Look Again, Think Slowly: Enhancing Visual Reflection in Vision-Language Models

Pu Jian^{1,2}, Junhong Wu^{1,2}, Wei Sun ^{1,2}, Chen Wang ^{1,2}, Shuo Ren¹, Jiajun Zhang^{1,2,3}*

¹Institute of Automation, Chinese Academy of Sciences

²School of Artificial Intelligence, University of Chinese Academy of Sciences
³Wuhan AI Research

{jianpu2023, wujunhong2021, sunwei2023, chenwang2020, shuo.ren}@ia.ac.cn jjzhang@nlpr.ia.ac.cn

Abstract

Recent advances in text-only "slow-thinking" reasoning have prompted efforts to transfer this capability to vision-language models (VLMs), for training visual reasoning models (VRMs). However, such transfer faces critical challenges: Effective "slow thinking" in VRMs requires vi**sual reflection**, the ability to check the reasoning process based on visual information. Through quantitative analysis, we observe that current VRMs exhibit limited visual reflection, as their attention to visual information diminishes rapidly with longer generated responses. To address this challenge, we propose a new VRM Reflection-V¹, which enhances visual reflection based on reasoning data construction for cold-start and reward design for reinforcement learning (RL). Firstly, we construct vision-centered reasoning data by leveraging an agent that interacts between VLMs and reasoning LLMs, enabling cold-start learning of visual reflection patterns. Secondly, a visual attention based reward model is employed during RL to encourage reasoning based on visual information. Therefore, Reflection-V demonstrates significant improvements across multiple visual reasoning benchmarks. Furthermore, Reflection-V maintains a stronger and more consistent reliance on visual information during visual reasoning, indicating effective enhancement in visual reflection capabilities.

1 Introduction

Recently, "slow-thinking" reasoning has emerged as a significant advancement in large language models (LLM) domain (Shao et al., 2024), demonstrating remarkable capabilities in solving complex reasoning tasks, such as OpenAI-o1 (Jaech et al., 2024) and DeepSeek-R1 (Guo et al., 2025a). The superior performance of "slow thinking" LLM primarily benefits from its ability to perform "reflec-

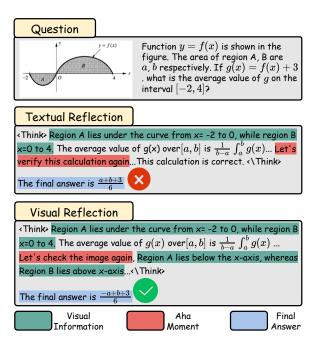


Figure 1: Existing "slow-thinking" VLMs claimed "aha moment" is often merely a textual reflection. We instead highlight visual reflection, where VLM actively verifies and refines its reasoning based on visual inputs.

tion" during reasoning (Yan et al., 2024). This reflection mechanism allows models to check and revise intermediate steps before generating the final answer, thereby avoiding errors that may arise from shortcut inference (Snell et al., 2024; Yang et al., 2025d; Cheng et al., 2024), which is also called "aha moment". Inspired by this success, some researchers attempt to integrate "slow thinking" into vision-language models (VLMs), enabling the trained visual reasoning models (VRMs) to generate more accurate and deliberate solutions (Wang et al., 2025; Chen et al., 2025; Tan et al., 2025; Huang et al., 2025). Specifically, they leverage "slow thinking" LLMs to reason based on image descriptions generated by VLMs, thereby introducing reflection patterns into reasoning data. These visual reasoning data are often used for supervised

^{*} Corresponding Author

¹The related codes are released in this URL: https://github.com/jian0805/ReflectionV

fine-tuning (SFT) (Thawakar et al., 2025; Xu et al., 2024), providing a cold-start initialization for subsequent reinforcement learning (RL) (Huang et al., 2025; Tan et al., 2025; Yang et al., 2025c).

In this paper, we propose that the true "aha moment" in visual reasoning arises when a model engages in **visual reflection**—that is, when it actively verifies and refines its reasoning based on the visual input, as shown in Figure 1. However, current distillation-based approaches to training VRMs often miss this crucial aspect. By transferring superficial reflective behaviors from LLMs trained solely on text, these methods encourage reasoning patterns that are detached from the visual modality. This is because the cold-start data for these VRMs still originates from text-only reasoning on visual descriptions, and the RL stage uses rewards based solely on textual outputs (Huang et al., 2025; Meng et al., 2025). Thus, instead of promoting visual insight, these VRMs risk reinforcing textual hallucinations and visual neglect (Zhong et al., 2024; Favero et al., 2024). As a result, VRMs may appear reflective while actually bypassing the visual content, undermining both the reliability and robustness of their reasoning.

Consistent with the previous discussion, we conducted a detailed analysis of existing VRMs and found that they struggle with visual reflection. Specifically, experiments observe that existing VRMs' attention to and reliance on visual information decline rapidly as the number of generated tokens increases. And VRMs trained by distilling text-only reflection data even exhibit lower reliance on visual prompts than their backbone VLMs. This indicates that existing VRMs struggle to attend to and leverage visual information during reflection, thereby degrading into text-only reflection models.

To address the aforementioned challenges, we propose a novel two-stage training strategy for training VRMs. In the cold-start stage, we focus on resolving the limitations of image description-based approaches in incorporating visual reflection patterns within training data. Specifically, we leverage a multi-modal agent, where LLMs interact with VLMs, to complete reasoning in an LLM-VLM interleaved way. This data construction paradigm ensures that visual information can be continuously accessed and repeatedly utilized during reasoning, thereby introducing a visual reflection pattern for VRMs to learn. In the RL stage, to further promote the visual reflection behavior learned from cold-start data, we introduce a visual attention based re-

ward for group relative policy optimization (GRPO) (Shao et al., 2024). This reward encourages VRMs to consistently attend to visual information.

Reflection-V, our VRM trained with the proposed strategy, achieves significant improvements on benchmarks focusing on mathematical (Lu et al., 2023; Wang et al., 2024a), multi-disciplinary (Yue et al., 2024a,b), and general reasoning (Chen et al., 2024a). At the 7B scale, it is comparable to or even surpasses several widely used very large VLMs, like GPT-40 (Hurst et al., 2024) and InternVL2.5-38B (Chen et al., 2024b). Notably, the aforementioned quantitative analyses and case study further show that compared to the base model, Reflection-V maintains more sustained attention to visual information and actively engages in visual reflection, representing the emergence of the true "aha moment" in visual reasoning.

2 VRMs Struggle with Visual Reflection

In this section, we claim that existing VRMs struggle to perform visual reflection. To support this claim, we analyze the visual tokens' role during reasoning. Specifically, we quantify visual tokens' effect using the following metrics: attention weight and a visual dependency measure. This observation later motivates our proposed methodology.

2.1 Visual Attention Weight

To capture how the contribution of visual tokens varies during reasoning, we track the attention weights from response tokens to visual tokens as more tokens are generated. Let $T_{\rm res}$ and $T_{\rm vis}$ denote the sets of response and visual tokens, respectively. For the h-th layer, let $a_{nj}^{(h)}$ represent the attention weight from the n-th response token to the j-th visual token. Thus the total attention from the n-th response token to $T_{\rm vis}$ is given by

$$Attn(n, T_{vis}) = \frac{\sum_{h} \sum_{j \in T_{vis}} a_{nj}^{(h)}}{\sum_{h} \sum_{j \in T_{vis}} \mathbf{1}_{a_{nj}^{(h)} > 0}}.$$
 (1)

2.2 Visual Dependency Measure

Furthermore, after generating several tokens, we drop the visual tokens and assess VRMs' reliance on the visual token during reasoning by measuring the divergence in subsequent generations, which is quantified based on the divergence between the next-token prediction distributions with and without visual tokens. We use the Hellinger distance

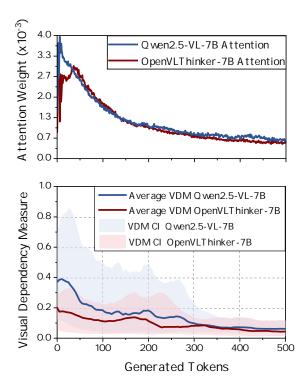


Figure 2: Attention weights on visual tokens and the visual dependency measure during reasoning on the MMMU dataset. Both metrics decline sharply as more tokens are generated, and RL-enhanced models (e.g., OpenVLThinker-7B) do not mitigate this decay.

