

# VisionLeaf: Entropy-Guided Leaf-First Reasoning for Efficient and Accurate Think-with-Image

## Supplementary Material

### Contents

<b>A Training Data</b>	<b>13</b>
<b>B Training Settings</b>	<b>13</b>
<b>C Prompt</b>	<b>14</b>
<b>D Evaluation Benchmarks</b>	<b>14</b>
<b>E Discussion on Training Methods</b>	<b>14</b>
<b>F. More Comparison Results</b>	<b>14</b>
F.1. Stability Analysis . . . . .	14
F.2. Efficiency Analysis . . . . .	14
F.3. Experiments on MiMO-VL . . . . .	16
F.4. Additional Ablation Study . . . . .	16
<b>G Background and Theorem Proof</b>	<b>16</b>
G.1. Background on Conditional Expectation and Variance . . . . .	16
G.2. Theorem Proof . . . . .	17
<b>H More Case Studies</b>	<b>20</b>
H.1. Success Cases . . . . .	20
H.2. Failed Cases . . . . .	20
<b>I. Limitations and Future Work</b>	<b>20</b>

### A. Training Data

As we adopt the same training dataset as DeepEyes, we conduct the following analysis to better illustrate the distribution and characteristics of the training data. Our analysis focuses on two key factors: task type and input resolution. The composition of the dataset is summarized below:

- **Visual Search [45]:** The training dataset consists of the Visual Search Dataset, which is derived from COCO2017 [21]. This dataset is essential for developing robust natural-image understanding, as accurate responses require recognizing subtle visual cues and making precise object-level distinctions.
- **ArxivQA [19]:** This dataset comprises samples with complex visual semantics, including scientific plots, diagrams, and schematic charts. Its inclusion enables the model to interpret abstract and symbolic visual representations that extend beyond conventional natural-image scenarios.

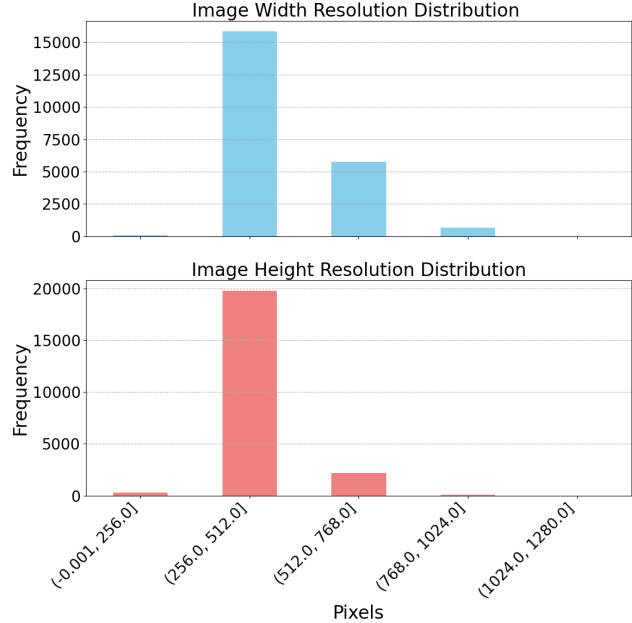


Figure 7. Distribution of image resolutions in the training dataset.

- **ThinkLite-VL [41]:** The training data also includes multimodal question–answering examples from ThinkLite-VL. This subset specifically targets tasks such as arithmetic reasoning, commonsense inference, and general problem-solving, which are critical for enhancing the model’s overall reasoning robustness and mitigating modality-specific overfitting.

Furthermore, we conducted an analysis of the training data’s resolution distribution, as illustrated in Fig. 7. The figure demonstrates that the majority of images are concentrated around resolutions such as  $512 \times 256$  pixels. This presents a significant distribution shift when compared to high-resolution benchmarks like the VStar dataset and HR-Bench, which commonly feature resolutions exceeding  $2K$ . The model’s strong performance across benchmarks with diverse input resolutions highlights the robust generalizability of our proposed methods.

### B. Training Settings

The model training was conducted on a cluster of  $8 \times \text{H800}$  GPUs, each equipped with 80 GB of memory. We employed a train batch size of 128. To accommodate long-context tasks, the maximum prompt length was set to 8192 tokens, and the maximum response length was con-

figured at 16384 tokens.

