

---

# Rethinking Visual Prompting for Multimodal Large Language Models with External Knowledge

---

Yuanze Lin<sup>✉\*</sup>   Yunsheng Li<sup>♣</sup>   Dongdong Chen<sup>♣</sup>  
Weijian Xu<sup>♣</sup>   Ronald Clark<sup>♣</sup>   Philip Torr<sup>♣</sup>   Lu Yuan<sup>♣</sup>  
♣ University of Oxford   ♣ Microsoft

## Abstract

In recent years, multimodal large language models (MLLMs) have made significant strides by training on vast high-quality image-text datasets, enabling them to generally understand images well. However, the inherent difficulty in explicitly conveying fine-grained or spatially dense information in text, such as masks, poses a challenge for MLLMs, limiting their ability to answer questions requiring an understanding of detailed or localized visual elements. Drawing inspiration from the Retrieval-Augmented Generation (RAG) concept, this paper proposes a new visual prompt approach to integrate fine-grained external knowledge, gleaned from specialized vision models (*e.g.*, instance segmentation/OCR models), into MLLMs. This is a promising yet underexplored direction for enhancing MLLMs' performance. Our approach diverges from concurrent works, which transform external knowledge into additional text prompts, necessitating the model to indirectly learn the correspondence between visual content and text coordinates. Instead, we propose embedding fine-grained knowledge information directly into a spatial embedding map as a visual prompt. This design can be effortlessly incorporated into various MLLMs, such as LLaVA and Mipha, considerably improving their visual understanding performance. Through rigorous experiments, we demonstrate that our method can enhance MLLM performance across nine benchmarks, amplifying their fine-grained context-aware capabilities.

## 1 Introduction

The advancement of large language models (LLMs) [56, 42, 43, 17] has revolutionized how machines process and generate human-like text, demonstrating remarkable abilities in reasoning, translation, and contextual understanding. The integration of language and vision into unified models, such as GPT-4V [41], represents a significant leap forward in enabling machines to understand and interact with the world in a manner akin to human cognition. As these models continue to evolve, they promise to further blur the lines between human and machine cognition, opening new frontiers in AI research and application [32, 51, 46, 33, 55, 31].

Despite their remarkable capabilities, most of the MLLMs (shown in Figure 1 (a)) trained with image-text pairs still often struggle in fine-grained multimodal comprehension capacities, *e.g.*, correctly count objects or output precise location of one specific object. This is partially because of the lack of high-quality data with exceptionally fine-grained text description. More importantly, text itself has the inherent difficulty in accurately conveying highly fine-grained or spatially dense information. As a result, current MLLMs often fail to accurately interpret pixel-level visual content of localized regions within an image, which in return harms the overall comprehension capacity for the image and thereby causes the notorious “hallucination” problem [25].

---

\*Work done during an internship at Microsoft Redmond.

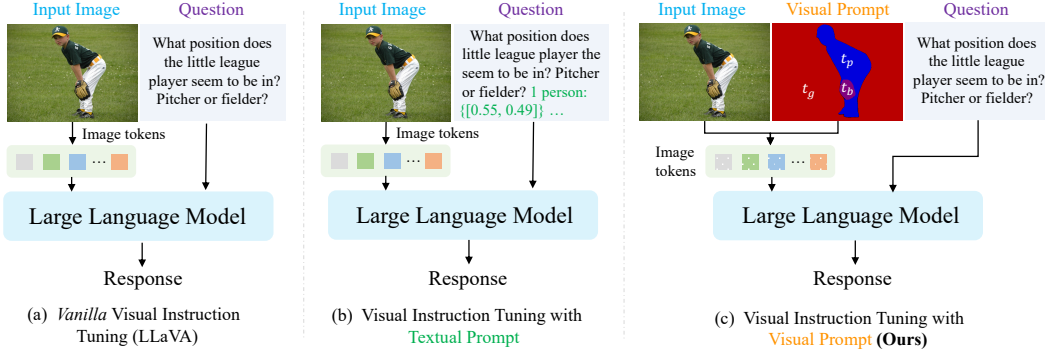


Figure 1: **Different training paradigms.** (a) means the original visual instruction tuning of LLaVA [34]. (b) denotes visual instruction tuning with external textual prompts [23] (e.g., 1 person and the center coordinates of its bounding box: [0.55,0.49]), note that we neglect the template prefix of textual prompts for visualization. (c) is the proposed auxiliary visual prompt, which is a feature map composed with different object regions. For each pixel, it is filled out with the textual embedding of the corresponding *categories* or *OCR text* ( $t_g$ ,  $t_p$  and  $t_b$  in the example visual prompt mean the textual embeddings of *grass*, *person* and *baseball glove*).

To tackle this challenge, one line of work [6, 62, 5] explicitly integrates region coordinates information into the text prompt and trains on specialized region-level chatting data. However, this still demands that the model implicitly learns to understand coordinates and establish connections with visual content, thereby increasing the learning complexity. Another line of work [50, 61, 30] proposes incorporating Region of Interest (ROI) features directly into model learning, necessitating bespoke model architectures. In contrast to these approaches, rather than starting from scratch to learn region information, this paper explores leveraging finely-grained recognition outcomes directly obtainable from existing vision models as external knowledge for MLLMs, inspired by the RAG concept. Concurrent with our work, one recent approach [23] introduces external knowledge, such as regional coordinates from object detection and Optical Character Recognition (OCR) technologies, into MLLMs (shown in Figure 1 (b)), helping understand localized multimodal content. However, this method still integrates external knowledge through the text prompt, mandating implicit learning of content-to-coordinate correspondence by the model. Furthermore, it lacks support for more nuanced external knowledge, such as instance masks.

In this paper, we propose a new visual prompt paradigm to insert external knowledge, e.g., localized information, into MLLMs addressing the challenge of fine-grained multimodal content correspondence. As illustrated in Figure 1 (c), the core idea is, rather than treating local context information as a part of text prompts, we embed them directly within the visual prompts. Specifically, we start by leveraging panoptic segmentation [60] and OCR detection [15] models, and a pre-trained text encoder to generate pixel-wise text embeddings, which are served as the local context information for MLLMs. Subsequently, we extend the original visual prompts by adding the newly generated context information in a spatial-wise manner. This integrated prompt is then assimilated into MLLMs, improving fine-grained visual content comprehension. Consequently, our approach is capable of enabling MLLMs to discern contexts in the pixel-level space and improve their performance.

With the proposed visual prompt paradigm, we train a bunch of MLLMs on the LLaVA-1.5 datasets [34]. The experimental results show that, even with 3 billion parameters, our method improves upon the leading open-source MLLMs such as LLaVA-1.5 [35, 34] and Qwen-VL [3], without needing additional training data. Remarkably, our models showcase superior performance across a wide array of benchmarks when compared to the 7-billion MLLM variants, including LLaVA-1.5, Qwen-VL, and InstructBLIP [13], and in some instances, even outperform their 13-billion MLLM counterparts. Our experimental results confirm the significance of integrating our proposed prompt approach with MLLMs to enhance cognitive capabilities.

The contributions can be summarized as follows:

- We systematically investigate integrating localized information into MLLMs. Empirical findings suggest that our proposed visual prompt significantly outperforms the previous prompt paradigm relying solely on textual prompts containing coordinates.

