# ZoomEye: Enhancing Multimodal LLMs with Human-Like Zooming Capabilities through Tree-Based Image Exploration

Anonymous ACL submission

#### Abstract

Multimodal Large Language Models (MLLMs) have demonstrated impressive capabilities in vision-language understanding. Recently, with the integration of test-time scaling techniques, these models have also shown strong potential in visual reasoning. However, most existing reasoning approaches remain text-level in nature: MLLMs are prompted to explore various combinations of textual tokens via their underlying language model, while the visual input remains fixed throughout the reasoning process. This paradigm limits the model's ability to fully exploit rich visual information, particularly when dealing with images containing numerous fine-grained elements. In such cases, vision-level reasoning becomes crucial-where models dynamically zoom into specific regions of the image to gather detailed visual cues necessary for accurate decision-making. In this paper, we propose Zoom Eye, a training-free, model-agnostic tree search algorithm tailored for vision-level reasoning. Zoom Eye treats an image as a hierarchical tree structure, where each child node represents a zoomed-in subregion of its parent, and the root corresponds to the full image. The algorithm enables MLLMs to simulate human-like zooming behavior by navigating from root to leaf nodes in search of task-relevant visual evidence. We experiment on a series of elaborate high-resolution benchmarks and the results demonstrate that Zoom Eye not only consistently improves the performance of a series of MLLMs with large margin (e.g., InternVL2.5-8B increases by 15.71% and 17.69% on HR-Bench) but also enables small 3-8B MLLMs to outperform strong large models such as GPT-4o.

011

014

022

027

034

041

042

## 1 Introduction

By integrating powerful language models (Touvron et al., 2023; Yang et al., 2024) with visual encoders (Radford et al., 2021; Sun et al., 2023; Zhai et al., 2023), Multimodal large language models (MLLMs) are able to jointly process textual and



Figure 1: **Top**: When dealing with a high-resolution image, MLLMs effectively perceive the dominant objects but often fail to recognize finer details, highlighting the need for vision-level reasoning. **Bottom**: Applied with Zoom Eye, MLLMs could perform vision-level reasoning, allowed to explore the image details until they can answer the question.

visual inputs, achieving impressive performance in vision-language understanding (Zhao et al., 2024a; Bai et al., 2023; Chen et al., 2024b; Li et al., 2024). Recently, drawing on test-time scaling techniques that enhance reasoning abilities in LLMs, such as OpenAI-o1 (Jaech et al., 2024) and DeepSeek-R1 (Guo et al., 2025), a series of literature tries to investigate these reasoning techniques in MLLMs to further improve the visual reasoning capabilities (Xu et al., 2024; Dong et al., 2024; Yao et al.,

045

051 052 053

073

081

090

091

098

100

101

102

103

104

106

#### 2024a; Shen et al., 2025; Meng et al., 2025)

However, these methods predominantly operate at the textual level, leveraging the generative capacity of the underlying language model without modifying the perception of the image itself. That is, the visual input remains static throughout the reasoning process, restricting the model's ability to process fine-grained visual content, especially on an elements-rich high-resolution image. As illustrated in the top of Figure 1, for the same image, the MLLM accurately recognizes the dominant object whereas it struggles to perceive the detailed one. This gap highlights the need for vision-level reasoning, where the model actively interacts with the image by zooming in and out to selectively attend to informative regions, as demonstrated in the bottom of Figure 1, much like how humans visually process complex scenes. A similar visionlevel zooming mechanism has been adopted in the closed-source OpenAI-o3 (OpenAI, 2025). In contrast, our goal is to develop an open-source visionlevel reasoning method, making this capability accessible to the broader research community.

When viewing a high-resolution image, humans typically start with a global scan, then gradually zoom into areas of interest for closer inspection (Figure 2(b)). If the desired information is not found, they zoom out and explore alternative regions (as shown in Figure 2 (c)). Inspired by this, structuring an image as a tree is highly logical for simulating similar actions in an MLLM: the root denotes the full image, each child node corresponds to a zoomed-in sub-region of its parent, and deeper nodes indicate higher zoom levels. This hierarchical representation, combined with a search algorithm, allows models to (1) explore fine-grained regions (node lookahead) and (2) return to the previous view to inspect other regions (node backtracking). Similar tree-based search strategies have shown strong performance in text-based LLM reasoning(Yao et al., 2024b; Hao et al., 2023; Feng et al., 2023; Zhu et al., 2023).

In this paper, we propose Zoom Eye, a tree search algorithm for vision-level reasoning, which navigates MLLMs in the dense image context by the hierarchical and visual nature of images (**contribution #1**). This method simulates the actions of zooming in and out to inspect image details and seek out crucial information. Given a question, the adopted MLLM first identifies the pertinent objects. We then introduce two types of *confidence values* by prompting the MLLM to rec-



Figure 2: Zoom Eye enables MLLMs to (a) answer the question directly when the visual information is adequate, (b) zoom in gradually for a closer examination, and (c) zoom out to the previous view and explore other regions if the desired information is not initially found.

ognize the presence of these relevant objects. These *confidence values* are used to prioritize each candidate node during the tree search, determining the sequence of node selection. The search concludes based on a stopping criterion when the MLLM can confidently answer the question. This process is illustrated in the bottom part of Figure 1. Finally, the MLLM formulates a final response based on the visual information gathered during the search.

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

130

131

132

133

134

135

136

138

We adapt Zoom Eye to a series of mainstream MLLMs, including Qwen2.5VL (Bai et al., 2025), LLaVA-v1.5 (Liu et al., 2024a), LLaVA-OneVision (Li et al., 2024), InternVL2.5 (Chen et al., 2024a), and evaluate them on a suite of elaborate high-resolution visual understanding benchmarks. Equipped with Zoom Eye, all evaluated models achieve substantial performance improvements compared to the baseline (**contribution #2**).

Additionally, our analysis also reveals certain deficiencies in visual understanding exhibited by these models, which we detail in §4.3 (**contribution #3**). Addressing these limitations is part of our future work. More importantly, as discussed in §4.4.1, we observe a visionlevel test-time scaling phenomenon analogous to what has been observed in text-based LLMs: performance consistently improves with an increasing number of search steps. This finding suggests that vision-level reasoning benefits from deeper exploratory search and opens new avenues for scaling MLLM inference beyond static image perception (**contribution #4**).

# 139

141

#### 2 Preliminary

In this section, we describe briefly the prevalently 140 adopted image preprocessing methods and image-142 text input ways of MLLMs.

Image preprocessing. For a given image I, a 143 **naive** processing style is to simply resize it to a 144 preset fixed resolution and then feed it into an 145 vision encoder to generate visual representations. 146 This could be formulated as:  $\mathbf{v} = \mathcal{F}(R(\mathbf{I})) =$ 147  $(v_1, v_2, \ldots, v_{L_v})$ , where  $\mathcal{F}$  is the vision encoder, 148 R is the resize operation, and  $L_v$  is the number of 149 visual representations. Due to the constraints of 150 the naive version's fixed and limited resolution, an-151 other method, known as AnyRes, was introduced. 152 It divides the original image into several equal-area 153 blocks and imposes a maximum limit, M, on the 154 155 number of divided blocks. The vision encoder then independently encodes each block and the overall 156 image. Finally, all the encoded visual representations are integrated together. This allows flexible 158 processing of various resolutions. Denoting  $\mathbf{I}^{(0)}$ 159 as the whole image and  $\{\mathbf{I}^{(1)}, \dots, \mathbf{I}^{(a)}\}\ (a \leq M)$ 160 as the blocks, the AnyRes could be formulate as: 161  $\mathbf{v} = \mathcal{F}(A(\mathbf{I})) = [\mathbf{v}_0, \mathbf{v}_1, \dots, \mathbf{v}_a],$  where A denotes 162 the AnyRes operation and  $\mathbf{v}_i = \mathcal{F}(R(\mathbf{I}^{(i)})) =$ 163  $(v_{(i,1)}, v_{(i,2)}, \ldots, v_{(i,L_v)}), i = 0, 1, \ldots, a.$  It is noteworthy that the naive method can be considered a special case of AnyRes when a = 0. 166

Imga-Text joint input for MLLM. Common 167 MLLMs link a vision encoder to the pre-trained 168 LLM via projection or alignment modules, allow-169 ing language generation through the autoregressive 170 capabilities of their LLM base. Specifically, given 171 an image I and an input prompt x, I is first encoded 172 into a set of visual representations as described 173 in the previous sub-section. 174 Subsequently, these visual representations, along with the text 175 input, are fed into the LLM base of the MLLM. 176 Assuming the length of the output sequence and text input are  $L_y$  and  $L_x$  respectively, the 178 probability for a MLLM  $\Phi_{\theta}$  to generate an output 179  $\mathbf{y} = (y_1, y_2, \dots, y_{L_y})$  conditioned on the visual input  $\mathcal{F}(\cdot(\mathbf{I})) = (v_{(0,1)}, \dots, v_{(a,L_v)})$  and the text 181 input  $\mathbf{x} = (x_1, x_2, \dots, x_{L_r})$  is:  $\Phi_{\theta}(\mathbf{y} | \mathcal{F}(\cdot(\mathbf{I})), \mathbf{x}) =$  $\prod_{i=1}^{L_y} \Phi_{\theta}(y_i | v_{(0,1):(a,L_v)}, x_{1:L_x}, y_{1:i-1}),$ where 183  $\mathcal{F}(\cdot)$  could represent  $\mathcal{F}(R)$  as naive resize or  $\mathcal{F}(A)$  as AnyRes.