For the optimization core, we utilized the vanilla GRPO algorithm to calculate the advantage during training. Consistent with prior work such as DAPO [50] and DeepEyes [60], we omitted the KL divergence term from the objective function. To further enhance model performance, we applied the standard data augmentation settings commonly used in segmental training approaches.

To significantly accelerate the sampling and inference procedure, we leveraged `sglang` as the inference engine. We adopted an asynchronous sampling method and restricted the maximum number of interactive turns to 5. This limitation effectively mitigates potential overlong thinking cycles, thereby preventing unnecessary increases in overall training time.

### C. Prompt

The following visualization illustrates the prompt instantiated by our tokenizer’s chat template. It embeds the `image_zoom_in_tool` schema definition and places it within a multimodal interaction context, as shown in Fig. 8.

### D. Evaluation Benchmarks

**HR-Bench** [40]. HR-Bench is a benchmark designed to evaluate the perceptual capabilities of multimodal large language models on high-resolution images. It addresses the limitations of existing datasets by providing images at 4K and 8K resolutions, offered in two variants: HR-Bench 8K and HR-Bench 4K.

**VStar** [45]. VStar is a dataset developed to assess the visual search and detailed perception capabilities of MLLMs on high-resolution images. It is constructed from 191 high-resolution images with an average resolution of approximately  $2246 \times 1582$  pixels. The benchmark comprises two sub-tasks: an attribute recognition task with 115 samples and a spatial relationship reasoning task with 76 samples. Each sample is posed as a multiple-choice question designed to test a model’s ability to process high-resolution inputs, focus on fine-grained visual details, and perform guided visual search.

**MME-RealWorld** [54]. MME-RealWorld is a large-scale benchmark tailored for real-world applications, featuring 13,366 high-resolution images with an average resolution of approximately  $2000 \times 1500$  pixels. It contains 29,429 manually annotated question-answer pairs covering 43 sub-tasks across 5 real-world domains: OCR, diagrams and tables, remote sensing, autonomous driving, and monitoring. This dataset is designed to be highly challenging, requiring models to identify small or densely distributed objects and to perform robust high-resolution perception and reasoning

in complex real-world scenarios. In our experiments, we adopt the official *MME-RealWorld-Lite* split, which preserves the original 5 domains and 43 subtasks while uniformly downsampling to at most 50 QA pairs per subtask (or all samples if fewer than 50), yielding a lightweight yet representative subset that significantly reduces evaluation cost for high-resolution benchmarks.

### E. Discussion on Training Methods

With the growing emergence of visual agents that make decisions based on histories of visual observations, language, and actions, a variety of training paradigms—most notably supervised fine-tuning and reinforcement learning—have been adopted. Many prior works [16, 25] employ an SFT warm-up stage followed by RL training. To more clearly isolate the contribution of RL within the overall training pipeline, we conduct all experiments using only the RL stage. Recent studies [5, 14] further suggest that RL-trained models often exhibit stronger generalization capabilities than their SFT counterparts, particularly under out-of-distribution conditions.

### F. More Comparison Results

#### F.1. Stability Analysis

To more clearly compare the performance characteristics of DeepEyes and VisionLeaf, we perform a statistical analysis of their total number of *crashed outputs* across all interaction turns. We define a crashed output as any case in which the model exhibits severe repetitive function calling or produces a substantially malformed output format. This metric serves as a direct indicator of model stability under diverse scenarios. As shown in Fig. 9, VisionLeaf exhibits a markedly lower crash rate than DeepEyes—almost reaching zero in our evaluation—indicating that VisionLeaf delivers more stable and reliable behavior.

Table 3. **Efficiency Analysis.** We report the training time and the inference time on the same settings.

Method	Training Time	Inference	
		Time	Tokens
DeepEyes	40.8h	22.5min	246
VisionLeaf	36.4h	12.2min	139

#### F.2. Efficiency Analysis

Furthermore, we assess the efficiency of VisionLeaf by comparing total inference time across diverse benchmarks. As depicted in Fig. 10, the proposed method achieves an approximate  $3 \times$  reduction in inference time under the full benchmark setting. To ensure a controlled and fair evaluation, all experimental configurations are held constant,