(Favero et al., 2024), defined as

$$H_{\text{dist}}(p,q) = 2^{-\frac{1}{2}} \sqrt{\sum_{i=1}^{k} (\sqrt{p_i} - \sqrt{q_i})^2}$$
 (2)

to quantify the divergence between two probability distributions. Therefore, when the number of generated response tokens is n, the visual dependency measure $\mathrm{VDM}(n|T_{\mathrm{vis}},T_q)$ for a given imagequestion pair (T_{vis},T_q) is given by

$$H_{\text{dist}}\left(p(\cdot|T_{res}^{< n}, T_{q}, T_{\text{vis}}), p(\cdot|T_{res}^{< n}, T_{q})\right). \tag{3}$$

2.3 Diminishing Visual Attention and Dependency

As Figure 2 illustrates, both the mean Visual Dependency Measure on MMMU (Yue et al., 2024a) and the layer-wise attention from response tokens to visual tokens in VRM, decline sharply as generation proceeds: after roughly 300 tokens, visual attention falls to only 20–30 % of its initial level. This analysis reveals that VRMs typically lack visual reflection ability in long-chain reasoning, as they rarely refer back to visual tokens when performing reflective checking of reasoning process.

Although widely used RL boosts VRMs' reasoning performance, it fails to equip them with visual reflection capability. Instead, it further reinforces over-reliance on previously generated text. As Figure 2 shows, OpenVLThinker, RL based on Qwen2.5-VL, exhibits even lower focus and dependence on visual tokens during reasoning than the base model.

3 Method

In the above analysis, we observe that as the length of reasoning process increases, VRMs rapidly reduce their reliance on and attention to visual information. This limitation hampers their ability to perform visual reflection and prevents them from fully benefiting from "slow thinking" paradigms like DeepSeek-R1 (Guo et al., 2025a). To address this problem, we propose a two-stage strategy for training VRMs. This strategy consists of: (1) cold-start initialization (Yang et al., 2025e) on reasoning data with visual reflection, and (2) reinforcement learning with a visual attention-based reward. In the first stage, we leverage a multimodal agent, where LLMs interact with VLMs, to construct visual reasoning data exhibiting visual reflection, and use it to perform supervised finetuning (SFT) on the base VLM. The second stage applies GRPO with the proposed reward function that explicitly encourages sustained attention to visual tokens. The related codes can be found at https://github.com/jian0805/ClearVQA

3.1 Reasoning Data with Visual Reflection Construction

Existing visual reasoning studies typically rely on LLMs to perform reasoning directly based on image captions, thereby constructing visual reasoning data (Liu et al., 2024; Yu et al., 2024; Huang et al., 2025). The absence of visual modality during reasoning makes it impossible for these reasoning data to exhibit the visual reflection we claimed. Inspired by recent advances in visual-language agents (Chen et al., 2023; Zhou et al., 2024), we employ crafted prompts to guide the interaction between the LLM and VLM in completing the reasoning task. This data construction paradigm ensures that visual information is continuously accessible and repeatedly utilized during reasoning, thereby introducing a visual reflection pattern. The data construction process is detailed in the following sections.

Visual Reflection Data Construction Based on

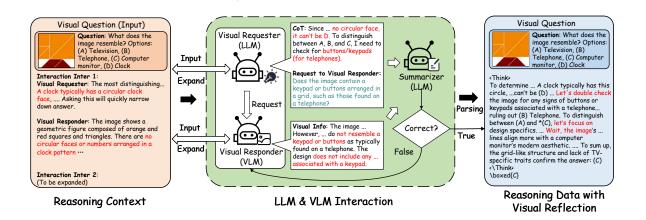


Figure 3: A workflow for constructing reasoning data with a visual reflection pattern. LLMs and VLMs perform reasoning through interaction, ensuring that visual information is continuously acquired and repeatedly utilized, thereby creating a visual reflection pattern in the reasoning process.

LLM-VLM Interaction. As shown in Figure 3, in our data construction process, the LLM and VLM interact by taking on the following three roles: 1) **Visual requester**, played by the LLM, which determines what visual information is needed to answer the question based on the existing reasoning context and proposes a request to the VLM. During this process, the CoT output from the VLM is added as a partial solution to the reasoning context. 2) Vi**sual responder**, played by the VLM, replies to the request made by the visual requester, revealing visual information related to the visual question. The visual description generated by the visual responder is added to the reasoning context. 3) Summarizer, played by the LLM, summarizes the existing reasoning context after each round of interaction, generating the final answer. If the generated answer does not match the ground truth, all outputs from the summarizer are discarded, and a new round of interaction begins. Otherwise, the output is added to the reasoning context, completing the reasoning process for the visual question.

Post-generation processing. After constructing visual reflection data through LLM and VLM interactions, we perform the following steps to ensure data quality: 1) Non-Reflection Filtering. We filter out the data where the summarizer produces the correct answer just after the first interaction. These samples lack sustained reliance on visual information, and the reasoning solution does not exhibit a visual reflection pattern. 2) Cohesion Enhancement. In the previous data construction process, the text generated across different VLM-LLM interaction rounds may lack coherence. We

employ the LLM to process and refine the text into a cohesive reasoning process. All detailed prompts are provided in Appendix B.

3.2 Visual Attention Based Reward

Following existing works (Wei et al., 2025; Xiao et al., 2025; Zhang et al., 2025a), we adopt GRPO, a rule-based reinforcement learning algorithm, to enhance the reasoning capabilities of VRMs. Building on the original reward function, we introduce a visual attention—based reward to encourage the model to maintain sustained attention to visual tokens throughout the reasoning process while preserving overall performance.

Specifically, based on the analyses in Section 2.1 that attention weights on visual tokens rapidly decrease as the number of generated tokens increases, our reward rule follows the principle: For a VRM reasoning process, VRMs receive a higher reward if relatively higher attention to visual tokens is maintained after generating several tokens. Therefore, the visual attention based reward is given by

$$r_{v} = \begin{cases} \frac{\sum\limits_{n>|T_{res}|/2}^{\sum} \text{Attn}(n,T_{\text{vis}})}{\sum\limits_{n<|T_{res}|/2}^{\sum} \text{Attn}(n,T_{\text{vis}})} & \text{if } r_{a} = 1\\ 0 & \text{if } r_{a} = 0 \end{cases}$$
 (4)

Here, r_a is the accuracy reward taking values from $\{0, 1\}$. Refer to function (2), $\operatorname{Attn}(n, T_{\text{vis}})$ represents the average attention weight of the n-th response token to the visual tokens (averaged over all attention heads). $|T_{res}|$ denotes the total number of tokens in a VRM's reasoning process.

