

Zoomer: Adaptive Image Focus Optimization for Black-box MLLM

Jiaxu Qian, Chendong Wang, Yifan Yang, Chaoyun Zhang, Huiqiang Jiang, Xufang Luo, Yu Kang, Qingwei Lin, Anlan Zhang, Shiqi Jiang, Ting Cao, Tianjun Mao, Suman Banerjee, Guyue Liu, Saravan Rajmohan, Dongmei Zhang, Yuqing Yang, Qi Zhang, Lili Qiu

Abstract

Recent advancements in multimodal large language models (MLLMs) have broadened the scope of vision-language tasks, excelling in applications like image captioning and interactive question-answering. However, these models struggle with accurately processing visual data, particularly in tasks requiring precise object recognition and fine visual details. Stringent token limits often result in the omission of critical information, hampering performance. To address these limitations, we introduce Zoomer, a novel visual prompting mechanism designed to enhance MLLM performance while preserving essential visual details within token limits. Zoomer features three key innovations: a prompt-aware strategy that dynamically highlights relevant image regions, a spatial-preserving orchestration schema that maintains object integrity, and a budget-aware prompting method that balances global context with crucial visual details. Comprehensive evaluations across multiple datasets demonstrate that Zoomer consistently outperforms baseline methods, achieving up to a 26.9% improvement in accuracy while significantly reducing token consumption.

1. Introduction

Vision-language understanding has witnessed remarkable progress with multimodal large language models (MLLMs) such as GPT-4o, Gemini Pro, and Claude 3.5 [8, 13]. While these black-box MLLMs have demonstrated impressive capabilities in controlled settings—even being touted as having PhD-level skills—our large-scale investigation reveals a critical hallucination problem: their inherent small object blindness and limited visibility lead to catastrophic failures in real-world visual reasoning tasks.

Through systematic analysis, we identify two fundamental limitations that create this problem. First, despite their sophisticated architectures, current black-box MLLMs process images with uniform attention, lacking the human-like ability to focus on relevant regions while maintaining spatial awareness. As shown in Figure 1, this limitation leads to surprising failures in entry-level tasks like object counting—a problem that persists across different models and

prompting strategies. Our investigation reveals that this is not merely an implementation issue but a fundamental limitation of how these models process visual information.

Second, we identify a critical tension between token budgets and visual fidelity in black-box MLLMs. While commodity black-box MLLM has treated token limitations, typically fixed token counts, as a simple engineering constraint to ensure computational efficiency and user fairness [4], our analysis reveals it as a fundamental resource allocation problem. As demonstrated in Figure 2, the standard approach of uniform downsampling leads to catastrophic loss of fine details—a problem that cannot be solved by better downsampling algorithms or prompt engineering alone.

This investigation leads us to formulate a novel problem in visual prompting: *How can we achieve region-aware visual processing within the strict token constraints of black-box MLLMs?* This question challenges the conventional view that visual prompting is merely a simple input formatting task. Instead, we show that effective visual prompting requires solving three interconnected challenges: **Region Selection**: How to identify and prioritize task-relevant regions without access to model internals; **Spatial Preservation**: How to maintain structural relationships while maximizing detail preservation; **Budget Optimization**: How to allocate limited tokens between global context and local details.

Based on these insights, we propose Zoomer, a visual prompting framework that addresses these challenges through three technical innovations:

- A prompt-aware visual emphasis mechanism that enables selective attention without modifying model architecture
- A spatial-preserving orchestration schema that maintains global context while preserving local details
- A budget-aware region selection strategy that optimally allocates tokens across image regions

In a comprehensive evaluation across datasets such as *Vstar* [28], *CVBench* [24], and *RealworldQA* [29], Zoomer consistently outperformed baseline methods. Notably, in the *Vstar* dataset, Zoomer-Patches achieved a 26.9% accuracy improvement over the baseline, while in *RealWorldQA*, Zoomer-Adaptive outperformed the baseline by

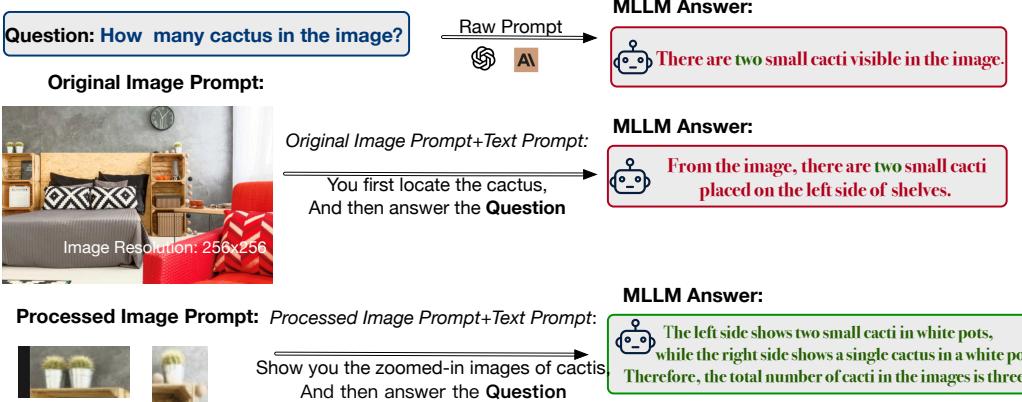


Figure 1. Illustration of a black-box MLLM’s approach to counting cacti in an image. The model identifies two small cacti on the left side and overlooks the single cactus on the right side of the image, arriving at a total of three cacti. The processed prompt highlights specific regions of interest to facilitate the correct object count.

