

Visually-Guided Policy Optimization for Multimodal Reasoning

Zengbin Wang^{1*}, Feng Xiong^{1*}, Liang Lin¹, Xuecai Hu^{1†}
Yong Wang^{1†}, Yanlin Wang², Man Zhang³, Xiangxiang Chu¹
¹AMAP, Alibaba Group ²SYSU ³BUPT

Abstract

Reinforcement learning with verifiable rewards (RLVR) has significantly advanced the reasoning ability of vision-language models (VLMs). However, the inherent text-dominated nature of VLMs often leads to insufficient visual faithfulness, characterized by sparse attention activation to visual tokens. More importantly, our empirical analysis reveals that temporal visual forgetting along reasoning steps exacerbates this deficiency. To bridge this gap, we propose Visually-Guided Policy Optimization (VGPO), a novel framework to reinforce visual focus during policy optimization. Specifically, VGPO initially introduces a *Visual Attention Compensation* mechanism that leverages visual similarity to localize and amplify visual cues, while progressively elevating visual expectations in later steps to counteract visual forgetting. Building on this mechanism, we implement a dual-grained advantage re-weighting strategy: the intra-trajectory level highlights tokens exhibiting relatively high visual activation, while the inter-trajectory level prioritizes trajectories demonstrating superior visual accumulation. Extensive experiments demonstrate that VGPO achieves better visual activation and superior performance in mathematical multimodal reasoning and visual-dependent tasks. The code has been released at <https://github.com/wzb-bupt/VGPO>.

1 Introduction

Vision-Language Models (VLMs) (Hurst et al., 2024; Bai et al., 2025b; Comanici et al., 2025), renowned for their multimodal reasoning over both image and text, have achieved remarkable progress across multimodal tasks (Lu et al., 2024; Zhang et al., 2024b; Yue et al., 2025; Zhou et al.; Qiao et al., 2025). Recently, the emerging Reinforcement Learning with Verifiable Rewards

* Equal contribution.

† Project leads and corresponding authors.

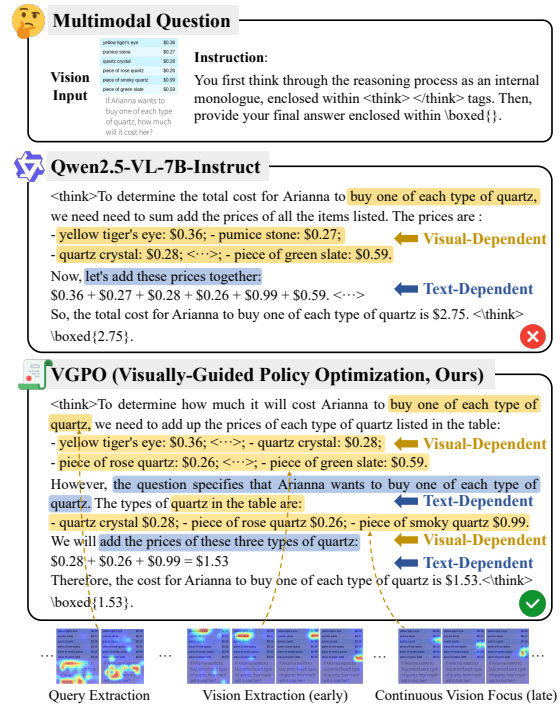


Figure 1: A multimodal reasoning example with visual input. Our VGPO not only attends to the query and visual elements in the early reasoning steps, like the base model (Bai et al., 2025b), but also exhibits better visual activation and sustains visual focus in the later steps to ensure consistent visual faithfulness.

(RLVR) (Shao et al., 2024; Yu et al., 2025; Zheng et al., 2025; Chu et al., 2025; Dai et al., 2026) has further advanced into complex step-by-step reasoning, yielding enhanced logical coherence and superior performance (Huang et al., 2025b; Meng et al., 2025; Wang et al., 2025a; Huang et al., 2025a; Yang et al., 2025c; Wang et al., 2025e).

Despite exhibiting better logical reasoning, the inherent text-dominated inference nature of VLMs often leads to *insufficient visual faithfulness* (Wang et al., 2024; Fu et al., 2024; Liu et al., 2025c,b). As discussed in recent literature (Jian et al., 2025; Yang et al., 2025b; Favero et al., 2024) and evidenced by Figures 1 and 2, models often assign

lower attention to input image tokens compared to query tokens and generated text tokens, resulting in sparse activation of input image. More critically, this deficit will be exacerbated by *temporal visual forgetting*: as reasoning chains extend, attention to visual inputs progressively decays. This over-reliance on textual priors rather than visual facts often leads to hallucinations or reasoning errors (Sun et al., 2025b; Wang et al., 2025e; Tian et al., 2025). Therefore, a natural question arises:

Q1: How can we effectively amplify visual activation and mitigate temporal visual forgetting to ensure visual faithfulness?

Recent works primarily focused on the former aspect “amplify visual response” through external interventions. For example, (1) Look-Back (Yang et al., 2025b) and latent reasoning methods (Li et al., 2025a; Sun et al., 2025a; Yang et al., 2025d) aimed to introduce external learnable special tokens (e.g., <back>, <latent_strat>, <latent_end>) that act as explicit triggers to revisit or reconstruct visual input through initial supervised fine-tuning and subsequent autonomous generation. (2) VAPO (Tian et al., 2025) explored an additional GPT-5 model to verify the visual faithfulness during intermediate steps and re-activate visual responses through probing visual-related questions. (3) PAPO (Wang et al., 2025e) and VPPO (Huang et al., 2025a) identified and highlighted visual tokens by contrasting the KL divergence between original versus noisy images through two forward processes. While effective, these methods inevitably introduce extra training tokens, auxiliary models, or additional forward passes. In this case, another question arises:

Q2: How can we leverage the model’s own internal states to localize visual activation, without any external dependencies?

To comprehensively answer these two questions, the key solution lies in (i) precisely localizing visual activations; (ii) amplifying visual activations; and (iii) sustaining visual expectation in the later reasoning steps. For the first part, our empirical experiment in Figure 2 reveals that the inherent hidden states similarity between generated tokens and input image tokens serves as a precise *Visual Focus Score*, enabling spontaneous localization of visually grounded tokens without external supervision. For the second and third parts, visual focus score serves as the basis for further designing a *Visual Attention Compensation* mechanism that

re-weights visually grounded tokens and progressively improves visual expectations in later reasoning steps to combat temporal visual forgetting.

Building on the above two insights, we propose **Visually-Guided Policy Optimization (VGPO)**. Relying solely on the model’s self-generated hidden states to localize visual tokens, VGPO integrates the Visual Attention Compensation mechanism into policy optimization through a *dual-grained advantage re-weighting* strategy. At *intra-trajectory* level, it dynamically re-weights token advantages to amplify specific visual cues, while applying progressive incentives to compensate for temporal visual decay in later steps. At *inter-trajectory* level, it globally prioritizes superior visual accumulation, effectively compensating for holistic visual negligence by favoring trajectories that sustain consistent grounding. Extensive experiments across mathematical and vision-dependent multimodal reasoning tasks demonstrate that VGPO achieves better visual activation and less visual forgetting.

In summary, our contributions are four-fold.

- We reveal that the inherent hidden states similarity between generated tokens and image tokens serves as a reliable *Visual Focus Score* to localize visual tokens, thereby facilitating targeted attention modulation without external supervision.
- We design a *Visual Attention Compensation* mechanism that leverages visual focus score to amplify visual cues, while progressively elevating visual expectations in the latter reasoning step to counteract temporal visual forgetting.
- We propose VGPO, which integrates the visual attention compensation mechanism into policy optimization using a *Dual-Grained Advantage Re-Weighting* strategy to amplify visual activation and mitigate visual forgetting across intra-trajectory and inter-trajectory levels.
- Extensive experiments demonstrate that VGPO achieves better visual activation and state-of-the-art performance on mathematical and vision-dependent multimodal reasoning tasks.

2 Preliminary

2.1 Group Relative Policy Optimization

Group Relative Policy Optimization (GRPO) (Shao et al., 2024) is a resource-efficient variant of reinforcement learning, which eliminates the value model (Schulman et al., 2017) by estimating advantages through group-based computation. Formally,

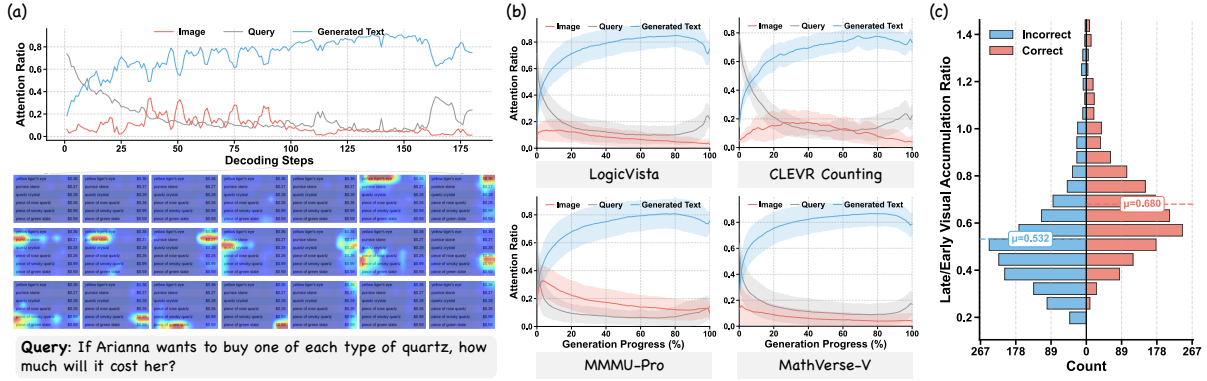


Figure 2: Analysis of the inference nature of multimodal reasoning trajectory (based on Qwen2.5-VL-7B (Bai et al., 2025b)). (a) An example of attention allocation across image, query, and generated text tokens (normalized to 1 at each step). (b) Average attention statistics on four visual-dependent benchmarks (Xiao et al., 2024; Li et al., 2023; Yue et al., 2025; Zhang et al., 2024b). (c) Distribution of late/early visual accumulation ratios for incorrect (left) vs. correct (right) samples of these four benchmarks. Incorrect samples often exhibit higher visual forgetting.

we consider visual input (I, q, a) , comprising a visual input I , a textual query q , and the corresponding ground truth answer a , the algorithm begins by sampling a group of \mathcal{G} candidate responses $\{o_i\}_{i=1}^{\mathcal{G}}$ from the current policy $\pi_{\theta_{\text{old}}}$. For each response o_i , a binary reward $r_i \in \{0, 1\}$ is assigned based on the exact match with the ground truth a . Subsequently, the advantage \hat{A}_i for each response is computed by normalizing the assigned rewards within the sampled group: $\hat{A}_i = \frac{r_i - \text{mean}(\{r_k\}_{k=1}^{\mathcal{G}})}{\text{std}(\{r_k\}_{k=1}^{\mathcal{G}})}$.