#### Methodology 3

188

In this section, we introduce the Zoom Eye algorithm. Firstly, we brief the general tree search

algorithm. Subsequently, we elaborate on our implementation by initializing the components of the tree search algorithms in detail.

189

190

191

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

210

211

212

213

214

215

216

217

218

219

220

221

223

224

225

226

227

# 3.1 Abstraction of Tree Search

Tree node. Typically, a node in the tree structure comprises the following attributes:(1) id: The unique identifier of the node. (2) depth: Represents the level of the node within the tree. (3) value: Used to store numeric or textual data in the node. (4) children: A list of references to the node's children nodes, which facilitates traversal of the tree structure. (5) Other custom attributes

Tree search. The abstraction of the tree search algorithm could be modeled as a tuple  $(T, Q, \mathcal{R}, \mathcal{S})$ , where T is the tree structure consisting of a set of nodes, Q is a container that holds all the nodes that might be accessed in the next search step,  $\mathcal{R}$  is a ranking function used to select the highest priority node based on the used search algorithm, and  $\mathcal{S}$  represents the stopping criterion. The abstract search process is shown in Algorithm 1.

| Alg.  | 1 Abstraction of Tree Search Algorithm |
|-------|--|
| Requi | re: $T, Q, \mathcal{R}, \mathcal{S}$   |
| 1: In | itialize $Q$ as the empty queue {}     |
| 2: Q  | .append( $T$ .root)                    |
| 3: wl | hile $Q$ is not empty <b>do</b>        |
| 4:    | $n_t \leftarrow Q.pop()$               |
| 5:    | if $\mathcal{S}(n_t) ==$ True then     |
| 6:    | break                                  |
| 7:    | $s \leftarrow n_t.children.size$       |
| 8:    | for $j = 1,, s$ do                     |
| 9:    | $Q$ .append( $n_t$ .children[j])       |
| 10:   | $Q.sort(\mathcal{R})$                  |
|       |  |

Consider the example of a DFS search for a node with a value of 5 in the tree, in this case,  $\mathcal{R}$  is a function that sorts the nodes in Q in descending order of depth, and in ascending order of id when depths are equal. Meanwhile, S is a function checking if a node's value equals 5.

A specific implementation of Zoom Eye search involves three key questions: 1. How to formulate the image as a tree T (§3.2). 2. How to set the ranking function  $\mathcal{R}$  (§3.3). 3. How to determine the stopping criterion S (§3.4). Finally, we provide a description of the overall algorithm in §3.5.

# 3.2 Tree Representation for Image

We model the overall image as a tree T. A specific *node*, denoted as  $n_t$ , represents an image patch view  $\{I, b_t\}$ , where I is the image and  $\mathbf{b}_t = (x_{1,t}, y_{1,t}, x_{2,t}, y_{2,t})$  is the normalized bounding box coordinates. If the size of  $n_t$ 's image patch



Figure 3: Two image input methods for MLLMs with distinct image processing.

exceeds the predefined resolution by the image encoder, it can be further divided into four equalsized sub-patches, serving as its children with size 4. Nodes are recursively divided until they meet the resolution limit. At the start of the search, the root node  $T.root = {I, (0, 0, 1, 1)}$  representing the overall image is visited.

However, due to the detailed nature of highresolution images and information loss from downsampling to the vision encoder's fixed resolution, MLLMs frequently struggle to accurately capture key parts of an image initially. Consequently, MLLMs should be be allowed to continuously scan and zoom into the current view (i.e., explore deeper nodes) for more focused information. In our implementation, we consider two image input methods to enable MLLMs to perceive the local patch represented by  $n_t$ : (1) Local Input: only the local patch is provided, suitable for earlier single-image input MLLMs with naive image preprocessing method (Li et al., 2023; Liu et al., 2024c,a). (2) Global+Local Input: both the global image and local patch are input, ideal for advanced MLLMs using AnyRes preprocessing method (Liu et al., 2024b; Li et al., 2024; Chen et al., 2024b). In this case, we use the visual prompt with a red rectangle to emphasize the local focus, applying naive processing to the global image and AnyRes to the local patch, as shown in Figure 3. Denoting  $\mathcal{V}(n_t)$  as the final image input, we have:

$$\mathcal{V}(n_t) = \begin{cases} [\mathcal{F}(R(\mathbf{I}.\operatorname{crop}(\mathbf{b}_t))] & \operatorname{Local} \\ [\mathcal{F}(R(\mathbf{I})), \mathcal{F}(A(\mathbf{I}.\operatorname{crop}(\mathbf{b}_t))] & \operatorname{Global+Local} \end{cases}$$
(1)

#### 3.3 Ranking Function

As shown in Algorithm 1,  $\mathcal{R}$  is used to rank the nodes with the priority value to determine which

| A 1     |            |   | D I | 1 '    | <b>T</b> |         | 0  | <b>C</b> . | •     | <b>n</b> · · | •     |
|---------|------------|---|-----|--------|----------|---------|----|------------|-------|--------------|-------|
| лі      | α_         | , | Ran | Vino.  | Hund     | nun iti | XT | Nton       | ning  | ( 'r1†       | erion |
| <b></b> | <b>×</b> . | 4 | ran | KIII Z | 1 unc    | JUOII   | œ  | SIUD       | UIII2 | UIII         | UIUII |
|         | <u> </u>   |   |     | C 3    |          |         |    |            | · 0   |              |       |

one to visit in the next step. A well-defined  $\mathcal{R}$ strategically steers the search process. In Zoom Eye, we adopt the MLLM to calculate the priority value and use  $\mathcal{R}$  to sort nodes by the value. Specifically, let o denote the visual cue that is crucial for answering the question, a MLLM should have the following capabilities: (1) It could perceive whether o exists within the visible view; (2) If ooccupies a small area and is not clearly visible, it can leverage the common sense knowledge to infer whether o might be discerned through further zooming. Thus, we query the MLLM with two prompts  $p_e(o)$  and  $p_l(o)$  (e.g., "Is there a *o* in the sub-patch?", "Is it possible to find a *o* by further zooming the sub-patch?") to trigger these two capabilities, and use the ratio of the next-word probability of the token "Yes" and "No" as priority values. We refer to these two values as existing confidence and latent confidence, denoted as  $c_e$  and  $c_l$ .

262

264

265

267

269

270

273

274

275

276

277

278

279

281

282

285

286

287

290

291

293

The overall priority value for a node is the weighted sum of  $c_e$  and  $c_l$ . We introduce a weight function  $\mathcal{W}(d)$  that is related to a node's depth. When the depth is shallow, indicating minimal zoom and the MLLM might not clearly perceive the cue, assign more weight  $c_l$ . As depth increases, shift more weight to  $c_e$ . Finally, ranking function  $\mathcal{R}$  is introduced to rank nodes by the overall priority value, as shown in Algorithm 2.

# 3.4 Stopping Criterion

Zoom Eye exits the search process when the MLLM provides feedback that the current view

256

257

261

is sufficient to answer the provided question, denoted as  $q_s$ . Specifically, we query the MLLM with a prompt  $p_a(q_s)$  (e.g., "Could you answer  $q_s$  now?") and use the same method as described in §3.3 to quantify the positive feedback. We refer to it as *answering confidence*, denoted as  $c_a$ . When  $c_a$  exceeds a predefined threshold  $\tau$ , the search terminates. The implementation of S is shown in Algorithm 2.