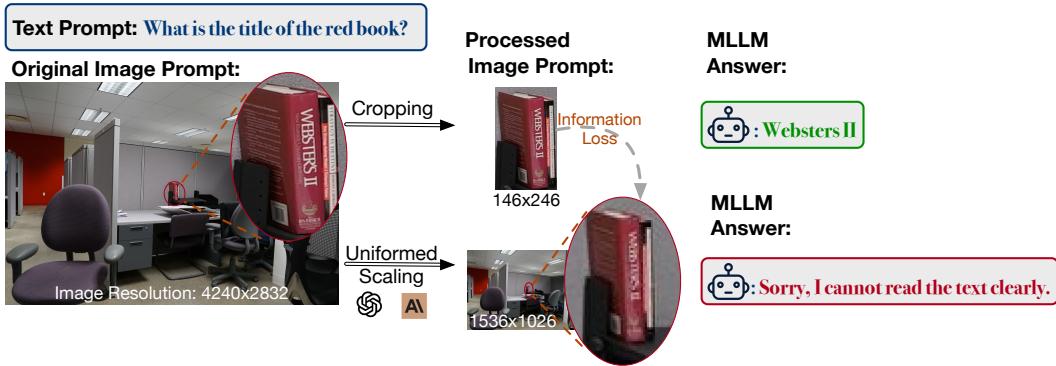


Figure 2. Illustration of information loss during image processing in black-box MLLMs. The original high-resolution image (4240x2832) is downscaled to meet token limits (1536x1026), leading to the loss of critical details. Cropping to focus on a region of interest (146x246) allows the model to correctly identify the book title as “Webster’s II”.

12.1%. Alongside these accuracy gains, Zoomer significantly reduced token usage. For example, on the TerraIncognita dataset, Zoomer achieved 6.4% higher accuracy with a 67% reduction in token consumption compared to the baseline. These results confirm that Zoomer not only addresses the limitations of visual processing in black-box MLLMs but also enhances efficiency. Moreover, across APIs like GPT-4o¹, Gemini-1.5Pro², and Claude-3.5-Sonnet³, Zoomer demonstrated consistent improvements in both accuracy and token efficiency, solidifying its potential to optimize MLLM performance in real-world, high-resolution applications.

The contribution of this work can be summarized as follows: (1) We conduct a detailed investigation of GPT-4o’s image prompting strategy, exposing key limitations in its handling of visual inputs. (2) We introduce Zoomer, a novel mechanism that addresses the challenge of preserv-

ing visual detail in black-box MLLMs while adhering to token constraints. (3) We present extensive experimental results across multiple datasets, demonstrating that Zoomer achieves significant improvements in both accuracy and token efficiency, offering valuable insights into enhancing multimodal processing in constrained environments.

2. Pilot Experiments

One of the primary challenges faced by black-box MLLMs is their inability to process visual inputs efficiently, leading to diminished accuracy in fine-grained visual tasks. Models like GPT-4o often struggle with recognizing detailed or occluded objects, particularly when dealing with complex images. These limitations are compounded by token constraints, which restrict the amount of image data that can be processed in a single prompt. Furthermore, according to the Vision pricing calculator⁴, GPT handle images by resizing and splitting them into basic units of 512×512 pixels. Each of these units corresponds to 170 tokens. This method of

¹<https://platform.openai.com/>

²<https://gemini.google.com/>

³<https://anthropic.com/>

⁴<https://openai.com/api/pricing/>

Method	Accuracy	Prompt Tokens
Unaltered Input	57%	955
Image Crop	58%	270
Zoomed Crop	64%	270

Table 1. Performance of different methods on Image from Vstar.

processing not only imposes a strict limit on the image resolution but also increases the computational overhead due to the additional tokens generated from splitting. As a result, vital visual details may be lost when images are downsampled or resized to fit within the token limits, leading to poor performance on tasks that require precise visual grounding.

To evaluate this issue, we conducted a series of pilot experiments using GPT-4o-0513 on the *Vstar-Bench* dataset. This dataset challenges MLLMs to accurately identify detailed objects within high-resolution images, making it an ideal test for the model’s capacity to handle fine-grained visual information. The experiments compared three different image processing strategies: (1) an unprocessed prompt (**Unaltered Input**), where the image is fed to the model in its original form; (2) a prompt where the image is cropped to focus on the target object (**Image Crop**); and (3) a prompt where the cropped image is further enlarged to emphasize the most relevant visual features (**Zoomed Crop**). Both the Image Crop and Zoomed Crop methods were constrained to fit within GPT-4o’s patch size limit of 512x512 pixels.

As shown in Table 1, the **Zoomed Crop** method significantly outperformed the others, achieving an accuracy of 0.76 with a token usage of 270. In comparison, the **Unaltered Input** method, despite processing the entire image, only achieved an accuracy of 0.64 while consuming 955 tokens. Similarly, the **Image Crop** method, although reducing the token count to 270, did not yield any improvement in accuracy compared to the unprocessed input.

These results highlight a fundamental problem in current black-box MLLMs: they fail to efficiently manage the trade-off between image resolution and token constraints. In cases like those presented by the *Vstar-Bench* dataset, where fine-grained visual information is critical, processing unaltered high-resolution images leads to excessive token consumption without improving accuracy. While the **Image Crop** method reduces token usage, it fails to improve performance because simply cropping an image without emphasizing key details does not provide sufficient context for the model to interpret the visual input accurately.

The superior performance of the **Zoomed Crop** method underscores the importance of vision enhancement techniques in black-box MLLMs. By focusing on the most relevant portions of the image, **Zoomed Crop** preserves criti-

cal details while remaining within token limits, enabling the model to interpret detailed visual inputs more effectively. This approach resolves a common issue faced by black-box MLLMs, where downscaling or cropping images to meet token requirements often leads to a loss of essential information, reducing the overall effectiveness of the model.

Our experiments reveal that without adaptive techniques like **Zoomed Crop**, black-box MLLMs struggle to process high-resolution images efficiently, limiting their performance on tasks that require precise visual recognition. These findings demonstrate the necessity of vision enhancement strategies to address the inherent limitations of token-constrained MLLMs.

3. Related Work

3.1. Multimodal LLMs: Open-Source and Black-Box Models

The integration of visual and textual modalities in large language models (LLMs) has led to significant advancements in multimodal models (MLLMs) like GPT-4o, Gemini Pro and Claude3-Sonnet. These models rely on effective visual encoding strategies to bridge the gap between language and vision. Approaches such as CLIP [31] align visual and language embeddings through contrastive learning, while models like Flamingo [1] and BLIP-2 [6] use cross-attention mechanisms or pretraining modules to link vision encoders with LLMs. However, these methods often rely on fixed low-resolution inputs (e.g., 224x224), limiting their ability to process high-resolution images or non-standard aspect ratios [16], which hampers performance on fine-grained tasks such as OCR and small object detection.