The policy π_{θ} is then updated to maximize the following surrogate objective:

$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E}_{\{o_i\}_{i=1}^{\mathcal{G}} \sim \pi_{\theta_{\text{old}}}(\cdot|I, q)} \left[\frac{1}{\mathcal{G}} \sum_{i=1}^{\mathcal{G}} \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \left(\min \left(r_{i,t}(\theta) \hat{A}_i, \text{clip} \left(r_{i,t}(\theta), 1 - \varepsilon, 1 + \varepsilon \right) \hat{A}_i \right) \right) \right], \quad (1)$$

where $r_{i,t}(\theta) = \frac{\pi_{\theta}(o_{i,t}|q, I, o_{i,<t})}{\pi_{\theta_{\text{old}}}(o_{i,t}|q, I, o_{i,<t})}$ represents probability ratio. Following DAPO Yu et al. (2025), we adopt clip-higher setting, asymmetric clipping hyper-parameters, and exclude the KL penalty.

2.2 Key Findings in Multimodal Reasoning

To investigate the underlying flaw in multimodal reasoning, we conduct an empirical analysis of how internal attentions are allocated and evolved throughout VLMs’ generation in Figure 2. This reveals three key findings that inspire our method. **(1) Text-dominated Inference and Sparse Visual Activation.** First, we analyze the attention allocation among three types of tokens: input image I , input query q , and current generated text $o_{<t}$. As illustrated in Figure 2(a)(b), we observe a significant dominance of textual priors: the model

heavily attends to both generated text history (blue line) and input text query (gray line), leaving visual attention (red line) lower than the sum of these text attentions. Furthermore, visual activation is often sparsely activated in the early stage, manifesting as a brief “glance” the image and followed by long periods of neglect in the later stage.

(2) Temporal Visual Forgetting. Second, we identify a critical correlation between reasoning length and visual activation, which we term *temporal visual forgetting*. The intensity of visual attention (red line) in Figure 2(b) across four benchmarks shows a similar trend: visual activation initially improves and progressively decays as the generation step increases. To further validate the impact of this decay, we statistically compared the attention allocation of correct versus incorrect reasoning samples. Figure 2(c) reveals that correct samples often exhibit higher late/early visual accumulation ratios (Average, 0.680 vs. 0.532), indicating better temporal visual forgetting mitigation.

(3) Inherent Visual Similarity as a Precise Visual Focus Score. Finally, despite the sparsity of visual activation, we observe a promising phenomenon in the bottom of Figure 2(a): when visual tokens are activated (*i.e.*, the peaks in the red line), their attention maps on the original image are highly accurate and semantically grounded based on current base models. This indicates the model can locate relevant information but lacks sustained focus in the following reasoning steps. Therefore, this visual similarity can serve as a precise, intrinsic *visual focus score* to detect and re-weight critical visual tokens, eliminating the need for auxiliary models.

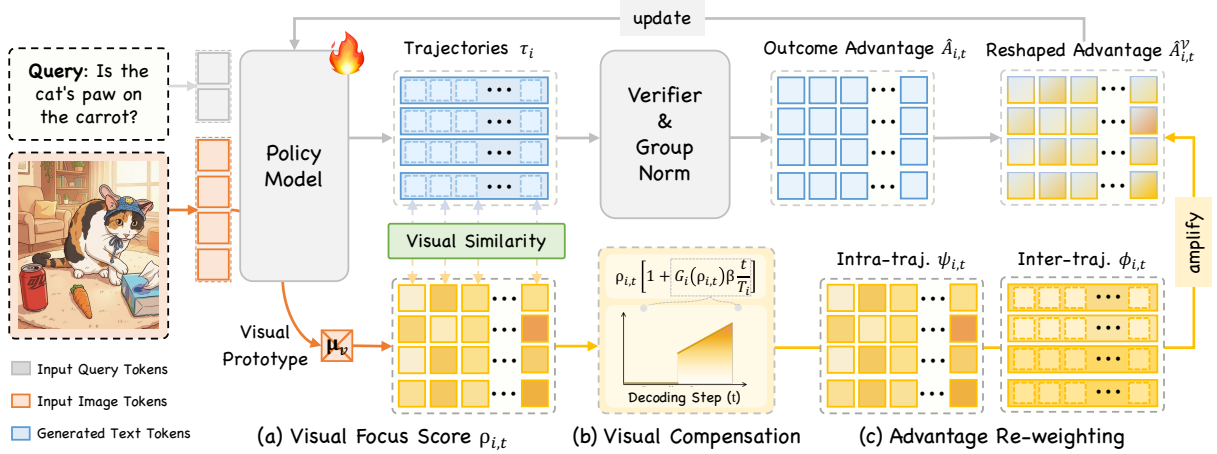


Figure 3: Overview of **Visually-Guided Policy Optimization** framework. Given query and image, (a) VGPO firstly utilizes the intrinsic hidden state similarity between generated tokens and visual prototype to derive a *Visual Focus Score* for visual token localization. (b) Then, *Visual Attention Compensation* (VAC) mechanism leverages this score to re-focus visual tokens, while progressively elevating visual expectations along decoding steps to counteract temporal visual forgetting. (c) Finally, *Dual-grained Advantage Re-weighting* strategy integrates VAC mechanism into intra- and inter-trajectory levels to explicitly incentivize sustained visual faithfulness during policy updates.

3 Methodology

In this section, we present Visually-Guided Policy Optimization (VGPO) framework. As in Figure 3, VGPO is structured into three key components, including Visual Focus Score (Section 3.1), Visual Attention Compensation mechanism (Section 3.2), and Dual-grained Advantage Re-weighting strategy (Section 3.3) to transform intrinsic visual attention into explicit signals to incentivize sustained visual focus throughout the reasoning process.

3.1 Visual Focus Score

To quantify the visual engagement of each generated token without external supervision, we introduce an intrinsic metric based on hidden state similarity, as we discussed in Section 2.2. Let $\mathcal{H}_v = \{\mathbf{h}_k^v\}_{k=1}^{N_v} \subset \mathbb{R}^d$ denote the hidden states of the N_v input image tokens. To capture the global visual semantics, we derive a visual prototype $\boldsymbol{\mu}_v$:

$$\boldsymbol{\mu}_v = \sum_{k=1}^{N_v} \alpha_k \mathbf{h}_k^v, \quad \text{s.t.} \quad \sum \alpha_k = 1, \quad (2)$$

where α_k represents the weight of the k -th image token. For the i -th trajectory τ_i at decoding step t , we measure the cosine semantic similarity $\mathcal{S}(\cdot, \cdot)$ between the current token hidden state $\mathbf{h}_{i,t}$ and the established visual prototype $\boldsymbol{\mu}_v$:

$$\mathcal{S}(\mathbf{h}_{i,t}, \boldsymbol{\mu}_v) = \frac{(\mathbf{h}_{i,t})^\top \boldsymbol{\mu}_v}{\|\mathbf{h}_{i,t}\|_2 \|\boldsymbol{\mu}_v\|_2 + \epsilon}, \quad (3)$$

where ϵ is a smoothing term. Finally, the visual

focus score $\rho_{i,t}$ for the t -th token in trajectory τ_i is obtained via a scaling function to normalize the values to the range $[0, 1]$:

$$\rho_{i,t} = \frac{1}{2} (\mathcal{S}(\mathbf{h}_{i,t}, \boldsymbol{\mu}_v) + 1). \quad (4)$$

3.2 Visual Attention Compensation

However, directly employing the raw visual focus score $\rho_{i,t}$ is suboptimal due to the *temporal visual forgetting* phenomenon observed in Section 2.2. As the reasoning chain extends, intrinsic visual attention naturally decays, causing visually grounded tokens in later steps to often exhibit suppressed focus scores and lead to insufficient optimization signals. To mitigate this, we introduce a *Visual Attention Compensation* mechanism that counteracts this decay by linearly elevating the expectations in the later reasoning step. Specifically, for i -th trajectory τ_i with token length T_i :

$$w_{i,t} = \rho_{i,t} \cdot \left[1 + G_i(\rho_{i,t}) \cdot \beta \cdot \frac{t}{T_i} \right], \quad (5)$$

where the ratio $\frac{t}{T_i}$ indicates a linear strategy to progressively amplify the visual optimization signal as the reasoning step extends. β is a hyper-parameter controlling visual compensation intensity. $G_i(\cdot)$ acts as a visual gate to selectively identify tokens with higher visual similarity (filtering out noisy text tokens) in the later reasoning step. Formally,

$$G_i(\rho_{i,t}) = \begin{cases} 1, & \text{if } t > (1-\gamma)T_i \text{ and } \rho_{i,t} \geq \mathcal{Q}_\kappa(\mathcal{P}_{\text{tail},i}), \\ 0, & \text{otherwise} \end{cases}, \quad (6)$$

where $\mathcal{P}_{\text{tail},i} = \{\rho_{i,k} \mid k > (1-\gamma)T_i\}$ represents the set of score in the tail of trajectory τ_i , and $\mathcal{Q}_\kappa(\cdot)$ returns threshold for the top κ -percent of this set.

3.3 Dual-grained Advantage Re-weighting

To effectively integrate visual attention compensation into the policy optimization, we introduce a Dual-grained Advantage Re-weighting mechanism. This approach modulates the standard advantage function by assessing visual relevance at both the *Intra-trajectory* and *Inter-trajectory* levels.

