

# BRIDGING PERCEPTION AND REASONING: TOKEN REWEIGHTING FOR RLVR IN MULTIMODAL LLMs

Anonymous authors  
 Paper under double-blind review

## ABSTRACT

Extending Reinforcement Learning with Verifiable Rewards (RLVR) to multi-modal large language models (MLLMs) faces a fundamental challenge: their responses inherently interleave **perception-related tokens**, which ground visual content, with **reasoning-related tokens**, which construct reasoning chains. These token types instantiate distinct yet interdependent capacities — *visual grounding and symbolic reasoning* — making isolated optimization insufficient. Through token-level empirical analysis, we demonstrate that optimizing either perception- or reasoning-only tokens consistently underperforms full optimization, underscoring their inherent coupling. To address this, we propose a *plug-and-play* **Token-Reweighting (ToR)** strategy that explicitly models this interdependence by identifying critical tokens of both types and dynamically reweighting them during RLVR training. Applied on top of existing methods (*e.g.*, GRPO and DAPO), ToR delivers consistent performance gains across multiple multi-modal reasoning benchmarks, achieving state-of-the-art performance with both accurate visual grounding and coherent reasoning.

## 1 INTRODUCTION

Reinforcement Learning with Verifiable Rewards (RLVR) has substantially advanced the reasoning ability of large language models (LLMs) on complex tasks (Lambert et al., 2025; Shao et al., 2024; Guo et al., 2025; Yang et al., 2025a). Extending RLVR to multimodal large language models (MLLMs), however, is non-trivial: generated responses *interleave* tokens that ground visual content (**perception**) with tokens that drive symbolic inference (**reasoning**), as illustrated in Figure 1. Existing RLVR variants for MLLMs typically optimize these capabilities in isolation — either via chain-of-thought objectives for reasoning (Huang et al., 2025; Wei et al., 2025) or perception-oriented rewards and augmentations for perception (Wang et al., 2025e; Xiao et al., 2025; ?) — leaving their interaction underexplored.

We hypothesize that this separated optimization is suboptimal because perception and reasoning are fundamentally interdependent at the token level. To empirically validate this claim, we conduct a controlled “selective optimization” study under Group Relative Policy Optimization (GRPO) (Shao et al., 2024). We identify reasoning-related tokens via high **next-token entropy** (following recent insights on reasoning forks (Wang et al., 2025c; Cheng et al., 2025)), and perception-related

tokens via **visual sensitivity**, measured as the change in token log-probability when conditioning on the image versus a text-only context (details in Section 3). We then train models while masking gradients on non-selected tokens, comparing three settings: optimizing only reasoning-related tokens, only perception-related tokens, and all tokens (vanilla GRPO).

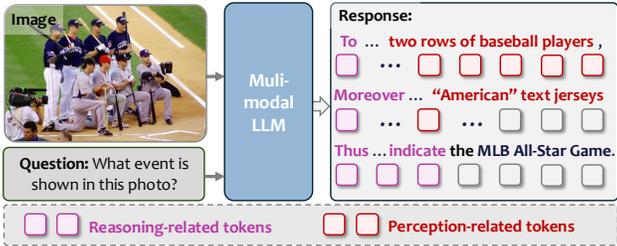


Figure 1: MLLM responses typically involve two types of critical tokens: (1) reasoning-related tokens to construct reasoning chains, and (2) perception-related tokens to ground and represent visual content.



Figure 2: Performance comparison over the wemath benchmark (Qiao et al., 2024) when optimizing different token types with GRPO (Shao et al., 2024). Results across selection ratios 20%, 30%, or 50% show that optimizing either reasoning-only or perception-only tokens underperforms all tokens. Qualitative examples are selected from the best-performing checkpoints.

Across selection ratios (20%, 30%, 50%), optimizing only reasoning tokens underperforms vanilla GRPO (e.g.,  $\sim 2\%$  absolute drop on we-math benchmark), and optimizing only perception tokens fares worse (e.g.,  $\sim 3\%$  drop, with 20% and 30% tokens even worse than the baseline model); neither matches training on all tokens (Figure 2). Qualitatively, reasoning-only models produce coherent chains of thought yet misinterpret key visual content, while perception-only models preserve low-level grounding but fail to integrate it into coherent reasoning. These results support our hypothesis: *perception and reasoning are coupled capabilities that demand joint optimization*.

Building on this insight, we propose **Token-Reweighting (ToR)**, a lightweight, plug-and-play module for RLVR that jointly optimizes perception and reasoning. Instead of treating all tokens equally or optimizing subsets in isolation, ToR strategically identifies the most critical perception- and reasoning-related tokens and adaptively reweights their importance in the policy gradient calculation. This mechanism explicitly models their interdependence, encouraging the MLLM to integrate visual grounding into its logical deliberations. As shown in Figure 2, applying ToR to GRPO (**ToR-GRPO**) not only recovers the performance lost from selective optimization but surpasses the standard GRPO baseline, confirming our effectiveness. Our contributions are threefold:

- We conduct the first systematic, token-level analysis to reveal the critical interdependence between perception and reasoning in MLLMs. Our controlled experiments quantitatively demonstrate that optimizing either capability in isolation is detrimental.
- We introduce **Token-Reweighting (ToR)**, a plug-and-play RLVR training strategy that explicitly models this interdependence by dynamically reweighting critical perception- and reasoning-related tokens during policy optimization.
- We empirically demonstrate that ToR delivers consistent and substantial gains when integrated with state-of-the-art RLVR algorithms (e.g., GRPO and DAPO), setting a new state-of-the-art across a diverse suite of both multi-modal reasoning and perception benchmarks.

## 2 PRELIMINARIES

In this section, we revisit Reinforcement Learning with Verifiable Rewards (RLVR) procedure and representative RLVR optimization strategies for Multi-modal Large Language Models (MLLMs).