Model	Math-Reasoning		Multi-Disciplinary		General	Hallucination	
	MathVision	MathVista	MMMU	MMMU-Pro	M3CoT	HallBench	
Cl	osed-Source \	Vision-Lang	uage Mode	els			
GPT-40 (Hurst et al., 2024)	30.4	60.0	69.1	51.9	†74.2	-	
GPT-4V (Yang et al., 2023)	22.8	49.9	56.8	41.2	62.6	65.3	
0	pen-Source V	ision-Langu	age Mode	ls			
QwenVL2.5-3B (Bai et al., 2025)	21.2	62.3	[†] 51.2	31.6	[†] 55.6	45.1	
QwenVL2.5-7B (Bai et al., 2025)	25.1	68.2	[†] 54.3	36.9	†60.5	49.5	
InternVL2.5-8B (Chen et al., 2024b)	19.7	63.6	56.0	30.5	[†] 41.5	49.0	
InternVL-2.5-38B (Chen et al., 2024b)	32.2	71.9	<u>57.6</u>	46.0	† <u>68.9</u>	-	
LLaVA-OneVision-72B (Li et al., 2024a)	30.1	67.5	56.8	31.0	[†] 61.5	47.9	
Kimi-VL-16B (Team et al., 2025)	21.4	68.7	55.7	-	-	-	
Open-S	ource Vision-	Language F	Reasoning	Models			
TVC-7B (Sun et al., 2025a)	22.7	68.1	-	-	-	-	
R1-VL-7B (Zhang et al., 2025a)	24.7	63.5	44.5	-	-	-	
MM-Eureka-7B ^{\dightarrow} (Meng et al., 2025)	26.9	<u>73.0</u>	[†] 51.3	†36.7	[†] 63.5	47.8	
R1-Onevision-7B ^{\dightarrow} (Yang et al., 2025e)	29.9	64.1	†48.7	21.6	[†] 53.1	41.7	
OpenVLThinker-7B ^{\dightarrow} (Deng et al., 2025)	25.3	70.2	52.5	37.3	†62.2	42.3	
Ours (Ti	raining strateg	y emphasizir	ıg visual re	flection)			
Reflection-V-3B	27.9	66.3	56.9	38.2	62.9	49.3	
Reflection-V-7B	33.9	73.3	61.3	<u>42.7</u>	71.1	<u>53.9</u>	

Table 1: Performance of Reflection-V across various visual reasoning benchmarks, compared to existing VLMs. † indicates results reproduced by us. \$\display\$ denotes vision-language reasoning Models also based on the Qwen2.5-7B series. Bold and underlined scores represent the best and second-best performance among open-source models for each benchmark.

Based on the observation shown in Figure 2, we calculate the visual attention based reward using the last layer where the attention to visual tokens is most significant. The overall reward r_o in GRPO is the weighted sum of the accuracy reward r_a , visual attention-based reward r_v , and format reward r_f (Shao et al., 2024), given by

$$r_o = r_a + \lambda_v r_v + \lambda_f r_f. \tag{5}$$

 λ_v and λ_f are scaling coefficients set to 0.5 and 0.1, respectively.

4 Experiments

4.1 Experimental Setup

Implementations. To construct the cold-start data, we use the open-source VLM Qwen-2.5-VL-72B (Bai et al., 2025) and a reasoning-capable LLM, QWQ-32B (Team, 2025), to interactively generate data. Our method is evaluated using the Qwen-2.5-VL-7B-Instruct as the base model. During the cold-start stage, we train for 3 epochs on 2 NVIDIA H100 GPUs. The model, after cold-start initialization, is subsequently trained using GRPO

with visual attention based reward for 12 epochs on 8 NVIDIA H100 GPUs, based on the Verl training framework (Sheng et al., 2024; Zheng et al., 2025b). For GRPO, 16K reasoning samples are collected from a diverse multimodal corpus. The detailed composition of training data is shown in Appendix C. Train details for cold-start initialization and GRPO stages is provided in Appendix A.

Benchmarks for Evaluation. We conduct a comprehensive experimental analysis to assess how our method improves visual reasoning. To ensure a well-rounded evaluation, we select widely recognized visual reasoning benchmarks that emphasize math, multi-disciplinary, and general reasoning skills. For evaluating math reasoning, we use Math-Vista (Lu et al., 2023) and MathVision (Wang et al., 2024a), which are standard tests for visual reasoning models. To evaluate performance across multiple disciplinary such as physics, chemistry, and computer science, we adopt MMMU and MMMU-Pro (Yue et al., 2024a,b). M3CoT (Chen et al., 2024a) is used to assess general reasoning ability, as it covers a broad range of knowledge-intensive and commonsense-based reasoning questions. Fur-

Model	MathVision	MathVista	MMMU	MMMU-Pro	МЗСоТ
Reflection-V-3B	27.94	66.30	56.89	38.17	62.95
w/o VAR	26.52	65.60	55.80	36.56	61.79
w/o Cold-Start	24.27	64.20	53.98	34.75	59.55
w/o Cold-Start + VAR	23.60	63.90	53.21	33.97	58.81
Reflection-V-7B	33.71	73.30	61.33	42.71	71.07
w/o VAR	32.47	72.40	60.10	41.95	69.28
w/o Cold-Start	29.01	70.40	58.81	39.06	65.87
w/o Cold-Start + VAR	28.53	69.80	58.03	38.24	64.63

Table 2: Ablation results for cold-start based on visual reflection data and visual attention based reward on performance improvement. VAR denotes visual attention based reward.

Model	MathVision	MathVista	MMMU	MMMU-Pro	МЗСоТ
Reflection-V-3B VR SFT \rightarrow Cap&R SFT	27.96	66.30	56.89	38.17	62.95
	25.04	63.90	54.22	33.59	60.41
Reflection-V-7B VR SFT \rightarrow Cap&R SFT	33.88	73.30	61.33	42.71	71.07
	29.31	69.00	58.41	37.95	66.25

Table 3: Comparative results of cold-start initialization using data with visual reflection pattern and mage caption-based reasoning data on visual reasoning performance improvement.

thermore, we use HallBench (Guan et al., 2024) to quantify the extent of visual hallucinations in VLMs or VRMs.

4.2 Main Result

We evaluate the performance of our model, Reflection-V, on visual reasoning benchmarks across three categories: math, multi-disciplinary, and general, as shown in Table 1. The results indicate that our model significantly outperforms Qwen2.5-VL (Bai et al., 2025) base model and other open-source models of similar scale in reasoning capability. Even compared to existing vision-language reasoning models based on RL, Reflection-V-7B achieves a notable margin of improvement.

Notably, Reflection-V-7B reaches comparable or even superior performance compared to some widely used, large-scale closed-source and open-source VLMs. For instance, on MathVision and MathVista, Reflection-V-7B outperforms GPT-4o and InternVL-2.5-38B (Chen et al., 2024b). On MMMU and M3CoT, Reflection-V-7B surpasses InternVL-2.5-38B and LLaVA-OneVision-72B (Li et al., 2024a), and is comparable to GPT-4o (Hurst et al., 2024). On MMMU-Pro, Reflection-V-7B outperforms LLaVA-OneVision-72B and GPT-4V (Yang et al., 2023), while being comparable to InternVL-2.5-38B. In contrast to existing vision-

language reasoning models, which show improved math reasoning but a decline in multi-disciplinary and general reasoning capabilities, Reflection-V demonstrates improvements across all three categories. Additionally, experimental results show that the proposed method is effective across models of different scales.

Surprisingly, thanks to "visual reflection" we proposed, Reflection-V exhibits significantly fewer visual hallucinations. Specifically, Compared to the Qwen2.5-VL base model, Reflection-V shows a significantly higher all accuracy (aAcc) (Guan et al., 2024) on HallBench. This result suggests that, due to sustained focus on visual information, the visual hallucinations commonly observed in VRMs are notably suppressed. In contrast, other VRMs experience even more severe visual hallucinations than the base model. This can be attributed to the fact that, as mentioned earlier, VRMs struggle to maintain sustained attention to visual information as more textual tokens are generated during reasoning.

4.3 Ablation Study

We ablate the cold-start and visual attention based reward components of our method to evaluate the impact of each design on enhancing visual reasoning capabilities of VRMs. Experimental results in Table 2 demonstrate that both components signifi-

Model	MathVision	MathVista	MMMU	MMMU-Pro	МЗСоТ
Reflection-V-3B (Qwen2.5-VL/QWQ cold start)	27.9	66.3	56.9	38.2	62.9
Reflection-V-3B (InternVL3/Qwen3 cold start)	27.1	67.6	58.0	36.4	64.2

Table 4: Performance comparison between cold-start data constructed with InternVL3/Qwen3 and Qwen2.5-VL/QWQ. In this experiment, InternVL3-38B, Qwen3-32B, and Qwen2.5-VL-72B are employed.