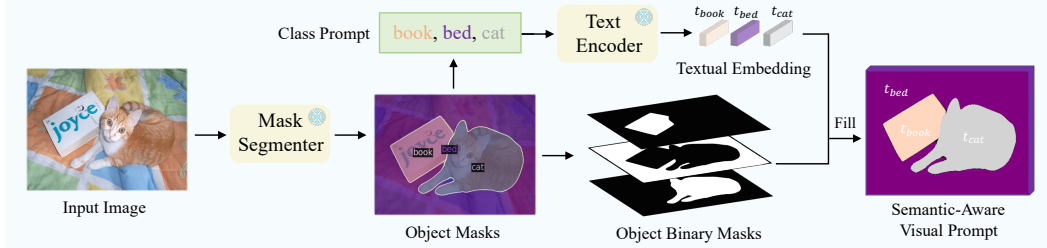


Figure 2: **Auxiliary visual prompt generation.** It firstly generates the panoptic segmentation masks [60] for the input image, there’s a class category for each mask region, then we can obtain the textual embeddings (e.g.,  $t_{\text{book}}$ ,  $t_{\text{bed}}$  and  $t_{\text{cat}}$ ) through a pre-trained text encoder for all the classes (e.g., **book**, **bed**, **cat**). Finally, the auxiliary visual prompt can be generated by concatenating these textual embeddings within the corresponding mask regions together. Note that we can also adopt the OCR model [15] to obtain the texts and the regions, we don’t display it here for clearer explanation.

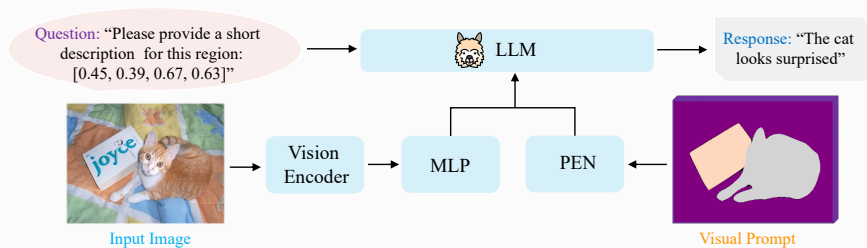


Figure 3: **The illustration of visual instruction tuning with the generated visual prompt.** Our proposed visual prompt can be easily combined with existing multimodal large language models (e.g., LLaVA [34]), note that *PEN* means prompt embedding network.

- We propose to integrate contextual embeddings within local contours (e.g., object masks) as the visual prompt, which facilitates the establishment of correlations between image pixels and contexts, thereby enhancing the fine-grained cognitive capabilities of various MLLMs across a spectrum of benchmarks.
- Based on our proposed approach, our model with 3B parameters surpasses both existing 7B and 13B models across diverse benchmarks, all without the need for extra training data.

## 2 Related Work

**Large Language Models.** The initial potential of large language models (LLMs) was showcased by foundational works like BERT [14] and GPT [47]. They sparked a wave of scaling efforts, leading to a range of influential projects, such as T5 [49], GPT-3 [4], Flan-T5 [12], and PaLM [9]. As the volume of training data expanded and the dimensions of model parameters grew, these scaling endeavors led to the creation of ChatGPT [40, 44]. Models like LLaMA [56] and GPT-4 [43] have been trained on extensive corpora and demonstrated remarkable capabilities in diverse cognitive tasks. Additionally, lightweight LLMs with fewer than 3B parameters, *i.e.*, Phi [2, 39] and StableLM-2 [53] have shown performance comparable to larger models [8]. In our work, we adopt Phi-2 [39] and Vicuna-7B [8] as our language backbone.

**Multimodal Large Language Models.** Influenced by the success of instruction tuning from LLM, LLaVA [35] and MiniGPT-4 [64] have adopted visual instruction tuning to improve LLMs’ interaction with visual data, yielding impressive outcomes. Kosmos-2 [45] and Shikra [7] have advanced MLLMs by enhancing visual comprehension capabilities. while works like LLaVA-Phi [66], MobileVLM [11] and Bunny [20] mainly focus on optimizing training recipes and architecture design for lightweight MLLMs. To solve the challenge of understanding fine-grained information in images, existing approaches propose to learn coordinate representations [6, 7, 62] and Region of Interest (ROI) features [45, 61], which use inflexible visual referral formats or necessitate the collection of

region-level training data. On the contrary, we focus on utilizing external knowledge to improve the fine-grained vision-language alignment for MLLMs without collecting extra chatting data.

**Prompting Multimodal Large Language Models.** Inspired by the ability of GPT-4V [41] to process diverse inputs, ViP-LLaVA [5] collects a visual prompt instruction dataset containing various visual prompts, *e.g.*, scribbles and arrows, for MLLM fine-tuning. Contemporary to our work, [23] has offered advanced insights in prompting MLLMs through external knowledge, which introduces bounding box and OCR coordinates into text prompt, however, it’s still challenging to interpret the pixel-level contexts. In this paper, we investigate how to efficiently utilize external knowledge to enhance multimodal fine-grained alignment of MLLMs and introduce a novel visual prompt paradigm incorporating pixel-level contextual information.

### 3 Proposed Method

In this section, we propose a new visual prompt paradigm that integrates local external information to enhance the capability of MLLMs. In section 3.1, we outline the design of the auxiliary visual prompt that contains local contextual information. Using the auxiliary visual prompt, in section 3.2, we further embed it into MLLMs by merging it with the original visual tokens. Finally, we briefly introduce the details of training in section 3.3.

#### 3.1 Auxiliary Visual Prompt with External Knowledge

In this section, we propose a method to generate local contextual external knowledge to assist MLLMs. In contrast to [23], which focuses solely on object detection and OCR information and integrates them as part of the text prompt, we enhance the granularity of local external knowledge by leveraging a panoptic segmentation model. Additionally, we continue to utilize an OCR model but transform both types of external knowledge into pixel-wise embeddings. Further details are provided below.

As shown in Figure 2, given the input image  $I \in \mathbb{R}^{3 \times H \times W}$ , we can obtain the fine-grained external knowledge by an off-the-shelf panoptic segmentation model [60] and an OCR model [15]. The generation of the external knowledge can be expressed as:

$$\{M_j, C_j\}_{j=1}^{N_s} = f_{\text{seg}}(I), \quad \{B_j, T_j\}_{j=1}^{N_o} = f_{\text{ocr}}(I), \quad (1)$$

where  $f_{\text{seg}}(\cdot)$  and  $f_{\text{ocr}}(\cdot)$  mean panoptic segmentation and optical character recognition (OCR) models,  $N_s$  and  $N_o$  are the numbers of detected mask regions and OCR bounding boxes.  $\{M_j, C_j\}_{j=1}^{N_s}$  is the set of mask regions and the corresponding classes, and  $\{B_j, T_j\}_{j=1}^{N_o}$  represents the set of detected OCR bounding boxes and texts.

With the detected classes  $\{C_j\}_{j=1}^{N_s}$  and OCR texts  $\{T_j\}_{j=1}^{N_o}$ , a pre-trained text encoder ( $f_{\text{text}}(\cdot)$ ) is leveraged to generate the texture embeddings as:

$$\begin{aligned} \mathcal{T}_s &= \{t_1, \dots, t_{N_s}\} = \{f_{\text{text}}(C_1), \dots, f_{\text{text}}(C_{N_s})\}, \\ \mathcal{T}_o &= \{\hat{t}_1, \dots, \hat{t}_{N_o}\} = \{(f_{\text{text}}(T_1), \dots, f_{\text{text}}(T_{N_o}))\}, \end{aligned} \quad (2)$$

where  $t_i \in \mathbb{R}^{1 \times d} (1 \leq i \leq N_s)$  and  $\hat{t}_i \in \mathbb{R}^{1 \times d} (1 \leq i \leq N_o)$  denote the  $i$ th textual embedding vector of the classes for the detected mask region and OCR texts respectively, while  $d$  is the embedding dimension.