## System Prompt

```
1 <|im_start|>system
2 You are a helpful assistant.
3
4 # Tools
5
6 You may call one or more functions to assist with the user query.
7
8 You are provided with function signatures within <tools></tools> XML tags:
9 <tools>
10 {
11   "type": "function",
12   "function": {
13     "name": "image_zoom_in_tool",
14     "description": "Zoom in on a specific region of an image by cropping it based on a
15       bounding box (bbox) and an optional object label.",
16     "parameters": {
17       "type": "object",
18       "properties": {
19         "bbox_2d": {
20           "type": "array",
21           "items": { "type": "number" },
22           "minItems": 4,
23           "maxItems": 4,
24           "description": "The bounding box of the region to zoom in, as [x1, y1, x2, y2
25             ], where (x1, y1) is the top-left corner and (x2, y2) is the bottom-right
26             corner."
27         },
28         "label": {
29           "type": "string",
30           "description": "The name or label of the object in the specified bounding box
31             (optional)."
```

Figure 8. The prompt used during the training procedure.

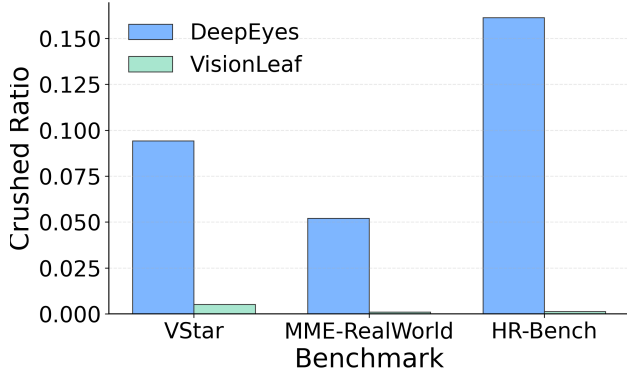


Figure 9. Comparison of the crushed ratio of VisionLeaf and DeepEyes.

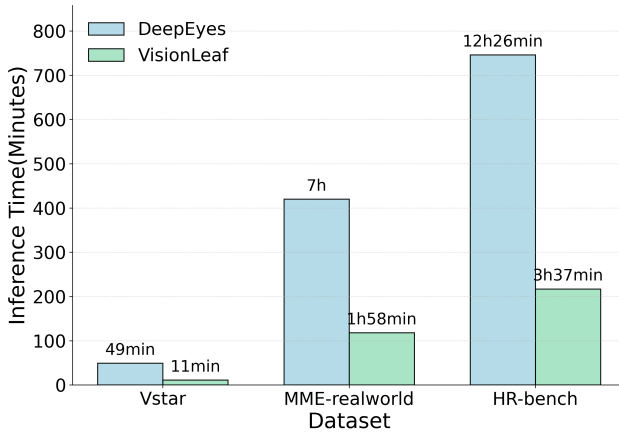


Figure 10. Comparison of the inference time of VisionLeaf and DeepEyes across multiple benchmarks.

and inference time is measured using a consistent evaluation framework. In addition, Table. 3 reports training and inference times on 8 GPUs, as well as the number of generated tokens. Collectively, these findings highlight that VisionLeaf not only improves task performance but also substantially reduces computational overhead in both training and inference phases.

### F.3. Experiments on MiMO-VL

To further assess the generalization performance of our method, we conduct experiments on MiMO-VL, a representative model gaining increasing traction in the VLM community. As shown in Fig. 11, our approach achieves a 1.4% improvement over the baseline method, DeepEyes, providing clear evidence of its strong generalization capability across architectures.

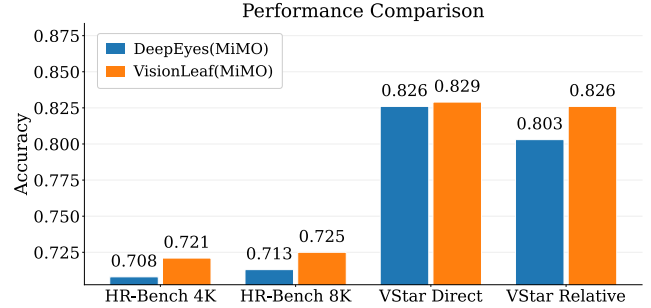


Figure 11. Performance Comparison on MiMO-VL. To further validate the generality of our method across different model architectures, we conduct additional experiments on MiMO-VL.

