>Chain of Evidence Clue: The yellow shoes is most likely to appear on the feet of the players.</p> <p>Proposal + &lt;seg&gt; → ✗ Ans: Yes.</p>
---	---	--	--

(d) The model fails to provide accurate pixel masks.

Figure 9. Failure cases of Our FAST system.



## 1296 F LIMITATION AND FUTURE WORK

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

While the FAST framework has demonstrated significant advancements in emulating human-like cognitive processes in visual AI through its fast and slow thinking mechanisms, several limitations warrant attention. Firstly, the system’s reliance on a predefined set of negative data for training the switch adapter may not encapsulate the full spectrum of real-world complexities. Firstly, this could lead to suboptimal performance when faced with novel or unexpected scenarios. Secondly, despite its fine-grained analysis capability, the pixel-level mask decoder might struggle with highly textured or patterned images where segmentation becomes challenging. Lastly, the generalizability of FAST across various domains and tasks necessitates further validation to ensure its robustness and reliability in diverse applications. We plan to develop advanced learning mechanisms that will allow the model to generalize more effectively beyond the predefined negative dataset. Additionally, we will focus on optimizing it for real-time applications to reduce computational overhead and response times.

For the recent large reasoning model like OpenAI o1 model, while these models leverage reinforcement learning and internal chains of thought to achieve scalability, they often require significant computational resources, making them less efficient.

In contrast, FaST is designed to prioritize multimodal reasoning with a clear focus on transparency and adaptability. The use of interpretable reasoning modes (System 1 and System 2) ensures that FaST provides insights into its decision-making processes, which is critical for applications requiring explainability. Additionally, FaST’s modular design allows it to balance computational efficiency and accuracy dynamically, making it suitable for diverse and resource-constrained environments. These strengths highlight FaST’s unique contributions and its complementary potential to scalable reasoning approaches.

## 1322 G SOCIAL IMPACTS

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

## 1333 H ETHICAL SAFEGUARDS

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

## 1344 I REPRODUCIBILITY

1345

1346

1347

1348

1349

Our FAST framework is implemented in PyTorch (Paszke et al., 2019). All the experiments are conducted on eight NVIDIA A100-80GB GPUs. Our full implementation shall be publicly released upon paper acceptance to guarantee reproducibility. The codes are available at the anonymous link <https://anonymous.4open.science/r/Sys2-LLaVA-8B0F/> for the review process.

1350 All Experiments (switch, proposal, and seg adapter) are conducted on eight NVIDIA A100-80GB  
1351 SXM GPUs<sup>5</sup>. Reproducing the fine-tuning process would require approximately 15 A100 GPU  
1352 days.

1353

1354

## 1355 J LICENSES FOR EXISTING ASSETS

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

All the methods we used for comparison are publicly available for academic usage. The switch adapter is implemented based on the released code (<https://github.com/penghao-wu/vstar>) with an MIT license. The proposal adapter is implemented on the released code (<https://github.com/dongyh20/Chain-of-Spot>) with an Apache-2.0 license. The seg adapter is implemented on the released code (<https://github.com/dvlab-research/LISA>) with an Apache-2.0 license.

---

<sup>5</sup><https://www.nvidia.com/en-sg/data-center/a100/>