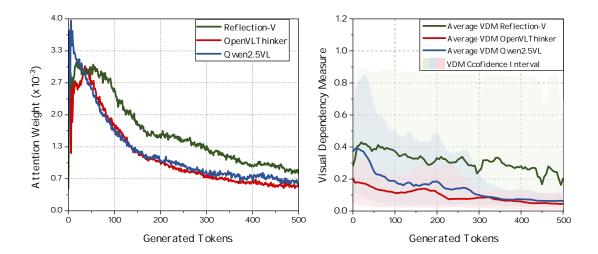


Figure 4: Attention weights (last layer) on visual tokens and visual dependency measure of Reflection-V-7B on MMMU benchmark, compared to OpenVLThinker-7B and Qwen2.5VL-7B. The shown attention weights represent the mean value across all samples. Visual dependency measure quantifies the difference in probability distributions for next token prediction based on generated tokens, before and after discarding visual tokens. The light-green, light-blue, and light-red bands represent the confidence intervals of the visual dependency measure for Reflection-V-7B, Qwen 2.5-VL-7B, and OpenVLThinker-7B, respectively.

cantly improve VRM's performance. Notably, coldstart yields a particularly substantial performance gain. This indicates that emphasizing the continuous reliance and repeated utilization of visual information in SFT data significantly improves visual reasoning performance. Furthermore, with coldstart, the performance improvement from visual attention based reward becomes more pronounced. We believe this occurs because cold-start, based on visual reflection data, guides VRMs in how to increase their attention to visual information.

We conduct a further ablation to validate the superiority of emphasizing visual reflection pattern in cold-start data. Specifically, we replace the reasoning data containing visual reflection patterns with image caption-based reasoning data (derived from the same origin data) during cold-start initialization, then compare their performance across benchmarks. As shown in Table 3, our method outperforms the "caption then reasoning" data construction paradigm by a significant margin. This result also illustrates that the improved visual reasoning performance originates not from distilling

high-quality data from larger models but rather from the intentional incorporation of visual reflection patterns in constructed data.

Besides, We conducted additional experiments to analyze whether the proposed cold-start data construction method exhibits any bias toward specific model families. Specifically, we replaced the VLM and LLM used for constructing cold-start data with InternVL3-38B (Zhu et al., 2025) and Qwen3-32B (Thinking mode) (Yang et al., 2025a), and compared the results with those from constructing cold-start data using Qwen2.5VL-72B and QWQ-32B. The results are shown in Table 4. The experimental findings indicate that the performance gap between InternVL3/Qwen3 and Qwen2.5-VL/QWQ in constructing cold-start data is minimal. This suggests that our cold-start data construction method does not exhibit any bias toward specific model families.

To further validate that the proposed method improves performance by achieving the claimed visual reflection, we present further analyses below.

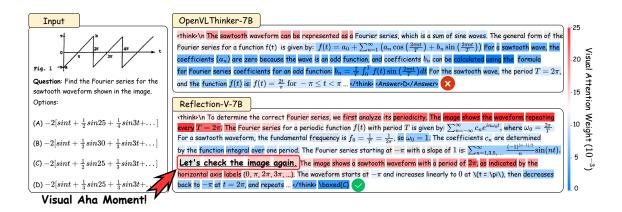


Figure 5: Our model, Reflection-V, exhibits the visual reflection capacity that we claim, in contrast to other RL-based visual reasoning models. The background color of tokens in the figure indicates the magnitude of the visual attention weight. This visual reflection capacity is demonstrated by the recheck and attention to visual tokens again that appear alongside "Aha moment", like "Let's check the image again".

4.4 Analyses

In Section 2, we demonstrate that existing VRMs struggle with visual reflection through three metrics: visual attention weight, and visual dependency measure. Based on these metrics, in this subsection, we analyze whether the performance improvement of Reflection-V genuinely stems from the training strategy that emphasizes visual reflection.

Our method leads to more sustained visual attention and dependence. We compare the attention weight of response tokens to visual tokens for Reflection-V and OpenVLThinker-7B of the same scale, at different generated tokens. We find that, in the middle and deep transformer layers, Reflection-V exhibits significantly higher attention weights to visual tokens than OpenVLThinker, the model also trained through SFT cold-start initialization and GRPO, based on Qwen2.5-VL, as shown in Figure 4. As the number of generated tokens increases, the decrease in attention weight to visual tokens is slower in Reflection-V-7B than in OpenVLThinker-7B. Furthermore, to investigate whether the proposed method enhances VRMs' reliance on visual information, we compare the visual dependency measure, as referred to in equation (3), of Reflection-V and OpenVLThinker at different generated tokens, as shown in Figure 4. The results indicate that, benefiting from the emphasis on visual reflection, Reflection-V significantly mitigates the diminishing of dependence on visual information as generated tokens increase, compared to OpenVLThinker. To sum up, Reflection-V exhibits more sustained reliance on visual information. Experimental results also show that this feature enhances visual reasoning accuracy.

Better performance, sustained visual attention, and reliance all originate from visual re**flection**. As discussed earlier, the proposed method improves visual reasoning performance while maintaining visual attention and dependency during reasoning. Figure 5 presents a comprehensive example demonstrating that these gains are indeed due to the model's visual reflection ability. In this example, Reflection-V actively verifies and refines its reasoning by rechecking the visual input. When textual "Aha moments" like "Let's check the image again" appear, the visual attention weight rises sharply during next-token prediction, representing the true "aha moment" in visual reasoning. As Figure 5 shows, visual reflection capacity enables Reflection-V to reveal critical visual information absent from the reasoning context, thereby inferring the correct answer. Additionally, Figure 4 shows that the upper bound of the confidence interval for visual dependency measure shows virtually no decline with increasing generated tokens. This suggests that during reasoning, as the number of generated tokens increases, Reflection-V maintains consistent dependency on visual tokens throughout the generation process. The observed decline in average visual dependency measure stems from the reduced frequency of visual reflection as the number of generated tokens increases. This phenomenon aligns with the re-emergent, image-focused attention derived from visual reflection, which is observed in Figure 4. These results demonstrate that, when Reflection-V engages in visual reflection, it maintains the same level of focus and reliance on

visual tokens as at the start of reasoning.

5 Related Works

Visual Reasoning model. Large VLMs typically project inputs from non-text modalities into textual representations that LLMs can process (Guo et al., 2025b; Li et al., 2025; Zhang et al., 2025c), achieving strong performance in vision understanding (Bai et al., 2025; Zhang et al., 2025b). The advancement of LLMs has redefined state-of-theart performance across a vast landscape of tasks (Yang et al., 2024; Wang et al., 2024b). Motivated by recent advances in LLM domain (Zhang et al., 2024b; Zheng et al., 2025a; Dai et al., 2025; Sun et al., 2025b), researchers enhance Large VLMs reasoning with step-level reasoning SFT datasets (Xu et al., 2024; Cheng et al., 2025) and RL (Yang et al., 2025e; Wang et al., 2025; Chu et al., 2025; Chung et al., 2025; Yang et al., 2025b). However, as discussed earlier, these trained VRMs typically struggle with visual reflection, leading VRMs to reason without visual grounding after many tokens are generated.

Visual forgetting alleviation. Consistent with visual reflection that we claim, some recent studies emphasize alleviating forgetting visual cues during long-chain inference. M3ID (Favero et al., 2024) employs mutual information decoding to amplify image influence while weakening linguistic priors, thereby promoting continuous reliance on visual cues. But diminishing linguistic priors lowers performance on complex reasoning tasks (Bitton Guetta et al., 2024; Zhang et al., 2024a). TVC (Sun et al., 2025a), a concurrent work, periodically replays visual tokens during inference to reuse visual cues, but it cannot flexibly invoke visual reflection when required. Distinct from these works, we embed visual reflection capability into VRMs based on data generated by LLM-VLM interaction, and reinforce this capability during RL. As a result, the trained VRMs can actively refine their reasoning based on the visual input when needed.