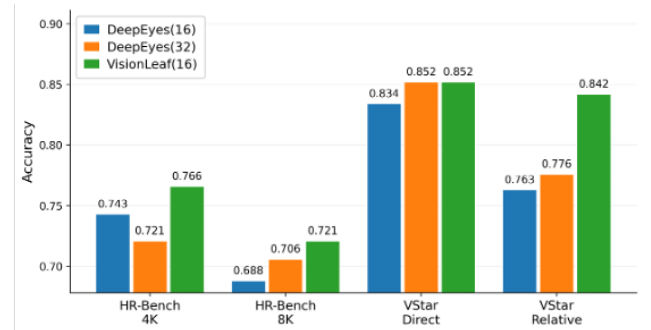


Figure 12. Larger Rollout Comparison. To provide a more thorough comparison, we evaluate VisionLeaf against the baseline using a larger number of rollouts.

### F.4. Additional Ablation Study

To further investigate the efficiency of our approach, we conduct a comparative study against the baseline using an increased GRPO rollout budget. Specifically, we evaluate DeepEyes with 32 rollouts to understand the effect of scaling root-level rollouts. As shown in Fig. 12, the performance improvement is marginal, while the associated training cost rises considerably. These findings indicate that increasing the GRPO rollout number alone yields diminishing returns and is substantially less efficient than our method.

## G. Background and Theorem Proof

### G.1. Background on Conditional Expectation and Variance

Let  $X$  and  $Y$  be square-integrable random variables on a common probability space. The conditional expectation  $\mathbb{E}[X | Y]$  is the (a.s. unique) random variable such that for every bounded measurable function  $\phi$ ,

$$\mathbb{E}[X \phi(Y)] = \mathbb{E}[\mathbb{E}[X | Y] \phi(Y)]. \quad (12)$$

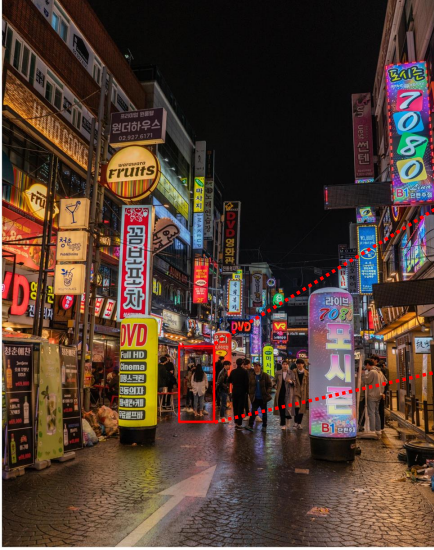
The conditional variance of  $X$  given  $Y$  is defined as

$$\text{Var}[X | Y] = \mathbb{E}[(X - \mathbb{E}[X | Y])^2 | Y]. \quad (13)$$



What color hat does the woman wearing the mask wear?

Original Image



VisionLeaf

**Think Step 1:**

I need to zoom on the woman wearing the mask to accurately determine the color of her hat.

**Tool Calling 1:**



**Answer:**

The woman wearing the mask wears a **black** hat. ✓

**Think Step 2:**

Now I can see that the woman wearing the mask is wearing a **black** hat. Therefore, I can provide the answer.

Figure 13. Fine-grained object perception case illustrating detailed visual reasoning. The correct match is highlighted in **green**, demonstrating the model’s strong ability to perform fine-grained object perception.

We will repeatedly use the *law of total expectation*

$$\mathbb{E}[X] = \mathbb{E}[\mathbb{E}[X | Y]], \quad (14)$$

and the *law of total variance*

$$\text{Var}[X] = \mathbb{E}[\text{Var}(X | Y)] + \text{Var}(\mathbb{E}[X | Y]). \quad (15)$$

All such relationships have analogous versions when we additionally condition on another random variable  $Z$ , e.g.

$$\text{Var}[X | Z] = \mathbb{E}[\text{Var}(X | Y, Z) | Z] + \text{Var}(\mathbb{E}[X | Y, Z] | Z). \quad (16)$$

## G.2. Theorem Proof

### Variance Analysis of Early Stage Response

We first recall a standard conditional variance identity that will be used throughout.

**Lemma G.1** (Conditional law of total variance). *Let  $X$  be a square-integrable random variable and let  $(Y, Z)$  be a pair of random variables. Then*

$$\text{Var}[X | Z] = \mathbb{E}[\text{Var}(X | Y, Z) | Z] + \text{Var}(\mathbb{E}[X | Y, Z] | Z). \quad (17)$$

This is the standard law of total variance applied conditionally on  $Z$ .