30

306

307

311

312

313

314

315

318

319

320

323

324

325

326

330

331

332

333

336

## 3.5 Overall Search Algorithm

With the above notations in place, we now describe how Zoom Eye works for a given image-question pair ( $\mathbf{I}$ , q). The complete algorithm workflow is shown in Appendix D.4.

**Generating visual cues to guide the search.** Before search, the MLLM has to predefine the visual cues essential for addressing q, enabling a targeted and guided search based on these cues. We utilize the in-context capability from the LLM base of the MLLM, using a sequence of contextual examples as prefixes to generate visual cues. Ultimately, the MLLM produces k visual cues  $\{o_1, \ldots, o_k\}$ pertinent to q. Each  $o_i$  ( $i \in \{1, \ldots, k\}$ ) can be categorized into two types: (type 1) those requiring a search for a single instance, and (type 2) those requiring identification of all instances in the image.

|   | Question   | Visual cues | Туре              |
|---|--|-------------|-------------------|
| 1 | What is the color of the dog?                        | dog         | type 1            |
| 2 | What is the relative position of the dog to the cat? | dog,<br>cat | type 1,<br>type 1 |
| 3 | How many dogs in the image?                          | all dogs    | type 2            |

Table 1: Examples of visual cues and their types.

Searching for cues. For each cue  $o_i$  ( $i \in$  $\{1, \ldots, k\}$ ), Zoom Eye explores the image tree to capture pertinent visual information. When searching for type 1 cues, the search is guided with  $\mathcal{R}$  and concludes as soon as it meets S, then the current node is recorded in a list L. For a single type 1 clue, as shown in line 1 of Table 1, the applied  $q_s$  for Sis the input question q. If multiple type 1 clues are generated as in line 2 of Table 1, we introduce a decomposed question template  $p_{dq}(o_i)$  such as "what is the location of the  $\{o_i\}$ ?" specific to each cue. In this case, the applied  $q_s$  of  $o_i$  is  $p_{da}(o_i)$ . If a type 2 cue is generated, as shown in line 3 of Table 1, S is not applied, and we search the whole tree to add all nodes with sufficient existing confidence to L.

Answering the question using the searched cues.

Given the searched nodes  $L = \{n_1^*, \ldots, n_K^*\}$ , the MLLM formulates a response to the input question q by synthesizing information of these nodes. Denoting  $\mathbf{b}_i^* = (x_{1,i}^*, y_{1,i}^*, x_{2,i}^*, y_{2,i}^*)$ as the bounding-box of  $n_i^*$  ( $i \in \{1, \ldots, K\}$ ), we union the bounding-box coordinates of all nodes in L to create a union bounding-box  $\mathbf{b}^* =$  $(\min_i x_{1,i}^*, \min_i y_{1,i}^*, \max_i x_{2,i}^*, \max_i y_{2,i}^*)$ . For the two distinct image input methods, we apply Eq. 1 to feed the focused region  $\mathbf{b}^*$  along with q into models and derive the final response. 337

338

339

341

342

343

344

346

347

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

369

370

371

372

373

374

375

376

377

378

379

380

381

382

384

385

## 4 **Experiments**

#### 4.1 Implementation Details

**Local input.** We select LLaVA-v1.5-7B (Liu et al., 2024a) as the base MLLM, with the naive image processing. We set  $\tau$  at 0.8 and define W as  $\frac{1-b}{D^2} \times d^2 + b$ , where D denotes the depth of the image tree, d is the depth of the visited node during the search, and b is a bias value, set here at 0.2.

**Global + Local input.** We select Qwen2.5VL-3B (Bai et al., 2025), LLaVA-ov(oneVision)-7B (Li et al., 2024), and InternVL2.5-8B (Chen et al., 2024a) as our MLLMs, with the AnyRes image processing. For LLaVA-ov and InternVL, we define the maximum AnyRes block as 12, and for QwenVL, we set the max pixels as 12, 845, 056. We set  $\tau$  at 0.6 and define W similarly to the above, except with *b* of 0.6.

For both input implementation, we set the maximum search depth at 2 when searching for type 2 cues to save costs. Additionally, the decomposed question template  $p_{dq}(o_i)$  is assigned as "What is the appearance of the  $\{o_i\}$ ?". More details are described in Appendix D.

#### 4.2 Results on High-Resolution Benchmark

**Evaluated benchmark.** We evaluate Zoom Eye on two meticulously curated high-resolution benchmarks. The first,  $V^*$  **Bench** (Wu and Xie, 2024), with an average resolution of 2246x1582, features sub-tasks in attribute recognition and spatial reasoning. The second, **HR-Bench 8K** (Wang et al., 2024) boasts average resolution of 7680, which consists of two sub-tasks: Fine-grained Single-instance Perception (FSP) and Fine-grained Cross-instance Perception (FCP). The 8K images are cropped around the objects in question to produce **HR-Bench 4K**. Both benchmarks are comprised of rich visual elements and required detailed perception to accurately respond. More results are displayed in Ta-

| $V^*$ Bench                              |        |         | H       | IR-Bench | 4K    | HR-Bench 8K |        |       |         |
|--|--------|---------|---------|----------|-------|-------------|--------|-------|---------|
| Model                                    | Attr   | Spatial | Overall | FSP      | FCP   | Overall     | FSP    | FCP   | Overall |
| Open-source MLLMs                        |        |         |         |          |       |             | 1      |       |         |
| minigptv2-7B (Chen et al., 2023a)        | -      | -       | -       | 25.75    | 25.25 | 25.50       | 26.0   | 26.25 | 26.13   |
| LLaVA-v1.6-7B (Liu et al., 2024b)        | 60.87  | 63.16   | 61.78   | 49.0     | 46.75 | 47.88       | 37.25  | 44.25 | 40.75   |
| LLaVA-v1.6-13B (Liu et al., 2024b)       | 60.0   | 64.47   | 61.78   | 49.75    | 41.25 | 45.50       | 38.0   | 38.25 | 38.13   |
| Yi-VL-34B (AI et al., 2024)              | -      | -       | -       | 46.0     | 42.75 | 44.38       | 39.50  | 38.50 | 39.0    |
| LLaVA-HR-X-7B (Luo et al., 2024)         | 51.30  | 64.47   | 56.54   | 57.75    | 46.25 | 52.0        | 42.0   | 41.25 | 41.63   |
| Closed-source MLLMs                      |        |         |         |          |       |             | 1      |       |         |
| OWen-VL-max (Bai et al., 2023)           | -      | -       | -       | 65.0     | 52.0  | 58.50       | 54.0   | 51.0  | 52.50   |
| GPT4o (Achiam et al., 2023)              | -      | -       | 66.0    | 70.0     | 48.0  | 59.0        | 62.0   | 49.0  | 55.5    |
| Baseline and Local Input Zoom Eye        |        |         |         |          |       |             |        |       |         |
| LLaVA-v1.5-7B (Liu et al., 2024a)        | 43.47  | 56.57   | 48.68   | 38.5     | 33.75 | 36.13       | 33.0   | 31.25 | 32.13   |
| LLaVA-v1.5-7B w/ Zoom Eye                | 83.45  | 82.89   | 83.25   | 67.75    | 38.75 | 53.25       | 65.50  | 36.0  | 50.75   |
| $\Delta$                                 | +40.48 | +26.32  | +34.57  | +29.25   | +5.0  | +17.12      | +32.50 | +4.75 | +18.62  |
| Baseline and Global+Local Input Zoom Eye |        |         |         |          |       |             | 1      |       |         |
| Qwen2.5VL-3B (Bai et al., 2025)          | 80.87  | 71.05   | 76.96   | 82.75    | 49.0  | 65.88       | 80.5   | 45.25 | 62.88   |
| Qwen2.5VL-3B w/ Zoom Eye                 | 88.70  | 89.47   | 89.01   | 86.75    | 53.50 | 70.13       | 84.75  | 52.0  | 68.38   |
| $\Delta$                                 | +7.83  | +18.42  | +12.05  | +4.0     | +4.50 | +4.25       | +4.25  | +6.75 | +5.50   |
| LLaVA-ov-7B (Li et al., 2024)            | 75.65  | 75.0    | 75.39   | 72.0     | 54.0  | 63.0        | 67.25  | 52.25 | 59.75   |
| LLaVA-ov-7B w/ Zoom Eye                  | 93.91  | 85.53   | 90.58   | 84.25    | 55.0  | 69.63       | 88.5   | 50.0  | 69.25   |
| $\Delta$                                 | +18.26 | +10.53  | +14.19  | +12.25   | +1.0  | +6.63       | +21.25 | -2.25 | +10.0   |
| InternVL2.5-8B (Chen et al., 2024a)      | 67.83  | 71.05   | 69.11   | 75.75    | 56.25 | 66.0        | 61.5   | 53.25 | 57.38   |
| InternVL2.5-8B w/ Zoom Eye               | 86.09  | 82.89   | 84.82   | 88.75    | 61.50 | 75.13       | 89.75  | 57.5  | 73.63   |
| $\Delta$                                 | +18.26 | +11.84  | +15.71  | +13.0    | +5.25 | +9.13       | +28.25 | +4.25 | +16.25  |