6 Conclusion

In this paper, we propose that the true "aha moment" in visual reasoning arises when a model engages in visual reflection—that is, when it actively verifies and refines its reasoning based on the visual input. Through quantitative studies, we reveal that existing VRMs struggle with such visual reflection. Therefore, to address this critical

challenge, we propose a two-stage training strategy combining LLM-VLM interaction-driven reflective reasoning patterns with visual attention based RL. This training strategy significantly improves performance across multiple benchmarks. Experiments confirm that such improvement is derived from sustained visual attention and reliance, demonstrating the effectiveness of visual reflection. This work establishes a foundation for integrating visual reflection into VRMs, narrowing the gap between visual and text-only reasoning on complex tasks.

Limitations

Firstly, due to computational constraints, we could not conduct experiments like GRPO with models larger than 7B parameters. Thus, we limited our exploration to the 3B and 7B parameter scales. Secondly, our cold-start initialization (based on constructed reasoning data with visual reflection pattern), reinforcement learning (with visual attention based reward), and evaluation presently involve relatively limited categories of visual-language datasets. In future work, we plan to include a wider range of visual-language datasets covering diverse problem types to further evaluate the generalization ability of the proposed method.

Acknowledgments

We thank the Wuhan AI Research for providing the computational resources. We also appreciate the outstanding open-source code repositories contributed by the LLM and multimodal communities. Finally, we thank all the reviewers for their detailed reviews and insightful comments.

This work is supported by National Key R&D Program of China 2022ZD0160602 and the Strategic Priority Research Program of Chinese Academy of Sciences under Grant XDA04080400.

References

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. 2025. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*.

Nitzan Bitton Guetta, Aviv Slobodkin, Aviya Maimon, Eliya Habba, Royi Rassin, Yonatan Bitton, Idan Szpektor, Amir Globerson, and Yuval Elovici. 2024. Visual riddles: a commonsense and world knowledge challenge for large vision and language models. *Advances in Neural Information Processing Systems*, 37:139561–139588.

- Hardy Chen, Haoqin Tu, Fali Wang, Hui Liu, Xianfeng Tang, Xinya Du, Yuyin Zhou, and Cihang Xie. 2025. Sft or rl? an early investigation into training rllike reasoning large vision-language models. *arXiv* preprint arXiv:2504.11468.
- Liangyu Chen, Bo Li, Sheng Shen, Jingkang Yang, Chunyuan Li, Kurt Keutzer, Trevor Darrell, and Ziwei Liu. 2023. Large language models are visual reasoning coordinators. *Advances in Neural Information Processing Systems*, 36:70115–70140.
- Qiguang Chen, Libo Qin, Jin Zhang, Zhi Chen, Xiao Xu, and Wanxiang Che. 2024a. M3cot: A novel benchmark for multi-domain multi-step multi-modal chain-of-thought. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8199–8221.
- Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. 2024b. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv* preprint arXiv:2412.05271.
- Kanzhi Cheng, Yantao Li, Fangzhi Xu, Jianbing Zhang, Hao Zhou, and Yang Liu. 2024. Vision-language models can self-improve reasoning via reflection. *arXiv preprint arXiv:2411.00855*.
- Kanzhi Cheng, Li YanTao, Fangzhi Xu, Jianbing Zhang, Hao Zhou, and Yang Liu. 2025. Vision-language models can self-improve reasoning via reflection. In Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 8876–8892.
- Xu Chu, Xinrong Chen, Guanyu Wang, Zhijie Tan, Kui Huang, Wenyu Lv, Tong Mo, and Weiping Li. 2025. Qwen look again: Guiding vision-language reasoning models to re-attention visual information. *arXiv* preprint arXiv:2505.23558.
- Jiwan Chung, Junhyeok Kim, Siyeol Kim, Jaeyoung Lee, Min Soo Kim, and Youngjae Yu. 2025. Don't look only once: Towards multimodal interactive reasoning with selective visual revisitation. *arXiv* preprint arXiv:2505.18842.
- Runpeng Dai, Linfeng Song, Haolin Liu, Zhenwen Liang, Dian Yu, Haitao Mi, Zhaopeng Tu, Rui Liu, Tong Zheng, Hongtu Zhu, et al. 2025. Cde: Curiosity-driven exploration for efficient reinforcement learning in large language models. *arXiv* preprint arXiv:2509.09675.
- Yihe Deng, Hritik Bansal, Fan Yin, Nanyun Peng, Wei Wang, and Kai-Wei Chang. 2025. Openvlthinker: An early exploration to complex vision-language reasoning via iterative self-improvement. *arXiv preprint arXiv:2503.17352*.
- Alessandro Favero, Luca Zancato, Matthew Trager, Siddharth Choudhary, Pramuditha Perera, Alessandro

- Achille, Ashwin Swaminathan, and Stefano Soatto. 2024. Multi-modal hallucination control by visual information grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14303–14312.
- Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, et al. 2024. Hallusionbench: an advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14375–14385.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025a. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Jiawei Guo, Feifei Zhai, Pu Jian, Qianrun Wei, and Yu Zhou. 2025b. Crop: Contextual region-oriented visual token pruning. *arXiv preprint arXiv:2505.21233*.
- Yunzhuo Hao, Jiawei Gu, Huichen Will Wang, Linjie Li, Zhengyuan Yang, Lijuan Wang, and Yu Cheng. 2025. Can mllms reason in multimodality? emma: An enhanced multimodal reasoning benchmark. *arXiv* preprint arXiv:2501.05444.
- Wenxuan Huang, Bohan Jia, Zijie Zhai, Shaosheng Cao, Zheyu Ye, Fei Zhao, Zhe Xu, Yao Hu, and Shaohui Lin. 2025. Vision-r1: Incentivizing reasoning capability in multimodal large language models. *arXiv* preprint arXiv:2503.06749.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. 2024. Openai o1 system card. *arXiv preprint arXiv:2412.16720*.
- Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. 2017. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2901–2910.
- Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. 2016. A diagram is worth a dozen images. In Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14, pages 235–251. Springer.

- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. 2024a. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*.
- Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong Feng, Lingpeng Kong, and Qi Liu. 2024b. Multimodal arxiv: A dataset for improving scientific comprehension of large vision-language models. *arXiv* preprint arXiv:2403.00231.
- Zhaowei Li, Wei Wang, YiQing Cai, Qi Xu, Pengyu Wang, Dong Zhang, Hang Song, Botian Jiang, Zhida Huang, and Tao Wang. 2025. Unifiedmllm: Enabling unified representation for multi-modal multi-tasks with large language model. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 334–344.
- Zhuowan Li, Xingrui Wang, Elias Stengel-Eskin, Adam Kortylewski, Wufei Ma, Benjamin Van Durme, and Alan L Yuille. 2023. Super-clevr: A virtual benchmark to diagnose domain robustness in visual reasoning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14963–14973.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2023. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*.
- Pan Lu, Ran Gong, Shibiao Jiang, Liang Qiu, Siyuan Huang, Xiaodan Liang, and Song-chun Zhu. 2021. Inter-gps: Interpretable geometry problem solving with formal language and symbolic reasoning. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 6774–6786.
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*, pages 3195–3204.
- Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2263–2279.

- Fanqing Meng, Lingxiao Du, Zongkai Liu, Zhixiang Zhou, Quanfeng Lu, Daocheng Fu, Tiancheng Han, Botian Shi, Wenhai Wang, Junjun He, et al. 2025. Mm-eureka: Exploring the frontiers of multimodal reasoning with rule-based reinforcement learning. arXiv preprint arXiv:2503.07365.
- Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. 2022. A-okvqa: A benchmark for visual question answering using world knowledge. In *European conference on computer vision*, pages 146–162. Springer.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.
- Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. 2024. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv:* 2409.19256.
- Wenhao Shi, Zhiqiang Hu, Yi Bin, Junhua Liu, Yang Yang, See Kiong Ng, Lidong Bing, and Roy Lee. 2024. Math-llava: Bootstrapping mathematical reasoning for multimodal large language models. In Findings of the Association for Computational Linguistics: EMNLP 2024, pages 4663–4680.
- Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. 2019. Towards vqa models that can read. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8317–8326.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2024. Scaling llm test-time compute optimally can be more effective than scaling model parameters. arXiv preprint arXiv:2408.03314.
- Hai-Long Sun, Zhun Sun, Houwen Peng, and Han-Jia Ye. 2025a. Mitigating visual forgetting via takealong visual conditioning for multi-modal long cot reasoning. *arXiv preprint arXiv:2503.13360*.
- Wei Sun, Qianlong Du, Fuwei Cui, and Jiajun Zhang. 2025b. An efficient and precise training data construction framework for process-supervised reward model in mathematical reasoning. *arXiv* preprint *arXiv*:2503.02382.
- Huajie Tan, Yuheng Ji, Xiaoshuai Hao, Minglan Lin, Pengwei Wang, Zhongyuan Wang, and Shanghang Zhang. 2025. Reason-rft: Reinforcement fine-tuning for visual reasoning. *arXiv preprint arXiv:2503.20752*.
- Kimi Team, Angang Du, Bohong Yin, Bowei Xing, Bowen Qu, Bowen Wang, Cheng Chen, Chenlin Zhang, Chenzhuang Du, Chu Wei, et al. 2025. Kimivl technical report. *arXiv preprint arXiv:2504.07491*.