**Lemma G.2** (Recursive conditional variance along a trajectory). *Let  $(s_t)_{t=1}^{T+1}$  be a sequence of context random variables in a rollout, and let  $R$  be a square-integrable terminal random variable. Assume that, for each  $t$ ,*

$$\mathbb{P}(R \in \cdot | s_t, s_{t+1}) = \mathbb{P}(R \in \cdot | s_{t+1}), \quad (18)$$

*i.e.,  $s_{t+1}$  is sufficient for the future return. Then*

$$\text{Var}[R | s_t] = \mathbb{E}[\text{Var}(R | s_{t+1}) | s_t] + \text{Var}(\mathbb{E}[R | s_{t+1}] | s_t). \quad (19)$$

*Proof.* Apply Lemma G.1 with  $X = R$ ,  $Y = s_{t+1}$  and  $Z = s_t$ :

$$\begin{aligned} \text{Var}[R | s_t] &= \mathbb{E}[\text{Var}(R | s_t, s_{t+1}) | s_t] \\ &\quad + \text{Var}(\mathbb{E}[R | s_t, s_{t+1}] | s_t). \end{aligned} \quad (20)$$

The sufficiency assumption implies that, for each  $t$ , the conditional distribution of  $R$  given  $(s_t, s_{t+1})$  coincides with that of  $R$  given  $s_{t+1}$ , hence

$$\begin{aligned} \text{Var}(R | s_t, s_{t+1}) &= \text{Var}(R | s_{t+1}), \\ \mathbb{E}[R | s_t, s_{t+1}] &= \mathbb{E}[R | s_{t+1}], \end{aligned} \quad (21)$$

and substituting these identities yields (19).  $\square$

What is the text before “Vision”?



Original Image



**Think Step 1:**  
I need to zoom on the text before “Vision” to accurately identify it.

**Tool Calling 1:**



**Answer:**

The text before “Vision” is “World”. ✓

**Think Step 2:**

Now I can see that the text before “Vision” is “World”.

Figure 14. Fine-grained example of word-level OCR reasoning. The correct match is highlighted in green, demonstrating the model’s strong ability to perform fine-grained word-level OCR reasoning.

### Analysis on Prefix Sharing Sampling

We now derive the variance identities for group-centered returns and for tree sampling.

**Lemma G.3** (Variance of group-centered returns). *Let  $R_1, \dots, R_G$  be conditionally independent and identically distributed copies of a square-integrable random variable  $R$  given  $s$ , with*

$$\mathbb{E}[R | s] = \mu_s, \quad \text{Var}[R | s] = \sigma_s^2 < \infty. \quad (22)$$

Define the group mean

$$\bar{R} := \frac{1}{G} \sum_{i=1}^G R_i, \quad (23)$$

and the centered quantities

$$\hat{A}_i := R_i - \bar{R}, \quad i = 1, \dots, G. \quad (24)$$

Then, for each  $i$ ,

$$\text{Var}[\hat{A}_i | s] = \left(1 - \frac{1}{G}\right) \text{Var}[R | s]. \quad (25)$$

*Proof.* Write  $\sigma_s^2 = \text{Var}[R | s]$ . Since the  $R_i$  are conditionally i.i.d. given  $s$ ,

$$\text{Var}[R_i | s] = \sigma_s^2, \quad \text{Cov}[R_i, R_j | s] = 0 \quad (i \neq j). \quad (26)$$

We have

$$\begin{aligned} \hat{A}_i &= R_i - \bar{R} \\ &= R_i - \frac{1}{G} \sum_{j=1}^G R_j \\ &= \left(1 - \frac{1}{G}\right) R_i - \frac{1}{G} \sum_{j \neq i} R_j. \end{aligned} \quad (27)$$

Thus

$$\begin{aligned} \text{Var}[\hat{A}_i | s] &= \text{Var}\left[\left(1 - \frac{1}{G}\right) R_i - \frac{1}{G} \sum_{j \neq i} R_j \mid s\right] \\ &= \left(1 - \frac{1}{G}\right)^2 \text{Var}[R_i | s] + \text{Var}\left[\frac{1}{G} \sum_{j \neq i} R_j \mid s\right] \\ &\quad - 2\left(1 - \frac{1}{G}\right) \text{Cov}\left(R_i, \frac{1}{G} \sum_{j \neq i} R_j \mid s\right). \end{aligned} \quad (28)$$

The covariance term vanishes because  $R_i$  is conditionally independent of each  $R_j$ ,  $j \neq i$ , given  $s$ . For the second variance term,

$$\begin{aligned} \text{Var}\left[\frac{1}{G} \sum_{j \neq i} R_j \mid s\right] &= \frac{1}{G^2} \sum_{j \neq i} \text{Var}[R_j | s] \\ &= \frac{G-1}{G^2} \sigma_s^2. \end{aligned} \quad (29)$$



On which hand does the woman wearing sunglasses wear a watch?