In order to generate a pixel-wise visual prompt for the external knowledge instead of a pure text description for the regions with coordinates and category names, the auxiliary visual prompt is initialized as a zero tensor  $\mathcal{P} \in \mathbb{R}^{H \times W \times d}$  and then filled with the newly generated texture embeddings for the external knowledge as:

$$\begin{aligned} \mathcal{P}_{j,k} &= \begin{cases} t_u & \text{if } (j,k) \in M_u \\ \mathcal{P}_{j,k} & \text{otherwise} \end{cases} \quad \forall u \in \{1, \dots, N_s\}, \\ \mathcal{P}_{j,k} &= \mathcal{P}_{j,k} + \begin{cases} \hat{t}_v & \text{if } (j,k) \in B_v \\ 0 & \text{otherwise} \end{cases} \quad \forall v \in \{1, \dots, N_o\}. \end{aligned} \quad (3)$$



Note, for some regions, if the confidence of the class prediction given by the segmentation model is low or the OCR model fails to detect any text, we leave the region area with zero values. For the regions that are occupied by both model, we simply add the text embeddings directly. We leave the investigation of more refined fusion techniques to future research.

With the auxiliary visual prompt containing pixel-level local contextual information from panoptic segmentation and OCR models, MLLMs can effectively capture finer-grained features. The next challenge is to establish a clearer connection between the newly generated external knowledge and the original image feature. This will help alleviate the model’s difficulties in learning their relationship effectively.

### 3.2 Visual Prompt Infusion

In this section, we introduce the visual prompt infusion that incorporates the proposed auxiliary visual prompts into the MLLMs. Previous methods [23] choose to append the external knowledge (embeddings for object category and its coordinates) to the text prompts, which requires the model to learn the correspondence of visual content within the specified coordinates encoded in the external knowledge and, as a result, increasing the difficulties of the learning process of the model. To address this challenge, we propose to merge the auxiliary visual prompt directly with the image features in a pixel-wise manner.

Specifically, as shown in Figure 3, the image tokens are first generated via an image encoder  $f_{\text{img}}(\cdot)$  and an MLP projector ( $f_{\text{MLP}}(\cdot)$ ):

$$\mathcal{F}_v = f_{\text{MLP}}(f_{\text{img}}(I)), \quad (4)$$

where  $\mathcal{F}_v \in \mathbb{R}^{N_v \times d_v}$ ,  $N_v$  and  $d_v$  represent the number of image tokens and the embedding dimension. Then, the auxiliary visual prompt is further processed by a prompt embedding network (PEN) as

$$\mathcal{F}_p = f_{\text{PEN}}(\mathcal{P}). \quad (5)$$

For the prompt embedding network, we employ three convolutional layers, with an activation layer (ReLU) inserted between each pair of them. This network primarily serves to align the feature space and spatial size between the image tokens and the auxiliary visual prompts.

When combining the image tokens and the processed auxiliary visual prompt, we mainly consider two options, both of which operate pixel-wise. **(1) feature fusion:**  $\hat{\mathcal{F}}_v = f(\text{Concat}(\mathcal{F}_v, \mathcal{F}_p))$ , where  $f$  is a linear layer that maps the embedding  $\mathbb{R}^{N_v \times d_{2v}} \rightarrow \mathbb{R}^{N_v \times d_v}$  to maintain the total number of image tokens unchanged; **(2) feature addition,**  $\hat{\mathcal{F}}_v = \mathcal{F}_v + \mathcal{F}_p$ , which sums the two types of features directly.

The advantages of the pixel-wise infusion for both options facilitate the model’s comprehension of the correspondence between external knowledge and original visual features. This explicit guidance enables the model to easily understand the pixel categories as well as the potential OCR text description it conveys. Consequently, it aids the model in disambiguating complex scenes, accentuating salient features, and distinguishing finer objects.

### 3.3 Training

Training MLLMs involves predicting responses based on multimodal inputs using an autoregressive approach. The objective is to maximize the probability of generating tokens that match the ground-truth answer  $Y_a$ . With the new visual embedding  $\hat{\mathcal{F}}_v$ , this can be mathematically expressed as follows:

$$P(Y_a | \hat{\mathcal{F}}_v, \mathcal{F}_t) = \prod_{i=1}^L P_{\theta}(y_i | \hat{\mathcal{F}}_v, \mathcal{F}_t, Y_{a, < i}). \quad (6)$$

Here,  $L$  represents the sequence length of the ground truth answer  $Y_a$ ,  $\theta$  means the trainable parameters.  $Y_{a, < i}$  represents all the answer tokens preceding the current prediction token  $x_i$ , where  $i$  denotes the step in the sequence of text token generation.  $\mathcal{F}_t \in \mathbb{R}^{N_t \times d_t}$  is the token embedding of the input question,  $N_t$  and  $d_t$  denote the number of text tokens and token embedding dimension. By infusing these enriched visual cues into the training pipeline, MLLMs can develop a more

Table 1: The ablation study of different visual prompts. *Mipha-3B* is the baseline with standard visual & text prompt. *Mipha-3B+LAF* denotes using textual prompting with LoRA Augmented Fine-tuning following (LAF) [23]. *feature fusion* and *feature addition* represent two prompt fusion methods we use to insert the auxiliary visual prompt to the original image features.

Method	VQAv2	GQA	SQA <sup>1</sup>	VQA <sup>T</sup>	MME-P	MME-C	MMB	MM-Vet	POPE
Mipha-3B	81.3	63.9	70.9	56.6	1488.9	295.0	69.7	32.1	86.7
Mipha-3B + LAF	81.6 $\uparrow$	62.6 $\downarrow$	71.4 $\uparrow$	57.8 $\uparrow$	1472.3 $\downarrow$	356.8 $\uparrow$	71.0 $\uparrow$	34.8 $\uparrow$	88.5 $\uparrow$
Ours (feature fusion)	81.9 $\uparrow$	64.8 $\uparrow$	71.6 $\uparrow$	57.6 $\uparrow$	1493.5 $\uparrow$	345.5 $\uparrow$	71.3 $\uparrow$	34.3 $\uparrow$	88.5 $\uparrow$
Ours (feature addition)	<b>82.4<math>\uparrow</math></b>	<b>65.3<math>\uparrow</math></b>	<b>71.8<math>\uparrow</math></b>	<b>57.8<math>\uparrow</math></b>	<b>1501.2<math>\uparrow</math></b>	<b>369.1<math>\uparrow</math></b>	<b>71.5<math>\uparrow</math></b>	<b>35.1<math>\uparrow</math></b>	<b>88.7<math>\uparrow</math></b>

Table 2: The ablation study of using different vision encoders, i.e., SigLIP v.s. CLIP. The results of *Mipha-3B* on CLIP are from [65].

Method	Vis Enc	VQAv2	GQA	SQA <sup>1</sup>	VQA <sup>T</sup>	MME-P	MME-C	MMB	MM-Vet	POPE
Mipha-3B	CLIP	78.6	62.3	68.2	53.0	-	-	68.4	31.0	86.9
Ours	CLIP	<b>79.7<math>\uparrow</math></b>	<b>63.7<math>\uparrow</math></b>	<b>70.1<math>\uparrow</math></b>	<b>54.8<math>\uparrow</math></b>	<b>1445.5</b>	<b>308.4</b>	<b>70.1<math>\uparrow</math></b>	<b>33.7<math>\uparrow</math></b>	<b>88.8<math>\uparrow</math></b>
Mipha-3B	SigLIP	81.3	63.9	70.9	56.6	1488.9	295.0	69.7	32.1	86.7
Ours	SigLIP	<b>82.4<math>\uparrow</math></b>	<b>65.3<math>\uparrow</math></b>	<b>71.8<math>\uparrow</math></b>	<b>57.8<math>\uparrow</math></b>	<b>1501.2<math>\uparrow</math></b>	<b>369.1<math>\uparrow</math></b>	<b>71.5<math>\uparrow</math></b>	<b>35.1<math>\uparrow</math></b>	<b>88.7<math>\uparrow</math></b>

comprehensive understanding of visual content, leading to better alignment between visual and textual representations. To accelerate the training process, we follow LoRA Augmented Training (LAF) strategy [23] to perform fine-tuning on Mipha-3B [65] and LLaVA-1.5 [34] using LoRA [21].