### 2.1 REINFORCEMENT LEARNING WITH VERIFIABLE REWARDS

Reinforcement Learning with Verifiable Rewards (RLVR) enhances multi-modal large language models (MLLMs) by aligning model outputs with verifiable answers. Given a multi-modal input with ground truth from a batch of  $B$  samples,  $(I^b, q^b) \in \{I^b, q^b\}_{b=1}^B$ , and  $y^b \in \{y^b\}_{b=1}^B$ , where  $I^b$ ,

$q^b$ , and  $y^b$  denotes the image, question, and ground-truth respectively, the model  $\pi_\theta$  generates output  $\mathbf{o}^b$  containing a reasoning process and a prediction. The prediction is enclosed in `\boxed{.}` while reasoning is delimited by `<think>...</think>`, enabling automated verification against the ground truth answers. RLVR employs a binary reward function  $\mathbf{R}(\cdot)$  to determine whether the answer is correct by comparing the model output  $\mathbf{o}^b$  with ground truth  $y^b$ . The goal of RLVR is to maximize the reward function, formalized as:

$$\mathcal{J}_{\text{RLVR}}(\theta) = \max_{\theta} \mathbb{E}_{\{(I^b, q^b) | y^b\}_{b=1}^B} \mathbb{E}_{\mathbf{o}^b \sim \pi_\theta(\cdot | I^b, q^b)} [\mathbf{R}(\mathbf{o}^b, y^b)]. \quad (1)$$

## 2.2 RLVR OPTIMIZATION ALGORITHMS

**Group Relative Policy Optimization (GRPO).** As a widely adopted RLVR optimization strategy, GRPO stabilizes training by computing advantages within response groups (Shao et al., 2024). Concretely, given a batch of samples  $\{(I^b, q^b) | y^b\}_{b=1}^B$ , the GRPO objective is:

$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E}_{\{(I^b, q^b)\}_{b=1}^B} \mathbb{E}_{\{o_i^b\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\mathbf{o}^b | I^b, q^b)} \left\{ \frac{1}{G} \sum_{i=1}^G \frac{1}{|o_i^b|} \sum_{t=1}^{|o_i^b|} \min \left[ \frac{\pi_\theta(o_{i,t}^b | I^b, q^b, o_{i,<t}^b)}{\pi_{\theta_{\text{old}}}(o_{i,t}^b | I^b, q^b, o_{i,<t}^b)} \hat{A}_{i,t}^b, \text{clip} \left( \frac{\pi_\theta(o_{i,t}^b | I^b, q^b, o_{i,<t}^b)}{\pi_{\theta_{\text{old}}}(o_{i,t}^b | I^b, q^b, o_{i,<t}^b)}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{i,t}^b \right] - \beta \mathbb{D}_{KL}[\pi_\theta \| \pi_{\text{ref}}] \right\}, \quad (2)$$

where  $G$  is the rollout group size for each input, and the clip function restricts the importance ratio within  $[1 - \epsilon, 1 + \epsilon]$ . Moreover, the advantage can be formulated as:

$$\hat{A}_{i,t}^b = \frac{\mathbf{R}_i^b - \text{mean}(\mathbf{R}^b)}{\text{std}(\mathbf{R}^b)}, \text{ where } \mathbf{R}_i^b = \mathbb{I}(\text{is\_equivalent}(o_i^b, y^b)), \quad (3)$$

where  $\mathbb{I}(\cdot)$  is the indicator function, `is_equivalent`( $\cdot$ ) extracts the predictions from the model output  $o_i^b$  and compares with the ground truth  $y^b$ .

**Decoupled Clip and Dynamic Sampling Policy Optimization (DAPO).** DAPO (Yu et al., 2025b) improves upon GRPO by removing KL regularization and introducing clip-higher, dynamic sampling, and token-level loss from the GRPO loss, achieving new state-of-the-art performance. Specifically, the DAPO objective for the batch of samples  $\{(I^b, q^b) | y^b\}_{b=1}^B$  can be formalized as:

$$\mathcal{J}_{\text{DAPO}}(\theta) = \mathbb{E}_{\{(I^b, q^b)\}_{b=1}^B} \mathbb{E}_{\{o_i^b\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\mathbf{o}^b | I^b, q^b)} \frac{1}{\sum_{i=1}^G |o_i|} \sum_{i=1}^G \sum_{t=1}^{|o_i^b|} \left\{ \min \left[ \frac{\pi_\theta(o_{i,t}^b | I^b, q^b, o_{i,<t}^b)}{\pi_{\theta_{\text{old}}}(o_{i,t}^b | I^b, q^b, o_{i,<t}^b)} \hat{A}_{i,t}^b, \text{clip} \left( \frac{\pi_\theta(o_{i,t}^b | I^b, q^b, o_{i,<t}^b)}{\pi_{\theta_{\text{old}}}(o_{i,t}^b | I^b, q^b, o_{i,<t}^b)}, 1 - \epsilon_{\text{low}}, 1 + \epsilon_{\text{high}} \right) \hat{A}_{i,t}^b \right] \right\}, \quad (4)$$

where  $\epsilon_{\text{low}}$  and  $\epsilon_{\text{high}}$  decouple the clipping thresholds for negative and positive advantages, respectively. In this work, we apply our token-reweighting strategy to both GRPO and DAPO, demonstrating its general applicability across various RLVR optimization strategies.

## 3 APPROACH

In this section, we elaborate on our token re-weighting strategy in detail. Firstly, we identify and analyze the effects of reasoning- and perception-related tokens; then, we illustrate our dynamic reweighting scheme based on the token types, and demonstrate the application to existing RLVR algorithms, such as GRPO and DAPO.

### 3.1 TOKEN IDENTIFICATION

In this subsection, we present our token identification strategy and analyze the effects of emphasizing different types of tokens during GRPO optimization.