- Qwen Team. 2025. Qwq-32b: Embracing the power of reinforcement learning. *URL: https://qwenlm. github.io/blog/qwq-32b*.
- Omkar Thawakar, Dinura Dissanayake, Ketan More, Ritesh Thawkar, Ahmed Heakl, Noor Ahsan, Yuhao Li, Mohammed Zumri, Jean Lahoud, Rao Muhammad Anwer, et al. 2025. Llamav-o1: Rethinking step-by-step visual reasoning in llms. *arXiv preprint arXiv:2501.06186*.
- Haozhe Wang, Chao Qu, Zuming Huang, Wei Chu, Fangzhen Lin, and Wenhu Chen. 2025. Vl-rethinker: Incentivizing self-reflection of vision-language models with reinforcement learning. *arXiv preprint arXiv:2504.08837*.
- Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Houxing Ren, Aojun Zhou, Mingjie Zhan, and Hongsheng Li. 2024a. Measuring multimodal mathematical reasoning with math-vision dataset. *Advances in Neural Information Processing Systems*, 37:95095–95169.
- Wei Wang, Zhaowei Li, Qi Xu, Yiqing Cai, Hang Song, Qi Qi, Ran Zhou, Zhida Huang, Tao Wang, and Li Xiao. 2024b. Qcrd: Quality-guided contrastive rationale distillation for large language models. *arXiv* preprint arXiv:2405.13014.
- Yichen Wei, Yi Peng, Xiaokun Wang, Weijie Qiu, Wei Shen, Tianyidan Xie, Jiangbo Pei, Jianhao Zhang, Yunzhuo Hao, Xuchen Song, et al. 2025. Skywork r1v2: Multimodal hybrid reinforcement learning for reasoning. *arXiv preprint arXiv:2504.16656*.
- Wenyi Xiao, Leilei Gan, Weilong Dai, Wanggui He, Ziwei Huang, Haoyuan Li, Fangxun Shu, Zhelun Yu, Peng Zhang, Hao Jiang, et al. 2025. Fast-slow thinking for large vision-language model reasoning. arXiv preprint arXiv:2504.18458.
- Guowei Xu, Peng Jin, Li Hao, Yibing Song, Lichao Sun, and Li Yuan. 2024. Llava-o1: Let vision language models reason step-by-step. *arXiv preprint arXiv:2411.10440*.
- Hanqi Yan, Qinglin Zhu, Xinyu Wang, Lin Gui, and Yulan He. 2024. Mirror: Multiple-perspective self-reflection method for knowledge-rich reasoning. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7086–7103.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. 2025a. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.
- Shuo Yang, Yuwei Niu, Yuyang Liu, Yang Ye, Bin Lin, and Li Yuan. 2025b. Look-back: Implicit visual re-focusing in mllm reasoning. *arXiv preprint arXiv:2507.03019*.
- Wen Yang, Junhong Wu, Chen Wang, Chengqing Zong, and Jiajun Zhang. 2024. Language imbalance driven rewarding for multilingual self-improving. *arXiv* preprint arXiv:2410.08964.

- Wen Yang, Junhong Wu, Chen Wang, Chengqing Zong, and Jiajun Zhang. 2025c. Implicit cross-lingual rewarding for efficient multilingual preference alignment. arXiv preprint arXiv:2503.04647.
- Wenkai Yang, Shuming Ma, Yankai Lin, and Furu Wei. 2025d. Towards thinking-optimal scaling of test-time compute for llm reasoning. *arXiv* preprint *arXiv*:2502.18080.
- Yi Yang, Xiaoxuan He, Hongkun Pan, Xiyan Jiang, Yan Deng, Xingtao Yang, Haoyu Lu, Dacheng Yin, Fengyun Rao, Minfeng Zhu, et al. 2025e. R1-onevision: Advancing generalized multimodal reasoning through cross-modal formalization. *arXiv* preprint arXiv:2503.10615.
- Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. 2023. The dawn of lmms: Preliminary explorations with gpt-4v (ision). *arXiv preprint arXiv:2309.17421*, 9(1):1.
- Qiying Yu, Quan Sun, Xiaosong Zhang, Yufeng Cui, Fan Zhang, Yue Cao, Xinlong Wang, and Jingjing Liu. 2024. Capsfusion: Rethinking image-text data at scale. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14022–14032.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. 2024a. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567.
- Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang, Kai Zhang, Shengbang Tong, Yuxuan Sun, Botao Yu, Ge Zhang, Huan Sun, et al. 2024b. Mmmu-pro: A more robust multi-discipline multimodal understanding benchmark. arXiv preprint arXiv:2409.02813.
- Jingyi Zhang, Jiaxing Huang, Huanjin Yao, Shunyu Liu, Xikun Zhang, Shijian Lu, and Dacheng Tao. 2025a. R1-vl: Learning to reason with multimodal large language models via step-wise group relative policy optimization. *arXiv preprint arXiv:2503.12937*.
- Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Yu Qiao, et al. 2024a. Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems? In *European Conference on Computer Vision*, pages 169–186. Springer.
- Wanyue Zhang, Yibin Huang, Yangbin Xu, JingJing Huang, Helu Zhi, Shuo Ren, Wang Xu, and Jiajun Zhang. 2025b. Why do mllms struggle with spatial understanding? a systematic analysis from data to architecture. *arXiv preprint arXiv:2509.02359*.

- Wanyue Zhang, Ziyong Li, Wen Yang, Chunlin Leng, Yinan Bai, Qianlong Du, Chengqing Zong, and Jiajun Zhang. 2024b. Chinesewebtext 2.0: Large-scale high-quality chinese web text with multi-dimensional and fine-grained information. *arXiv* preprint arXiv:2411.19668.
- Yunhao Zhang, Xiaohan Zhang, Chong Li, Shaonan Wang, and Chengqing Zong. 2025c. Mulcogbench: a multi-modal cognitive benchmark dataset for evaluating chinese and english computational language models: Y. zhang et al. *Language Resources and Evaluation*, pages 1–24.
- Yuze Zhao, Jintao Huang, Jinghan Hu, Xingjun Wang, Yunlin Mao, Daoze Zhang, Zeyinzi Jiang, Zhikai Wu, Baole Ai, Ang Wang, Wenmeng Zhou, and Yingda Chen. 2024. Swift:a scalable lightweight infrastructure for fine-tuning. *Preprint*, arXiv:2408.05517.
- Tong Zheng, Hongming Zhang, Wenhao Yu, Xiaoyang Wang, Xinyu Yang, Runpeng Dai, Rui Liu, Huiwen Bao, Chengsong Huang, Heng Huang, et al. 2025a. Parallel-r1: Towards parallel thinking via reinforcement learning. *arXiv preprint arXiv:2509.07980*.
- Yaowei Zheng, Junting Lu, Shenzhi Wang, Zhangchi Feng, Dongdong Kuang, and Yuwen Xiong. 2025b. Easyr1: An efficient, scalable, multi-modality rl training framework.
- Weihong Zhong, Xiaocheng Feng, Liang Zhao, Qiming Li, Lei Huang, Yuxuan Gu, Weitao Ma, Yuan Xu, and Bing Qin. 2024. Investigating and mitigating the multimodal hallucination snowballing in large vision-language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11991–12011.
- Jingqi Zhou, Sheng Wang, Jingwei Dong, Lei Li, Jiahui Gao, Jiyue Jiang, Lingpeng Kong, and Chuan Wu. 2024. Proreason: Multi-modal proactive reasoning with decoupled eyesight and wisdom. *arXiv preprint arXiv:2410.14138*.
- Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, et al. 2025. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv preprint arXiv:2504.10479*.