## 4 Experiment

In this section, we conduct a comprehensive comparison of our method with existing state-of-the-art (SOTA) multimodal models. Additionally, we perform a series of ablation studies to further validate the proposed method. Finally, we provide visualization examples for in-depth analysis.

**Models.** For the vision encoder, we adopt SigLIP-384px [59] for experiments. We leverage Phi-2-2.7B [39] and Vicuna-7B [8] model as the language decoder. For the multimodal projector, same as LLaVA [34], we adopt a two-layer MLP. We use OpenSeed [60] and PaddleOCRv2 [15] to generate the per-pixel externally knowledge for pixel class and OCR text, and leverage UAE-Large-V1 [27] to extract the textual embedding.

**Training Setting.** We fine-tune the models on LLaVA-Instruct-150K dataset [34] using LoRA [21] for 1 epoch, at a learning rate of  $2e-4$  and a batch size of 256 on  $32 \times V100$  32GB GPUs. For the setting of LoRA, we set LoRA rank to be 128 and LoRA’s hyperparameter  $\alpha$  as 256. Note that we fix all the weights of pre-trained modules, i.e., vision encoder, language encoder and MLP, during training. Our models’ weights are initialized from Mipha-3B [65] and LLava-7B [34].

**Benchmarks and Baselines.** We evaluate our approach using 9 popular benchmarks to comprehensively assess its multimodal capabilities. These benchmarks include: VQA-v2 test-dev split [18], GQA test-dev-balanced split [22], ScienceQA-IMG test split [38], MME perception [16], MME cognition [16], MMBench test split [37], MM-Vet test split [58], TextVQA [52], and POPE [28].

We compare our results with a bunch of state-of-the-art multimodal large language models (MLLMs): BLIP-2 [26], InstructBLIP [13], Shikra-13B [7], IDEFICS80/9B [24], Qwen-VL [3], mPLUG-Owl2 [57], LLaVA-v1.5-13/7B [34], LAF-7B [23], and multimodal small language models (MSLMs) [65]: MobileVLM [11], LLaVA-Phi [66], MC-LLaVA [1], Imp-v1 [54], MoE-LLaVA-3.6B [29], TinyLLaVA-share-Sig-Phi [63], Bunny [20] and Mipha [65].

### 4.1 Ablation Studies

In this section, we conduct an ablation study to assess the effectiveness of the proposed approach. By default, the experiments are conducted using Mipha-3B [65] with Phi-2 [39] as the language backbone unless otherwise specified.

**Visual input example, context-related questions:**



What is funny about this image?



Why are the persons wearing black clothes?

LLaVA-7B	The funny aspect of this image is the juxtaposition of two women pointing fingers at a cat, which is sitting at a table with a plate of food. The cat appears to be looking at the women with an angry expression, as if it is reacting to their pointing fingers. This creates a humorous and unexpected scene, as it is not common to see a cat reacting to human gestures in such a manner. The image captures a light-hearted and amusing moment between the women and the cat.	The persons in the image are wearing black clothes, which is a common attire for religious ceremonies or events. The group of people, including a man and a woman, are walking down the street, possibly participating in a religious procession or a funeral. The black clothing signifies their respect and adherence to the religious or cultural traditions associated with the event they are attending.
<b>Ours</b>	The funny aspect of this image is the contrast between the two women pointing at the cat and the cat’s reaction to the situation. <i>One woman is making a funny face, while the cat is sitting on a table with a plate of food, seemingly unbothered by the attention. The cat’s calm demeanor in the face of the women’s playful behavior creates a humorous scene, as the women try to entertain themselves and the cat seems to be enjoying the attention without reacting to their attempts.</i>	The persons in the image are wearing black clothes, <i>which are typically associated with mourning or a period of grieving. This suggests that they might be attending a funeral or a memorial service for someone who has passed away. The black attire signifies respect and solidarity with the bereaved family and friends during this difficult time. It is a common practice in many cultures to wear black clothes to express sympathy and support for those persons who have experienced a loss.</i>

Table 3: The challenging examples on LLaVA-1.5, our approach can provide reasonable reponses for the tricky questions which need to fully understand the contexts of the images.

Table 4: The ablation study of introducing OCR information into the visual prompt.

Method	VQAv2	GQA	SQA <sup>1</sup>	VQA <sup>T</sup>	MME-P	MME-C	MMB	MM-Vet	POPE
Mipha-3B	81.3	63.9	70.9	56.6	1488.9	295.0	69.7	32.1	86.7
Ours (w/o OCR)	81.9 $\uparrow$	64.7 $\uparrow$	71.3 $\uparrow$	57.1 $\uparrow$	1498.3 $\uparrow$	355.2 $\uparrow$	70.8 $\uparrow$	34.0 $\uparrow$	87.9 $\uparrow$
Ours	<b>82.4<math>\uparrow</math></b>	<b>65.3<math>\uparrow</math></b>	<b>71.8<math>\uparrow</math></b>	<b>57.8<math>\uparrow</math></b>	<b>1501.2<math>\uparrow</math></b>	<b>369.1<math>\uparrow</math></b>	<b>71.5<math>\uparrow</math></b>	<b>35.1<math>\uparrow</math></b>	<b>88.7<math>\uparrow</math></b>

**Prompting MLLMs with Different Approaches.** In Table 1, we present the results of the ablation study for four different prompting strategies: (1) Mihpa-3B baselines with vanilla text prompt, as used by LLaVA-1.5 [34]. (2) Mihpa-3B + LAF proposed in [23] that appends external local contextual knowledge to the text prompts. (3) The proposed auxiliary visual prompt inserted via feature fusion. (4) The proposed auxiliary visual prompt added via feature addition.

From Table 1, we note that compared to the baseline (1) with vanilla prompts, both proposed fusion strategies (3) and (4) exhibit a significant improvement. This suggests that external knowledge is indeed beneficial in enhancing the capabilities of MLLMs. In comparison to Mihpa-3B+LAF (2), which inserts external local contextual knowledge into the text prompt, (4) outperforms it in 8 out of 9 benchmarks, notably for GQA [22] and MME-P [16]. This implies that explicitly linking external local knowledge to the original visual features reduces the model’s learning burden in establishing spatial relationships, consequently enhancing performance. Furthermore, we empirically observe that directly adding auxiliary visual prompts yields slightly better results than concatenation. Therefore, we adopt feature addition as our default setting for subsequent experiments.

**The Effect of Using Different Vision Encoders.** In Table 2, we further ablate the effectiveness brought by different vision encoders, i.e., CLIP [48] v.s. SigLIP [59]. From the results, we can draw two conclusions. First, for both vision encoders, our methods have consistent improvement compared

Table 5: The comprehensive multi-modal evaluation across 9 distinct benchmarks to thoroughly assess model performance: VQAv2 [18], GQA [22], SQA<sup>1</sup>: ScienceQA-IMG [38], VQA<sup>T</sup>: TextVQA [52], MME-P: MME Perception [16], MME-C: MME Cognition [16], MMB: MM-Bench [37], MM-Vet [58], and POPE [28]. The included proprietary in-house data not publicly accessible, denoted as †. The image resolution used by the visual backbone is indicated in the column labeled *Res.*, while the columns *PT* and *IT* represent the data sizes in the pretraining and visual instruction tuning stages, respectively.