### 3.1.1 REASONING-RELATED TOKENS

**Identification of reasoning-related tokens.** Recent works (Wang et al., 2025c; Cheng et al., 2025) have shown that high-entropy tokens often correspond to critical “forking” points in reasoning chains, directly reflecting the model’s *decision uncertainty*. Moreover, retaining gradients for only a small subset of such tokens has been shown to benefit optimization, highlighting their potential importance for reasoning.

Motivated by this, we identify reasoning-related tokens using an entropy-based criterion. Given the  $i$ -th generated response  $\mathbf{o}_i^b = \{o_{i,1}^b, o_{i,2}^b, \dots, o_{i,L_i}^b\}$  conditioned on the image  $I^b$ , and then question  $q^b$  from the input batch  $\{(I^b, q^b)\}_{b=1}^B$ , the prediction entropy at position  $t$  is computed as:

$$H_{i,t}^b = - \sum_{v \in \mathcal{V}_{\text{top-}p}} P_{\theta}(o_{i,t}^b = v \mid \mathbf{o}_{i,<t}^b, I^b, q^b) \cdot \log P_{\theta}(o_{i,t}^b = v \mid \mathbf{o}_{i,<t}^b, I^b, q_i^b), \quad (5)$$

where  $\mathcal{V}_{\text{top-}p}$  denotes the set of vocabulary tokens within the top- $p$  cumulative probability ( $p = 0.95$ ). This top- $p$  truncation avoids the influence of extremely low-probability tokens and is consistent with the rollout sampling process.

To aggregate across the rollout batch, we collect all token entropies into a set:

$$\mathcal{H} = \{H_{i,t}^b \mid b = 1, \dots, B; i = 1, \dots, G; t = 1, \dots, L_i^b\}. \quad (6)$$

The reasoning-token set is then defined by selecting the top- $\alpha_r$  fraction of tokens with the highest entropy across the batch:

$$\mathcal{T}_r = \{(b, i, t) \mid H_{i,t}^b \geq \text{Percentile}_{1-\alpha_r}(\mathcal{H})\}, \quad (7)$$

where  $\alpha_r$  controls the fraction of selected tokens. Intuitively, these high-uncertainty positions correspond to pivotal points where reasoning chains are shaped.

**Influences of reasoning-related tokens.** We evaluate the practical effect of the reasoning-token set  $\mathcal{T}_r$  by conducting GRPO training constrained to tokens in this set with  $\alpha_r = 20\%, 30\%, 50\%$ . As shown in Figure 3, all selection ratios underperform the baseline that optimizes over all tokens. Notably, the 30% and 50% settings converge to similar outcomes, both falling short of the full-token GRPO. These results indicate that while  $\mathcal{T}_r$  captures critical decision points, optimizing solely on it fails to preserve sufficient contextual information — *particularly when errors originate from perception* — motivating the investigation of perception-related tokens.

Acc w/ Reasoning tokens over steps

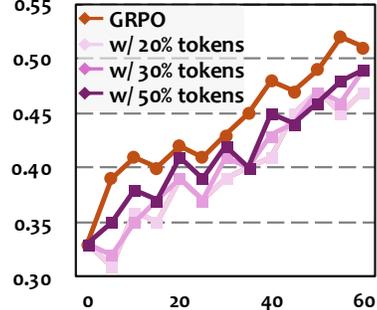


Figure 3: Performance on the Geo3K validation set for different ratios of reasoning-related tokens during GRPO training. Training data: Geo3K training set, Model: Qwen-VL-2.5 7B.

### 3.1.2 PERCEPTION-RELATED TOKENS

**Identification of perception-related tokens.** While reasoning-related tokens capture decision uncertainty, perception-related tokens highlight positions whose predictions strongly depend on visual inputs. To quantify this dependence, we compare the token log-probabilities under two conditions: (i) *with image*, conditioned on both  $I^b$  and  $q^b$ ; and (ii) *without image*, where the image channel is replaced by an empty placeholder  $\emptyset$  and the model is conditioned only on  $q^b$ .

Given the  $i$ -th generated response  $\mathbf{o}_i^b = \{o_{i,1}^b, \dots, o_{i,T_i}^b\}$ , the visual-sensitivity score at position  $t$  is computed as:

$$S_{i,t}^b = |\log \pi_{\theta}(o_{i,t}^b \mid \mathbf{o}_{i,<t}^b, I^b, q_i^b) - \log \pi_{\theta}(o_{i,t}^b \mid \mathbf{o}_{i,<t}^b, \emptyset, q_i^b)|, \quad (8)$$

where  $\pi_{\theta}(o_t \mid \cdot)$  is the token-level policy. A large  $S_{i,t}$  indicates a strong visual influence on token  $o_{i,t}$ . Aggregating across the rollout batch, the perception-token set is defined as the top- $\alpha_p$  fraction of tokens with the highest visual-sensitivity:

$$\mathcal{T}_p = \{(b, i, t) \mid S_{i,t}^b \geq \text{Percentile}_{1-\alpha_p}(\{S_{b,i,t}^b \mid b = 1, \dots, B; i = 1, \dots, G; t = 1, \dots, L_i^b\})\}, \quad (9)$$

where  $\alpha_p$  controls the selection ratio. This batch-level percentile selection ensures consistent token selection thresholds across rollout batches.

**Influences of perception-related tokens.** We evaluate the practical effect of the perception-token set  $\mathcal{T}_p$  by constraining GRPO optimization to tokens in this set with  $\alpha_p = 20\%, 30\%, 50\%$ . As shown in Figure 4, all selection ratios underperform the full-token baseline. In particular, 20% and 30% result in performance drops exceeding 5%, whereas 50% shows the smallest performance gap but still underperforms the baseline. These results indicate that although perception-related tokens capture visually sensitive positions, optimizing solely on them is insufficient. *Effective GRPO training requires attending to both reasoning- and perception-related tokens to capture critical decision points while leveraging visual context.*

Acc w/ Perception tokens over steps

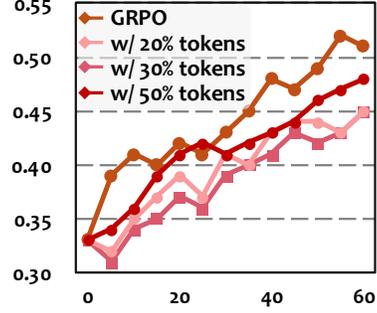


Figure 4: Performance on the Geo3K validation set for different ratios of perception-related tokens during GRPO training, Training data: Geo3K training set, Model: Qwen-VL-2.5 7B.