In contrast, open-source multimodal models [12, 15, 30, 36, 37] allow for architectural modifications and fine-tuning to accommodate any-resolution inputs. However, black-box MLLMs such as GPT-4o and Gemini Pro, which impose strict token limits for computational efficiency, require alternative solutions. The need to downsample or crop images to meet these constraints often results in the loss of crucial visual details, particularly in tasks requiring detailed visual understanding. While position embedding interpolation [2, 5, 9, 19, 26] and patch-based cropping [12, 30] widely adopted in open-source models offer promising directions for any aspect ratio and any-resolution image processing, they are not applicable to black-box models, where architectural changes and extra training/fine-tuning are not permitted.

3.2. Object Detection

Traditional object detection models, such as Faster R-CNN [22] and YOLO [21], effectively identify and localize objects within predefined categories. However, they struggle with open-set scenarios, where novel objects not seen during training need to be detected.

Recent advances address this limitation through open-set detection models that leverage natural language processing. For instance, OV-DETR [34] integrates CLIP with object detection to generate category-specific bounding boxes from textual prompts, enabling detection in open-world settings. Similarly, GLIP [14] reframes detection as a grounding problem, improving alignment between visual regions and textual descriptions. DetCLIP [32] extends this further using pseudo labels from large-scale captioning datasets, enhancing generalization. Grounding DINO [17], built on the DETR framework [3], also advances open-set detection through natural language integration.

In addition, SAM [11] and SAM-2 [20] offer zero-prompt or minimal-prompt segmentation for arbitrary objects but lack robust text-prompt handling. EVF-SAM [35] overcomes this by extending SAM’s capabilities to better manage complex text-based object segmentation.

By incorporating these models, Zoomer enhances its ability to dynamically detect and emphasize regions of interest (RoIs), enabling black-box MLLMs to focus on the most relevant visual content without losing critical details, which is essential for maintaining high performance across varied resolutions.

4. Method Overview

Inspired by the observation derived from our pilot experiments, we propose Zoomer, a comprehensive visual prompting mechanism designed to effectively address the loss of detail in images that occurs during the naive resizing process in current black box multimodal LLMs, such as GPT-4 and Gemini 1.5. As illustrated in Figure 3, our mechanism comprises three key components: (1) A prompt-aware visual emphazizer that allocates high-fidelity image slices based on prompt texts to facilitate efficient and focused visual encoding; (2) A spatial-preserving encoding schema that consolidates the collected image slices while maintaining their relative spatial positions to create a condensed visual input; (3) A budget-aware prompting strategy that maximizes the accuracy of results obtained from the black box models while fits the budget requirement from users.

4.1. Prompt-aware Visual Emphasizer

The prompt-aware visual emphazizer utilizes a multi-scale emphasizing strategy to prioritize image slices that are most relevant to the input prompts. By analyzing the semantic content of the prompts, this component dynamically selects and enhances specific regions of the image at varying resolutions. This approach not only enriches the contextual information available to the model but also mitigates the adverse effects of losing critical details during the resizing process.

Prompt Tokenization Prompt tokenization is a critical first

step in which input prompts are parsed into meaningful tokens. This process segments the prompt into components that can be easily analyzed for semantic relevance. Specifically, the prompt is divided into structural components, and our focus is on processing the relevant sections that contribute directly to visual emphasis.

To enhance the extraction of semantically relevant tokens, we apply advanced natural language processing (NLP) techniques. First, we use the NLTK library⁵ to remove stopwords, reducing noise and ensuring that the model’s attention remains on the most critical visual elements. By eliminating these non-essential words, we concentrate on key terms that directly influence the visual emphasis.

In addition to basic stopword removal, we utilize dependency parsing [7, 23] to analyze the syntactic structure of the prompt. This deeper analysis identifies core entities and relationships, such as subject-object pairs and action verbs, which are crucial for interpreting the user’s intent. By focusing on these core semantic elements, we ensure that the visual emphasis aligns precisely with the underlying meaning of the prompt.

Finally, we strip away any irrelevant formatting or non-content-related details, allowing the visual emphazizer to focus solely on the essential information. This multi-layered tokenization approach ensures an optimal match between the tokenized prompt and the image features selected for emphasis.

Multi-Scale Emphasizing Algorithm: Given a key object term extracted from the text prompt, the Multi-Scale Emphasizing Algorithm 1 utilizes a state-of-the-art object detection model to localize the corresponding object in the image prompt. In our experiments, we primarily employ GroundingDINO [17] as our localization model.

The encoder in such models typically downsamples the input image to a resolution of 224×224 or 336×336 , potentially resulting in information loss when localizing the target object at a coarse granularity. To address this limitation, we propose a Multi-Scale Emphasizing Algorithm that processes the original image at multiple resolutions. The algorithm divides the input image into patches at various granularities, e.g. 2×2 , 3×3 , and beyond. For each generated patch, we apply the object detection model to localize the target object. The algorithm retains bounding boxes returned by the model that exceed a predefined confidence threshold. These high-confidence bounding boxes collectively form the output of our algorithm, providing a comprehensive multi-scale representation of the target object’s location.