Method	LM	Res.	PT	IT	VQAv2	GQA	SQA <sup>1</sup>	VQA <sup>T</sup>	MME-P	MME-C	MMB	MM-Vet	POPE
Multimodal Large Language Models													
BLIP-2 [26]	Vicuna (13B)	224	129M	-	65.0	41.0	61.0	42.5	1293.8	290.0	-	22.4	85.3
InstructBLIP [13]	Vicuna (7B)	224	129M	1.2M	-	49.2	60.5	50.1	-	-	36	26.2	-
InstructBLIP [13]	Vicuna (13B)	224	129M	1.2M	-	49.5	63.1	50.7	1212.8	291.8	-	25.6	78.9
Shikra [7]	Vicuna (13B)	224	600K	5.5M	77.4	-	-	-	-	-	58.8	-	-
IDEFICS-9B [24]	LLaMA (7B)	224	353M	1M	50.9	38.4	-	25.9	-	-	48.2	-	-
IDEFICS-80B [24]	LLaMA (65B)	224	353M	1M	60.0	45.2	-	30.9	-	-	54.5	-	-
Qwen-VL [3]	Qwen (7B)	448	1.4B <sup>†</sup>	50M <sup>†</sup>	78.8	59.3	67.1	<b>63.8</b>	-	-	38.2	-	-
Qwen-VL-Chat [3]	Qwen (7B)	448	1.4B <sup>†</sup>	50M <sup>†</sup>	78.2	57.5	68.2	61.5	1487.5	360.7	60.6	-	-
mPLUG-Owl2 [57]	LLaMA (7B)	448	400M	1.23M	79.4	56.1	68.7	58.2	1450.2	313.2	64.5	<b>36.2</b>	85.8
LLaVA-1.5 [34]	Vicuna (7B)	336	558K	665K	78.5	62.0	66.8	58.2	1510.7	316.1	64.3	30.5	85.9
LAF-7B [23]	Vicuna (7B)	336	558K	665K	79.0	60.5	-	60.1	1482.7	397.9	67.3	35.2	88.9
<b>LLaVA-1.5<sup>+</sup>(Ours)</b>	Vicuna(7B)	336	558K	665K	<b>79.8<sup>†</sup></b>	<b>63.3<sup>†</sup></b>	<b>69.5<sup>†</sup></b>	<b>59.8<sup>†</sup></b>	<b>1515.3<sup>†</sup></b>	<b>399.5<sup>†</sup></b>	<b>67.6<sup>†</sup></b>	<b>34.9<sup>†</sup></b>	<b>88.9<sup>†</sup></b>
Multimodal Small Language Models													
MobileVLM-1.7B [10]	M-LLaMA (1.4B)	336	558K	665K	-	56.1	57.3	41.5	1196.2	-	53.2	-	84.5
MobileVLM-3B [10]	M-LLaMA (2.7B)	336	558K	665K	-	59.0	61.2	47.5	1288.9	-	59.6	-	84.9
MobileVLM-v2-1.7B [11]	M-LLaMA (1.4B)	336	1.2M	2.4M	-	59.3	66.7	52.1	1302.8	-	57.7	-	84.3
MobileVLM-v2-3B [11]	M-LLaMA (2.7B)	336	1.2M	2.4M	-	61.1	70.0	57.5	1440.5	-	63.2	-	84.7
LLaVA-Phi [66]	Phi-2 (2.7B)	336	558k	665K	71.4	-	68.4	48.6	1335.1	-	59.8	28.9	85.0
MC-LLaVA [1]	Phi-2 (2.7B)	384	558k	665K	64.2	49.6	-	38.6	-	-	-	-	80.6
Imp-v1 [54]	Phi-2 (2.7B)	384	558k	665K	79.5	58.6	70.0	59.4	1434.0	-	66.5	33.1	88.0
MoE-LLaVA-3.6B [29]	Phi-2 (2.7B)	384	558k	1.59M	79.9	62.6	70.3	57.0	1431.3	-	68.0	35.9	85.7
TinyLLaVA [63]	Phi-2 (2.7B)	384	1.2M	665k	79.9	62.0	69.1	59.1	1464.9	-	66.9	32.0	86.4
Bunny-3B [20]	Phi-2 (2.7B)	384	2M	695K	79.8	62.5	70.9	-	1488.8	289.3	68.6	-	86.8
Mipha-3B [65]	Phi-2 (2.7B)	384	558K	665K	81.3	63.9	70.9	56.6	1488.9	295.0	69.7	32.1	86.7
<b>Mipha-3B<sup>+</sup>(Ours)</b>	Phi-2 (2.7B)	384	558K	665K	<b>82.4<sup>†</sup></b>	<b>65.3<sup>†</sup></b>	<b>71.8<sup>†</sup></b>	<b>57.8<sup>†</sup></b>	<b>1501.2<sup>†</sup></b>	<b>369.1<sup>†</sup></b>	<b>71.5<sup>†</sup></b>	<b>35.1<sup>†</sup></b>	<b>88.7<sup>†</sup></b>

to the baselines, which validates the stability of our methods. Second, SigLIP emerges as the stronger vision encoder when compared to CLIP. Therefore, we opt to utilize SigLIP as the default vision encoder in subsequent sections.

**The Effect of Introducing OCR Information Into Visual Prompt.** In Table 4, we perform the ablation of using OCR information or not, we can conclude that adopting the information from OCR can further improve the model’s overall performance, especially, when incorporating OCR information for some text-specific tasks (*e.g.*, TextVQA [52] and MM-Vet [58]), it can achieve remarkable performance boosts.

## 4.2 Main Results

In Table 5, we compare our methods with other state-of-the-art (SOTA) models. We divide the table into sections for language models smaller than 3B and those beyond 7B to provide a clearer comparison. From the results, we observe that our model achieves the best performance on 7 out of 9 benchmarks for larger language models (>7B) and attains the highest accuracy on 7 out of 9 benchmarks for relatively smaller language models (<3B). Note that, in Table 5, some models, *e.g.* Shikra-13B [7], Qwen-VL [3], are trained with million or billion level data, while our model is only trained on the dataset used by LLaVA-1.5 without any extra data for neither pre-training nor fine-tuning, which highlights the exceptional multimodal understanding and reasoning capabilities of our models. In addition, on top of the LLaVA-1.5 framework, our approach can bring more remarkable and consistent improvement on all benchmarks compared with LAF [23]. It justifies the proposed infusion strategy, which involves inserting external knowledge in a pixel-wise manner directly into the visual features, as being more effective than appending it to the text prompt [23].

## 4.3 Quantitative Result Analysis

We present visualization results in Table 3 and 6 to further illustrate the improvement of our model in terms of both global image understanding and local object and text recognition. Table 3 demonstrates that compared to LLaVA-1.5 7B [34], our approach generates more detailed and contextually relevant responses, *e.g.*, “The cat’s calm demeanor in the face of the women’s playful behavior” for the left

---

**Visual input example, relationship-aware and text-related questions:**

---



On the right desk, what is to the left of the laptop?



What are all the scene text in the image?

LLaVA-1.5 7B	To the left of the laptop, there is a stack of books on the desk.	The scene text in the image is "Shakespeare's Dramas, Sonnets, & Poems".
<b>Ours</b>	There is a <b>desk lamp</b> to the left of the laptop on the right desk.	The scene text in the image includes the title " <b>Shakespeare's Comedies, Histories, and Tragedies</b> ".