### 3.2 TOKEN REWEIGHTING

We unify reasoning- and perception-critical tokens into a single optimization framework by introducing *token reweighting*, which selectively amplifies their contributions while ignoring irrelevant tokens. Specifically, we integrate token-specific reweighting directly into the RL objectives. For the given input batch  $\{(I^b, q^b)\}_{b=1}^B$ , we constructed the reasoning-related token set  $\mathcal{T}_r$  and the perception-related token set  $\mathcal{T}_p$ , and thus the GRPO objective with the token reweighting strategy (**ToR-GRPO**) is:

$$\begin{aligned} \mathcal{J}_{\text{ToR-GRPO}}(\theta) = & \mathbb{E}_{\{(I^b, q^b)\}_{b=1}^B} \mathbb{E}_{\{o_i^b\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(o^b | I^b, q^b)} \left\{ \frac{1}{G} \sum_{i=1}^G \frac{1}{|o_i^b|} \sum_{t=1}^{|o_i^b|} [\gamma_r \cdot \mathbb{I}[(b, i, t) \in \mathcal{T}_r] \right. \\ & \left. + \gamma_p \cdot \mathbb{I}[(b, i, t) \in \mathcal{T}_p]] \cdot \min \left[ r_{\theta}(o_{i,t}^b) \cdot \hat{A}_{i,t}^b, \text{clip} \left[ r_{\theta}(o_{i,t}^b), 1 - \epsilon, 1 + \epsilon \right] \cdot \hat{A}_{i,t}^b \right] - \beta \mathbb{D}_{KL}[\pi_{\theta} \| \pi_{\text{ref}}] \right\}. \end{aligned} \quad (10)$$

Similarly, the DAPO objective with the token-weighting strategy (**ToR-DAPO**) is:

$$\begin{aligned} \mathcal{J}_{\text{ToR-DAPO}}(\theta) = & \mathbb{E}_{\{(I^b, q^b)\}_{b=1}^B} \mathbb{E}_{\{o_i^b\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot | I^b, q^b)} \frac{1}{\sum_{i=1}^G |o_i^b|} \sum_{i=1}^G \sum_{t=1}^{|o_i^b|} \left\{ [\gamma_r \cdot \mathbb{I}[(b, i, t) \in \mathcal{T}_r] \right. \\ & \left. + \gamma_p \cdot \mathbb{I}[(b, i, t) \in \mathcal{T}_p]] \cdot \min \left[ r_{\theta}(o_{i,t}^b) \cdot \hat{A}_{i,t}^b, \text{clip} \left[ r_{\theta}(o_{i,t}^b), 1 - \epsilon_{\text{low}}, 1 + \epsilon_{\text{high}} \right] \cdot \hat{A}_{i,t}^b \right] \right\}. \end{aligned} \quad (11)$$

$r_{\theta}(o_{i,t}^b)$  is the importance sampling ratio, formalized as:

$$r_{\theta}(o_{i,t}^b) = \frac{\pi_{\theta}(o_{i,t}^b | I_i^b, q_i^b, o_{i,<t}^b)}{\pi_{\theta_{\text{old}}}(o_{i,t}^b | I_i^b, q_i^b, o_{i,<t}^b)}. \quad (12)$$

Here,  $\gamma_r$  weights reasoning-related tokens, emphasizing critical decision points in reasoning chains;  $\gamma_p$  weights perception-related tokens, emphasizing the integration of visual context. Tokens outside these sets are excluded from optimization (*i.e.*, assigned a weight of zero). This reweighting ensures gradients focus on tokens essential for both reasoning and perception, improving both training efficiency and effectiveness.

**Discussion** Compared with existing approaches, ToR offers several key advantages: **● Plug-and-play**. A simple reweighting mask integrates seamlessly into standard RLVR objectives, without

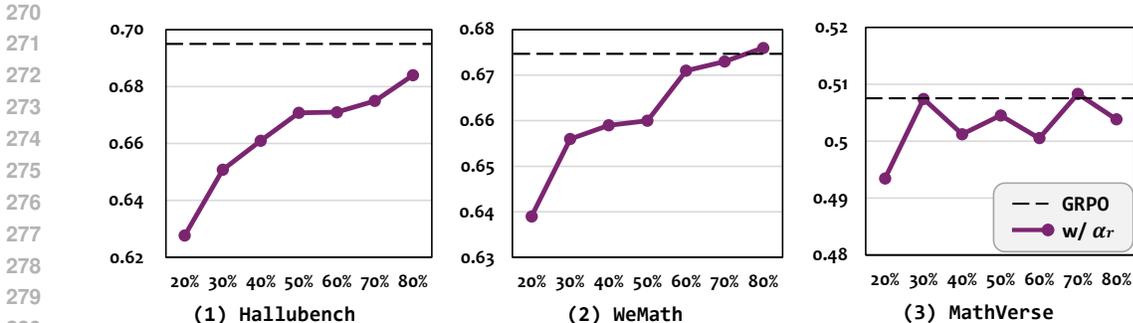


Figure 5: Performance comparison with different ratios of reasoning-related tokens ( $\alpha_r$ ).

introducing extra pipeline modifications. **2 Self-contained.** Critical tokens are identified purely from the model’s intrinsic uncertainty and visual sensitivity, eliminating the need for external priors. **3 Joint optimization.** By explicitly modeling the interdependence between reasoning and perception tokens, ToR enables balanced and simultaneous enhancement of both capabilities.