⁵<https://www.nltk.org/>

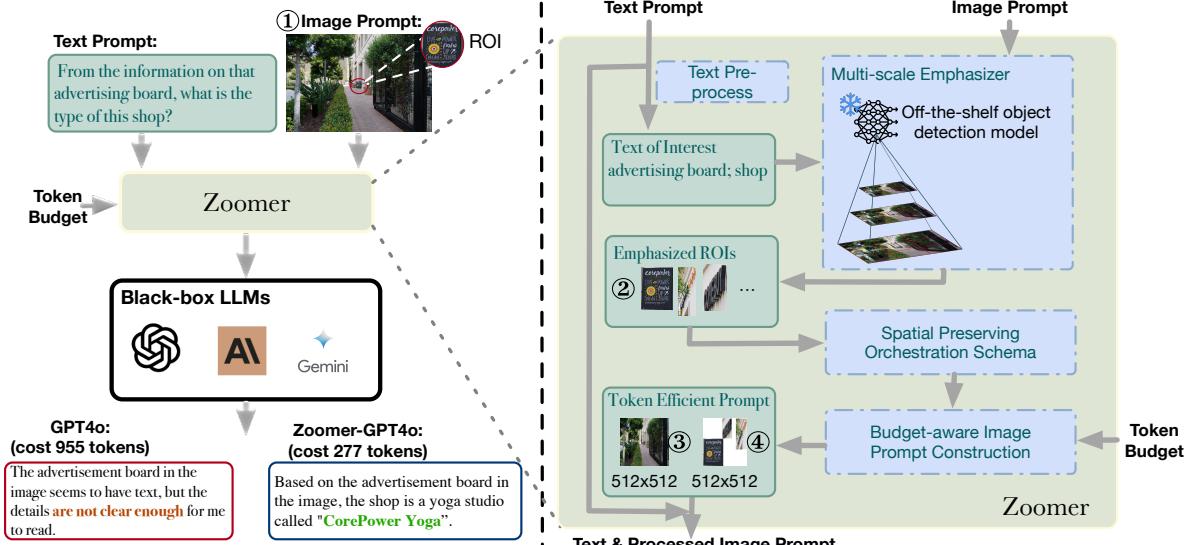


Figure 3. The Zoomer framework. Left: Raw Input image (1) and text prompt are processed by Zoomer and then fed into a black-box LLM (e.g., GPT-4o) for analysis, resulting in more accurate and detailed responses compared to standard input methods with even token saving. Right: Zoomer processes the text to extract key terms and uses a multi-scale emphasizes module (§4.1) with an off-the-shelf object detection model to identify regions of interest (ROIs). The identified ROIs (2) are then processed through a spatial preserving orchestration schema (§4.2) for a filtered emphasized patch (4) and a budget-aware image prompt construction module (§4.3) to create a token-efficient prompt within the specified budget. A scaled global view (3) is also generated for potential prompting.

Algorithm 1 Multi-Scale Emphasizing Algorithm

Require: I : input image, k : key object term, M : object detection model, T : confidence threshold
Ensure: B : set of bounding boxes

- 1: $B \leftarrow \emptyset$
- 2: $S \leftarrow \{2, 3, \dots, S_{\max}\}$ \triangleright Set of scaling factors
- 3: **for** each $s \in S$ **do**
- 4: $P_s \leftarrow \text{DivideIntoPatches}(I, s \times s)$
- 5: **for** each patch $p \in P_s$ **do**
- 6: $b, c \leftarrow M(p, k)$ \triangleright Get bounding box and confidence
- 7: **if** $c \geq T$ **then**
- 8: $B \leftarrow B \cup \{b\}$
- 9: **end if**
- 10: **end for**
- 11: **end for**
- 12: **return** B

4.2. Spatial-preserving Orchestration Schema

Building upon the Multi-Scale Emphasizing Algorithm, we introduce a Spatial-preserving Orchestration Schema to maintain the structural integrity of the image during the encoding process. This schema filters the bounding boxes obtained from the Multi-Scale Emphasizing Algorithm and ensures that the relative positions of the selected image slices are preserved, facilitating a more faithful representation of the original image layout and enabling coherent reconstruction when processed by the multimodal LLM. To refine the selection of bounding boxes, we implement a Non-Maximum Suppression (NMS) based slice filtering method. NMS is employed to eliminate redundant and overlapping slices, retaining only the most salient features that align with the prompt. The process works as described in Algorithm 2.

By setting an appropriate threshold T for the Intersection of Union (IoU) of bounding boxes around the selected regions, we ensure that only the highest-quality slices are retained for the encoding process. This filtering step enhances computational efficiency by reducing the number of slices to be processed and improves the clarity and relevance of the visual information provided to subsequent stages of the model.

The resulting set of filtered slices are then orchestrated to preserve their original relative positions within the image. This orchestration process involves the following steps: **Slice Extraction**: For each bounding box b_i in the filtered set F , we extract the corresponding image slice from the original image. **Blank Image Creation**: We create a new blank image with the same dimensions as the original image. **Slice Placement**: We place each extracted slice onto the blank image at its original position, leaving the rest of the image blank. **Image Shrinking**: The resulting image, containing only the selected slices in their original positions

The resulting set of filtered slices are then orchestrated to preserve their original relative positions within the image. This orchestration process involves the following steps: **Slice Extraction**: For each bounding box b_i in the filtered set F , we extract the corresponding image slice from the original image. **Blank Image Creation**: We create a new blank image with the same dimensions as the original image. **Slice Placement**: We place each extracted slice onto the blank image at its original position, leaving the rest of the image blank. **Image Shrinking**: The resulting image, containing only the selected slices in their original positions

with the rest left blank, is then shrunk to a predetermined size while maintaining its aspect ratio.

4.3. Budget-aware Prompting Strategy:

Our approach incorporates a sophisticated budget-aware prompting strategy that optimizes the allocation of token budget for image processing. This strategy begins with a user-specified total token budget B_{total} , allowing for customization based on specific task requirements or computational constraints. We propose four varieties of Zoomer to accommodate different budget scenarios and task requirements:

- **Zoomer-Local(④):** This variant utilizes only the spatial-preserving schema to consolidate all focused image slices into a single image patch(④ in Figure 3). It is optimal for scenarios with very limited token budgets, prioritizing the most relevant visual information.
- **Zoomer-Adaptive (④ + ◇ ③):** This approach dynamically includes a global view of the original image if the cropped portion falls below a certain threshold T_A . This allows the MLLM to better understand the overall scene context when the budget permits, while still focusing on key areas of interest.
- **Zoomer-Global (④ + ③):** This variant assigns a global view to all images, regardless of the specific regions of interest. It is suitable for tasks that require consistent overall context and when the token budget is sufficient to include both global and local information.
- **Zoomer-Patches(② + ③):** This is the most token-intensive approach, assigning each image slice its own patch without spatial preservation, along with a global view. It provides the most detailed information but requires the largest token budget.