Table 6: The challenging examples on LLaVA-1.5. Our approach can generate accurate responses for text-related questions.

example; “mourning or a period of grieving” and “express sympathy and support for those persons who have experienced a loss” for the right example, which all need a deeper understanding of the global image context. Meanwhile, Table 6 highlights our model’s ability to correctly recognize objects’ spatial relationships, such as between a “desk lamp” and a “laptop” from the left image, and exhibit stronger OCR capability in detecting words written on a book from the right image, compared to LLaVA-1.5 7B [34]. These visualizations validate the effectiveness of our proposed methods and support the conclusion that incorporating external local contextual information in a spatial-wise manner improves the model’s fine-grained recognition capability and enhances its overall ability for global image understanding. Note that we’ve shown more ablation study experiments and visualization result analysis in the Appendix.

## 5 Limitations and Broader Impact

Our method relies on pre-trained models for panoptic segmentation and OCR detection in a zero-shot manner. The performance of these models will significantly impact the performance of our proposed method, particularly when there is a substantial domain gap between the images from specific benchmarks and the training set of the segmentation or OCR models. While the proposed approach holds promise for significantly enhancing the cognitive capabilities of multimodal models and may inspire new methodologies and techniques in the development of robust multimodal AI systems, users must be aware of potential negative societal impacts. For instance, biases may manifest in various forms; for example, biased responses may be generated by the model if the training data of MLLMs, panoptic segmentation, and OCR detection models contain certain biases.

## 6 Conclusion

In this paper, we have proposed a method for leveraging external knowledge, such as localized contextual information, to enhance the capabilities of multimodal language models (MLLMs). To accomplish this objective, we propose extracting pixel-wise contextual information using a panoptic segmentation and OCR model, and then directly integrate this with the visual features. This enables the model to better understand both fine-grained objects and the overall global image context. Experimental results from ablations and comparisons with state-of-the-art methods demonstrate the effectiveness of our approach. We hope this paper can shed light on the importance of external knowledge for MLLMs and an effective way to leverage such knowledge.



## References

- [1] Multi-crop llava-3b, 2023. URL <https://huggingface.co/visheratin/MC-LLaVA-3b>.
- [2] Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.
- [3] 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.
- [4] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [5] Mu Cai, Haotian Liu, Siva Karthik Mustikovela, Gregory P Meyer, Yuning Chai, Dennis Park, and Yong Jae Lee. Making large multimodal models understand arbitrary visual prompts. *arXiv preprint arXiv:2312.00784*, 2023.
- [6] Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechu 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*, 2023.
- [7] Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multimodal llm’s referential dialogue magic. *arXiv preprint arXiv:2306.15195*, 2023.
- [8] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023), 2023.
- [9] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- [10] Xiangxiang Chu, Limeng Qiao, Xinyang Lin, Shuang Xu, Yang Yang, Yiming Hu, Fei Wei, Xinyu Zhang, Bo Zhang, Xiaolin Wei, et al. Mobilevlm: A fast, reproducible and strong vision language assistant for mobile devices. *arXiv preprint arXiv:2312.16886*, 2023.
- [11] Xiangxiang Chu, Limeng Qiao, Xinyu Zhang, Shuang Xu, Fei Wei, Yang Yang, Xiaofei Sun, Yiming Hu, Xinyang Lin, Bo Zhang, et al. Mobilevlm v2: Faster and stronger baseline for vision language model. *arXiv preprint arXiv:2402.03766*, 2024.
- [12] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022.
- [13] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. *Advances in Neural Information Processing Systems*, 36, 2024.
- [14] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [15] Yuning Du, Chenxia Li, Ruoyu Guo, Cheng Cui, Weiwei Liu, Jun Zhou, Bin Lu, Yehua Yang, Qiwen Liu, Xiaoguang Hu, et al. Pp-ocrv2: Bag of tricks for ultra lightweight ocr system. *arXiv preprint arXiv:2109.03144*, 2021.



- [16] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiaowu Zheng, Ke Li, Xing Sun, et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023.
- [17] Google. Google bard. <https://bard.google.com/chat/>, 2023.
- [18] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913, 2017.
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.
- [20] Muyang He, Yexin Liu, Boya Wu, Jianhao Yuan, Yueze Wang, Tiejun Huang, and Bo Zhao. Efficient multimodal learning from data-centric perspective. *arXiv preprint arXiv:2402.11530*, 2024.
- [21] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [22] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709, 2019.
- [23] Qirui Jiao, Daoyuan Chen, Yilun Huang, Yaliang Li, and Ying Shen. Enhancing multimodal large language models with vision detection models: An empirical study. *arXiv preprint arXiv:2401.17981*, 2024.
- [24] Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander Rush, Douwe Kiela, et al. Obelics: An open web-scale filtered dataset of interleaved image-text documents. *Advances in Neural Information Processing Systems*, 36, 2024.
- [25] Chunyuan Li, Zhe Gan, Zhengyuan Yang, Jianwei Yang, Linjie Li, Lijuan Wang, Jianfeng Gao, et al. Multimodal foundation models: From specialists to general-purpose assistants. *Foundations and Trends® in Computer Graphics and Vision*, 16(1-2):1–214, 2024.
- [26] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023.
- [27] Xianming Li and Jing Li. Angle-optimized text embeddings. *arXiv preprint arXiv:2309.12871*, 2023.
- [28] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023.
- [29] Bin Lin, Zhenyu Tang, Yang Ye, Jiayi Cui, Bin Zhu, Peng Jin, Junwu Zhang, Munan Ning, and Li Yuan. Moe-llava: Mixture of experts for large vision-language models. *arXiv preprint arXiv:2401.15947*, 2024.
- [30] Yuanze Lin, Yujia Xie, Dongdong Chen, Yichong Xu, Chenguang Zhu, and Lu Yuan. Revive: Regional visual representation matters in knowledge-based visual question answering. *Advances in Neural Information Processing Systems*, 35:10560–10571, 2022.
- [31] Yuanze Lin, Chen Wei, Huiyu Wang, Alan Yuille, and Cihang Xie. Smaug: Sparse masked autoencoder for efficient video-language pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2459–2469, 2023.
- [32] Yuanze Lin, Yi-Wen Chen, Yi-Hsuan Tsai, Lu Jiang, and Ming-Hsuan Yang. Text-driven image editing via learnable regions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7059–7068, 2024.

- [33] Yuanze Lin, Ronald Clark, and Philip Torr. Dreampolisher: Towards high-quality text-to-3d generation via geometric diffusion. *arXiv preprint arXiv:2403.17237*, 2024.
- [34] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023.
- [35] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024.
- [36] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023.
- [37] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? *arXiv preprint arXiv:2307.06281*, 2023.
- [38] Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. *Advances in Neural Information Processing Systems*, 35: 2507–2521, 2022.
- [39] Microsoft. Phi-2: The surprising power of small language models, 2023. URL <https://www.microsoft.com/en-us/research/blog/phi-2-the-surprising-power-of-small-language-models>.
- [40] OpenAI. Chatgpt. <https://openai.com/blog/chatgpt>, 2022.
- [41] OpenAI. Gpt-4v(ision) system card. [https://cdn.openai.com/papers/GPTV\\_System\\_Card.pdf](https://cdn.openai.com/papers/GPTV_System_Card.pdf), 2023.
- [42] OpenAI. Chatgpt. <https://openai.com/blog/chatgpt/>, 2023.
- [43] OpenAI. Gpt-4 technical report. 2023.
- [44] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- [45] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*, 2023.
- [46] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022.
- [47] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. *OpenAI*, 2018.
- [48] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [49] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- [50] Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abdelrahman Shaker, Salman Khan, Hisham Cholakkal, Rao M Anwer, Erix Xing, Ming-Hsuan Yang, and Fahad S Khan. Glamm: Pixel grounding large multimodal model. *arXiv preprint arXiv:2311.03356*, 2023.