## 4 EXPERIMENT

In this section, we elaborate on the effectiveness of our token reweighting strategy. Specifically, we first introduce the details of our experimental settings. Next, we present the ablation studies, and finally, we compare our results with those of state-of-the-art methods across various benchmarks.

### 4.1 EXPERIMENTAL SETTINGS

We adopt the Geometry3K (Lu et al., 2021) dataset for training and validation, following existing methods (Xiao et al., 2025), which consists of 2,100 training samples and 300 validation samples. Moreover, we employ the multi-modal framework EasyR1 (Yaowei et al., 2025) for reinforcement learning training, and following the works in (?) for evaluation. Evaluation is conducted with five benchmarks, including four for visual reasoning: MathVerse (Zhang et al., 2024), MathVision (Wang et al., 2024), MathVista (Lu et al., 2024), and WeMath (Qiao et al., 2024), and one for visual perception: HallusionBench (Guan et al., 2024).

**Implementation details.** We adopt Qwen2.5-VL-7B (Bai et al., 2025) as our baseline model by following existing works in (Meng et al., 2025; ?). Specifically, all experiments are conducted using 8 NVIDIA H800 GPUs (80 GB memory for each), with the default settings in EasyR1: a learning rate of  $1e^{-6}$ , a global batch size of 128, a rollout batch size of 512, a rollout  $n$  as 12, and a rollout temperature of 0.95.

### 4.2 ABLATION STUDIES

In this section, we analyze the influences of different components on our token re-weighting strategy. Specifically, we utilize the Qwen-2.5-VL 7B (Bai et al., 2025) as the backbone, and conduct GRPO on the training set of Geometry3K. With the aim to evaluate the reasoning and perception ability, we select the Hallusion Bench (Guan et al., 2024), the Wemath (Qiao et al., 2024), and the MathVerse (Zhang et al., 2024) for evaluation, where Hallusion Bench focuses on the perception abilities, and MathVerse targets reasoning capabilities, whereas Wemath focuses on both.

**1 The effects of different ratios for reasoning and perception tokens.** In this section, we study the effect of varying token ratios (*e.g.*, the proportion of reasoning-related tokens  $\alpha_r$  and perception-related tokens  $\alpha_p$ ). Results for reasoning-related tokens are presented in Figure 5, while those for perception-related tokens are shown in Figure 6. Specifically, we vary  $\alpha_r$  and  $\alpha_p$  from 20% to 80% over 60 GRPO steps, and report the performance of full-token GRPO as a dashed reference line.

From Figure 5, two observations emerge: (1) On perception-intensive benchmarks such as HallusionBench, none of the ratios matches the performance of full-token GRPO. (2) On reasoning-

324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377

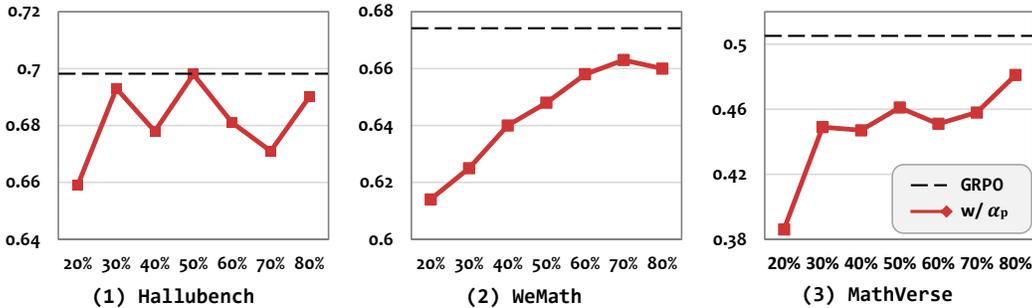


Figure 6: Performance comparison with different ratios of perception-related tokens ( $\alpha_p$ ).

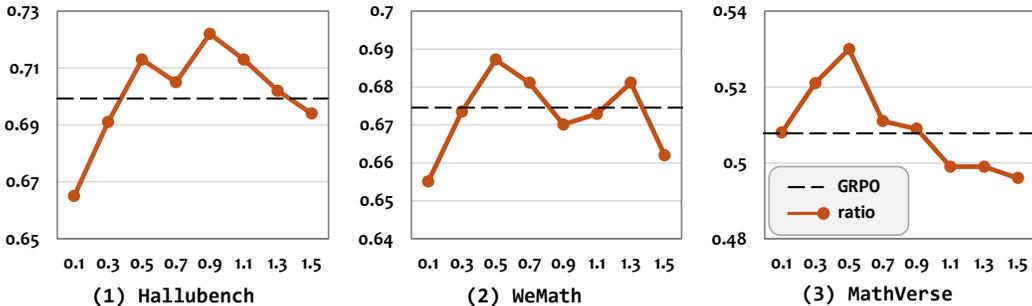


Figure 7: Performance comparison with varying combination ratios of reasoning- and perception-related tokens (reasoning weight  $\gamma_r$  fixed at 1, perception weight  $\gamma_p$  varying from 0.1 to 1.5).

focused benchmarks (e.g., WeMath and MathVerse), adjusting token ratios produces results comparable to direct GRPO, but yields little to no further improvement.

Similarly, Figure 6 shows that while perception-only reweighting achieves performance close to full-token GRPO on perception-oriented benchmarks, it performs substantially worse on reasoning-demanding tasks, particularly on MathVerse, indicating a notable gap from full-token GRPO. These results collectively confirm that focusing exclusively on either reasoning- or perception-related tokens is insufficient for robust multimodal reasoning.