The selection among these varieties depends on the user-specified budget and the nature of the task. For each variant, the number of high-resolution slices or patches N is calculated based on the available budget and the token cost per slice or patch. These slices are selected from the output of our Multi-Scale Emphasizing Algorithm, prioritizing based on their relevance to the key term of text prompts. To present the methods more clearly and vividly, we refer to Figure 4, which outlines the methodology, and Figure 5, which showcases a specific case study.

5. Experiments

In this section, we evaluate the performance of Zoomer through a series of experiments designed to test its ability to improve token efficiency and preserve visual fidelity across different black-box MLLMs. Specifically, we aim to answer the following questions: **(i) Accuracy:** Does Zoomer improve accuracy across different black-box MLLMs on image-related tasks? **(ii) Efficiency:** How does Zoomer perform compared to baseline methods in terms of both ac-

curacy and token efficiency? **(iii) Component Contribution:** What is the impact of key components in Zoomer, such as multi-scale vision emphasize and the budget-aware prompt strategy?

5.1. Setup

Assessment and Datasets We evaluated our system on a series of challenging multimodal tasks, using commercial black-box MLLMs for applications ranging from visual-language reasoning to image understanding and question answering. The experiments were conducted on a variety of different public datasets, including:

- 1) *Vstar* [28]: A benchmark dataset focused on image classification, used to evaluate fine-grained visual recognition capabilities in object detection and classification tasks.
- 2) *CVBench* [24]: Contains 2 sub-category, *CVBench_{2D}* and *CVBench_{3D}*, respectively, representing two-dimensional and three-dimensional visual image, respectively, to evaluate the performance of the model when processing images of different dimensions, especially the understanding ability in complex scenes.
- 3) *RealworldQA* [29]: Used to test the multimodal question answering performance of the model in real-world scenarios, involving cross-language and cross-image information processing.
- 4) *MMVP* [25]: A validation set for multimodal visual processing, designed to evaluate the comprehensive understanding of models for complex visual scenes.
- 5) *ScienceQA* [18]: A multimodal scientific question-answering dataset featuring multiple-choice questions across a diverse range of science topics.
- 6) *MMMU* [33]: The validation part of a new benchmark, which designed to evaluate the performance of multimodal models on multidisciplinary tasks that require university-level subject knowledge and deliberate reasoning.
- 7) *HR* [27]: A high-resolution multimodal benchmark consisting of 4K and 8K images and corresponding questions.

Models We employed three black-box MLLMs—GPT-4o-0513, Claude-v3-Sonnet, and Gemini-Pro—accessed via their respective APIs (OpenAI, Claude, Google). Across all experiments, we set the temperature to 0 and used greedy decoding for consistency, optimizing the stability of outputs. NMS was applied with a confidence score threshold of 0.8 to filter irrelevant regions from high-resolution images.

Metrics We used classification accuracy across all examples as the primary evaluation metric. Additionally, we compared token usage for each model configuration to evaluate the efficiency improvements offered by Zoomer.

Baselines We compare Zoomer against the following baseline methods:

Acc./Tokens \ Bench	Vstar	CVBench-2D	CVBench-3D	RealworldQA	SQA-I	MMVP	MMMU	HR-4K	HR-8K
Method									
Raw	56.5%/955	68.5%/428	78.2%/895	67.6%/998	87.3%/353	83.3%/270	68.4%/608	50.6%/1105	46.8%/1105
Resize	41.9%/270	66.3%/270	75.2%/270	61.1%/270	86.8%/270	83.3%/270	62.9%/270	35.8%/270	33.4%/270
Zoomer-Local	67.1%/270	72.4%/270	86.2%/270	72.4%/270	88.3%/270	87.1%/270	59.8%/270	58.8%/270	57.7%/270
Zoomer-Adaptive	67.5%/419	72.9%/374	87.9%/408	74.7%/362	91.1%/308	88.7%/351	61.6%/312	60.8%/331	59.3%/324
Zoomer-Global	67.6%/540	73.1%/540	88.3%/540	75.3%/540	92.3%/540	88.9%/540	67.3%/540	61.3%/540	59.8%/540
Zoomer-Patches	71.7% /1029	74.6% /709	85.8%/1113	75.8% /997	92.8% /727	88.4%/726	68.9% /841	60.4%/713	58.9%/875

Table 2. Performance of GPT-4o across different datasets using various image prompt processing methods, focusing on accuracy and token consumption. Among these approaches: **Local**: Only the extracted RoIs are used. **Adaptive**: Selectively provides the MLLM with a global view of the image based on the prompt strategy. **Global**: Every request includes the global view of the image. **Patches**: Does not use the Spatial-Preserving Orchestration Schema; instead, each possible ROI is independently provided to the MLLM, including the global view.

Method	Accuracy		Tokens		Latency		Money Cost(\$10-e3)	
	Zero-Shot	15-Shot	Zero-Shot	15-Shot	Zero-Shot	15-Shot	Zero-Shot	15-Shot
Raw	78%	84%	963	13488	4.8s	18.7s	4.815	67.44
Resize	61%	74%	255	4080	2.9s	7.5s	1.275	20.4
Low-Detail	60%	70%	85	1360	2.1s	6.5s	0.425	6.8
Zoomer-Adaptive	83%	88%	315	5112	3.1s	9.8s	1.575	25.56

Table 3. Performance in terms of accuracy, latency, and image token cost on TerraIncognita under ICL conditions—specifically with 15 examples per question—and under zero-shot conditions.

API	Method	Vstar	CVBench-2D	RealworldQA	MMVP
GPT-4o	Raw	56.5%	68.5%	67.6%	83.3%
	Zoomer	71.7%	74.6%	75.8%	88.9%
Gemini-1.5Pro	Raw	53.1%	65.4%	64.0%	79.8%
	Zoomer	70.4%	73.2%	73.9%	87.8%
Claude-3.5-Sonnet	Raw	51.8%	66.7%	61.0%	80.2%
	Zoomer	69.7%	72.8%	74.1%	87.2%

Table 4. Accuracy of Different Black-box MLLM APIs. For *Vstar*, *CVBench-2D*, and *RealworldQA*, we used the Patches version of SysName. For *MMVP*, inspired by Table 2, we employed the Global version.