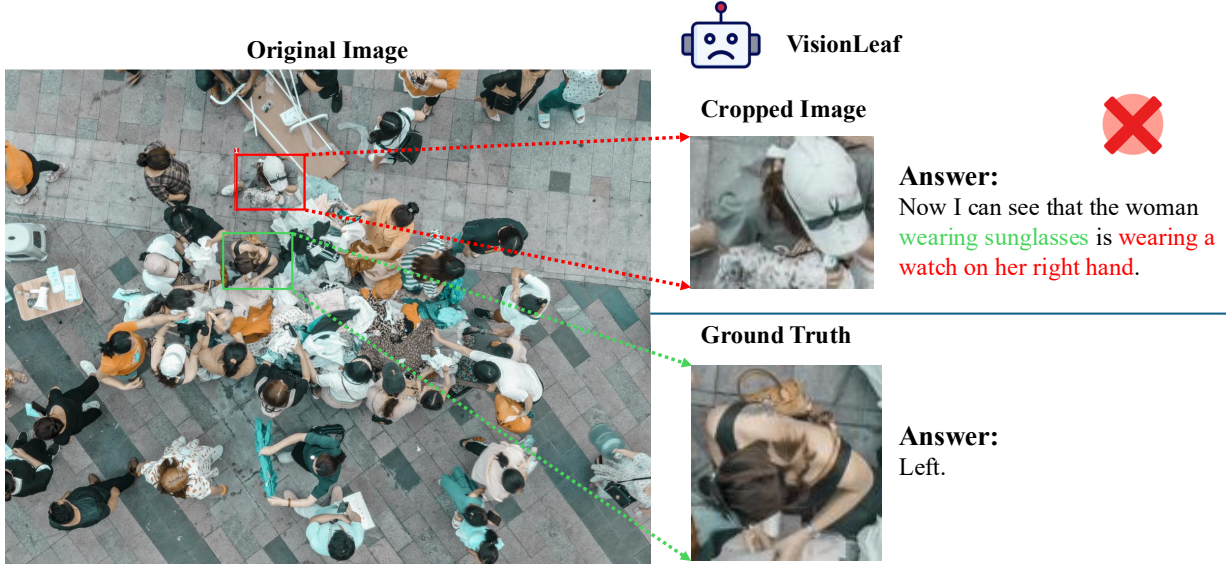


Figure 15. Failure case in which the reasoning procedure fails to capture all required details. We highlight the correct match in green and the mismatched prediction in red, illustrating that the model still encounters difficulties in complex reasoning scenarios.

Therefore

$$\begin{aligned}
 \text{Var}[\hat{A}_i | s] &= \left(1 - \frac{1}{G}\right)^2 \sigma_s^2 + \frac{G-1}{G^2} \sigma_s^2 \\
 &= \left(1 - \frac{2}{G} + \frac{1}{G^2} + \frac{G-1}{G^2}\right) \sigma_s^2 \\
 &= \left(1 - \frac{2}{G} + \frac{1}{G}\right) \sigma_s^2 \\
 &= \left(1 - \frac{1}{G}\right) \sigma_s^2 \\
 &= \left(1 - \frac{1}{G}\right) \text{Var}[R | s].
 \end{aligned} \tag{30}$$

□

**Lemma G.4** (Node-level group variance in a tree). *Let  $n$  be a (random) node extending  $s$ . Suppose that, conditional on  $n$ ,*

$$R_1^{(n)}, \dots, R_G^{(n)} \tag{31}$$

*are i.i.d. copies of a square-integrable random variable  $R^{(n)}$  with finite variance. Define*

$$\bar{R}^{(n)} := \frac{1}{G} \sum_{j=1}^G R_j^{(n)}, \quad \hat{A}_i^{(n)} := R_i^{(n)} - \bar{R}^{(n)}. \tag{32}$$

*Then, for each  $i$ ,*

$$\text{Var}[\hat{A}_i^{(n)} | n] = \left(1 - \frac{1}{G}\right) \text{Var}[R^{(n)} | n]. \tag{33}$$

*Proof.* Fix  $n$  and apply Lemma G.3 with  $s$  replaced by  $n$  and  $R$  replaced by  $R^{(n)}$ . □

To relate  $\text{Var}[R^{(n)} | n]$  to  $\text{Var}[R | s]$ , we use the conditional law of total variance once more.