By jointly considering both token types and their dynamics during training, we adopt a balanced ratio of 30% reasoning tokens and 30% perception tokens for re-weighting. Notably, in the initial stage of GRPO training on Geometry3K, about 12% of tokens overlap between the two categories, and this proportion evolves during optimization. For these overlapping tokens, our experiments show that while both reasoning- and perception-related tokens experience performance drops, the decrease is smaller for reasoning tokens; therefore, we assign them the weight of reasoning tokens during re-weighting.

**🔗 The effects of different combination ratios for different tokens.** In this section, we evaluate the impact of varying combination ratios between reasoning- and perception-related tokens. The results are shown in Figure 7, where the weight of reasoning tokens is fixed at 1, and the weight of perception tokens is varied from 0.1 to 1.5 in steps of 0.2. The performance of full-token GRPO is shown in a dashed reference line.

From Figure 7, we make the following observations: (1) Combining reasoning- and perception-related tokens consistently improves performance compared to using either type alone. (2) For reasoning-focused tasks, a relatively low proportion of perception tokens is preferable, whereas perception-oriented tasks benefit from a higher proportion. (3) Extremely low or high perception token weights tend to degrade performance. Overall, a weight of 0.5 for perception tokens provides a good balance across tasks, therefore, we adopt this ratio as a reference for subsequent experiments.

Table 1: Performance comparison of 7B-sized Multi-modal LLMs on different benchmarks. Following existing works in (Xiao et al., 2025), we highlight the data size with blue and red for SFT and RL, respectively. The best value in each column is shown in bold, and the second-best is underlined.

Model	Data Size	MathVerse	MathVision	MathVista	WeMath	HallusionBench
<i>Open-source Models</i>						
InternVL-2.5-8B (Chen et al., 2024)	-	39.5	19.7	64.4	-	67.3
InternVL-3-8B (Zhu et al., 2025)	-	39.5	29.3	71.6	37.1	-
LLaVA-OneVision-7B (Li et al., 2025a)	-	26.2	-	63.2	-	48.4
Qwen2.5-VL-7B-Instruct (Bai et al., 2025)	-	46.2	25.0	67.5	63.1	64.6
<i>reinforcement learning with verifiable reward strategies</i>						
R1-VL-7B (Zhang et al., 2025a)	260K+10K	40.0	24.7	63.5	-	-
Vision-R1-7B (Huang et al., 2025)	200K+10K	52.4	-	<b>73.5</b>	-	-
R1-OneVision-7B (Yang et al., 2025b)	155K+10K	46.1	22.5	63.9	62.1	65.6
OpenVLThinker-7B (Deng et al., 2025b)	35K+15K	48.0	25.0	71.5	67.8	70.8
MM-Eureka-Qwen-7B (Meng et al., 2025)	15K	50.5	28.3	71.5	65.5	68.3
ThinkLite-7B-VL (Wang et al., 2025d)	11K	50.2	27.6	<u>72.7</u>	69.2	71.0
VLAA-Thinker-7B (Chen et al., 2025a)	25K	49.9	26.9	68.8	67.9	68.6
NoisyRollout-7B (Liu et al., 2025a)	2.1K	<u>53.2</u>	<u>28.5</u>	72.6	<u>69.6</u>	<u>72.1</u>
GRPO (Shao et al., 2024)	2.1K (Geometry3K)	50.8	27.3	70.5	67.4	69.8
<b>ToR-GRPO</b>	2.1K (Geometry3K)	53.0	<b>28.6</b>	71.9	68.9	<b>72.4</b>
DAPO (Yu et al., 2025b)	2.1K (Geometry3K)	50.6	26.5	70.3	69.3	67.9
<b>ToR-DAPO</b>	2.1K (Geometry3K)	<b>53.4</b>	27.9	72.6	<b>72.1</b>	71.8

### 4.3 COMPARISON WITH STATE-OF-THE-ART APPROACHES

In this section, we compare our token re-weighting strategies with state-of-the-art approaches, including open-source multi-modal LLMs: InternVL-2.5 (Chen et al., 2024), InternVL-3 (Chen et al., 2024), LLaVA-OneVision (Li et al., 2025a), and Qwen-2.5-VL 7B (Bai et al., 2025), as well as RLVR-based methods: R1-VL (Zhang et al., 2025a), Vision-R1 (Huang et al., 2025), R1-OneVision (Yang et al., 2025b), Open-VLThinker (Deng et al., 2025b), MM-Eureka (Meng et al., 2025), ThinkLite (Wang et al., 2025d), VLAA-Thinker (Chen et al., 2025a), and NoisyRollout (?). The results are summarized in Table 1.

Unlike existing approaches that often rely on large-scale data augmentation or extensive chain-of-thought distillation, our token re-weighting method achieves substantial improvements with only 2.1K samples. Specifically, it yields an average gain of more than 1.5% over GRPO across benchmarks, with notable improvements of 2.2% on MathVerse and 2.6% on HallusionBench.

Furthermore, our approach establishes new state-of-the-art results on multiple benchmarks. For example, on the WeMath benchmark, applying token re-weighting to DAPO (**ToR-DAPO**) achieves an improvement of about 3%, significantly surpassing state-of-the-art methods.

## 5 RELATED WORK

In this section, we first briefly summarize methods that focus on enhancing reasoning capabilities in Multi-modal LLMs, and then, we illustrate related methods for reinforcement learning with verifiable rewards. Finally, we enumerate the differences between our approach and related methods.

### 5.1 REASONING IN MULTI-MODAL LLMs

Existing approaches that focus on reasoning in multi-modal LLMs can be broadly categorized into:

❶ **Extending reasoning LLMs with visual understanding.** Building upon the strong reasoning capabilities of recent LLMs, one research branch explores incorporating visual content into LLMs. Typical strategies include: (1) integrating visual encoders into LLMs to directly extend them to multimodal scenarios (Peng et al., 2025; Wang et al., 2025b); and (2) transforming images into captions and feeding them into LLMs, enhancing chain-of-thought generation to bridge perception and reasoning (Huang et al., 2025; Wang et al., 2025e).