Table 2: Results of different models on high-resolution benchmarks. FSP: Fine-grained Single-instance Perception; FCP: Finegrained Cross-instance Perception. More results are displayed in Table 5.

|             |           |           | МО       |             |                                |  |                        | AD        |                    | 1 | ]     | RS       |
|-------------|-----------|-----------|----------|-------------|--------------------------------|--|------------------------|-----------|--------------------|---|-------|----------|
| Method      | Calculate | Intention | Property | Orientation | $\operatorname{Color}^\dagger$ |  | Intention <sup>†</sup> | Attention | $Motion^{\dagger}$ |   | Count | Position |
| LLaVA-ov-7B | 36.33     | 27.55     | 55.0     | 14.94       | 34.19                          |  | 37.32                  | 71.89     | 30.61              | 1 | 32.95 | 61.40    |
| w/ Zoom Eye | 38.67     | 38.78     | 60.0     | 14.62       | 47.09                          |  | 38.56                  | 68.66     | 42.71              |   | 35.56 | 48.45    |
| $\Delta$    | +2.34     | +11.23    | +5.0     | -0.32       | +12.90                         |  | +1.24                  | -3.23     | +12.10             |   | +2.61 | -12.95   |

Table 3: Performance comparison on MME-RealWorld benchmark. This benchmark comprises numerous sub-tasks, and we only list those that exhibit obvious performance changes of Zoom Eye against the baseline. MO (Monitoring), AD (Autonomous Driving), and RS (Remote Sensing) are data categories within this benchmark. <sup>†</sup>This result is an average derived from multiple similar sub-tasks (e.g., Color is the average of Vehicle Color and Person Color).

#### ble 5.

386

387

389

396

397

400

401

402

403

404

405

406

407

**Main results.** As shown in Table 2, all evaluated models exhibit significant performance gains after incorporating Zoom Eye, highlighting its model-agnostic applicability. For instance, LLaVA-ov-7B achieves performance improvements of 14.19%, 6.63%, and 10.00% on  $V^*$  Bench, HR-Bench 4K, and HR-Bench 8K, respectively. In conjunction with the case studies presented in Figure 5, these results demonstrate that vision-level reasoning enables MLLMs to more effectively capture fine-grained and task-relevant visual information in complex scenes, thereby enhancing their overall visual understanding capabilities.

#### 4.3 Results on Real-World Benchmark

**Evaluated benchmark.** We further evaluate Zoom Eye on MME-RealWorld (Zhang et al., 2024), a manually annotated benchmark tailored for realworld applications, featuring an average resolution of  $2000 \times 1500$ . It includes 5 data categories and 43 sub-class tasks. Due to the page limit, we report on only 13 sub-tasks that show significant performance changes with Zoom Eye. These sub-tasks span 3 data categories, with similar types merged (e.g., Vehicle Color and Person Color into Color) to present average scores. Detailed results are provided in Appendix C. 408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

Results. As shown in Table 3, Zoom Eye improves the performance of LLaVA-ov-7B on most sub-tasks, especially on MO/Intention (+11.23%), MO/Color (+12.9%), and AD/Motion (+12.1%). However, we also notice that the model's performance with Zoom Eye decline on some sub-tasks. We selecte one error example each from MO/Orientation and RS/Position and display them in Figure 5. For MO/Orientation, the low direct response scores for LLaVA-ov, as seen in the Table 3, along with error example in the figure, suggest a probable deficiency of orientation data during training, negatively impacting model performance in this aspect. For RS/Position, despite Zoom Eye locates the target, the final response was incorrect, suggesting the model struggles to link positional relationships between the full image and sub-images, resulting in a marked decline in performance on this sub-task.



Figure 4: The relationship between the number of search steps and the performance of the MLLM. The experimental statistics are derived from LLaVA-ov-7B's results on  $V^*$  Bench.

These error examples reveal the model's deficiencies, by which we will guide the direction of improvements in the model's capabilities in our future work.

#### 4.4 Ablation Studies

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462 463

464

465

466

467

# 4.4.1 Vision-level test-time scaling

We progressively reduce the answering confidence threshold  $\tau$  and analyze the relationship between the number of search steps and the performance of the MLLM, as illustrated in Figure 4.

From the figure, it can be seen that as the number of search steps increases, the model performance improves and eventually stabilizes. This behavior is analogous to the test-time scaling in text-level reasoning, where the accuracy of the final answer improves with more CoT tokens being explored. This finding could be viewed as a form of **vision-level test-time scaling**, where exploring more detailed zoomed information instead of the static image could enhance the ability of MLLM to generate more accurate responses.

When deploying Zoom Eye in real-world scenarios, we can adjust the confidence threshold or the maximum number of search steps based on specific needs to achieve the best trade-off between performance and efficiency.

# 4.4.2 Does the Zoom operation contribute to the improvement of the MLLM?

By comparing the answer accuracy of MLLM when Zoom is successful versus when it fails, we investigate the contribution of the Zoom operation to the model. As shown in Table 4, the accuracy sees a remarkable improvement (from 54.55% to 93.45%) when Zoom is successfully applied. This substantial gain highlights the critical role of the Zoom operation. By effectively refining the model's focus on relevant visual details, it contributes to more

| Used MLLM   | Zoom Successfully | Performance |
|-------------|-------------------|-------------|
| LLaVA-ov-7B | ✓                 | 93.45       |
| LLaVA-ov-7B | ×                 | 54.55       |

Table 4: Comparison of MLLM performance conditioned on whether zoom is successful. A zoom is considered successful when the searched box covers at least 50% of the target object. The experimental statistics are derived from  $V^*$  Bench.

accurate and reliable responses, reinforcing its importance as a key mechanism for optimizing visual understanding.

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

#### 4.5 Case Study

We visualize some cases in Figure 5, along with error examples mentioned in §4.3. We present cases for single type 1 cue, multiple type 1 cues, and type 2 cue, which is corresponding to the examples in Table 1. From the figure, it can be observed that Zoom Eye accurately seeks out cues, enabling the MLLM to focus on the crucial visual information and respond to queries precisely.

#### 5 Related Work

Multimodal LLMs. Since the advent of large language models (LLMs), they have achieved success across various linguistic applications, such as incontext learning (Dong et al., 2022; Zhang et al., 2022) and retrieval augmented generation (Liu et al., 2024d; Zhao et al., 2024b,c), which facilitated the emergence of Multimodal LLMs, with pioneering works including (Alayrac et al., 2022; Li et al., 2023; Koh et al., 2023). Following these, LLaVA (Liu et al., 2024c) employed GPT-4 (Achiam et al., 2023) to develop training data, inspiring a series of works focused on visual instruction data (Liu et al., 2024a; Dai et al., 2023; Chen et al., 2023b). Since these models utilize pretrained vision encoders (Radford et al., 2021; Zhai et al., 2023) to process image information, the resolution that MLLMs can handle is limited by the input resolution of these encoders. To address this limitation, AnyRes was developed to flexibly manage images of varying resolutions (Liu et al., 2024b; Chen et al., 2024b). Additionally, there are efforts focused on utilizing high-resolution encoders to accommodate high-resolution images (Lu et al., 2024; Wei et al., 2025). However, despite these efforts, the perception of the image by the MLLM remains as the original image itself. We hope to enable MLLMs to explore the varying hierarchical features of images to capture key information.



Figure 5: Examples of Zoom Eye. The resolution of the image is displayed. Red rectangles are patches searched by Zoom Eye.