Intra-trajectory Re-weighting. This component aims to capture the distinctions of individual tokens based on their local visual saliency. To highlight visual significance within a specific reasoning path, we derive an intra-trajectory scaling factor $\psi_{i,t}$. Specifically, we first normalize the raw visual focus scores $w_{i,t}$ via Min-Max scaling to ensure numerical stability. Subsequently, we zero-center these scores relative to the trajectory’s mean, ensuring that only tokens with above-average visual activation are positively incentivized, while effectively suppressing non-visual tokens. Formally,

$$\widehat{w}_{i,t} = \frac{w_{i,t} - \min_k w_{i,k}}{\max_k w_{i,k} - \min_k w_{i,k} + \epsilon}, \quad (7)$$

$$\psi_{i,t} = \widehat{w}_{i,t} - \frac{1}{T_i} \sum_{k=1}^{T_i} \widehat{w}_{i,k}. \quad (8)$$

Inter-trajectory Re-weighting. While the intra-trajectory adjustment addresses local granularity, it is equally crucial to evaluate the global visual focus of generated sequences. Correspondingly, we design the inter-trajectory re-weighting strategy to incentivize the model to prioritize entire generations that exhibit superior aggregate visual score. For each trajectory τ_i , we compute a cumulative visual score s_i , which aggregates the compensated weights over all timesteps. As shown in Eq. 9, this score explicitly accounts for both the base relevance and the time-dependent compensation:

$$s_i = \sum_{t=1}^{T_i} w_{i,t} = \underbrace{\sum_{t=1}^{T_i} \rho_{i,t}}_{\text{Visual Relevance}} + \underbrace{\sum_{t=1}^{T_i} \rho_{i,t} \cdot G_i(\rho_{i,t}) \cdot \left(\beta \cdot \frac{t}{T_i}\right)}_{\text{Late-stage Compensation}}. \quad (9)$$

Finally, to derive the final inter-trajectory scaling factor ϕ_i , we apply group-wise normalization and zero-centering across the rollout group \mathcal{G} :

$$\widehat{s}_i = \frac{s_i - \min_{\tau_j \in \mathcal{G}} s_j}{\max_{\tau_j \in \mathcal{G}} s_j - \min_{\tau_j \in \mathcal{G}} s_j + \epsilon}, \quad (10)$$

$$\phi_i = \widehat{s}_i - \frac{1}{\mathcal{G}} \sum_{j=1}^{\mathcal{G}} \widehat{s}_j. \quad (11)$$

Overall Integration. By synergizing the intra-trajectory re-weighting with the inter-trajectory re-weighting, we formulate the final visual-focus advantage $\widehat{A}_{i,t}$. This mechanism essentially modulates the standard outcome-based advantage \widehat{A}_i (Eq. 1) to explicitly incorporate visual focus signals into the policy update. The final advantage is computed as a multiplicative integration of the base advantage and the dual-grained visual factors:

$$\widehat{A}_{i,t}^{\mathcal{V}} = \widehat{A}_i \cdot (1 + \psi_{i,t}) \cdot (1 + \phi_i). \quad (12)$$

Here, the modulation terms $(1 + \psi_{i,t})$ and $(1 + \phi_i)$ serve as fine-grained and coarse-grained scaling factors, respectively. This composite advantage function reshapes the optimization landscape, ensuring that the policy is driven not merely by the correctness of the final answer, but also by the faithful utilization of visual focus.

4 Experiments

In this section, we present the experimental methodology and the analysis of the results. Further details and results are provided in the Appendix A.

4.1 Experimental Settings

Models and Datasets. Consistent with prior studies (Wang et al., 2025e), we adopt Qwen2.5-VL-series (3B, 7B, and 32B) (Bai et al., 2025b) as our backbone models. We train these models in ViRL39K (Wang et al., 2025a), Geo3K (Lu et al., 2021), and MMK12 (Meng et al., 2025), while validating in MMK12-val (Meng et al., 2025).

Baselines and Evaluations. We conduct a comprehensive evaluation of our approach against a diverse set of leading VLMs at 7B scale. Our baselines include ThinkLite-VL-7B (Wang et al., 2025d), VL-Rethinker-7B (Wang et al., 2025a), MMEureka-7B (Meng et al., 2025), NoisyRollout-7B (Liu et al., 2025a), PAPO_D-7B (Wang et al., 2025e), and VPPO-RL-7B (Huang et al., 2025a). To assess our VGPO, we employ a suite of benchmarks categorized as follows: (1) *General Mathematical & Geometric Reasoning.* We utilize MathVista (Lu et al., 2024), MathVerse (Zhang et al., 2024b), We-Math (Qiao et al., 2025), MMK12 (Meng et al., 2025), GeoMath (Zhou et al.), and Geometry3K (Lu et al., 2021). (2) *Vision-dependent Multimodal Reasoning.* We adopt

Models	General Mathematical & Geometric Reasoning						Avg-Math	Vision-dependent Multimodal Reasoning				Avg-Vision
	MathVista	MathVerse	WeMath	MMK12	GeoMath	Geo3k		LogicVista	Counting	MMMU-Pro	MathVerse _V	
<i>Open-Source General MLLMs (32B-72B)</i>												
Qwen2.5-VL-72B	74.8	65.2 [†]	69.4 [†]	65.9 [†]	53.5 [†]	54.2 [†]	63.8	57.9 [†]	80.5 [†]	51.1	57.6	61.8
Qwen3-VL-32B	83.8	71.5 [†]	75.2	63.1 [†]	45.7 [†]	65.4 [†]	67.5	62.2	95.5 [†]	65.3	76.8	75.0
<i>Open-Source Multimodal Reasoning Models (~7B)</i>												
ThinkLite-VL-7B [†]	73.6	33.7	43.7	50.5	44.6	45.3	48.6	44.3	88.5	32.3	30.9	49.0
VL-Rethinker-7B [†]	65.2	68.8	67.5	66.0	50.1	44.4	60.3	47.0	76.5	37.6	64.8	56.5
MM-Eureka-7B [†]	66.3	66.7	66.6	61.7	49.6	41.1	58.7	47.9	79.0	40.6	63.6	57.8
NoisyRollout-7B [†]	72.9	58.8	64.4	51.2	51.6	55.2*	59.0	47.0	88.5	38.2	51.7	56.3
PAPOD _V -7B [†]	72.3	69.5	69.4	80.8	52.3	<u>48.8</u>	65.5	45.9	<u>89.0</u>	40.2	66.6	60.4
VPPO-RL-7B [†]	70.5	<u>71.1</u>	70.6	81.8	<u>53.2</u>	46.9	<u>65.7</u>	48.8	88.5	40.2	67.6	<u>61.3</u>
Qwen2.5-VL-7B	68.5	40.2	47.8	49.4	51.2	42.9	50.0	45.2	76.5	36.4	36.6	48.7
+ GRPO	70.1	66.7	69.8	72.9	52.8	43.4	62.6 Δ 25.2%	49.4	85.5	38.6	61.8	58.8 Δ 20.7%
+ DAPO	68.7	69.6	<u>70.8</u>	77.0	51.3	45.6	63.8 Δ 27.6%	47.4	85.5	39.0	66.6	59.6 Δ 22.4%
+ VGPO (Ours)	74.1	71.6	72.5	<u>81.5</u>	54.3	45.8	66.6 Δ 33.2%	49.4	95.5	<u>40.5</u>	67.6	63.3 Δ 30.0%

Table 1: Performance comparisons across mathematical and vision-dependent multimodal reasoning benchmarks. “**” indicates NoisyRollout is trained on Geo3K. “†” indicates our reproduction with official checkpoints. Other values are sourced from the official technical report. **Bold** and underlined indicate best and second-best results.

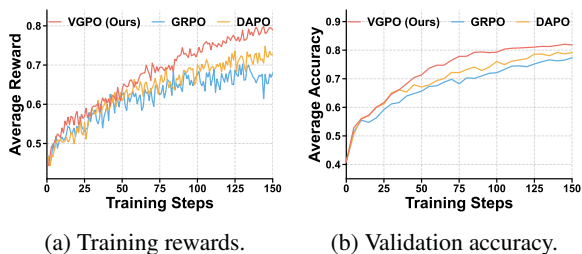


Figure 4: Training dynamics based on Qwen2.5-VL-7B: (a) training rewards and (b) validation accuracy on MMK12 (Meng et al., 2025) across GRPO (Shao et al., 2024), DAPO (Yu et al., 2025), and our VGPO.

LogicVista (Xiao et al., 2024), SuperClevr Counting (Li et al., 2023), MMMU-Pro (Yue et al., 2025), and MathVerse-V (Zhang et al., 2024b). Specifically, we adopt these datasets from PAPO-Eval (Wang et al., 2025e), which provide filtered subsets enabling robust verification without reliance on LLM-as-a-judge (Chen et al., 2024a).

Implementation Details. Following most settings in current methods (Wang et al., 2025e), we set the training epoch to 2, with learning rate of 1×10^{-6} , rollout batch size of 512, and maximum response length of 2,048. During the evaluation phase, we set temperature to 0.0 for all experiments. Regarding the specific hyper-parameters of our VGPO, we employ mean-pooling strategy for visual prototype and set compensation intensity $\beta = 0.3$, position $\gamma = 0.5$, and κ -percent $\kappa = 0.2$.

4.2 Main Results

Superior Performance on Multimodal Reasoning. As shown in Table 1, our proposed VGPO

achieves significant improvements over the base Qwen2.5-VL-7B model, delivering substantial relative gains of 33.2% on general mathematical reasoning and 30.0% on vision-dependent multimodal reasoning tasks. Furthermore, when compared to other advanced reasoning models initialized from the same backbone, VGPO consistently secures the top performance, surpassing previous state-of-the-art results with superior average accuracies of 66.6% in general mathematical benchmarks and 63.3% in vision-dependent tasks. Notably, despite the significant disparity in model scale, our 7B model exhibits highly competitive performance comparable to the much larger Qwen2.5-VL-72B, highlighting the efficiency of our approach in maximizing the potential of smaller-scale models.

Training Dynamics and Generalization. We further analyze the training stability and generalization capabilities in Figure 4. VGPO demonstrates a more stable and efficient learning trajectory, achieving consistently higher training rewards compared to baselines like GRPO and DAPO. This superiority can also extend to generalization, where VGPO maintains a clear lead in validation accuracy throughout the training steps. Unlike competitive methods that may suffer from instability or slower convergence, VGPO establishes a robust optimization path, ensuring steady performance gains.

Scalability across Model Scales and Training Data. To validate the scalability, as summarized in Table 2, experiments ranging from 3B to 32B parameters (Qwen2.5-VL-based) demonstrate that VGPO consistently outperforms baselines, achieving significant relative gains (e.g., 13.8% on math-

Models	General Mathematical & Geometric Reasoning						Avg-Math	Vision-dependent Multimodal Reasoning				Avg-Vision
	MathVista	MathVerse	WeMath	MMK12	GeoMath	Geo3k		LogicVista	Counting	MMMU-Pro	MathVerse _V	
<i>Exploring More Scalable Models (7B → 3B, 32B)</i>												
Qwen2.5-VL-3B	54.0	44.1	50.2	44.1	44.3	27.6	44.1	39.2	58.5	28.6	37.5	41.0
+ DAPO	63.9	57.1	62.9	64.1	47.6	36.3	55.3 $\Delta_{25.4\%}$	43.0	65.5	31.5	53.1	48.3 $\Delta_{17.8\%}$
+ VGPO	65.0	61.4	62.1	71.6	50.2	36.1	57.7 $\Delta_{30.8\%}$	45.9	78.0	32.5	58.1	53.6 $\Delta_{30.7\%}$
Qwen2.5-VL-32B	76.8	62.4	68.8	60.0	53.7	50.8	62.1	55.7	80.5	46.0	58.4	60.2
+ DAPO	69.7	73.0	78.7	82.0	55.2	51.6	68.4 $\Delta_{10.1\%}$	59.1	85.5	47.3	67.4	64.8 $\Delta_{7.6\%}$
+ VGPO	75.3	73.8	78.7	87.3	56.5	52.3	70.7 $\Delta_{13.8\%}$	59.3	90.0	47.3	70.1	66.7 $\Delta_{10.8\%}$
<i>Exploring More Training DataSets (39K → 2.1K, 6.4K)</i>												
Qwen2.5-VL-7B	68.5	40.2	47.8	49.4	51.2	42.9	50.0	45.2	76.5	36.4	36.6	48.7
+ DAPO w/ Geo3K ^{2.1K}	71.6	56.3	63.6	51.1	49.4	52.3	57.4 $\Delta_{14.8\%}$	46.5	85.5	36.8	50.3	54.8 $\Delta_{12.5\%}$
+ VGPO w/ Geo3K ^{2.1K}	72.2	62.5	66.6	53.0	52.6	55.6	60.4 $\Delta_{20.8\%}$	47.7	80.5	37.0	57.8	55.8 $\Delta_{14.6\%}$
+ DAPO w/ MMK12 ^{6.4K}	72.6	66.2	67.5	64.9	51.2	42.3	60.8 $\Delta_{21.6\%}$	47.4	89.5	37.1	61.2	58.8 $\Delta_{20.7\%}$
+ VGPO w/ MMK12 ^{6.4K}	74.1	68.5	68.9	66.0	53.2	43.4	62.4 $\Delta_{24.8\%}$	48.8	90.0	38.6	63.6	60.3 $\Delta_{23.8\%}$

Table 2: Performance comparisons across different model scales and training datasets.