❷ **Enhancing reasoning abilities within MLLMs.** Another research branch seeks to endow existing MLLMs with reasoning skills. Representative approaches include: (1) transferring reasoning

priors from reasoning LLMs into MLLMs through model merging, thereby leveraging the complementary strengths of both models (Chen et al., 2025b); and (2) adapting RLVR algorithms, such as Group Relative Policy Optimization (GRPO) (Shao et al., 2024), to enhance reasoning capabilities in multimodal settings (Huang et al., 2025; Liu et al., 2025c).

## 5.2 REINFORCEMENT LEARNING WITH VERIFIABLE REWARDS

RLVR has recently demonstrated its effectiveness in aligning LLMs with verifiable reasoning outcomes. Among the implementations, GRPO (Shao et al., 2024) has shown great success with more stable advantage estimation. Subsequent refinements have further improved its efficiency and effectiveness (Yu et al., 2025b; Liu et al., 2025b; Wang et al., 2025c). When extending RLVR to multimodal settings, current methods primarily focus on two distinct branches:

❶ **Emphasizing the importance of visual understanding within RLVR.** Specifically, works like (Xiao et al., 2025; Yu et al., 2025a) introduce perception-oriented rewards to incentivize accurate visual grounding and understanding. Moreover, works such as (Wang et al., 2025e; Li et al., 2025b) employ data augmentation strategies to improve model robustness and sensitivity against visual variations. Works like (Zhang et al., 2025c; Liu et al., 2025d) incorporate image manipulation tools like cropping and zooming in to focus on critical regions in the image during the reasoning process.

❷ **Constructing coherent reasoning chains with RLVR.** Representative works include: (1) distilling chain-of-thought reasoning patterns from stronger reasoning models to improve reasoning coherence in multimodal tasks (Huang et al., 2025; Wei et al., 2025); and (2) modifying different components of GRPO (e.g., clip ratios or advantage estimation) to emphasize reasoning-critical tokens better and stabilize training (Zhang et al., 2025b; Meng et al., 2025; Wang et al., 2025a).

**Differences.** Unlike existing methods that focus on optimizing either perception or reasoning abilities in isolation, our work systematically investigates and addresses the interdependence between these two capabilities. Specifically, through comprehensive token-level analysis, we demonstrate that perception and reasoning tokens exhibit complex interactions during training, where optimizing one type can inadvertently impair the other. To address this challenge, we propose a simple token-reweighting strategy that explicitly balances the optimization of both perception and reasoning tokens, leading to significant performance improvements across both capabilities.

## 6 CONCLUSION

In this work, through systematic token-wise analysis, we uncover a fundamental challenge in extending RLVR to multimodal LLMs: **the intrinsic interdependence between perception and reasoning**. We show that overemphasizing either capability inevitably impairs the other, yet current approaches overlook this issue and optimize them in isolation.

To address this, we proposed a simple yet effective **Token-Reweighting (ToR)** strategy that identifies and reweights perception- and reasoning-related tokens during RLVR training. We apply ToR over current RLVR algorithms (e.g., GRPO and DAPO), and achieves significant performance gains across diverse benchmarks, consistently enhancing both perception and reasoning.

**Limitations and Future Work.** While ToR establishes a foundation for addressing perception-reasoning interdependence, several promising directions remain open: ❶ **Fine-grained token identification strategies:** Precisely localizing critical regions in images using models like SAM (Kirillov et al., 2023) to identify more fine-grained perception-critical tokens. ❷ **Dynamic token reweighting:** Dynamically assigning weights to tokens based on their gradient contributions or connections to final outcomes (Yu et al., 2025c). ❸ **Extending beyond tokens:** As tokens derive meaning from context, future work could explore perception-reasoning interdependence at broader contextual levels — optimizing tokens within their semantic context, preserving their contextual relationships. ❹ **Exploring broader applications:** Extending this interdependence framework to more complex scenarios, e.g., unified multi-modal generation and understanding tasks (Wu et al., 2025; Deng et al., 2025a), where visual tokens participate in reasoning processes.

## ETHICS STATEMENT

Our token re-weighting (ToR) approach significantly enhances the reasoning capabilities of multi-modal LLMs. As a simple, plug-and-play strategy, it can be readily integrated into existing RLVR frameworks, making it a practical tool for improving multi-modal reasoning performance. Given its focus on model optimization, this method does not introduce new ethical risks beyond those already inherent to large language models.

## REPRODUCIBILITY STATEMENT

Our ToR strategy is easy to implement. We provide implementation details, including the framework, token selection criteria, and token weights used in our experiments. The code will be made publicly available upon acceptance to facilitate reproducibility and further research.

## REFERENCES

- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
- Hardy Chen, Haoqin Tu, Fali Wang, Hui Liu, Xianfeng Tang, Xinya Du, Yuyin Zhou, and Cihang Xie. SFT or RL? an early investigation into training rl-like reasoning large vision-language models. *Transactions on Machine Learning Research*, 2025a. ISSN 2835-8856. URL <https://openreview.net/forum?id=wZI5qkQeDF>.
- Shiqi Chen, Jinghan Zhang, Tongyao Zhu, Wei Liu, Siyang Gao, Miao Xiong, Manling Li, and Junxian He. Bring reason to vision: Understanding perception and reasoning through model merging. In *Forty-second International Conference on Machine Learning*, 2025b. URL <https://openreview.net/forum?id=ntCAP6tMoX>.
- Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024.
- Daixuan Cheng, Shaohan Huang, Xuekai Zhu, Bo Dai, Wayne Xin Zhao, Zhenliang Zhang, and Furu Wei. Reasoning with exploration: An entropy perspective. *arXiv preprint arXiv:2506.14758*, 2025.
- Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Weihao Yu, Xiaonan Nie, Ziang Song, et al. Emerging properties in unified multimodal pretraining. *arXiv preprint arXiv:2505.14683*, 2025a.
- Yihe Deng, Hritik Bansal, Fan Yin, Nanyun Peng, Wei Wang, and Kai-Wei Chang. OpenVLThinker: Complex vision-language reasoning via iterative SFT-RL cycles. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025b. URL <https://openreview.net/forum?id=gfX1nqBKtu>.
- Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, et al. 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*, pp. 14375–14385, 2024.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang, Shirong Ma, Xiao Bi, et al. Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. *Nature*, 645(8081):633–638, 2025.
- Wenxuan Huang, Bohan Jia, Zijie Zhai, Shaosheng Cao, Zheyu Ye, Fei Zhao, Zhe Xu, Yao Hu, and Shaohui Lin. Vision-r1: Incentivizing reasoning capability in multimodal large language models. *arXiv preprint arXiv:2503.06749*, 2025.