Tree-based search. Tree-based search algorithms 509 have been applied in text-only LLM reasoning and 510 511 have demonstrated superior performance. Early works such as (Wei et al., 2022; Wang et al., 2022) relied on chain reasoning, a method susceptible to 513 errors in one step propagating through subsequent 514 steps. Consequently, ToT (Yao et al., 2024b) pro-515 516 posed a tree-based reasoning method that leverages the expansiveness of tree structures to widen the 517 reasoning space. Simultaneously, several similar 518 studies were also introduced, which define a de-519 composed question step as a node and utilize beam 520 search (Xie et al., 2023) and Monte-Carlo Tree 521 Search (Hao et al., 2023) to uncover optimal solu-522 tions. Subsequently, TS-LLM (Feng et al., 2023) utilized reinforcement learning to increase search 524 depth, further enhancing reasoning performance. In our work, we conceptualize an image as a tree to search for crucial visual information using a spe-527 cific algorithm. A close-related work is  $V^*$  (Wu 528 and Xie, 2024), and we describe the detailed comparison with it in Appendix B. 530

#### 6 Limitations

531

Although Zoom Eye offers several advantages, such as strong interpretability, model-agnostic, and training-free, it also comes with certain limitations. First, the current search procedure relies on heuristic strategies, including manually defined ranking functions and stopping criteria. While these designs are effective in many settings, they may not generalize optimally across all image types or task conditions. Second, the image is partitioned into fixed-size patches to construct the hierarchical tree structure, which may not align well with the semantic regions of the image. As a result, some visual cues may be fragmented or overlooked during traversal. Lastly, Zoom Eye is primarily tailored for natural images with spatially distributed visual elements. It is less applicable to document understanding tasks, where layout, reading order, and structured information (e.g., tables, forms) are central. Addressing these challenges—such as by integrating learnable search strategies or adaptive patch partitioning—will be an important direction for future work. 542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

## 7 Conclusion

To address the limitations of text-level visual reasoning, we propose Zoom Eye, a type of visionlevel reasoning method, a tree search algorithm designed to navigate the hierarchical and visual nature of images to capture detailed crucial information. Through prompts guiding MLLMs, we develop a ranking function and stopping criterion for Zoom Eye, which steers models to efficiently search along the image tree, seek out pertinent information, and accurately respond to related queries. Experiments show the broad-applicability and effectiveness of Zoom Eye, which substantially improves MLLMs' performance. Notably, Zoom Eye exhibits a testtime scaling phenomenon analogous to that observed in text-level reasoning. Meanwhile, through the analysis of failure cases, we identify several inherent limitations in current MLLMs' visual reasoning capabilities, which we aim to address in future work.

#### References

574

576

579

580

581

584

590

591

592

593

594

595

596

597

598

603

611

612

613

614

615

616

617

618

619

623

629

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- 01. AI, :, Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang, Shiming Yang, Tao Yu, and 13 others. 2024. Yi: Open foundation models by 01.ai. *Preprint*, arXiv:2403.04652.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, and 1 others. 2022. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716– 23736.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 1(2):3.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, and 1 others. 2025. Qwen2. 5-vl technical report. arXiv preprint arXiv:2502.13923.
- Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. 2023a. Minigpt-v2: large language model as a unified interface for vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*.
- Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2023b. Sharegpt4v: Improving large multimodal models with better captions. *arXiv preprint arXiv:2311.12793*.
- Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, and 1 others. 2024a. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*.
- Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, and 1 others. 2024b. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi.

2023. Instructblip: Towards general-purpose visionlanguage models with instruction tuning. *Preprint*, arXiv:2305.06500.

- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Tianyu Liu, and 1 others. 2022. A survey on incontext learning. *arXiv preprint arXiv:2301.00234*.
- Yuhao Dong, Zuyan Liu, Hai-Long Sun, Jingkang Yang, Winston Hu, Yongming Rao, and Ziwei Liu. 2024. Insight-v: Exploring long-chain visual reasoning with multimodal large language models. *arXiv preprint arXiv:2411.14432*.
- Xidong Feng, Ziyu Wan, Muning Wen, Stephen Marcus McAleer, Ying Wen, Weinan Zhang, and Jun Wang. 2023. Alphazero-like tree-search can guide large language model decoding and training. *arXiv preprint arXiv:2309.17179*.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shi-rong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Shibo Hao, Yi Gu, Haodi Ma, Joshua Hong, Zhen Wang, Daisy Wang, and Zhiting Hu. 2023. Reasoning with language model is planning with world model. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 8154–8173.
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, and 1 others. 2024. Openai o1 system card. *arXiv preprint arXiv:2412.16720*.
- Jing Yu Koh, Ruslan Salakhutdinov, and Daniel Fried. 2023. Grounding language models to images for multimodal inputs and outputs. In *International Conference on Machine Learning*, pages 17283–17300. PMLR.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024. Llavaonevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024a. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306.

630

631

632

646 647 648

643

644

645

653

654

655

656

657

658

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

- 690 694 710 711 712 713 715 716 719 722 723 724 725 730 731 732 733 734 736
- 737
- 739

- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024b. Llavanext: Improved reasoning, ocr, and world knowledge.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024c. Visual instruction tuning. Advances in neural information processing systems, 36.
- Jingyu Liu, Jiaen Lin, and Yong Liu. 2024d. How much can rag help the reasoning of llm? arXiv preprint arXiv:2410.02338.
- Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, and 1 others. 2024. Deepseek-vl: towards real-world vision-language understanding. arXiv preprint arXiv:2403.05525.
- Gen Luo, Yiyi Zhou, Yuxin Zhang, Xiawu Zheng, Xiaoshuai Sun, and Rongrong Ji. 2024. Feast your eyes: Mixture-of-resolution adaptation for multimodal large language models. arXiv preprint arXiv:2403.03003.
- Fanqing Meng, Lingxiao Du, Zongkai Liu, Zhixiang Zhou, Quanfeng Lu, Daocheng Fu, Botian Shi, Wenhai Wang, Junjun He, Kaipeng Zhang, and 1 others. 2025. Mm-eureka: Exploring visual aha moment with rule-based large-scale reinforcement learning. arXiv preprint arXiv:2503.07365.
- M Minderer, A Gritsenko, A Stone, M Neumann, D Weissenborn, A Dosovitskiy, A Mahendran, A Arnab, M Dehghani, Z Shen, and 1 others. 2022. Simple open-vocabulary object detection with vision transformers. arxiv 2022. arXiv preprint arXiv:2205.06230, 2.
- OpenAI. 2025. o3/o4 mini system card. https:// openai.com/index/o3-o4-mini-system-card/.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, and 1 others. 2021. Learning transferable visual models from natural language supervision. In International conference on machine learning, pages 8748–8763. PMLR.
- Haozhan Shen, Peng Liu, Jingcheng Li, Chunxin Fang, Yibo Ma, Jiajia Liao, Qiaoli Shen, Zilun Zhang, Kangjia Zhao, Qianqian Zhang, Ruochen Xu, and Tiancheng Zhao. 2025. Vlm-r1: A stable and generalizable r1-style large vision-language model. arXiv preprint arXiv:2504.07615.
- Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. 2023. Eva-clip: Improved training techniques for clip at scale. arXiv preprint arXiv:2303.15389.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, and 1 others. 2023. Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288.

Wenbin Wang, Liang Ding, Minyan Zeng, Xiabin Zhou, Li Shen, Yong Luo, and Dacheng Tao. 2024. Divide, conquer and combine: A training-free framework for high-resolution image perception in multimodal large language models. arXiv preprint arXiv:2408.15556.

740

741

742

743

744

745

746

747

748

749

750

751

753

756

757

758

759

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. arXiv preprint arXiv:2203.11171.
- Haoran Wei, Lingyu Kong, Jinyue Chen, Liang Zhao, Zheng Ge, Jinrong Yang, Jianjian Sun, Chunrui Han, and Xiangyu Zhang. 2025. Vary: Scaling up the vision vocabulary for large vision-language model. In European Conference on Computer Vision, pages 408-424. Springer.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824-24837.
- Penghao Wu and Saining Xie. 2024. V?: Guided visual search as a core mechanism in multimodal llms. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 13084-13094.
- Yuxi Xie, Kenji Kawaguchi, Yiran Zhao, Xu Zhao, Min-Yen Kan, Junxian He, and Qizhe Xie. 2023. Decomposition enhances reasoning via self-evaluation guided decoding. Preprint, arXiv:2305.00633.
- Guowei Xu, Peng Jin, Hao Li, Yibing Song, Lichao Sun, and Li Yuan. 2024. Llava-cot: Let vision language models reason step-by-step. Preprint, arXiv:2411.10440.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, and 1 others. 2024. Qwen2 technical report. arXiv preprint arXiv:2407.10671.
- Huanjin Yao, Jiaxing Huang, Wenhao Wu, Jingyi Zhang, Yibo Wang, Shunyu Liu, Yingjie Wang, Yuxin Song, Haocheng Feng, Li Shen, and 1 others. 2024a. Mulberry: Empowering mllm with o1-like reasoning and reflection via collective monte carlo tree search. arXiv preprint arXiv:2412.18319.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2024b. Tree of thoughts: Deliberate problem solving with large language models. Advances in Neural Information Processing Systems, 36.
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. Sigmoid loss for language image pre-training. Preprint, arXiv:2303.15343.