Method	Avg-Math	Avg-Vision	Overall
DAPO Baseline	63.8	59.6	62.2
+ Intra-trajectory	66.1	62.5	64.6
+ Inter-trajectory	65.3	62.0	64.0
+ Intra- & Inter-trajectory	66.6	63.3	65.3

Table 3: Ablation study on the impact of Intra- and Inter-trajectory re-weighting strategies.

Method	Avg-Math	Avg-Vision	Overall
Qwen2.5-VL-7B	50.0	48.7	49.5
+ DAPO	63.8	59.6	62.2
+ DAPO w/ Entropy	65.6	61.9	64.1
+ DAPO w/ KL _{perception}	65.7	61.3	63.9
+ VGPO (Ours)	66.6	63.3	65.3

Table 4: Comparisons with other advantage shaping strategies, including the Entropy-based (Cheng et al., 2025) and KL-based (Huang et al., 2025a) methods.

emathical tasks, 10.8% on visual-dependent tasks) based on 32B model. Additionally, evaluations on different training datasets with varying sizes reveal superior data efficiency, where our method surpasses the DAPO baseline even in data-constrained regimes. These results suggest that incorporating visual focus into the training process is a generalizable strategy to enhance multimodal reasoning.

4.3 Quantitative Analysis

Ablation Study on Re-weighting Strategy. As detailed in Table 3, while the baseline model provides a solid foundation, the addition of both intra- and inter-trajectory mechanisms provides extra gains in reliability and specific tasks. Their combination achieves peak average performance across all metrics, confirming that these two strategies are highly complementary and work synergistically to enhance multimodal reasoning capability.

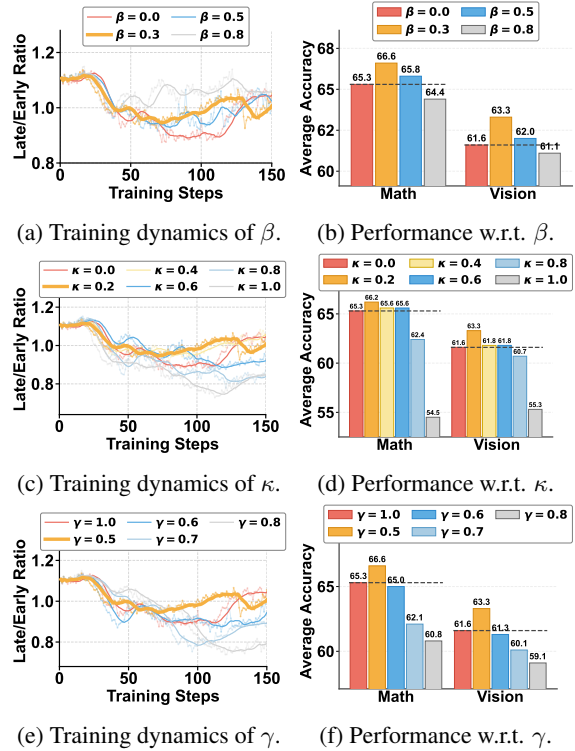


Figure 5: Ablation study of training dynamics of the late/early ratio and performance on hyperparameters β , κ , and γ . The Late/Early Ratio is calculated by dividing the visual attention score of the late reasoning stage by that of the early stage.

Comparison with Existing Advantage Shaping Methods. We evaluate the efficacy of our VGPO against established advantage shaping strategies, specifically regularized derivatives of DAPO. As shown in Table 4, while entropy and KL regularization effectively regulate advantage estimation relative to vanilla DAPO, VGPO achieves superior empirical performance with a peak average accuracy of 63.3%. This advantage is particularly pronounced in vision-centric tasks like LogicVista

Method	Avg-Math	Avg-Vision	Overall
Baseline (DAPO)	63.8	59.6	62.2
+ Step-Function	64.7	60.7	63.1
+ Exponential	65.1	61.0	63.5
+ Linear (Ours)	66.6	63.3	65.3

Table 5: Ablation study on the impact of different compensation schedules.

Method	Avg-Math	Avg-Vision	Overall
Baseline (DAPO)	63.8	59.6	62.2
+ Full-trajectory	53.0	54.2	53.5
+ Late-trajectory (Ours)	66.6	63.3	65.3

Table 6: Ablation study on the impact of Full-trajectory versus Late-trajectory compensation.

and Counting, confirming that strengthening visual focus is a more effective strategy for enhancing multimodal reasoning than generic regularization. **Sensitivity to Hyperparameters.** To evaluate the impact of different parameters on our method and uncover the underlying relation between training dynamics and performance, we conduct a sensitivity analysis on three key hyperparameters: compensation intensity β , threshold κ , and tail ratio γ . As in Figure 5, while the model exhibits performance fluctuations across different settings, it achieves optimal accuracy with specific configurations (*i.e.*, $\beta = 0.3$, $\kappa = 0.2$, and $\gamma = 0.5$). Crucially, by correlating these performance peaks with the evolution of the Late/Early Ratio, we identify a decisive factor for training stability: the model yields the best performance when this ratio converges to or stabilizes near 1. This convergence implies that the optimal hyperparameter configuration facilitates a necessary equilibrium between late and early visual focus. By harmonizing the contributions of these two stages, our method ensures robust performance across diverse reasoning tasks.

Comparison of Different Compensation Schedules. To explore the impact of different visual compensation schedules, we compare our Linear strategy against Step-Function and Exponential schedules, as summarized in Table 5. The Linear strategy consistently achieves the best performance. This is because the Exponential schedule tends to overcorrect by placing too much emphasis on the final tokens (often formatting or calculation), while the Step-Function schedule introduces training instability by abruptly changing the compensation value. The Linear schedule, however, aligns well with the progressive and continuous decay of visual atten-

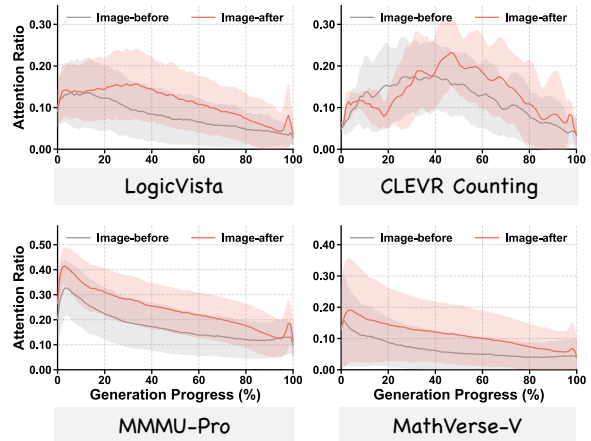


Figure 6: Comparison of the vision attention ratio distribution before and after our VGPO across four visual-dependent multimodal reasoning benchmarks.

tion observed in our empirical analysis. Detailed results and analyses are provided in Appendix B.4. **Comparison with Full-trajectory Compensation.** We justify our choice of the late-trajectory compensation strategy by comparing it with a full-trajectory approach, as shown in Table 6. The full-trajectory compensation significantly drops performance across most benchmarks. This is because current VLMs naturally exhibit high visual attention in the early stages, and enforcing additional visual compensation early can distract the model from parsing the textual query or overemphasize early visual attention. Our late-trajectory design specifically targets the later stages where visual decay occurs, achieving superior results. Detailed comparisons are provided in Appendix B.6. **Analysis of Visual Attention Allocation after our VGPO.** To investigate whether our VGPO achieves better visual activation and temporal visual forgetting mitigation, we compare the visual attention ratio of input image before and after our VGPO in Figure 6. Our VGPO exhibits higher visual attention allocation throughout the entire generation process and sustains better temporal visual forgetting mitigation compared with the baseline. This prolonged visual grounding is critical for accurate long-chain reasoning, ensuring the model does not lose track of visual evidence in later stages.

5 Related Work

Multimodal Reasoning Challenges. Following the milestone of Large Language Models (LLMs) in complex step-by-step reasoning (Wei et al., 2022; Xia et al., 2025; Yang et al., 2025a), the research

focus has naturally shifted toward extending these abilities to Vision-Language Models (VLMs) for broader real-world applications through integrating vision encoders (Vaswani et al., 2017) and large language models (Zhang et al., 2024a; Li et al., 2025d; Zhu et al., 2025; Bai et al., 2025a). However, despite significant architectural advancements, current VLMs largely inherit the text-dominated inductive biases of their LLM backbones, frequently manifesting as insufficient visual faithfulness and severe visual hallucinations remain the primary bottleneck (Bai et al., 2024; Wu et al., 2024; Zhong et al., 2024; Liu et al., 2025c; He et al., 2025a).

Mainstream Strategies for Multimodal Reasoning. To enhance multimodal reasoning, existing methods primarily focus on an RL paradigm that enables the autonomous refinement of reasoning trajectories via rollout sampling (Shao et al., 2024; Huang et al., 2025b; Shen et al., 2025; Li et al., 2025c; He et al., 2026). Following this paradigm, a series of works have explored two kinds of improvements: including (1) training strategies like improving rollout diversity via mixing normal image and noisy augmentation (Liu et al., 2025a) or entropy-based regulation (Cheng et al., 2025) to balance exploration and exploitation; (2) Visual-centric refinement like fine-grained visual enhancement via KL divergence comparison between normal and noisy images (Huang et al., 2025a; Wang et al., 2025e), re-activating specific reasoning paths via introducing specific tokens (Sun et al., 2025a; Li et al., 2025a; Yang et al., 2025d), or verifying intermediate processes as the reward via auxiliary model (Tian et al., 2025). While effective, how to mitigate the need for external models and solely utilize the inherent states of modal abilities to achieve visual enhancement remains an open question.