Method	Model	Prompt Strategy	Vstar	CVBench-2D	CVBench-3D	RealworldQA	SQA-I	MMVP	MMMU
Default	EVF-SAM	Local	57.1%	67.6%	79.9%	68.5%	87.8%	84.0%	55.3%
		Adaptive	57.5%	71.3%	82.5%	72.1%	88.3%	85.3%	55.9%
		Global	57.8%	72.1%	83.0%	72.4%	88.8%	87.3%	56.8%
		Patches	57.1%	72.7%	83.8%	72.1%	86.8%	84.7%	56.1%
	Ground Dino	Local	58.1%	70.6%	82.5%	71.5%	85.3%	84.3%	55.6%
		Adaptive	58.3%	71.5%	83.9%	73.1%	90.1%	85.6%	56.3%
		Global	58.3%	71.8%	84.8%	73.4%	90.3%	87.8%	57.0%
		Patches	58.8%	70.6%	83.1%	72.6%	90.9%	87.7%	56.6%
	EVF-SAM	Local	58.4%	69.2%	84.0%	72.5%	88.3%	84.7%	57.1%
		Adaptive	58.5%	71.3%	85.3%	73.1%	91.1%	86.8%	58.7%
		Global	58.4%	72.1%	85.8%	73.4%	91.8%	87.8%	59.6%
		Patches	60.2%	71.3%	85.7%	73.1%	91.3%	86.8%	59.2%
Multi-Resolution	EVF-SAM	Local	63.6%	71.8%	82.6%	70.2%	88.8%	85.1%	56.8%
		Adaptive	64.4%	71.8%	84.3%	70.6%	90.3%	86.5%	58.5%
		Global	66.4%	72.2%	84.7%	71.4%	91.8%	86.9%	59.3%
		Patches	66.2%	72.6%	83.6%	70.2%	92.1%	86.9%	58.8%
	Ground Dino	Local	63.7%	72.1%	85.2%	70.4%	85.3%	85.7%	57.5%
		Adaptive	64.3%	72.4%	86.1%	70.6%	90.1%	87.6%	58.5%
		Global	64.3%	72.8%	87.9%	71.7%	90.3%	88.8%	59.9%
		Patches	67.2%	73.7%	87.1%	73.0%	90.9%	88.0%	59.7%
	EVF-SAM	Local	67.1%	72.4%	86.2%	72.4%	88.3%	87.1%	57.7%
		Adaptive	67.5%	72.9%	87.9%	74.7%	91.1%	88.7%	59.3%
		Global	67.6%	73.1%	88.3%	75.3%	92.3%	88.9%	59.8%
		Patches	71.7%	74.6%	85.8%	75.8%	92.8%	88.4%	58.9%

Table 5. Performance of Zoomer Across Datasets for Different Emphasis Methods, Models, and Prompt Strategies.

1) *Raw*: This baseline feeds MLLM the unmodified prompt, with no adjustments made to the image.

2) *Resize*: Here, images larger than 512x512 pixels are resized to fit within the GPT-4o’s patch limit, while smaller images remain unchanged.

5.2. Main results

Table 2 compares the performance of Zoomer against baseline methods across various datasets using GPT-4o. The results show that Zoomer, particularly in its Patches and Adaptive versions, consistently outperformed baseline approaches in terms of accuracy while maintaining lower token usage than the Raw method. Zoomer consistently outperformed baseline methods, showing accuracy improvements up to 26% across multiple tasks. For example, in the *Vstar* dataset, Zoomer-Patches achieved an accuracy of 0.717, compared to 0.565 using the Raw baseline, marking a 26.9% improvement. In *RealworldQA*, which demands complex multimodal reasoning, Zoomer-Adaptive achieved 0.758 accuracy, outperforming the 0.676 accuracy of the Raw method by 12.1%. These results highlight that Zoomer effectively preserves fine-grained visual details, enabling improved object recognition and image understanding across real-world tasks, where precise detail retention is crucial.

We further evaluated Zoomer on Claude-3.5-Sonnet and Gemini-1.5Pro to assess its generalizability across different black-box MLLMs. Table 4 shows that Zoomer demonstrated robust performance across different black-box MLLMs, for example, Zoomer achieved an accuracy of 0.704 on *Vstar*, compared to 0.531 with the Raw baseline, marking a 32.6% improvement. Similarly, in Claude-

3.5-Sonnet, Zoomer outperformed the baseline by 34.5% on RealworldQA, improving from 0.610 to 0.741. These results suggest that Zoomer can consistently enhance performance across different architectures, making it a versatile tool for various MLLM-based applications.

A key contribution of Zoomer is its ability to reduce token consumption while maintaining or improving accuracy. As shown in Table 3, Zoomer consistently delivers both token efficiency and performance improvements across various benchmark datasets. For instance, on the TerraIncognita dataset, guided by ManyICL [10], Zoomer achieves 0.83 accuracy using 315 tokens, compared to the Raw baseline’s 0.78 accuracy with 963 tokens—a 67% reduction in token usage while improving performance by 6.4%. Additionally, Zoomer reduces latency, making it practical for real-time applications. In the TerraIncognita zero-shot setting, Zoomer lowered latency from 4.8s to 3.1s, a 35.4% reduction without sacrificing accuracy. This makes Zoomer highly suitable for tasks like autonomous driving, and real-time visual analytics, where both token efficiency and reduced latency are critical.

5.3. Findings

Here we analyze why the Zoomer-Patches version underperforms compared to Zoomer-Global and even the Local version on certain datasets. For example, on the *CVBench_{3D}* dataset, the accuracy of the Patches version is 0.025 lower than the Global version and 0.04 lower than the Local version. Similarly, on the *MMVP* dataset, the Patches version falls short by 0.005 compared to the Global version and by 0.003 compared to the Adaptive version. Given that these results are averaged over multiple measurements, and accounting for model fluctuations, we hypothesize that this performance drop occurs because the Patches version treats each ROI as an independent image and provides them separately to the MLLM. When there are too many ROIs, the model may fail to capture or integrate some of them, leading to a drop in accuracy.