- 794 795
- 797

- 803 804

806

- 808 809 810
- 811
- 812 813 814
- 815
- 816

- 817

819 820 821

822

823

824

825

826

818

arXiv:2410.10293.

Xinyu Zhu, Junjie Wang, Lin Zhang, Yuxiang Zhang, Yongfeng Huang, Jiaxing Zhang, Yujiu Yang, and 1 others. 2023. Solving math word problems via cooperative reasoning induced language models. In The 61st Annual Meeting Of The Association For Computational Linguistics.

Yi-Fan Zhang, Huanyu Zhang, Haochen Tian, Chaoyou

Fu, Shuangqing Zhang, Junfei Wu, Feng Li, Kun

Wang, Qingsong Wen, Zhang Zhang, and 1 oth-

ers. 2024. Mme-realworld: Could your multimodal

llm challenge high-resolution real-world scenarios

Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex

Tiancheng Zhao, Qianqian Zhang, Kyusong Lee, Peng

Liu, Lu Zhang, Chunxin Fang, Jiajia Liao, Kelei

Jiang, Yibo Ma, and Ruochen Xu. 2024a. Omchat:

A recipe to train multimodal language models with

strong long context and video understanding. arXiv

Xinping Zhao, Dongfang Li, Yan Zhong, Boren

Hu, Yibin Chen, Baotian Hu, and Min Zhang.

for retrieval-augmented generation. arXiv preprint

Xinping Zhao, Yan Zhong, Zetian Sun, Xinshuo Hu,

Zhenyu Liu, Dongfang Li, Baotian Hu, and Min

Zhang. 2024c. Funnelrag: A coarse-to-fine pro-

gressive retrieval paradigm for rag. arXiv preprint

Seer: Self-aligned evidence extraction

Smola. 2022. Automatic chain of thought prompt-

arXiv preprint

arXiv preprint

that are difficult for humans?

ing in large language models.

preprint arXiv:2407.04923.

arXiv:2408.13257.

arXiv:2210.03493.

2024b.

arXiv:2410.11315.

#### Α **Results of More MLLMs on High-Resolution Benchmark**

We present the results of additional MLLMs on high-resolution benchmarks in Table 5, including models of smaller or larger scale. Consistent with the findings in the main paper, all evaluated models exhibit improved performance after being adapted to Zoom Eye, further demonstrating the effectiveness of vision-level reasoning in handling complex visual scenarios.

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

#### **Compared with Other HR Processing** B **Methods**

Compared method. We choose two high-resolution processing methods to compare— $V^*$  (Wu and Xie, 2024) (published in CVPR 2024) and  $DC^2$  (Wang et al., 2024) (published in AAAI 2025).  $V^*$  is a LLM-guided search pipeline for MLLMs. To match the input resolution of the  $V^*$  model, we specifically trained a 224px version of the LLaVA-v1.5 model for a fair comparison. Apart from using CLIP-224 (Radford et al., 2021) as the vision encoder, all other settings were identical to those of LLaVA-v1.5.

 $DC^2$  (Divide, Conquer, and Combine) is a framework that supplements visual information using text for high-resolution images understanding. Like our approach, it recursively splits an image into patches. The MLLM then generates textual descriptions for each leaf patch. These descriptions are then relayed to the parent nodes, which create combined descriptions by synthesizing the contents from their child nodes with their own. This process continues up to the root node.

# **B.1** Zoom Eye vs. V<sup>\*</sup>.

From Table 6, it is evident that compared to  $V^*$ , our method offers several advantages: (1) The  $V^*$ pipeline requires specifically targeted training data, making zero-shot searches impossible, whereas our method utilizes the native capabilities of MLLMs, allowing adaptation to any MLLM without additional training; (2)  $V^*$ 's search process necessitates the integration of another specially trained MLLM to guide the search, along with an extra high-resolution image encoder (Minderer et al., 2022)(768px), while our approach operates at the native resolution of MLLMs and conducts searches independently; (3) Our method demonstrates superior performance.

|  | $V^*$ Bench |         | HR-Bench 4K |        |       | HR-Bench 8K |        |       |         |
|--|-------------|---------|-------------|--------|-------|-------------|--------|-------|---------|
| Model                                    | Attr        | Spatial | Overall     | FSP    | FCP   | Overall     | FSP    | FCP   | Overall |
| Baseline and Local Input Zoom Eye        |             |         |             |        |       |             |        |       |         |
| LLaVA-v1.5-13B (Liu et al., 2024a)       | 41.74       | 55.26   | 47.12       | 45.25  | 41.25 | 43.25       | 37.50  | 38.0  | 37.75   |
| LLaVA-v1.5-13B w/ Zoom Eye               | 87.83       | 81.58   | 85.34       | 73.0   | 43.25 | 58.13       | 67.25  | 45.50 | 56.38   |
| $\Delta$                                 | +46.09      | +26.32  | +38.22      | +27.75 | +2.00 | +14.88      | +29.75 | +7.50 | +18.63  |
| Baseline and Global+Local Input Zoom Eye |             |         |             |        |       |             |        |       |         |
| LLaVA-ov-0.5B (Li et al., 2024)          | 63.48       | 64.47   | 63.87       | 63.50  | 39.50 | 51.50       | 47.25  | 38.25 | 42.75   |
| LLaVA-ov-0.5B w/ Zoom Eye                | 85.22       | 73.68   | 80.62       | 75.50  | 39.75 | 57.63       | 68.50  | 38.25 | 53.38   |
| $\Delta$                                 | +21.74      | +9.21   | +16.75      | +12.00 | +0.25 | +6.13       | +21.25 | +0.00 | +10.63  |
| InternVL2.5-4B (Chen et al., 2024a)      | 69.57       | 71.05   | 70.16       | 77.50  | 53.75 | 65.63       | 63.00  | 49.25 | 56.13   |
| InternVL2.5-4B w/ Zoom Eye               | 85.22       | 77.63   | 82.20       | 81.25  | 56.75 | 69.00       | 80.00  | 52.25 | 66.13   |
| $\Delta$                                 | +15.65      | +6.58   | +12.04      | +3.75  | +3.00 | +3.37       | +17.00 | +3.00 | +10.00  |
| InternVL2.5-26B (Chen et al., 2024a)     | 73.91       | 72.37   | 73.30       | 82.00  | 66.25 | 74.13       | 73.00  | 61.75 | 67.38   |
| InternVL2.5-26B w/ Zoom Eye              | 91.30       | 86.84   | 89.53       | 89.75  | 68.25 | 79.00       | 89.25  | 63.00 | 76.13   |
| $\Delta$                                 | +17.39      | +14.47  | +16.23      | +7.75  | +2.00 | +4.87       | +16.25 | +1.25 | +8.75   |

Table 5: Results of more models on high-resolution benchmarks.

| Method                         | Input | Search | Zero | Indep. | $V^*$        | HR           |
|--------------------------------|-------|--------|------|--------|--------------|--------------|
|                                | Res.  | Res.   | shot | search | Bench        | Bench        |
| $V^*$ search(Wu and Xie, 2024) | 224   | 768    | ×    | ×      | 75.39        | 37.81        |
| Zoom Eye                       | 224   | 224    | ✓    |        | <b>81.58</b> | <b>47.63</b> |

Table 6: Performance comparison between Zoom Eye and  $V^*$  Search. Input Res.: The input resolution of the model generating the final response; Search Res.: The resolution required during the search process; Zero shot: Whether the method could be adapted for models without the need for specialized additional training; Indep. search: Whether the method could be applied directly to an MLLM instead of requiring an additional search model.

### **B.2** Zoom Eye vs. $DC^2$ .

Our approach differs from  $DC^2$  in two key ways: (1)  $DC^2$  uses textual modalities to supplement the missing visual information at high resolutions, whereas Zoom Eye employs simulated zooming operations, allowing the MLLM to actively discover missing visual details and provide more pertinent information for multimodal tasks; (2)  $DC^2$ is question-agnostic, generating descriptions consistently across different questions which may lead to unfocused textual content. In contrast, Zoom Eye is question-driven in its visual cues searching, yielding more precise visual information that is instrumental in answering the input question. Table 7 shows that Zoom Eye performs better than  $DC^2$ using the same MLLM.