Visual Perception Methods. Recent concurrent works have also explored enhancing visual perception in multimodal reasoning. PEARL (Zhang et al., 2025) and ViCrit (Wang et al., 2025c) enhance perception through external checklists or synthetic hallucination detection, but they require high computational or data construction costs. On the other hand, SSL4RL (Guo et al., 2025) and VisPlay (He et al., 2025b) utilize unlabeled data yet fail to maintain consistent visual attention during complex reasoning. VGPO overcomes these trade-offs by deriving a Visual Focus Score directly from the internal hidden state. By enforcing continuous visual guidance via an attention compensation mechanism, VGPO

mitigates visual decay and sparse activation. This ensures robust reasoning and cross-domain generalizability, eliminating the reliance on external models or outcome-based proxy tasks found in prior work.

6 Conclusion

In this paper, we explore the critical limitation of insufficient visual activation and its temporal decay during multimodal reasoning. To this end, we propose VGPO, a novel framework designed to alleviate this issue by leveraging intrinsic hidden states to autonomously ground visual focus. By strategically coupling a Visual Attention Compensation mechanism with a dual-grained advantage re-weighting strategy, our method ensures consistent and sustained visual reasoning capability. Comprehensive experiments validate that VGPO significantly enhances visual activation and delivers state-of-the-art performance.

Limitations

While VGPO effectively elevates visual expectations to amplify activation and mitigate forgetting during reasoning, we identify several avenues for future improvement:

- The strategy of progressively elevating visual expectations acts as a robust heuristic that may not represent the globally optimal solution for all reasoning patterns. For instance, in scenarios where the final reasoning steps rely strictly on logical deduction or calculation independent of visual cues, a mandated high visual focus might be less critical. Nevertheless, we regard this work as a pivotal milestone in identifying and addressing the overlooked issue of visual decay. We hope this study serves as a foundation for future research to explore more adaptive mechanisms that can dynamically rectify visual reliance according to the specific context of each reasoning step.
- As VGPO is designed to elicit and sustain vision intrinsic to the VLMs, its performance upper bound is inherently constrained by the representational quality of the visual encoder and projector. If the base model fails to encode critical visual features into the hidden states initially, our re-weighting strategy may not effectively recover this missing information. Consequently, our framework focuses on optimizing the utilization of visual cues rather than enhancing the raw perceptual capabilities of visual encoder itself.

References

- Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhifang Guo, Qidong Huang, Jie Huang, Fei Huang, Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng Li, and 45 others. 2025a. [Qwen3-vl technical report](#).
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, and 1 others. 2025b. [Qwen2.5-vl technical report](#). *arXiv preprint arXiv:2502.13923*.
- Zechen Bai, Pichao Wang, Tianjun Xiao, Tong He, Zongbo Han, Zheng Zhang, and Mike Zheng Shou. 2024. [Hallucination of multimodal large language models: A survey](#). *arXiv preprint arXiv:2404.18930*.
- Dongping Chen, Ruoxi Chen, Shilin Zhang, Yaochen Wang, YINUO Liu, Huichi Zhou, Qihui Zhang, Yao Wan, Pan Zhou, and Lichao Sun. 2024a. [Mllm-as-a-judge: Assessing multimodal llm-as-a-judge with vision-language benchmark](#). In *Forty-first International Conference on Machine Learning*.
- Qiguang Chen, Libo Qin, Jin Zhang, Zhi Chen, Xiao Xu, and Wanxiang Che. 2024b. [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*, pages 8199–8221, Bangkok, Thailand. Association for Computational Linguistics.
- Daixuan Cheng, Shaohan Huang, Xuekai Zhu, Bo Dai, Wayne Xin Zhao, Zhenliang Zhang, and Furu Wei. 2025. [Reasoning with exploration: An entropy perspective](#). *arXiv preprint arXiv:2506.14758*.
- Xiangxiang Chu, Hailang Huang, Xiao Zhang, Fei Wei, and Yong Wang. 2025. [Gpg: A simple and strong reinforcement learning baseline for model reasoning](#). *arXiv preprint arXiv:2504.02546*.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 1 others. 2025. [Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities](#). *arXiv preprint arXiv:2507.06261*.
- Yanqi Dai, Yuxiang Ji, Xiao Zhang, Yong Wang, Xiangxiang Chu, and Zhiwu Lu. 2026. [Harder is better: Boosting mathematical reasoning via difficulty-aware grpo and multi-aspect question reformulation](#). *arXiv preprint arXiv:2601.20614*.
- Alessandro Favero, Luca Zancato, Matthew Trager, Sidharth 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.
- Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A Smith, Wei-Chiu Ma, and Ranjay Krishna. 2024. [Blink: Multimodal large language models can see but not perceive](#). In *European Conference on Computer Vision*, pages 148–166. Springer.
- Xiaojun Guo, Runyu Zhou, Yifei Wang, Qi Zhang, Chenheng Zhang, Stefanie Jegelka, Xiaohan Wang, Jiajun Chai, Guojun Yin, Wei Lin, and 1 others. 2025. [Ssl4rl: Revisiting self-supervised learning as intrinsic reward for visual-language reasoning](#). *arXiv preprint arXiv:2510.16416*.
- Jinghan He, Junfeng Fang, Feng Xiong, Zijun Yao, Fei Shen, Haiyun Guo, Jinqiao Wang, and Tat-Seng Chua. 2026. [Active zero: Self-evolving vision-language models through active environment exploration](#). *arXiv preprint arXiv:2602.11241*.
- Jinghan He, Kuan Zhu, Haiyun Guo, Junfeng Fang, Zhenglin Hua, Yuheng Jia, Ming Tang, Tat-Seng Chua, and Jinqiao Wang. 2025a. [Cracking the code of hallucination in LVLMs with vision-aware head divergence](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3488–3501, Vienna, Austria. Association for Computational Linguistics.
- Yicheng He, Chengsong Huang, Zongxia Li, Jiaxin Huang, and Yonghui Yang. 2025b. [Visplay: Self-evolving vision-language models from images](#). *arXiv preprint arXiv:2511.15661*.
- Siyuan Huang, Xiaoye Qu, Yafu Li, Yun Luo, Zefeng He, Daizong Liu, and Yu Cheng. 2025a. [Spotlight on token perception for multimodal reinforcement learning](#). *arXiv preprint arXiv:2510.09285*.
- Wenxuan Huang, Bohan Jia, Zijie Zhai, Shaosheng Cao, Zheyu Ye, Fei Zhao, Zhe Xu, Yao Hu, and Shaohui Lin. 2025b. [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, and 1 others. 2024. [Gpt-4o system card](#). *arXiv preprint arXiv:2410.21276*.
- Pu Jian, Junhong Wu, Wei Sun, Chen Wang, Shuo Ren, and Jiajun Zhang. 2025. [Look again, think slowly: Enhancing visual reflection in vision-language models](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 9262–9281.
- Bangzheng Li, Ximeng Sun, Jiang Liu, Ze Wang, Jialian Wu, Xiaodong Yu, Hao Chen, Emad Barsoum, Muhao Chen, and Zicheng Liu. 2025a. [Latent visual reasoning](#). *arXiv preprint arXiv:2509.24251*.

- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2025b. LLaVA-onevision: Easy visual task transfer. *Transactions on Machine Learning Research*.
- Renda Li, Hailang Huang, Fei Wei, Feng Xiong, Yong Wang, and Xiangxiang Chu. 2025c. Adacurl: Adaptive curriculum reinforcement learning with invalid sample mitigation and historical revisiting. *arXiv preprint arXiv:2511.09478*.
- Yunxin Li, Zhenyu Liu, Zitao Li, Xuanyu Zhang, Zhenran Xu, Xinyu Chen, Haoyuan Shi, Shenyuan Jiang, Xintong Wang, Jifang Wang, and 1 others. 2025d. Perception, reason, think, and plan: A survey on large multimodal reasoning models. *arXiv preprint arXiv:2505.04921*.
- 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.
- Xiangyan Liu, Jinjie Ni, Zijian Wu, Chao Du, Longxu Dou, Haonan Wang, Tianyu Pang, and Michael Qizhe Shieh. 2025a. Noisyrollout: Reinforcing visual reasoning with data augmentation. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.
- Zhining Liu, Ziyi Chen, Hui Liu, Chen Luo, Xianfeng Tang, Suhang Wang, Joy Zeng, Zhenwei Dai, Zhan Shi, Tianxin Wei, and 1 others. 2025b. Seeing but not believing: Probing the disconnect between visual attention and answer correctness in vlms. *arXiv preprint arXiv:2510.17771*.
- Zujing Liu, Junwen Pan, Qi She, Yuan Gao, and Guisong Xia. 2025c. On the faithfulness of visual thinking: Measurement and enhancement. *arXiv preprint arXiv:2510.23482*.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2024. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *The Twelfth International Conference on Learning Representations*.
- 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*, pages 6774–6786, Online. Association for Computational Linguistics.
- Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Y. Tang, Manan Roongta, Colin Cai, Jeffrey Luo, Li Erran Li, Raluca Ada Popa, and Ion Stoica. 2025. Deepscaler: Surpassing o1-preview with a 1.5b model by scaling rl.
- Fanqing Meng, Lingxiao Du, Zongkai Liu, Zhixiang Zhou, Quanfeng Lu, Daocheng Fu, Tiancheng Han, Botian Shi, Wenhai Wang, Junjun He, and 1 others. 2025. Mm-eureka: Exploring the frontiers of multimodal reasoning with rule-based reinforcement learning. *arXiv preprint arXiv:2503.07365*.
- Runqi Qiao, Qiuna Tan, Guanting Dong, MinhuiWu MinhuiWu, Chong Sun, Xiaoshuai Song, Jiapeng Wang, Zhuoma GongQue, Shanglin Lei, YiFan Zhang, Zhe Wei, Miaoquan Zhang, Runfeng Qiao, Xiao Zong, Yida Xu, Peiqing Yang, Zhimin Bao, Muxi Diao, Chen Li, and Honggang Zhang. 2025. [We-math: Does your large multimodal model achieve human-like mathematical reasoning?](#) In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 20023–20070, Vienna, Austria. Association for Computational Linguistics.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Jay Shah, Ganesh Bikshandi, Ying Zhang, Vijay Thakkar, Pradeep Ramani, and Tri Dao. 2024. Flashattention-3: Fast and accurate attention with asynchrony and low-precision. *Advances in Neural Information Processing Systems*, 37:68658–68685.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, and 1 others. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.
- Haozhan Shen, Peng Liu, Jingcheng Li, Chunxin Fang, Yibo Ma, Jiajia Liao, Qiaoli Shen, Zilun Zhang, Kangjia Zhao, Qianqian Zhang, and 1 others. 2025. Vlm-r1: A stable and generalizable r1-style large vision-language model. *arXiv preprint arXiv:2504.07615*.
- Guohao Sun, Hang Hua, Jian Wang, Jiebo Luo, So-hail Dianat, Majid Rabbani, Raghuvver Rao, and Zhiqiang Tao. 2025a. Latent chain-of-thought for visual reasoning. *arXiv preprint arXiv:2510.23925*.
- Hai-Long Sun, Zhun Sun, Houwen Peng, and Han-Jia Ye. 2025b. Mitigating visual forgetting via take-along visual conditioning for multi-modal long cot reasoning. *arXiv preprint arXiv:2503.13360*.
- Kai Sun, Yushi Bai, Ji Qi, Lei Hou, and Juanzi Li. 2024. [MM-MATH: Advancing multimodal math evaluation with process evaluation and fine-grained classification](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 1358–1375, Miami, Florida, USA. Association for Computational Linguistics.