- [51] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- [52] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8317–8326, 2019.
- [53] Stability AI. Introducing stable lm 2, 2024. URL <https://stability.ai/news/introducing-stable-lm-2>.
- [54] Yi-Lin Sung, Linjie Li, Kevin Lin, Zhe Gan, Mohit Bansal, and Lijuan Wang. An empirical study of multimodal model merging. *arXiv preprint arXiv:2304.14933*, 2023.
- [55] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [56] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [57] Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Haowei Liu, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration. *arXiv preprint arXiv:2311.04257*, 2023.
- [58] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*, 2023.
- [59] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11975–11986, 2023.
- [60] Hao Zhang, Feng Li, Xueyan Zou, Shilong Liu, Chunyuan Li, Jianwei Yang, and Lei Zhang. A simple framework for open-vocabulary segmentation and detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1020–1031, 2023.
- [61] Shilong Zhang, Peize Sun, Shoufa Chen, Min Xiao, Wenqi Shao, Wenwei Zhang, Kai Chen, and Ping Luo. Gpt4roi: Instruction tuning large language model on region-of-interest. *arXiv preprint arXiv:2307.03601*, 2023.
- [62] Liang Zhao, En Yu, Zheng Ge, Jinrong Yang, Haoran Wei, Hongyu Zhou, Jianjian Sun, Yuang Peng, Runpei Dong, Chunrui Han, et al. Chatspot: Bootstrapping multimodal llms via precise referring instruction tuning. *arXiv preprint arXiv:2307.09474*, 2023.
- [63] Baichuan Zhou, Ying Hu, Xi Weng, Junlong Jia, Jie Luo, Xien Liu, Ji Wu, and Lei Huang. Tinyllava: A framework of small-scale large multimodal models. *arXiv preprint arXiv:2402.14289*, 2024.
- [64] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. MiniGPT-4: Enhancing vision-language understanding with advanced large language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=1tZbq88f27>.
- [65] Minjie Zhu, Yichen Zhu, Xin Liu, Ning Liu, Zhiyuan Xu, Chaomin Shen, Yaxin Peng, Zhicai Ou, Feifei Feng, and Jian Tang. A comprehensive overhaul of multimodal assistant with small language models. *arXiv preprint arXiv:2403.06199*, 2024.
- [66] Yichen Zhu, Minjie Zhu, Ning Liu, Zhicai Ou, Xiaofeng Mou, and Jian Tang. Llava-phi: Efficient multi-modal assistant with small language model. *arXiv preprint arXiv:2401.02330*, 2024.

## A Appendix

In the supplementary materials, we provide the following sections:

- (a) More implementation details in Section B.
- (b) Ablation study experiments in Section C.
- (c) Visualization result analysis in Section D.

## B Implementation Details

The training time for LLaVA-1.5 7B [34] and Mipha-3B [65] is approximately 14 hours and 9 hours, respectively, with a batch size of 256 on  $32 \times$  NVIDIA V100 32GB GPUs. For the initialization of the proposed prompt embedding network (PEN), we use Kaiming initialization technology [19]. The UAE-Large-V1<sup>2</sup> model is adopted as the pre-trained textual encoder to extract textual embeddings for the visual prompt.

## C Ablation Study

Next, we conduct more ablation study experiments to provide deeper insight into the components of our proposed approach.

**The Effect of Using Different Pre-trained Textual Encoders.** In Table 7, we perform an ablation study using different textual encoders, *i.e.*, CLIP [48] vs. UAE [27], to extract textual embeddings for the proposed visual prompt. We draw two conclusions from Table 7: (1) Using different textual encoders, the proposed approach consistently outperforms the baseline, demonstrating the robustness of our method. (2) Adopting UAE as the pre-trained textual encoder achieves significantly better performance. Therefore, we choose UAE as the default pre-trained textual encoder in our experiments.

**Object Detector v.s. Segmentation Model.** To determine the effect of using an object detector or segmentation model to incorporate pixel-level semantics into the proposed visual prompt, we conduct an ablation study with the popular object detector GroundingDINO [36] and the segmentation model OpenSeed [60]. The results are shown in Table 8. We observe that both GroundingDINO and OpenSeed significantly boost performance across all benchmarks. However, utilizing OpenSeed achieves better performance gains due to its fine-grained mask regions. Thus, we adopt OpenSeed by default to generate object regions.

**The Effect of Fine-Tuning with the Visual Prompt.** As displayed in Table 9, the model fine-tuned with the proposed visual prompt (*i.e.*, the third row) achieves remarkably better performance than the one fine-tuned without our visual prompt (*i.e.*, the second row) across all benchmarks. Specifically, without using our visual prompt for fine-tuning, the model even shows performance degradation on Text-VQA benchmark [52] and has negligible gains on Science-QA [38], VQAv2 [18], MME-P [16], and MME-C [16] benchmarks. All these results demonstrate the superiority of the proposed method.

## D Visualization Result Analysis

We’ve provided more visualization results in Table 10, 11, 12, and 13. Compared to LLaVA-1.5 7B [34], our method generates more reasonable and accurate responses to the questions.

As shown in Table 10, our approach can generate accurate movie titles, such as “The Godfather”, and the two actors’ names, such as “Al Pacino” and “Robert De Niro”. Additionally, it provides a corresponding introduction, such as “The movie is a classic crime drama film directed by Francis Ford Coppola, known for its iconic characters, storytelling, and memorable scenes” for the left example. In the right example, our method generates the precise title “The Lord of the Rings: The Fellowship of the Ring” and provides an accurate introduction, such as “The movie is an epic fantasy adventure that

---

<sup>2</sup><https://huggingface.co/WhereIsAI/UAE-Large-V1>

Table 7: The ablation study of using different textual encoders, i.e., CLIP v.s. UAE, to extract textual embeddings for the proposed visual prompt.

Method	Text Enc	VQAv2	GQA	SQA <sup>I</sup>	VQA <sup>T</sup>	MME-P	MME-C	MMB	MM-Vet	POPE
Mipha-3B	-	81.3	63.9	70.9	56.6	1488.9	295.0	69.7	32.1	86.7
Mipha-3B	CLIP	82.1 $\uparrow$	64.9 $\uparrow$	71.3 $\uparrow$	57.4 $\uparrow$	1497.2 $\uparrow$	361.5 $\uparrow$	71.1 $\uparrow$	34.6 $\uparrow$	88.5 $\uparrow$
Ours	UAE	<b>82.4<math>\uparrow</math></b>	<b>65.3<math>\uparrow</math></b>	<b>71.8<math>\uparrow</math></b>	<b>57.8<math>\uparrow</math></b>	<b>1501.2<math>\uparrow</math></b>	<b>369.1<math>\uparrow</math></b>	<b>71.5<math>\uparrow</math></b>	<b>35.1<math>\uparrow</math></b>	<b>88.7<math>\uparrow</math></b>

Table 8: The ablation study of using an object detector or a panoptic segmentation model to extract object regions for pixel-level textual embeddings.

Method	Region Generator	VQAv2	GQA	SQA <sup>I</sup>	VQA <sup>T</sup>	MME-P	MME-C	MMB	MM-Vet	POPE
Mipha-3B	-	81.3	63.9	70.9	56.6	1488.9	295.0	69.7	32.1	86.7
Mipha-3B	GroundingDINO	82.0 $\uparrow$	64.9 $\uparrow$	71.4 $\uparrow$	57.2 $\uparrow$	1491.7 $\uparrow$	350.2 $\uparrow$	71.0 $\uparrow$	34.5 $\uparrow$	88.4 $\uparrow$
Ours	OpenSeed	<b>82.4<math>\uparrow</math></b>	<b>65.3<math>\uparrow</math></b>	<b>71.8<math>\uparrow</math></b>	<b>57.8<math>\uparrow</math></b>	<b>1501.2<math>\uparrow</math></b>	<b>369.1<math>\uparrow</math></b>	<b>71.5<math>\uparrow</math></b>	<b>35.1<math>\uparrow</math></b>	<b>88.7<math>\uparrow</math></b>

Table 9: The ablation study of fine-tuning with and without the proposed visual prompt. The first, second and third rows mean Mipha-3B baseline, fine-tuning on Mipha-3B without and with the proposed visual prompt using LoRA [21].