Hyper-parameters	Value
Epoches	3
Batch size	8
Warmup ratio	0.1
Gradient accumulation	4
Learning rate scheduler	Cosine
GPUs	2
Optimizer	AdamW

Table 5: The hyperparameters used during cold-start initialization using the constructed data with visual reflection pattern.

Hyper-parameters	Value
Hyper-parameters	12
Batch size	512
Micro Batch size	8
Warmup	False
Rollout	16
Rollout Temperature	1.0
Rollout Top-P	0.99
Freeze Vision Tower	True
KL divergence coefficient	1×10^{-2}
Learning rate	5×10^{-6}
Weight Decay	1×10^{-2}
GPUs	8
Optimizer	AdamW
Framework	Verl
	l .

Table 6: The hyper-parameters used during GRPO with visual attention based reward.

A Implementation Details and Hyperparameters

During supervised fine-tuning (SFT), we set the learning rate to $1e^{-5}$, apply a cosine scheduler with a 0.1 warm-up ratio, use BF16 mixed precision, adopt a batch size of 8, and train for 3 epochs. In the reinforcement-learning (RL) phase, we limit both prompts and responses to 2048 tokens and apply KL divergence with a coefficient of $1e^{-2}$. Each training step processes 512 questions with 16 rollouts per question; rollout sampling uses a temperature of 1.0 and a top-p of 0.99. For validation, we evaluate with the pass@1 metric and set the sampling temperature to 0.5. Detailed hyperparameters are shown in Table 5 and Table 5.

B Prompts

Prompt Templates of Visual Requester

You currently need to address the following question: <question> The information you need is in an image, but you can't see the image right now. At the same time, you're not capable of complex reasoning.

However, you can can consult the following two Vision Expert for help. You can ask him a single question for information in the picture, for example, you could ask him, "What color is the bird in the picture?"

Use the following format: {'Thought': 'analyze the problem here.', 'Question':'Questions you want to ask the Vision EXPERT'}

<split>
And the information you know
currently is as follows:
<info>

Prompt Templates of Visual Responder

Please answer my question in a tone that provides a concise description of the image. If it is a yes/no question, focus on describing the relevant visual information, avoiding answering with yes/no.

Question:
<question>

Prompt Templates of Summarizer

The following is the available information: <info>

Please solve the following problems step by step: <question>

Use the following format: Thought: Conduct an analysis before you give me an answer.

Final Answer: "The final answer you get when you have finished reasoning."

Prompt Templates of Cohesion Enhancement

Below is the reasoning steps for the question <Question>, but there are some disjointed parts marked with "...". Please fill in the gaps to improve coherence. You can use some connecting phrases such as "Let's double check," "Let's check the image again," and "To sum up," and "Wait".

Use the following format:
'Thought': 'Reasoning steps',
'Final answer':'\boxed{...}'
The final answer (only choice
like A, B, C, D) MUST BE put in

The reasoning steps is:

<Reasoning>

\boxed{}.

Prompt Templates of RL Training and Evaluation

You FIRST think about the reasoning process as an internal monologue and then provide the final answer.

The reasoning process MUST BE enclosed within <think> </think> tags. The final answer MUST BE put in \boxed{}.

Qustion:

Table 7: Detailed composition of the datasets used to construct reasoning data with visual-reflection pattern for cold-start initialization.

Datasets	Samples
AI2D (Kembhavi et al., 2016)	~ 0.5K
A-OKVQA (Marino et al., 2019)	$\sim 0.5 \mathrm{K}$
M3CoT (train set) (Chen et al., 2024a)	$\sim 1.0 \text{K}$
CLEVR-Math (Johnson et al., 2017)	$\sim 0.5 \mathrm{K}$
ScienceQA (Masry et al., 2022)	$\sim 0.5 \mathrm{K}$
TextVQA (Singh et al., 2019)	$\sim 0.2 \mathrm{K}$

Table 8: Detailed composition of the datasets used to conduct GRPO.

Datasets	Samples
Geo3K (Lu et al., 2021)	~ 2.1K
AI2D (Kembhavi et al., 2016)	$\sim 1.5 \text{K}$
TextVQA (Singh et al., 2019)	$\sim 0.8 \mathrm{K}$
M3CoT (train set) (Chen et al., 2024a)	$\sim 3.0 \text{K}$
MathVerse (Zhang et al., 2024a)	$\sim 2.5 \text{K}$
Super-CLEVR (Li et al., 2023)	$\sim 0.5 \text{K}$
MathV360K (Shi et al., 2024)	$\sim 1.0 \text{K}$
A-OKVQA (Marino et al., 2019)	$\sim 0.5 \text{K}$
ScienceQA (Schwenk et al., 2022)	$\sim 0.5 \text{K}$
ChartQA (Masry et al., 2022)	$\sim 1.0 \text{K}$
ArxivQA (Li et al., 2024b)	$\sim 1.0 \text{K}$
EMMA (Hao et al., 2025)	$\sim 1.6 \text{K}$

C Data Resources

We collect data from a large multimodal corpus for (1) constructing reasoning data with visual reflection pattern (cold-start initialization stage) and (2) GRPO training, as summarized in Tables X and Y.

D Supplementary Experiments

Broader evaluation of Reflection-V's capability of sustained visual attention. Beyond the MMMU results reported in the main text, Figures 6 and 7 evaluate Reflection-V-7B on three additional visual-reasoning benchmarks spanning mathematical, multi-disciplinary, and general domains. Figure 6 traces last-layer attention from response tokens to visual tokens over 500 generated tokens: Reflection-V-7B consistently launches with higher visual-attention strength than Qwen2.5VL-7B and—crucially—decays far more slowly, retaining about 30%–40% of its initial level where baselines sink below 10%. This advantage extends to other VRMs such as OpenVLThinker-7B, whose

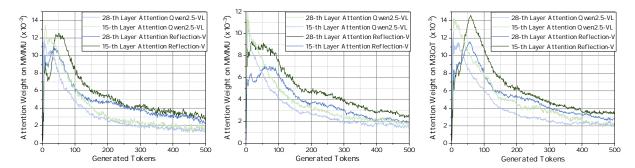


Figure 6: Attention weights (last layer) on visual tokens of Reflection-V-7B on multiple benchmarks, compared to OpenVLThinker-7B and Qwen2.5VL-7B. The shown attention weights represent the mean value across all samples.

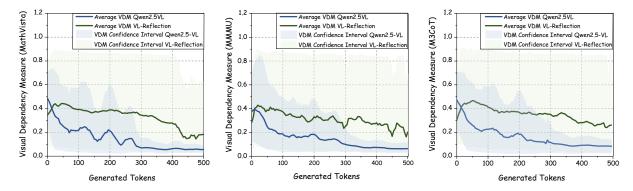


Figure 7: Visual dependency measur of Reflection-V-7B on multiple benchmarks, compared to OpenVLThinker-7B and Qwen2.5VL-7B. Visual dependency measure quantifies the difference in probability distributions for next token prediction based on generated tokens, before and after discarding visual tokens.

cold-start and RL stages, as noted earlier, further erode visual attention and dependency; Reflection-V therefore surpasses these models as well. Figure 7 reports the Visual Dependency Measure (VDM): the upper bound of Reflection-V's confidence interval remains nearly flat across all tasks, indicating sustained reliance on visual evidence, while the baselines exhibit a pronounced downward trend. These supplementary findings demonstrate that Reflection-V maintains robust visual attention and dependency across diverse reasoning scenarios, substantiating its superior visual-reflection capability.