- Xinyu Tian, Shu Zou, Zhaoyuan Yang, Mengqi He, Fabian Waschkowski, Lukas Wesemann, Peter Tu, and Jing Zhang. 2025. More thought, less accuracy? on the dual nature of reasoning in vision-language models. *arXiv preprint arXiv:2509.25848*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in Neural Information Processing Systems*, 30.
- Haozhe Wang, Chao Qu, Zuming Huang, Wei Chu, Fangzhen Lin, and Wenhui Chen. 2025a. V1-rethinker: Incentivizing self-reflection of vision-language models with reinforcement learning. *arXiv preprint arXiv:2504.08837*.
- Jiayu Wang, Yifei Ming, Zhenmei Shi, Vibhav Vineet, Xin Wang, Sharon Li, and Neel Joshi. 2024. Is a picture worth a thousand words? delving into spatial reasoning for vision language models. *Advances in Neural Information Processing Systems*, 37:75392–75421.
- Peijie Wang, Zhong-Zhi Li, Fei Yin, Dekang Ran, and Cheng-Lin Liu. 2025b. Mv-math: Evaluating multimodal math reasoning in multi-visual contexts. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 19541–19551.
- Xiyao Wang, Zhengyuan Yang, Chao Feng, Yongyuan Liang, Yuhang Zhou, Xiaoyu Liu, Ziyi Zang, Ming Li, Chung-Ching Lin, Kevin Lin, and 1 others. 2025c. Vicrit: A verifiable reinforcement learning proxy task for visual perception in vlms. *arXiv preprint arXiv:2506.10128*.
- Xiyao Wang, Zhengyuan Yang, Chao Feng, Hongjin Lu, Linjie Li, Chung-Ching Lin, Kevin Lin, Furong Huang, and Lijuan Wang. 2025d. Sota with less: Mcts-guided sample selection for data-efficient visual reasoning self-improvement. *arXiv preprint arXiv:2504.07934*.
- Zhenhailong Wang, Xuehang Guo, Sofia Stoica, Haiyang Xu, Hongru Wang, Hyeonjeong Ha, Xiushi Chen, Yangyi Chen, Ming Yan, Fei Huang, and 1 others. 2025e. Perception-aware policy optimization for multimodal reasoning. *arXiv preprint arXiv:2507.06448*.
- 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.
- Junyi Wu, Weitai Kang, Hao Tang, Yuan Hong, and Yan Yan. 2024. On the faithfulness of vision transformer explanations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10936–10945.
- Yu Xia, Rui Wang, Xu Liu, Mingyan Li, Tong Yu, Xiang Chen, Julian McAuley, and Shuai Li. 2025. Beyond chain-of-thought: A survey of chain-of-x paradigms for llms. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 10795–10809.
- Yijia Xiao, Edward Sun, Tianyu Liu, and Wei Wang. 2024. Logicvista: Multimodal llm logical reasoning benchmark in visual contexts. *arXiv preprint arXiv:2407.04973*.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 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*.
- Yi Yang, Xiaoxuan He, Hongkun Pan, Xiyang Jiang, Yan Deng, Xingtao Yang, Haoyu Lu, Dacheng Yin, Fengyun Rao, Minfeng Zhu, Bo Zhang, and Wei Chen. 2025c. R1-onevision: Advancing generalized multimodal reasoning through cross-modal formalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2376–2385.
- Zeyuan Yang, Xueyang Yu, Delin Chen, Maohao Shen, and Chuang Gan. 2025d. Machine mental imagery: Empower multimodal reasoning with latent visual tokens. *arXiv preprint arXiv:2506.17218*.
- Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gao-hong Liu, Juncai Liu, LingJun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, and 17 others. 2025. DAPO: An open-source LLM reinforcement learning system at scale. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.
- Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang, Kai Zhang, Shengbang Tong, Yuxuan Sun, Botao Yu, Ge Zhang, Huan Sun, Yu Su, Wenhui Chen, and Graham Neubig. 2025. MMMU-pro: A more robust multi-discipline multimodal understanding benchmark. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15134–15186, Vienna, Austria. Association for Computational Linguistics.
- Chi Zhang, Haibo Qiu, Qiming Zhang, Yufei Xu, Zhixiong Zeng, Siqi Yang, Peng Shi, Lin Ma, and Jing Zhang. 2025. Perceptual-evidence anchored reinforced learning for multimodal reasoning. *arXiv preprint arXiv:2511.18437*.
- Jingyi Zhang, Jiaxing Huang, Sheng Jin, and Shijian Lu. 2024a. Vision-language models for vision tasks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(8):5625–5644.

Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Yu Qiao, Peng Gao, and Hongsheng Li. 2024b. Multi-modal llm truly see the diagrams in visual math problems? In *Computer Vision – ECCV 2024: 18th European Conference, Milan, Italy, September 29–October 4, 2024, Proceedings, Part VIII*, page 169–186, Berlin, Heidelberg. Springer-Verlag.

Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong Liu, Rui Men, An Yang, and 1 others. 2025. Group sequence policy optimization. *arXiv preprint arXiv:2507.18071*.

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. *arXiv preprint arXiv:2407.00569*.

Yue Zhou, Litong Feng, Mengcheng Lan, Yiping Ke, Xue Jiang, and Wayne Zhang. Geomath: A benchmark for multimodal mathematical reasoning in remote sensing.

Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, and 1 others. 2025. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv preprint arXiv:2504.10479*.

A Experimental Settings

A.1 Training Datasets

We mainly conduct our training experiments on the ViRL39K (Wang et al., 2025a) dataset, a verifiable instruction-tuning benchmark designed for vision-language reasoning. This dataset comprises approximately 39,000 high-quality samples synthesized from proprietary collections and various upstream datasets, including Llava-OneVision (Li et al., 2025b), R1-OneVision (Yang et al., 2025c), MM-Eureka (Meng et al., 2025), MM-Math (Sun et al., 2024), M3CoT (Chen et al., 2024b), DeepScaleR (Luo et al., 2025), and MV-Math (Wang et al., 2025b). ViRL39K is distinguished by a rigorous filtering pipeline that removes unverifiable queries, thereby ensuring extensive coverage of topics ranging from general chart interpretation to complex STEM problem-solving.

A.2 Baselines

- **ThinkLite-VL** (Wang et al., 2025d). ThinkLite-VL bypasses knowledge distillation in favor of an MCTS-guided selection strategy that curates a compact, high-quality dataset based on sample difficulty. By applying reinforcement fine-tuning to these challenging instances, the method attains superior visual reasoning results while reducing data requirements by an order of magnitude.
- **VL-Rethinker** (Wang et al., 2025a). VL-Rethinker enhances VLMs’ slow-thinking capabilities through a distillation-free reinforcement learning framework that integrates Selective Sample Replay (SSR) to stabilize training and Forced Rethinking to incentivize self-reflection. By rehearsing high-value experiences and triggering explicit verification processes, the method effectively mitigates GRPO limitations and fosters the internalization of deliberate reasoning patterns.
- **MM-Eureka** (Meng et al., 2025). MMEureka advances multimodal reasoning by combining the MMK12 dataset with a Qwen2.5-VL-based GRPO pipeline that employs online filtering to ensure gradient efficacy. To address training instability in large models, it utilizes a two-stage strategy involving initial training on MMK12 without KL divergence followed by fine-tuning on Geo3k with KL regularization.
- **NoisyRollout** (Liu et al., 2025a). NoisyRollout enhances visual reasoning and robustness by

computing GRPO advantages from hybrid clean and distorted trajectories, restricting policy optimization to uncorrupted inputs. Furthermore, it incorporates a noise annealing schedule to gradually diminish distortion, thereby effectively balancing exploration and stability.

- **PAPO (Wang et al., 2025e)**. PAPO addresses perception errors by introducing an Implicit Perception Loss that maximizes the KL divergence between policy outputs on original and masked inputs to enforce visual grounding. To ensure training stability, it employs a Double Entropy Loss regularization, offering a model-agnostic framework that jointly optimizes perception and reasoning without requiring additional annotations or teacher models.
- **VPPO (Huang et al., 2025a)**. VPPO integrates visual perception into RL by quantifying token-level visual dependency via KL divergence between policy distributions on original and perturbed image inputs. Based on this metric, it employs advantage modulation and sparse gradient masking to concentrate optimization exclusively on pivotal visual tokens.

A.3 Evaluation Benchmarks

Mathematical & Geometric Reasoning:

- **MathVista (Lu et al., 2024)**. MathVista is designed to integrate challenges across diverse mathematical and visual domains, necessitating both fine-grained visual perception and compositional reasoning. The benchmark consists of 6,141 examples amassed from a wide array of sources, including 28 existing multimodal mathematics datasets and 3 constructed datasets.
- **MathVerse (Zhang et al., 2024b)**. MathVerse is designed to mitigate the textual bias in existing benchmarks by curating 2,612 high-quality problems that necessitate genuine visual interpretation. Each problem is manually reformulated into six distinct versions with varying multimodal information density, resulting in a comprehensive corpus of approximately 15,000 test samples.
- **We-Math (Qiao et al., 2025)**. We-Math transcends conventional accuracy metrics to scrutinize the underlying problem-solving mechanisms of LMMs through a hierarchical framework of

Parameter	Configuration
General Settings	
Rollout Number	8
Learning Rate	1e-6
Global Batch Size	128
Rollout Batch Size	512
Val Batch Size	1024
Max Prompt Length	4096
Max Response Length	2048
Reward	Binary Accuracy
GPU Usage	8×H20, 96G Memory
Qwen2.5-VL-3B/7B on ViRL39K (39K samples)	
Training Episodes	2
Total Optimization Steps	150
Qwen2.5-VL-7B on Geo3K (2.1K samples)	
Training Episodes	15
Total Optimization Steps	60
Qwen2.5-VL-7B on MMK12 (6.4K samples)	
Training Episodes	12
Total Optimization Steps	120
Qwen2.5-VL-32B on ViRL39K (39K samples)	
Training Episodes	2
Total Optimization Steps	150
GPU Usage	32×H20, 96G Memory

Table 7: Experimental hyperparameter configurations of our VGPO across 3B, 7B, and 32B-based backbones and different scalable training datasets.