5.4. Ablation study

To further understand the impact of the individual components within Zoomer, we conducted an ablation study focusing on two key variants: 1) Zoomer with multi-scale emphasize: Compared with the commonly used multi-resolution and directly use, multi-scale visual emphasis is used to identify and emphasize ROI in the image. 2) Zoomer with different models.

To further investigate the contributions of Zoomer’s components, we conducted an ablation study (Table 5). The results demonstrate that the combination of multi-scale visual emphasis and Patches prompt strategies delivers the best performance across all most datasets. Comparing different vision emphasis models, such as EVF-SAM and Ground Dino, further highlights the ef-

fectiveness of Zoomer. Despite differences in model capabilities, both show accuracy improvements across datasets. Additionally, when comparing different emphasis methods—Default, Multi-Resolution, and Multi-Scale—the Multi-Scale method consistently outperformed Multi-Resolution. We hypothesize that, while Multi-Scale crops images and may split objects, its pyramid-shaped multi-recall strategy compensates for this by enhancing recall. In contrast, although Multi-Resolution maintains object integrity, adjusting resolution disrupts the model’s performance, likely because most models are trained on fixed-size inputs, and changing the resolution weakens their inherent capabilities.

6. Conclusion

In this paper, we introduced Zoomer, a novel visual prompting mechanism designed to overcome the limitations of black-box MLLMs in processing images while adhering to token constraints. Our approach effectively balances the need to capture essential visual details without exceeding token budgets, a challenge commonly encountered in existing models like GPT-4o and Gemini Pro.

Through a comprehensive evaluation across datasets such as Vstar and RealWorldQA, Zoomer demonstrated significant improvements, particularly in fine-grained visual tasks. Our results show that Zoomer-Patches achieved a 26.9% accuracy gain over baseline methods in Vstar, and Zoomer-Adaptive provided a 12.1% improvement in RealWorldQA. These gains were achieved while drastically reducing token usage, with Zoomer delivering 6.4% higher accuracy in the TerraIncognita using 67% fewer tokens.

Although this work primarily focuses on improving the efficiency of visual processing in black-box MLLMs, another potential issue that arises in real-world applications is communication cost. For example, transferring large images from edge devices (e.g., wearable cameras or glasses) to cloud servers can be expensive in terms of bandwidth, latency, and energy consumption. While Zoomer is designed to reduce token usage, its application in minimizing data transmission costs is an area that could be explored in future work.

As part of future research, we plan to investigate how Zoomer could be adapted for edge ML applications, enabling local processing on devices such as wearable cameras. This would allow for more efficient handling of visual inputs at the edge, reducing the need for extensive data transfers to the cloud. We aim to measure latency, power consumption, and the overall impact on system performance to assess the feasibility of applying Zoomer in these scenarios.

In summary, Zoomer offers a practical solution for enhancing visual processing in constrained MLLMs and opens up new directions for future exploration.

References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: A Visual Language Model for Few-Shot Learning. 3
- [2] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-VL: A Versatile Vision-Language Model for Understanding, Localization, Text Reading, and Beyond. 3
- [3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-End Object Detection with Transformers. 4
- [4] Boyuan Chen, Zhuo Xu, Sean Kirmani, Brian Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14455–14465, 2024. 1
- [5] 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. 3
- [6] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. InstructBLIP: Towards General-purpose Vision-Language Models with Instruction Tuning. 3
- [7] Marie-Catherine De Marneffe and Christopher D Manning. The stanford typed dependencies representation. In *Coling 2008: proceedings of the workshop on cross-framework and cross-domain parser evaluation*, pages 1–8, 2008. 4
- [8] Zhaopeng Gu, Bingke Zhu, Guibo Zhu, Yingying Chen, Ming Tang, and Jinqiao Wang. Anomalygpt: Detecting industrial anomalies using large vision-language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1932–1940, 2024. 1
- [9] Wenyi Hong, Weihan Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxuan Zhang, Juanzi Li, Bin Xu, Yuxiao Dong, Ming Ding, and Jie Tang. CogAgent: A Visual Language Model for GUI Agents. 3
- [10] Yixing Jiang, Jeremy Irvin, Ji Hun Wang, Muhammad Ahmed Chaudhry, Jonathan H. Chen, and Andrew Y. Ng. Many-shot in-context learning in multimodal foundation models, 2024. 8
- [11] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything, 2023. 4
- [12] Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. LLaVA-NeXT-Interleave: Tackling Multi-image, Video, and 3D in Large Multimodal Models. . 3
- [13] Jiaxuan Li, Duc Minh Vo, Akihiro Sugimoto, and Hideki Nakayama. Evcap: Retrieval-augmented image captioning with external visual-name memory for open-world comprehension. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13733–13742, 2024. 1
- [14] Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, Kai-Wei Chang, and Jianfeng Gao. Grounded Language-Image Pre-training. . 4
- [15] Yanwei Li, Chengyao Wang, and Jiaya Jia. LLaMA-VID: An Image is Worth 2 Tokens in Large Language Models. . 3
- [16] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual Instruction Tuning. . 3
- [17] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, and Lei Zhang. Grounding DINO: Marrying DINO with Grounded Pre-Training for Open-Set Object Detection. . 4
- [18] Pan Lu, Swaroop Mishra, Tony 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. In *The 36th Conference on Neural Information Processing Systems (NeurIPS)*, 2022. 6
- [19] Gen Luo, Yiyi Zhou, Yuxin Zhang, Xiawu Zheng, Xiaoshuai Sun, and Rongrong Ji. Feast Your Eyes: Mixture-of-Resolution Adaptation for Multimodal Large Language Models. 3
- [20] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chaoyuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feichtenhofer. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. 4
- [21] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You Only Look Once: Unified, Real-Time Object Detection. 3
- [22] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. 3
- [23] Parth Sarthi, Salman Abdulla, Aditi Tuli, Shubh Khanna, Anna Goldie, and Christopher D. Manning. Raptor: Recursive abstractive processing for tree-organized retrieval. In *International Conference on Learning Representations (ICLR)*, 2024. 4
- [24] Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula, Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, Austin Wang, Rob Fergus, Yann LeCun, and Saining Xie. Cambrian-1: A fully open, vision-centric exploration of multimodal llms, 2024. 1, 6