Method	Visual Prompt	VQAv2	GQA	SQA <sup>I</sup>	VQA <sup>T</sup>	MME-P	MME-C	MMB	MM-Vet	POPE
Mipha-3B	-	81.3	63.9	70.9	56.6	1488.9	295.0	69.7	32.1	86.7
Mipha-3B <sup>+</sup>	$\times$	81.4 $\uparrow$	64.3 $\uparrow$	71.0 $\uparrow$	56.5 $\downarrow$	1489.2 $\uparrow$	303.2 $\uparrow$	70.4 $\uparrow$	33.5 $\uparrow$	87.4 $\uparrow$
Ours	$\checkmark$	<b>82.4<math>\uparrow</math></b>	<b>65.3<math>\uparrow</math></b>	<b>71.8<math>\uparrow</math></b>	<b>57.8<math>\uparrow</math></b>	<b>1501.2<math>\uparrow</math></b>	<b>369.1<math>\uparrow</math></b>	<b>71.5<math>\uparrow</math></b>	<b>35.1<math>\uparrow</math></b>	<b>88.7<math>\uparrow</math></b>

follows the journey of a fellowship of characters". In contrast, LLaVA’s responses are very general without fully understanding the global contexts within the images.

In Table 11, our method not only recognizes the text “Me: I’ll do it at 8. Time: 8.05. Looks like I gotta wait till 9 now”, but also understands its implication, such as “the puppet is aware of the time and is intentionally delaying any work or task it might be assigned”, in the left image. For the right image, our method generates more potential options, such as a “private school”, “a religious school”, etc. In contrast, LLaVA’s answers are either inappropriate or too limited.

Finally, as shown in Table 12 and 13, our method understands spatial relationships and accurately recognizes the text within the images. For example, in the left example of Table 12, our approach accurately names the person on the left as “Keira Knightley”. It also recognizes the plate number “S3302CD” in the left example of Table 13, while LLaVA’s answers are all incorrect. This indicates the strong fine-grained multimodal understanding capacity of our proposed method.

---

**Visual input example, context-related questions:**

---



Can you give a short introduction to this movie?



Can you give a short introduction to this movie?

---

**LLaVA-7B** The image features a scene from a movie where two men are standing close to each other, with one man placing his hand on the other's shoulder. They are both wearing tuxedos, which suggests that they are attending a formal event or a special occasion. The men appear to be engaged in a conversation or interacting with each other in some way, possibly discussing the event or sharing a moment of camaraderie.

The image features a young man wearing a brown shirt and holding a ring in his hand. He appears to be in a dark setting, possibly a cave or a dungeon. The man is looking at the camera, possibly engaging in a conversation or sharing a story. The ring he is holding suggests that the scene might be related to a fantasy or adventure movie, where the ring could represent a magical or powerful object.

---

**Ours** The image features a scene from the movie "The Godfather", where two men played by actors Al Pacino and Robert De Niro, are standing close to each other. One man has his hand on the other's shoulder, and they appear to be engaged in a conversation or a confrontation. The movie is a classic crime drama film directed by Francis Ford Coppola, known for its iconic characters, storytelling, and memorable scenes.

The image features a scene from the movie "The Lord of the Rings: The Fellowship of the Ring", where the character Gandalf is holding out his hand to the hobbit Frodo, offering him the One Ring. This iconic moment captures the essence of the story, where Frodo is entrusted with the task of destroying the ring to prevent its power from falling into the wrong hands. The movie is an epic fantasy adventure that follows the journey of a fellowship of characters as they embark on a perilous quest to rid Middle-earth of the dark lord Sauron and his army of orcs.

---

Table 10: The challenging examples on LLaVA-1.5, our approach can provide reasonable responses for the tricky questions which need to fully understand the contexts of the images.



---

**Visual input example, context-related questions:**

---

**Me:** I'll do it at 8  
**Time:** 8.05  
**Me:** looks like I gotta wait till 9 now



Can you explain this meme?



What kind of school does the boy go to? Answer it and give the rationale.

---

**LLaVA-7B** The meme features a green Kermit the Frog puppet lying on a bed, with a caption that reads, "Me looks like I gotta wait till 9" The puppet appears to be relaxing on the bed, possibly waiting for someone or something. The image is a playful and lighthearted representation of a common situation where people might have to wait for a certain time before they can proceed with their plans.

The boy likely goes to a private school, as he is wearing a suit and tie, which is a common attire for students in private schools. The presence of a tie suggests that the school has a formal dress code, and the boy is dressed appropriately for the occasion.

---

**Ours** The image features a green frog, which is a popular character from the television show "The Muppet Show", the frog is sitting on a bed, appearing to be in a relaxed and comfortable position. The puppet is holding its belly and seems to be enjoying its time. The bed has a white sheet, and the scene appears to be set in a cozy bedroom. The meme is captioned with a humorous message which reads "Me: I'll do it at 8. Time: 8.05. Looks like I gotta wait till 9 now", this implies that the puppet is aware of the time and is intentionally delaying any work or task it might be assigned. The meme is likely meant to be amusing and entertaining, as it combines the idea of a puppet with a human-like sense of humor and awareness of time.

The boy in the image is wearing a tie, which suggests that he might be attending a school with a formal dress code or a school that requires students to wear ties as part of their uniform. This could be a private school, a religious school, or a school with a strong emphasis on discipline and professionalism. The presence of a woman, who is wearing a dress, in the image, possibly his mother, indicates that the boy might be going to school with her support and guidance.

---

Table 11: The challenging examples on LLaVA-1.5, our approach can provide reasonable responses for the tricky questions which need to fully understand the contexts of the images.

---

**Visual input example, relationship-aware questions:**



Who is the person on the left?

Who is the person in the bottom right?

---

LLaVA-7B    The person on the left is Reese Witherspoon.

The person in the bottom right is a man named Tom Hiddleston.

---

**Ours**    The person on the left is **Keira Knightley**.

The person in the bottom right is **Jamie Dornan**.

---

Table 12: The challenging examples on LLaVA-1.5, our approach can provide reasonable responses for the tricky questions which need to fully understand the relationships among objects within images.

---

**Visual input example, text-related questions:**

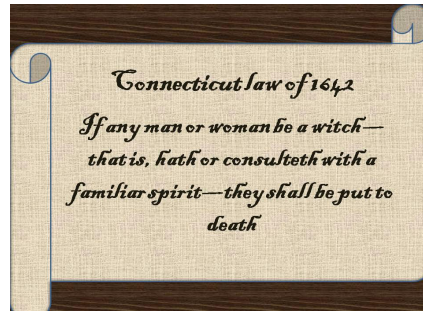
---



What are all the scene text in the image?

LLaVA-7B The scene text in the image is "S33020CD".

**Ours** The scene text in the image includes the license plate number "S3302CD" which is written in white on a black background.



What are all the scene text in the image?

The scene text in the image is "Connecticut Law of 1652".

The scene text in the image includes "Connecticut Law of 1642", which states "If any man or woman be a witch - that is hath or consulteth with a familiar spirit - they shall be put to death".

---

Table 13: The challenging examples on LLaVA-1.5, our approach can provide reasonable responses for the tricky questions which need to accurately recognize the texts within the images.