Scaling experiment of the proposed method.

We further conduct a scaling experiment on the MS-SWIFT framework (Zhao et al., 2024), leveraging LoRA as an efficient training technique on InternVL3-14B. The number of GRPO training epochs is also set to 12. The evaluation results on several visual reasoning benchmarks are shown in Table A. These results demonstrate that our method significantly outperforms GRPO, which solely relies on textual output-based reward. This experiment demonstrates that our approach can effectively scale to larger models and enhance their vi-

sual reasoning capabilities.

E Case Study

Model	Math-Reasoning		Multi-Disciplinary		General	
	MathVision	MathVista	MMMU	MMMU-Pro	M3CoT	
Closed-So	urce Vision-L	anguage Mo	odels			
GPT-40 (Hurst et al., 2024)	30.4	60.0	69.1	51.9	†74.2	
GPT-4V (Yang et al., 2023)	22.8	49.9	56.8	41.2	62.6	
Open-Sou	ırce Vision-La	anguage Mo	dels			
QwenVL2.5-3B (Bai et al., 2025)	21.2	62.3	†51.2	31.6	[†] 55.6	
QwenVL2.5-7B (Bai et al., 2025)	25.1	68.2	[†] 54.3	36.9	†60.5	
InternVL-2.5-38B (Chen et al., 2024b)	32.2	71.9	57.6	46.0	†68.9	
InternVL3-14B (Chen et al., 2024b)	35.9	73.8	64.1	48.9	70.1	
Open-Source Vision-Language Reasoning Models						
MM-Eureka-7B ^{\(\phi\)} (Meng et al., 2025)	26.9	73.0	[†] 51.3	†36.7	†63.5	
R1-Onevision-7B [⋄] (Yang et al., 2025e)	29.9	64.1	†48.7	21.6	[†] 53.1	
OpenVLThinker-7B ^{\dightarrow} (Deng et al., 2025)	25.3	70.2	52.5	37.3	†62.2	
Ours (Training strategy emphasizing visual reflection)						
InternVL3-14B (GRPO)	38.3	<u>75.6</u>	66.9	51.3	73.4	
Reflection-V (InternVL3-14B)	39.8	76.8	68.7	53.5	78.1	

Table 9: Performance of the proposed approach when scaled to 14B parameters (InternVL3-14B).

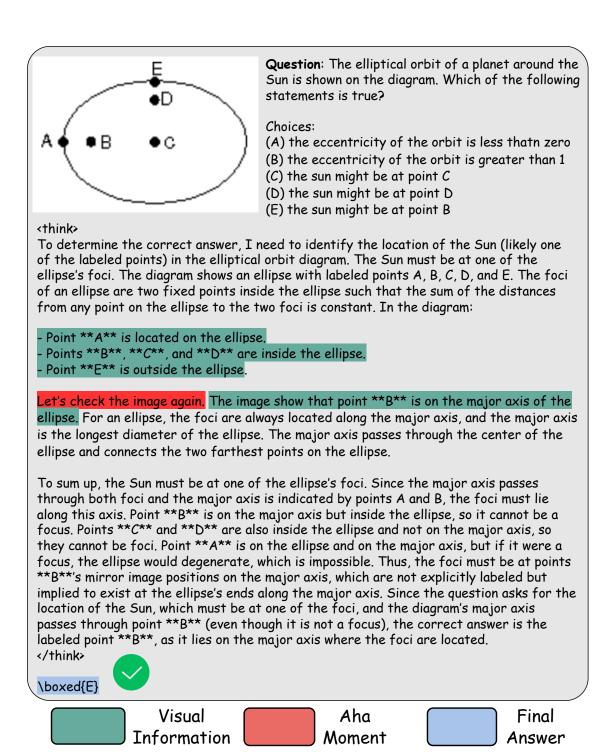
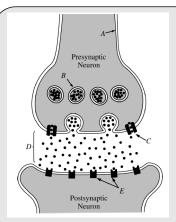


Figure 8: Case study 1 for our model, Reflection-V, performs visual reflection during visual reasoning. In this case, based on rechecking the image, Reflection-V reveals key visual information which is not in the generated reasoning context, therefore finally infers the correct answer.



Question: In the figure above, where does the neurotransmitter dopamine have its effect?

Choices: (A) A, (B) B, (C) C, (D) D, (E) E

<think>

To determine where dopamine has its effect, I need to identify the receptor sites on the postsynaptic neuron, which are likely labeled in the figure. In the image:

- **E** is shown as black structures on the surface of the postsynaptic neuron.
- These structures are likely dopamine receptors, as they are the sites where neurotransmitters bind to initiate the effect

Thus, the label **E** represents the receptor sites on the postsynaptic neuron.

Wait. let me confirm the labels again to ensure accuracy. To determine where dopamine has its effect, I need to confirm which label corresponds to the receptor sites on the postsynaptic neuron where neurotransmitters bind.

The label **E** in the figure is shown on the surface of the postsynaptic neuron. These labels are depicted as black structures, which likely represent receptor sites. The synaptic cleft, labeled **D**, is filled with neurotransmitters (indicated by the black dots), ...

To sum up, the receptor sites (E) are the final point of action for dopamine. The neurotransmitter dopamine binds to receptor sites on the postsynaptic neuron to produce its effect. From the information provided, label E is described as black structures on the postsynaptic neuron that face the synaptic cleft (labeled D). Neurotransmitters (likely dopamine) are shown in the synaptic cleft (D), and their binding to the receptor sites on the postsynaptic neuron (labeled E) is implied. Since dopamine binds to these receptor sites to initiate the effect, E represents the site of action for dopamine. The other labels (A,B,C,D) are associated with structural or functional elements of the synaptic region but are not explicitly described as receptor sites.

Final Answer: The neurotransmitter dopamine has its effect at label **E**. </think> \boxed{E}



Figure 9: Case study 2 for our model, Reflection-V, performs visual reflection during visual reasoning. In this case, based on rechecking the image, Reflection-V reveals key visual information which is not in the generated reasoning context, therefore finally infers the correct answer.

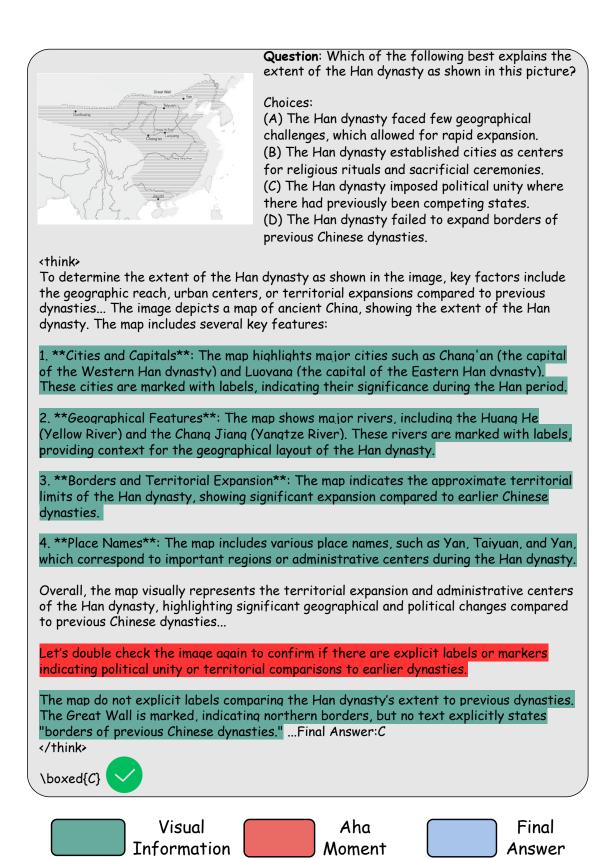


Figure 10: Case study 3 for our model, Reflection-V, performs visual reflection during visual reasoning. In this case, based on rechecking the image, Reflection-V reveals key visual information which is not in the generated reasoning context, therefore finally infers the correct answer.