- [25] Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide shut? exploring the visual shortcomings of multimodal llms, 2024. 6
- [26] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-VL: Enhancing Vision-Language Model’s Perception of the World at Any Resolution. 3
- [27] Wenbin Wang, Liang Ding, Minyan Zeng, Xiabin Zhou, Li Shen, Yong Luo, and Dacheng Tao. Divide, conquer and combine: A training-free framework for high-resolution image perception in multimodal large language models, 2024. 6
- [28] Penghao Wu and Saining Xie. V*: Guided visual search as a core mechanism in multimodal llms. *arXiv preprint arXiv:2312.14135*, 2023. 1, 6
- [29] xAI. Grok. 2024. 1, 6
- [30] Ruyi Xu, Yuan Yao, Zonghao Guo, Junbo Cui, Zanlin Ni, Chunjiang Ge, Tat-Seng Chua, Zhiyuan Liu, Maosong Sun, and Gao Huang. LLaVA-UHD: An LMM Perceiving Any Aspect Ratio and High-Resolution Images. 3
- [31] Yifan Yang, Weiquan Huang, Yixuan Wei, Houwen Peng, Xinyang Jiang, Huiqiang Jiang, Fangyun Wei, Yin Wang, Han Hu, Lili Qiu, and Yuqing Yang. Attentive Mask CLIP. 3
- [32] Lewei Yao, Jianhua Han, Youpeng Wen, Xiaodan Liang, Dan Xu, Wei Zhang, Zhenguo Li, Chunjing Xu, and Hang Xu. DetCLIP: Dictionary-Enriched Visual-Concept Parallelized Pre-training for Open-world Detection. 4
- [33] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of CVPR*, 2024. 6
- [34] Yuhang Zang, Wei Li, Kaiyang Zhou, Chen Huang, and Chen Change Loy. Open-Vocabulary DETR with Conditional Matching. pages 106–122. 4
- [35] Yuxuan Zhang, Tianheng Cheng, Rui Hu, Lei Liu, Heng Liu, Longjin Ran, and Xiaoxin Chen. EVF-SAM: Early Vision-Language Fusion for Text-Prompted Segment Anything Model. 14(8), . 4
- [36] Yi-Fan Zhang, Qingsong Wen, Chaoyou Fu, Xue Wang, Zhang Zhang, Liang Wang, and Rong Jin. Beyond LLaVA-HD: Diving into High-Resolution Large Multimodal Models, . 3
- [37] Xiangyu Zhao, Xiangtai Li, Haodong Duan, Haian Huang, Yining Li, Kai Chen, and Hua Yang. MG-LLaVA: Towards Multi-Granularity Visual Instruction Tuning. 3

7. Appendix

7.1. Details of the method

Figure 4 is an example of the Zoomer, and Figure 5 is the output of various versions of the Zoomer.

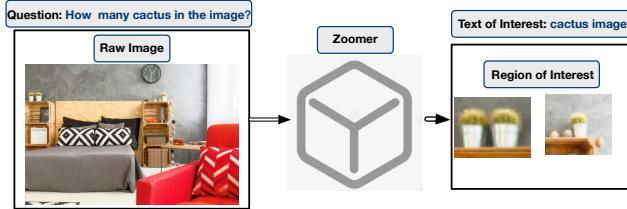


Figure 4. The example of applying Zoomer

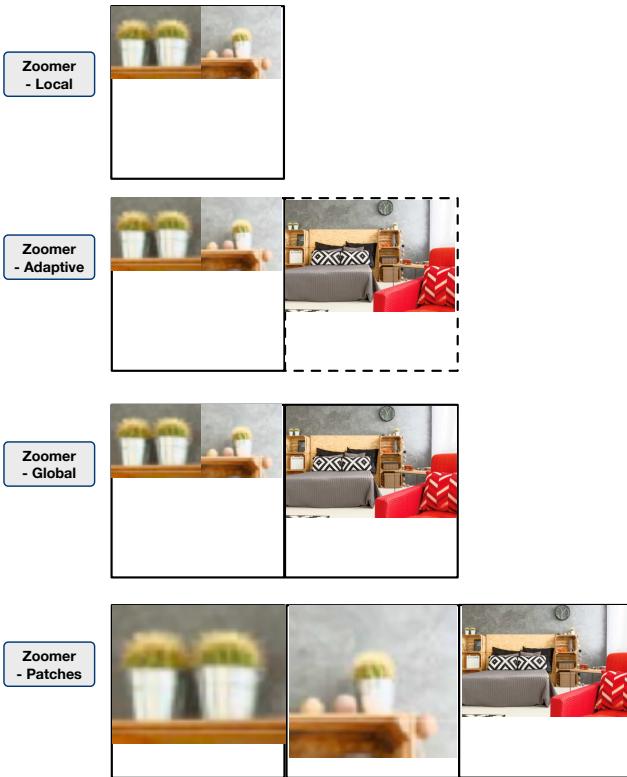


Figure 5. The example of different settings of Zoomer

7.2. Algorithm of the NMS

Algorithm 2 NMS-based Slice Filtering

Require: B : set of bounding boxes, T : IoU threshold

Ensure: F : set of filtered bounding boxes

```

1:  $F \leftarrow \emptyset$ 
2: Sort  $B$  in descending order of confidence scores
3: while  $B \neq \emptyset$  do
4:    $b_{\max} \leftarrow \arg \max_{b \in B} \text{score}(b)$ 
5:    $F \leftarrow F \cup b_{\max}$ 
6:    $B \leftarrow B \setminus b_{\max}$ 
7:   for each  $b \in B$  do
8:     if  $\text{IoU}(b_{\max}, b) \geq T$  then
9:        $B \leftarrow B \setminus b$ 
10:      end if
11:   end for
12: end while
13: return  $F$ 
```