**Lemma G.5** (Average comparison between root and node variances). *Let  $N$  be a random node extending  $s$ , and suppose that the conditional distribution of  $R$  given  $(s, N)$  coincides with that of  $R^{(N)}$  given  $N$ . Then*

$$\mathbb{E}[\text{Var}[R^{(N)} | N] | s] \leq \text{Var}[R | s]. \tag{34}$$

*Proof.* Apply Lemma G.1 with  $X = R$ ,  $Y = N$  and  $Z = s$ :

$$\text{Var}[R | s] = \mathbb{E}[\text{Var}(R | s, N) | s] + \text{Var}(\mathbb{E}[R | s, N] | s). \tag{35}$$

By the assumption on the conditional distributions,

$$\text{Var}(R | s, N) = \text{Var}[R^{(N)} | N], \tag{36}$$

so

$$\begin{aligned}
 \text{Var}[R | s] &= \mathbb{E}[\text{Var}[R^{(N)} | N] | s] + \text{Var}(\mathbb{E}[R^{(N)} | N] | s) \\
 &\geq \mathbb{E}[\text{Var}[R^{(N)} | N] | s],
 \end{aligned} \tag{37}$$

since the last term is a conditional variance and hence non-negative. □

Combining Lemma G.4 with Lemma G.5, we obtain that, on average over nodes  $N$  extending  $s$ ,

$$\begin{aligned} \mathbb{E}[\text{Var}[\widehat{A}_i^{(N)} \mid N \mid s]] &= \left(1 - \frac{1}{G}\right) \mathbb{E}[\text{Var}[R^{(N)} \mid N \mid s]] \\ &\leq \left(1 - \frac{1}{G}\right) \text{Var}[R \mid s]. \end{aligned} \tag{38}$$

This shows how the variance formulas for group-centered returns and their shared-prefix refinements follow from basic properties of conditional expectation and variance.

## H. More Case Studies

### H.1. Success Cases

To further illustrate the efficiency of VisionLeaf, we provide additional qualitative examples beyond those presented in the main paper. These cases are shown in Fig. 13 and Fig. 14. From these examples, we observe that the function-call mechanism plays a crucial role: the initial reasoning step determines which region of the image requires closer examination, and the subsequent zoom-in operation enables more focused, image-grounded analysis. This leads to a coherent “think-with-image” process that ultimately produces the correct answer. Moreover, this paradigm demonstrates strong performance in both OCR-based reasoning and general visual reasoning tasks.

### H.2. Failed Cases

Despite the effectiveness of VisionLeaf, certain limitations remain when handling complex reasoning tasks. As illustrated in Fig. 15, the model focuses on the woman wearing sunglasses while overlooking an additional constraint—that the target individual must also be wearing a watch. This leads to an incorrect localization. For comparison, we show the correct bounding box corresponding to the question alongside the model’s predicted crop region.

More specifically, this case is difficult even for humans: the black sunglasses blend closely with the subject’s dark hair, resulting in low-contrast visual cues that are challenging to distinguish. Given the limited capacity of the 7B model used in our experiments, such fine-grained visual reasoning remains particularly difficult, highlighting the constraints imposed by model size.

## I. Limitations and Future Work

**Limitations.** While the proposed reinforcement-learning pipeline effectively improves image-grounded reasoning, several limitations remain. All experiments in this work are conducted using a 7B model, whose relatively small capacity constrains both fine-grained visual perception and higher-level reasoning. Consequently, the model may underperform on tasks requiring accurate bounding-box construction or subtle visual discrimination, as discussed in

Sec. H.2. We attribute these challenges to the restricted model size and its associated representational limitations.

**Future Work.** To mitigate these issues, future work will explore scaling the approach to larger backbone models with stronger visual and linguistic reasoning capabilities. Increasing model capacity is expected to significantly alleviate difficulties in detecting low-contrast regions, localizing subtle visual attributes, and maintaining stable reasoning under complex multimodal constraints. More broadly, integrating our method with more capable foundation models may unlock substantially improved performance in multi-turn, tool-augmented visual reasoning tasks.