# C Complete Results on MME-RealWorld Benchmark

We provide the complete results of Zoom Eye on MME-RealWorld Benchmark (Zhang et al., 2024), as show in Table 8. This benchmark includes 5 data categories: Monitoring (MO), Autonomous Driving (AD), OCR, emote Sensing (RS), and Diagram

| Method                                 | $V^*$ Bench | HR-Bench 8K |
|--|-------------|-------------|
| LLaVA-v1.5-7B                          | 48.7        | 32.1        |
| w/ DC <sup>2</sup> (Wang et al., 2024) | 57.3        | 39.5        |
| w/ Zoom Eye                            | 83.3        | 50.8        |
| LLaVA-v1.5-13B                         | 47.1        | 37.8        |
| w/ $DC^2$                              | 57.3        | 40.5        |
| w/ Zoom Eye                            | 85.3        | 56.4        |

Table 7: Performance comparison between Zoom Eye and  $DC^2$ .

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

and Table (TD). Since Zoom Eye is not applicable to the TD task, we do not conduct tests on it. It could be observed that Zoom Eye improves the performance of LLaVA-ov-7B across most subtasks, with particularly significant improvements in certain tasks. For instance, it achieves a 20.22% improvement in the Person<sub>color</sub> task, a 29.11% improvement in the Motion<sub>vehicle</sub> task, and a 12.93% improvement in the Visual<sub>trafficsignal</sub> task, demonstrating the effectiveness of Zoom Eye. However, performance declines were observed in some subtasks when using Zoom Eye. We have analyzed these cases in the main paper, revealing certain limitations of the employed MLLM. Addressing these issues will be a focus of our future work.

### **D** Implementation Details

Due to the page limit of the main paper, we provide more implementation details here. In §D.1 and §D.2, we detail the implementation of Local Input and Global+Local Input, respectively. §D.3 describes the implementations common to both. Finally, based on the introductions in the first three subsections, we present the complete algorithm workflow of Zoom Eye in §D.4.

# D.1 Local Input

We select LLaVA-v1.5-7B (Liu et al., 2024a) and 13B as our MLLMs, with the vision encoder's input

875

- 881
- 883

885

- 887
- 888 889

89

891 892

893

894

| Task |  | LLaVA <sub>ov</sub> -7B | +ZoomEye | $\Delta\uparrow$ |
|------|--|-------------------------|----------|------------------|
|      | Calculate                                | 36.33                   | 38.67    | +2.34            |
|      | Intention                                | 27.55                   | 38.78    | +11.23           |
|      | Property                                 | 55.0                    | 60.0     | +5.0             |
| мо   | Vehicle <sub>counting</sub>              | 59.89                   | 61.14    | +1.25            |
|      | Personcounting                           | 61.35                   | 61.87    | +0.52            |
|      | Vehiclelocation                          | 33.82                   | 33.82    | -                |
|      | Vehicleorientation                       | 19.35                   | 18.71    | -0.64            |
|      | Vehicle <sub>color</sub>                 | 43.65                   | 49.24    | +5.59            |
|      | Person <sub>color</sub>                  | 24.72                   | 44.94    | +20.22           |
|      | Personorientation                        | 10.53                   | 10.53    | -                |
|      | Intentionego                             | 28.62                   | 28.95    | +0.33            |
|      | Intention <sub>pedestrian</sub>          | 52.43                   | 53.40    | +0.97            |
|      | Intention <sub>vehicle</sub>             | 30.92                   | 33.33    | +2.41            |
|      | Interaction <sub>other2other</sub>       | 12.94                   | 13.43    | +0.49            |
|      | Attention <sub>trafficsignal</sub>       | 71.89                   | 68.66    | -3.23            |
|      | Interaction <sub>ego2pedestrain</sub>    | 27.36                   | 28.30    | +0.94            |
|      | Interaction <sub>ego2trafficsignal</sub> | 22.86                   | 25.71    | +2.85            |
| AD   | Interaction <sub>ego2vehicle</sub>       | 20.79                   | 19.80    | -0.99            |
|      | Objectsidentify                          | 64.40                   | 64.85    | +0.45            |
|      | Motion <sub>vehicle</sub>                | 23.42                   | 52.53    | +29.11           |
|      | Motion <sub>multivehicles</sub>          | 34.26                   | 34.75    | +0.49            |
|      | Visual <sub>trafficsignal</sub>          | 60.20                   | 73.13    | +12.93           |
|      | Motion <sub>pedestrain</sub>             | 34.15                   | 40.85    | +6.70            |
|      | Object <sub>count</sub>                  | 37.92                   | 39.86    | +1.94            |
|      | Motion <sub>multipedestrians</sub>       | 31.24                   | 31.64    | +0.40            |
|      | Scene understanding                      | 64.80                   | 64.80    | -                |
|      | Character identification                 | 57.60                   | 56.40    | -1.20            |
|      | Adver & product                          | 76.64                   | 78.37    | +1.73            |
| OCR  | Book map poster                          | 77.17                   | 75.24    | -1.93            |
|      | License                                  | 80.16                   | 82.39    | +2.23            |
|      | Phone & address                          | 77.82                   | 81.28    | +3.46            |
|      | Text recog                               | 74.87                   | 77.13    | +2.26            |
|      | Color                                    | 59.60                   | 60.56    | +0.96            |
| RS   | Count                                    | 32.95                   | 35.56    | +2.61            |
|      | Position                                 | 61.40                   | 48.45    | -12.95           |

Table 8: Performance comparison between Zoom Eye and the baseline model on MME-RealWorld benchmark. MO (Monitoring), AD (Autonomous Driving), OCR and RS (Remote Sensing) are data categories within this benchmark.

#### **Prompt Templates of Local Input**

- $p_e(o)$ : <local patch> Is there a {o} in the image? Answer Yes or No.
- *p<sub>l</sub>(o):* <local patch> According to your common sense knowledge and the content of the image, is it possible to find a {o} in the image? Answer Yes or No and tell the reason.
- *p<sub>a</sub>(q):* <local patch> Question: {*q*}
   \nCould you answer the question based on the the available visual information? Answer Yes or No.

where the o and q are the input visual cue and question, which could be referred to §3.5.

As mentioned in §3.5, the final visual input uses the union of all searched patches. However, when multiple distant patches are combined, they may form a large image. For MLLMs using naive resize processing, information can still be lost during downsampling. Therefore, for the Local Input Zoom Eye with naive resize processing, when the area of  $b^*$  is relatively large (with the longer side exceeding 1000px), we skip the Union operation. Instead, we paste the searched patches onto a blank image according to their relative positions in the original image, and then feed it to the MLLMs. An example is shown in Figure 6.



resolution as 336px and naive processing. We set the threshold of the stopping criterion at  $\tau = 0.8$  and define the weighted function as  $W = \frac{1-b}{D^2} \times d^2 + b$ , where *D* denotes the depth of the image tree, *d* is the depth of the visited node during the search, and *b* is a bias value, set here at 0.2. The prompt templates for calculating *existing confidence*, *latent confidence*, and *answering confidence* (please refer to §3.3 and §3.4 for the discussion on these three confidence values) are set as:

Figure 6: If the area of the union bounding box is too large, we paste the searched patches onto a blank image according to their relative positions in the original image, and then feed it to the MLLMs.It is notable that this operation is only applied to Local Input, while for Local+Global, we consistently provide the MLLMs with the full union patch as input.

939

940

941

942

943

944

945

946

947

948

949

950

#### D.2 Global + Local Input

We select LLaVA-ov(oneVision)-0.5B (Li et al., 2024) and 7B as our MLLMs, with the vision encoder's input resolution as 384px and AnyRes processing. We define the maximum AnyRes block as 12, set  $\tau$  at 0.6 and define W as  $\frac{1-b}{D^2} \times d^2 + b$ , where D denotes the depth of the image tree, d is the depth of the visited node during the search, and b is a bias value, set here at 0.6. The prompt templates for calculating *existing confidence, latent confidence,* and *answering confidence* are set as:

**Prompt Templates of Global + Local Input** 

- *p<sub>e</sub>(o):* <*global image><local patch>Is* there a {*o*} in the zoomed-in view? Answer Yes or No.
- p<sub>l</sub>(o): <global image><local patch> According to your common sense knowledge and the content of the zoomed-in view, along with its location in the image, is it possible to find a {o} by further zooming in the current view? Answer Yes or No and tell the reason.
- *p<sub>a</sub>(q):* <*global* image><*local* patch>Question: {*q*} \nCould you answer the question based on the the available visual information? Answer Yes or No.