6,500 visual tasks spanning 67 distinct knowledge concepts. By employing a decompositional methodology that reduces composite problems into atomic sub-tasks, the benchmark facilitates a granular four-dimensional evaluation distinguishing among Insufficient Knowledge, Inadequate Generalization, Complete Mastery, and Rote Memorization.

- **MMK12 (Meng et al., 2025)**. MMK12 benchmark constitutes a comprehensive evaluation framework designed to assess multimodal reasoning capabilities at the K-12 educational level. Comprising 2,000 high-quality instances, the benchmark is stratified across four scientific disciplines, specifically Mathematics, Physics, Chemistry, and Biology, with 500 multimodal multiple-choice questions dedicated to each field.
- **GeoMath (Zhou et al.)**. Designed to address scarcity of reasoning data in Earth observation, GeoMath comprises 3,773 high-quality queries

derived from aerial imagery, encompassing six distinct mathematical subjects across 20 sub-topics. The visual data are acquired via proprietary drone flights capturing a wide range of altitudes and viewing angles.

- **Geometry3K (Lu et al., 2021).** Geometry3K comprises 3,002 geometry problems enriched with dense annotations in formal language, serving as a challenging task for abstract problem understanding and symbolic reasoning, requiring axiomatic knowledge.

Vision-Dependent Multimodal Reasoning:

- **LogicVista (Xiao et al., 2024).** LogicVista is an evaluation benchmark designed to assess the integrated logical reasoning abilities of VLMs in visual contexts. Addressing the limitations of prior work in systematically evaluating logical proficiency. The dataset consists of 448 multiple-choice questions, each densely annotated with the correct answer and a human-written explanation.
- **Super-CLEVR (Li et al., 2023).** Super-CLEVR constitutes a virtual benchmark designed to rigorously assess the out-of-distribution robustness and domain generalization of VQA models. Through controllable data synthesis, it decouples inherent multi-modal factors to facilitate the isolated analysis of four specific domain shifts: visual complexity, question redundancy, concept distribution, and concept compositionality.
- **MMMU-Pro (Yue et al., 2025).** MMMU-Pro advances the MMMU benchmark by employing a stringent curation process that eliminates text-solvable queries and expands candidate options to rigorously evaluate intrinsic multimodal reasoning. Additionally, it introduces a vision-only paradigm where inquiries are embedded within images, thereby necessitating integrated visual and textual interpretation.

B More Experimental Results

B.1 Detailed Results of Re-weighting Strategy

Table 8 provides the whole performance breakdown across all datasets, complementing the aggregated results in the main text.

Effectiveness of Intra-trajectory Strategy. This mechanism significantly enhances general mathematical reasoning. It yields substantial gains on

MathVista (+5.4%) and MMK12 (+4.7%), demonstrating that reinforcing self-consistency within individual reasoning paths is critical to re-focus visual activation for multimodal reasoning

Impact of Inter-trajectory Strategy. The Inter-trajectory module excels in geometry and visual counting tasks. It achieves a peak score of 54.9% on GeoMath and matches the top performance on the Counting dataset. These results suggest that cross-path comparison effectively mitigates visual hallucinations in object-centric scenarios.

Synergy of the Two Strategies. Combining both modules yields the most robust performance across mathematical and vision-dependent categories. Notably on WeMath, the integrated approach reverses a slight decline caused by Intra-module alone, achieving a superior score of 72.5%. Furthermore, the combined strategy secures top results on challenging benchmarks like LogicVista and MMMU-Pro. This confirms a synergistic effect, where the Intra-module refines logical consistency while the Inter-module broadens the reasoning scope.

B.2 Detailed Results of Advantage Shaping Methods

Table 9 presents a fine-grained comparison of VGPO against established advantage shaping strategies. Consistent with the main results, VGPO exhibits superior performance across both mathematical and vision-centric benchmarks.

Robustness in General Reasoning. In contrast to generic regularization methods like Entropy and KL, which rely on implicit constraints on policy divergence, VGPO explicitly leverages visual feedback. This mechanism translates into substantial gains in General Mathematical & Geometric Reasoning. VGPO secures the highest scores on 4 out of 6 datasets, establishing significant leads on MathVista and MathVerse. This demonstrates that visually-guided optimization not only enhances reasoning stability but does so while preserving the general problem-solving proficiency.

Efficacy in Visual Perception. Our detailed analysis highlights the pivotal role of VGPO in Vision-dependent Multimodal Reasoning. The performance gap is particularly evident in the *Counting* dataset, a task requiring rigorous visual grounding. VGPO achieves 95.5% accuracy, surpassing the Entropy and KL baselines by a significant margin. Similar dominance is observed on LogicVista. These results validate our hypothesis that explicitly reinforcing visual attention through VGPO is far

Models	General Mathematical & Geometric Reasoning						Avg-Math	Vision-dependent Multimodal Reasoning				Avg-Vision
	MathVista	MathVerse	WeMath	MMK12	GeoMath	Geo3k		LogicVista	Counting	MMMU-Pro	MathVerse _V	
DAPO Baseline	68.7	69.6	70.8	77.0	51.3	45.6	63.8	47.4	85.5	39.0	66.6	59.6
+ Intra-trajectory	74.1	71.1	69.8	81.7	54.3	45.6	66.1	48.1	95.0	40.2	66.5	62.5
+ Inter-trajectory	72.1	68.8	69.9	80.7	54.9	45.3	65.3	47.9	95.5	40.2	64.6	62.0
+ Intra- & Inter-trajectory	74.1	71.6	72.5	81.5	54.3	45.8	66.6	49.4	95.5	40.5	67.6	63.3

Table 8: The detailed ablation study on the impact of Intra- and Inter-trajectory re-weighting strategies.

Models	General Mathematical & Geometric Reasoning						Avg-Math	Vision-dependent Multimodal Reasoning				Avg-Vision
	MathVista	MathVerse	WeMath	MMK12	GeoMath	Geo3k		LogicVista	Counting	MMMU-Pro	MathVerse _V	
Qwen2.5-VL-7B	68.5	40.2	47.8	49.4	51.2	42.9	50.0	45.2	76.5	36.4	36.6	48.7
+ DAPO	68.7	69.6	70.8	77.0	51.3	45.6	63.8	47.4	85.5	39.0	66.6	59.6
+ DAPO w/ Entropy	70.2	72.1	70.8	80.9	53.9	45.6	65.6	48.3	90.5	40.1	68.7	61.9
+ DAPO w/ KL-perception	70.5	71.1	70.6	81.8	53.2	46.9	65.7	48.8	88.5	40.2	67.6	61.3
+ VGPO (Ours)	74.1	71.6	72.5	81.5	54.3	45.8	66.6	49.4	95.5	40.5	67.6	63.3

Table 9: The detailed comparisons with other advantage shaping strategies, including the Entropy-based (Cheng et al., 2025) and KL-based (Huang et al., 2025a) methods.

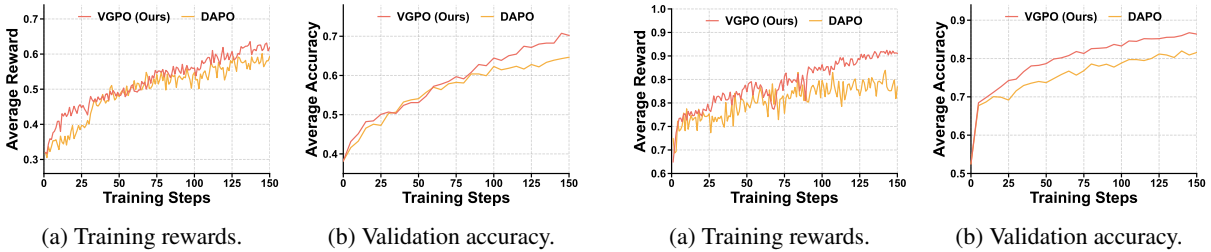


Figure 7: Training dynamics of Qwen2.5-VL-3B: (a) training rewards and (b) validation accuracy on MMK12 (Meng et al., 2025) across DAPO (Yu et al., 2025), and our VGPO.

Figure 8: Training dynamics of Qwen2.5-VL-32B: (a) training rewards and (b) validation accuracy on MMK12 (Meng et al., 2025) across DAPO (Yu et al., 2025), and our VGPO.

more effective for complex multimodal tasks than relying on implicit regularization techniques.

B.3 Analysis of Training Dynamics

To further investigate the convergence, stability, and scalability of our VGPO, we visualize the training dynamics of training rewards and validation accuracy on MMK12 dataset for both Qwen2.5-VL-3B (Figure 7) and Qwen2.5-VL-32B (Figure 8).

Dynamics on Qwen2.5-VL-3B. As shown in Figure 7, VGPO demonstrates a clear performance advantage. In terms of training rewards, VGPO maintains a consistently higher reward trajectory compared to DAPO throughout the training steps. More importantly, regarding validation accuracy, VGPO achieves a higher final accuracy, which indicates superior generalization capabilities.

Dynamics on Qwen2.5-VL-32B. The advantages of VGPO are even more pronounced on the larger 32B model, particularly regarding stability. Figure 8(a) highlights that DAPO suffers from severe oscillation during training, whereas VGPO exhibits a smooth and steady increase in rewards. This

suggests that our visual-guided strategy effectively reduces variance and stabilizes the optimization process for large-scale models. Consistent with the rewards, Figure 8(b) shows that VGPO consistently outperforms DAPO in validation accuracy, establishing a significant margin by the end of training.

B.4 Detailed Results of Different Compensation Schedules

Table 10 provides the detailed performance breakdown across all datasets, complementing the aggregated results in the main text.

We adopt the Linear Compensation Strategy since it aligns well with the progressive, continuous, and nearly linear decay of visual attention observed in our empirical analysis across four visual-dependent tasks (LogicVista, CLEVR Counting, MMMU-Pro, MathVerse-V). To explore potential alternatives, we compare our Linear strategy against Step-Function and Exponential schedules. As shown in Table 10, all these schedules exhibit higher performance compared to the baseline. However, the Linear strategy consistently achieves

Compensation Schedules	General Mathematical & Geometric Reasoning						Vision-dependent Multimodal Reasoning					
	MathVista	MathVerse	WeMath	MMK12	GeoMath	Geo3k	Avg-Math	LogicVista	Counting	MMMU-Pro	MathVerse _V	Avg-Vision
Baseline (DAPO)	68.7	69.6	70.8	77.0	51.3	45.6	63.8	47.4	85.5	39.0	66.6	59.6
w/ Step-Function	69.1	69.2	69.4	80.3	52.3	47.6	64.7	47.5	91.5	39.2	64.7	60.7
w/ Exponential	70.6	69.9	69.3	81.4	54.2	44.9	65.1	49.0	90.5	38.8	65.5	61.0
w/ Linear (Ours)	74.1	71.6	72.5	81.5	54.3	45.8	66.6	49.4	95.5	40.5	67.6	63.3

Table 10: Results of General Mathematical & Geometric Reasoning and Vision-dependent Multimodal Reasoning tasks with different compensation schedules. *Note:* For the Exponential schedule, we test powers of 1.0 and 2.0 (reporting the best result from 2.0). We test 1.0 for the Step-wise schedule.