#### **D.3** Additional Settings

For both input implementation, we set the maximum search depth at 2 when searching for type 2 cues to save costs. In  $\S3.5$ , we state that we search the whole tree to add all nodes with sufficient existing confidence to L if type 2 cue is generated. Thus, we introduce an additional threshold  $\tau_2$  for this condition, which is set at 0.8 for both implementation. The decomposed question template  $p_{da}(o_i)$  is assigned as "What is the appearance of the  $\{o_i\}$ ?". For type 1 search, a key aspect is determining the value of  $\tau$ . If it is set too low, an incorrect patch, which probably lead to erroneous guidance for MLLMs, may be selected. Conversely, setting  $\tau$  too high, surpassing the  $c_a$  values of all nodes in the tree, would compel MLLMs to search the entire tree unnecessarily, thus wasting time. Therefore, we adopt a strategy where  $\tau$  is progressively reduced as the number of search steps increases. Specifically, if the number of search steps exceeds the step threshold C, we reduce the value of  $\tau$  by

0.1. This reduction occurs every  $\delta$  steps, until the 984  $c_a$  value of a node having been visited surpasses  $\tau$ 985 or  $\tau$  falls below a predefined minimum limit  $\tau_{min}$ . 986 For both implementation, we set  $\delta$  at 2,  $\tau_{min}$  at 0, 987 and C as  $D \times 3$ . Finally, the in-context examples 988 we utilized to generate visual cues are denote as 989  $(q^{(1)}, \mathbf{o}^{(1)}, \dots, q^{(m)}, \mathbf{o}^{(m)})$  and are presented at the 990 end of this document. 991

## D.4 Complete Algorithm Workflow

With the aforementioned notation and description993in place, we provide the complete algorithm work-994flow in Algorithm 3, where the Zoom Eye search995method is shown in Algorithm 4.996

 $\leftarrow$ 

# Algorithm 3 Complete Algorithm Workflow of Zoom Eye

- **Require:** Multimodal LLM  $\Phi_{\theta}$ , input question-image pair (**I**, q), decomposed question template  $p_{dq}$ , in-context examples  $(q^{(1)}, \mathbf{0}^{(1)}, \dots, q^{(m)}, \mathbf{0}^{(m)}, q)$
- 1:  $\{o_1, \ldots, o_k\}$   $\Phi_{\theta}$ .generate $(q^{(1)}, \mathbf{o}^{(1)}, \ldots, q^{(m)}, \mathbf{o}^{(m)}, q)$ 2: Initialize *L* as the empty list 3: Build I as a tree *T* 4: for  $i = 1, \ldots, k$  do 5: if k == 1 then
- 6:  $q_s \leftarrow q$
- 7: else
- 8:  $q_s \leftarrow \mathbf{p}_{dq}(o_i)$
- 9:  $L.extend(ZOOM EYE(T, o_i, q_s))$
- 10:  $\mathbf{b}^* \leftarrow$  Union bounding-boxes of all nodes in L
- 11:  $n^* \leftarrow {\mathbf{I}, \mathbf{b}^*}$
- 12: Final response  $\leftarrow \Phi_{\theta}$ .generate $(\mathcal{V}(n^*), q)$

963 964 965

966

967

968

970

972

975

976

977

979

983

962

951

952

953

955

956

960

#### Algorithm 4 Zoom Eye Search

```
Require: Threshold of type 1 cue and type 2 cue (\tau, \tau_2),
     minimum limit \tau_{min}, interval \delta
 1: function ZOOM EYE(T, o_i, q_s)
 2:
         Initialize Q as the empty queue \{\}
 3:
         Q.append(T.root)
 4:
         Initialize L_i as the empty list
 5:
         search all \leftarrow o_i.startswith("all")
 6:
         if not search all then
             ZOOM EYE TYPE 1(T, Q, L_i, q_s, \tau)
 7:
 8:
         else
 9:
             ZOOM EYE TYPE 2(Q, L_i, \tau_2)
10:
         return L_i
11:
12: function ZOOM EYE TYPE 1(T, Q, L_i, q_s, \tau)
         import \mathcal{R} and \mathcal{S} from Algorithm 2
13:
         count \leftarrow 0
14:
15:
         C \leftarrow T.depth \times 3
16:
         Initialize n_m as T.root to record the node with the
     best c_a
17:
         while Q is not empty do
18:
             n_t \leftarrow Q.\mathsf{pop}()
19:
             N.append(n_t)
20:
             count \leftarrow count + 1
21:
             if count \ge C then
22:
                  \tau \leftarrow \tau - 0.1
23:
                  C \leftarrow C + \delta
                 if \tau < \tau_{min} then
24:
25:
                      break
26:
             if \mathcal{S}(n_t, q_s, \tau) == True then
27:
                  L_i.append(n_t)
28:
                  break
29:
             else if \mathcal{S}(n_m, q_s, \tau) == True then
30:
                  L_i.append(n_m)
31:
                  break
32:
             if n_t.c_a \ge n_m.c_a then
33:
                  n_m \leftarrow n_t
34.
              s \leftarrow n_t.children.size
35:
             for j = 1, ..., s do
36:
                  Q.append(n_t.children[j])
37:
             Q.sort(\mathcal{R}(o_i))
38:
39: function ZOOM EYE TYPE 2(Q, L_i, \tau_2)
40:
         while Q is not empty do
41:
             if n_t.depth \geq 2 then
42:
                 break
43:
             c_e \leftarrow calculate the existing confidence of n_t
44:
             if c_e \geq \tau_2 then
45:
                  L_i.append(n_t)
46:
             s \leftarrow n_t.children.size
47:
             for j = 1, ..., s do
48:
                  Q.append(n_t.children[j])
```

#### **In-Context Examples of** V\* **Bench**

User: Question: What is the color of the boy's bag? If you want to answer the question, which objects' information do you need?

Assistant: To answer the question, I need know the location of the boy with a bag so that I can determine the color of the bag. So I need the information about the following objects: boy with a bag.

*User: Is the yellow car on the left or right side of the white car?* 

Assistant: To answer the question, I need know the location of the yellow car and the white car so that I can determine the positional relationship between the two of them. So I need the information about the following objects: white car and yellow car.

User: Tell me the number on the black board.

Assistant: To answer the question, I need know the location of the black board so that I can determine the number on it. So I need the information about the following objects: black board

**User:** Is the girl with pink hair on the left or right side of the man with backpack?

Assistant: To answer the question, I need know the location of the girl with pink hair and the man with backpack so that I can determine the positional relationship between the two of them. So I need the information about the following objects: girl with pink hair and man with backpack.

User: What kind of animal is on the red sign?

Assistant: To answer the question, I need know the location of the red sign so that I can determine the kind of animal on it. So I need the information about the following objects: red sign.

*User:* From the information on that advertising board, what is the type of this shop?

Assistant: To answer the question, I need know the location of the advertising board so that I can determine the type of the shop. So I need the information about the following objects: advertising board.

#### **In-Context Examples of HR-Bench**

User: Question: What is the color of the boy's bag? If you want to answer the question, which objects' information do you need?

Assistant: To answer the question, I need know the location of the boy with a bag so that I can determine the color of the bag. So I need the information about the following objects: boy with a bag.

*User: Is the yellow car on the left or right side of the white car?* 

Assistant: To answer the question, I need know the location of the yellow car and the white car so that I can determine the positional relationship between the two of them. So I need the information about the following objects: white car and yellow car.

*User: Tell me the number on the black board above the dog.* 

Assistant: To answer the question, I need know the location of the black board above the dog so that I can determine the number on it. So I need the information about the following objects: black board above the dog.

*User: Is the girl with pink hair on the left or right side of the man with backpack?* 

Assistant: To answer the question, I need know the location of the girl with pink hair and the man with backpack so that I can determine the positional relationship between the two of them. So I need the information about the following objects: girl with pink hair and man with backpack.

User: What kind of animal is on the red sign?

Assistant: To answer the question, I need know the location of the red sign so that I can determine the kind of animal on it. So I need the information about the following objects: red sign.

User: How many cars in the image?

Assistant: To answer the question, I need know the location of all cars so that I can determine the number of cars. So I need the information about the following objects: all cars.