Models	General Mathematical & Geometric Reasoning						Vision-dependent Multimodal Reasoning					
	MathVista	MathVerse	WeMath	MMK12	GeoMath	Geo3k	Avg-Math	LogicVista	Counting	MMMU-Pro	MathVerse _V	Avg-Vision
Baseline (DAPO, Qwen2.5-VL-7B)	68.7	69.6	70.8	77.0	51.3	45.6	63.8	47.4	85.5	39.0	66.6	59.6
+VGPO (Ours)	74.1	71.6	72.5	81.5	54.3	45.8	66.6 ^{△4.4%}	49.4	95.5	40.5	67.6	63.3 ^{△6.2%}
Baseline (DAPO, Qwen2.5-VL-3B)	63.9	57.1	62.9	64.1	47.6	36.3	55.3	43.0	65.5	31.5	53.1	48.3
+VGPO (Ours)	65.0	61.4	62.1	71.6	50.2	36.1	57.7 ^{△4.3%}	45.9	78.0	32.5	58.1	53.6 ^{△11.0%}
Baseline (DAPO, Qwen2-VL-2B)	50.8	33.2	36.3	42.3	34.8	8.7	34.4	30.0	66.5	19.8	31.6	37.0
+VGPO (Ours)	51.1	37.8	42.5	44.2	37.0	16.8	38.2 ^{△11.1%}	30.2	83.5	21.9	36.4	43.0 ^{△16.2%}

Table 11: Performance comparison of stronger to weaker visual encoders on General Mathematical & Geometric Reasoning and Vision-dependent Multimodal Reasoning tasks.

the best performance. We attribute these results to the following reasons:

- **Exponential schedule (over-correction):** This schedule puts most of the bonus on the very last tokens, which often correspond to calculation or formatting rather than visual grounding. Overweighting those tokens dilutes the benefit of compensating truly visual-dependent steps and can hurt accuracy compared to the more balanced linear schedule.
- **Step-Function schedule (training instability):** This schedule makes position compensation jump from zero to full compensation at a single threshold. However, it does not match our observed progressive forgetting: visual attention decays gradually over the trajectory, so we use a linear schedule that increases compensation smoothly rather than a single step.

B.5 Robustness on Weaker Visual Encoders

To investigate the robustness of VGPO across different visual encoder capacities, we evaluate our method using Qwen2-VL-2B-Instruct (Bai et al., 2025b), which employs a significantly smaller and weaker visual encoder compared to the Qwen2.5-VL series. As shown in Table 11, despite the inherently constrained perceptual capacity of the 2B baseline, applying VGPO achieves even better relative improvements (+11.1% on Math and +16.2% on vision-dependent tasks). This empirical evidence confirms that VGPO demonstrates strong

generalizability and robustness, effectively resolving the policy-level bottleneck of temporal visual forgetting regardless of the base visual encoder’s strength.

B.6 Detailed Results of Full-trajectory Compensation

Table 12 provides the detailed performance breakdown across all datasets, complementing the aggregated results in the main text. As shown in Table 12, we find that the full-trajectory compensation strategy significantly drops performance on the majority of benchmarks. Moreover, the Late/Early Visual Activation Ratio exhibits a lower value (0.8591) than our VGPO (1.0696). This indicates that excessive elevation of early-trajectory compensation often overemphasizes early visual attention and leads to higher visual forgetting with lower performance.

B.7 Prompt Template

Template for Multi-modal Reasoning

SYSTEM

You are a helpful assistant.

USER

<question>

You first think through the reasoning process as an internal monologue, enclosed within <think> </think> tags. Then, provide your final answer enclosed within \boxed{ }.

Compensation Schedules	General Mathematical & Geometric Reasoning						Vision-dependent Multimodal Reasoning					
	MathVista	MathVerse	WeMath	MMK12	GeoMath	Geo3k	Avg-Math	LogicVista	Counting	MMMU-Pro	MathVerse _V	Avg-Vision
Baseline (DAPO)	68.7	69.6	70.8	77.0	51.3	45.6	63.8	47.4	85.5	39.0	66.6	59.6
Full-trajectory	71.4	50.3	53.8	60.1	44.9	37.3	53.0	40.5	92.0	36.4	47.8	54.2
Late-trajectory (Ours)	74.1	71.6	72.5	81.5	54.3	45.8	66.6	49.4	95.5	40.5	67.6	63.3

Table 12: Performance comparison of Full-trajectory compensation versus our Late-trajectory compensation strategy.

C Discussion about Visual Focus Score vs. Attention Weights

To validate the theoretical grounding of our proposed Visual Focus Score, we conduct an empirical correlation analysis between the Visual Focus Score and actual attention weights.

Empirical validation of correlation with Attention Weights. We select the validation dataset (MMK12-val (Meng et al., 2025)) for this analysis. At each generation step of each sample:

- For the **Visual Focus Score**, we compute the cosine similarity between the final hidden state of the t -th generated token and the mean-pooled visual prototype.
- For the **Attention Weights**, we compute the sum of the attention probabilities from the t -th generated token to all image token positions, averaged across all attention heads in the final layer.

We then calculate the Pearson correlation coefficient between these two step-wise sequences ($t = 1, 2, \dots, T$). We observe a strong positive correlation ($r \approx 0.67$), confirming that our similarity-based Visual Focus Score closely aligns with the actual visual attention mechanism.

Why use hidden-state similarity rather than actual attention weights? During implementation, extracting actual attention weights requires setting `output_attentions=True`. This causes the model to fall back to the less efficient eager attention implementation rather than Flash Attention (Shah et al., 2024), losing Flash Attention’s speed and memory benefits (typically incurring a 20%–30% additional computational cost). In contrast, our hidden-state-based implementation only requires `output_hidden_states=True`, which adds approximately 5%–10% overhead and remains fully compatible with Flash Attention, making it highly efficient for large-scale training.

D The Use of LLMs

This manuscript utilizes LLMs exclusively for the purpose of linguistic refinement and stylistic improvement to enhance clarity. The LLMs serve

strictly as a writing aid and contribute no intellectual content to the research conceptualization, methodological design, data execution, analysis, or interpretation. These core scientific components remain the sole work of the authors. All suggestions provided by the LLMs are critically reviewed and selectively adopted to ensure the accuracy and integrity of the text. The authors retain full accountability for the content and findings presented herein, and the LLM does not qualify for authorship or contributorship.

E Failure Case Analysis

To provide deeper insights into the boundaries of our proposed VGPO, we present both qualitative and quantitative analyses of failure cases where increased visual focus does not translate to improved accuracy.

Qualitative insight (hallucination): Our error analysis reveals that if the base visual encoder fundamentally misinterprets a feature (*e.g.*, misrepresenting a blurry “6” as an “8” due to low resolution), VGPO can actually reinforce this error. By incentivizing the LLM to strongly attend to a flawed visual embedding, the model produces confidently incorrect answers, rather than falling back on common-sense priors that might have occasionally guessed correctly.

Quantitative insight: As shown in our original analysis in Figure 2(c), while correct samples generally show higher visual accumulation, there is a subset of incorrect samples that also exhibit high visual focus scores. This quantitatively demonstrates that “looking closely” is a necessary but not sufficient condition for correctness. The model must also possess the correct world knowledge and reasoning logic to interpret the visual cues properly. Here we provide two typical failure cases to illustrate these points.

Failure Case 1: Visual Perception Error (GeoMath, id: geomath-2884)



Question: How many vertical bars are visible in the image, and what is the total number of elements (cow and bars) present in the image?

Ground Truth: 8

Model Prediction: To determine the number of vertical bars and the total number of elements in the image, let's break it down step by step.

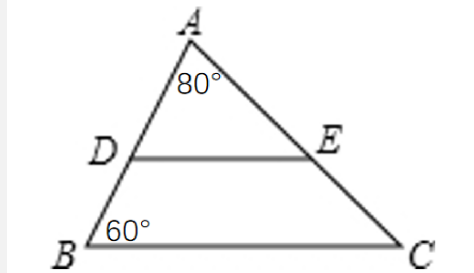
1. Count the vertical bars: The image shows a series of vertical bars that form a fence. By visually inspecting the image, it appears there are at least 10 vertical bars. (*visual counting error*)
2. Count the elements: There is one cow visible in the image. There are 10 vertical bars.
3. Calculate the total number of elements: Total elements = 10 + 1 = 11.

Therefore, the total number of elements is: $\boxed{11}$.

Failure Analysis: The model's prediction fails due to an Object Counting Error. It correctly identifies the task of counting bars and adding the cow, but falters on the initial visual perception step, incorrectly counting 10 vertical bars instead of the actual visible bars (7, not been occluded by other bars or the cow). While its mathematical logic ($10 + 1 = 11$) is internally consistent, the entire reasoning is built on a flawed premise. The root cause is a failure in basic visual perception, not a flaw in logical deduction.

Failure Case 2: Logical Reasoning Error (MathVerse-V, id: 2)

As shown in the figure, $DE \parallel BC$, then the size of angle CED is
()
Choices: A: 40° B: 60° C: 120° D: 140°



Question: .

Ground Truth: D

Model Prediction: To determine the size of angle CED, we start by analyzing the given information. First, we know that $\angle B = 60^\circ$ and $\angle A = 80^\circ$. Therefore, we can find $\angle C$ using the fact that the sum of the angles in a triangle is 180° :

$$\begin{aligned}\angle C &= 180^\circ - \angle A - \angle B \\ &= 180^\circ - 80^\circ - 60^\circ = 40^\circ\end{aligned}$$

Since $DE \parallel BC$, by the corresponding angles property (*wrong geometric theorem*), we have:

$$\angle CED = \angle C = 40^\circ$$

Therefore, the size of angle CED is 40° . Thus, the correct answer is: \boxed{A} .

Failure Analysis: The model correctly perceives the geometric setup and calculates $\angle C$ as 40° . The error stems from flawed knowledge: it misapplies the parallel line theorem, incorrectly equating $\angle CED$ with $\angle C$. It confuses supplementary consecutive interior angles with equal corresponding angles, showing that while the model "looks" correct, it lacks the proper logic to interpret what it sees.