- 540 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete  
541 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceed-*  
542 *ings of the IEEE/CVF international conference on computer vision*, pp. 4015–4026, 2023.
- 543 Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brah-  
544 man, Lester James Validad Miranda, Alisa Liu, Nouha Dziri, Xixi Lyu, Yuling Gu, Saumya  
545 Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Christo-  
546 pher Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Ha-  
547 jishirzi. Tulu 3: Pushing frontiers in open language model post-training. In *Second Conference on*  
548 *Language Modeling*, 2025. URL <https://openreview.net/forum?id=iluGbfHHpH>.
- 549 Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan  
550 Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. LLaVA-onevision: Easy visual task trans-  
551 fer. *Transactions on Machine Learning Research*, 2025a. ISSN 2835-8856. URL <https://openreview.net/forum?id=zKv8qULV6n>.
- 552 Yuting Li, Lai Wei, Kaipeng Zheng, Jingyuan Huang, Linghe Kong, Lichao Sun, and Weiran Huang.  
553 Vision matters: Simple visual perturbations can boost multimodal math reasoning. *arXiv preprint*  
554 *arXiv:2506.09736*, 2025b.
- 555 Xiangyan Liu, Jinjie Ni, Zijian Wu, Chao Du, Longxu Dou, Haonan Wang, Tianyu Pang, and  
556 Michael Qizhe Shieh. Noisyrollout: Reinforcing visual reasoning with data augmentation. In *2nd*  
557 *AI for Math Workshop @ ICML 2025*, 2025a. URL [https://openreview.net/forum?](https://openreview.net/forum?id=GdNm9PEAei)  
558 [id=GdNm9PEAei](https://openreview.net/forum?id=GdNm9PEAei).
- 559 Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min  
560 Lin. Understanding r1-zero-like training: A critical perspective. In *2nd AI for Math Workshop @*  
561 *ICML 2025*, 2025b. URL <https://openreview.net/forum?id=jLpClzavzn>.
- 562 Ziyu Liu, Zeyi Sun, Yuhang Zang, Xiaoyi Dong, Yuhang Cao, Haodong Duan, Dahua Lin, and Jiaqi  
563 Wang. Visual-rft: Visual reinforcement fine-tuning. In *Proceedings of the IEEE/CVF Interna-*  
564 *tional Conference on Computer Vision (ICCV)*, pp. 2034–2044, October 2025c.
- 565 Ziyu Liu, Yuhang Zang, Yushan Zou, Zijian Liang, Xiaoyi Dong, Yuhang Cao, Haodong  
566 Duan, Dahua Lin, and Jiaqi Wang. Visual agentic reinforcement fine-tuning. *arXiv preprint*  
567 *arXiv:2505.14246*, 2025d.
- 568 Pan Lu, Ran Gong, Shibiao Jiang, Liang Qiu, Siyuan Huang, Xiaodan Liang, and Song-chun Zhu.  
569 Inter-gps: Interpretable geometry problem solving with formal language and symbolic reasoning.  
570 In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics*  
571 *and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long*  
572 *Papers)*, pp. 6774–6786, 2021.
- 573 Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-  
574 Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning  
575 of foundation models in visual contexts. In *The Twelfth International Conference on Learning*  
576 *Representations*, 2024. URL <https://openreview.net/forum?id=KUNzEQMWU7>.
- 577 Fanqing Meng, Lingxiao Du, Zongkai Liu, Zhixiang Zhou, Quanfeng Lu, Daocheng Fu, Botian Shi,  
578 Wenhai Wang, Junjun He, Kaipeng Zhang, et al. Mm-eureka: Exploring visual aha moment with  
579 rule-based large-scale reinforcement learning. *CoRR*, 2025.
- 580 Yi Peng, Peiyu Wang, Xiaokun Wang, Yichen Wei, Jiangbo Pei, Weijie Qiu, Ai Jian, Yunzhuo Hao,  
581 Jiachun Pan, Tianyidan Xie, et al. Skywork r1v: Pioneering multimodal reasoning with chain-of-  
582 thought. *arXiv preprint arXiv:2504.05599*, 2025.
- 583 Runqi Qiao, Qiuna Tan, Guanting Dong, Minhui Wu, Chong Sun, Xiaoshuai Song, Zhuoma  
584 GongQue, Shanglin Lei, Zhe Wei, Miaoxuan Zhang, et al. We-math: Does your large multi-  
585 modal model achieve human-like mathematical reasoning? *arXiv preprint arXiv:2407.01284*,  
586 2024.
- 587 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,  
588 Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathemati-  
589 cal reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- 590  
591  
592  
593

- 594 Haozhe Wang, Chao Qu, Zuming Huang, Wei Chu, Fangzhen Lin, and Wenhui Chen. VL-  
595 rethiker: Incentivizing self-reflection of vision-language models with reinforcement learning.  
596 In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025a. URL  
597 <https://openreview.net/forum?id=4oYxzssbVg>.
- 598 Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Houxing Ren, Aojun Zhou, Mingjie Zhan, and Hong-  
599 sheng Li. Measuring multimodal mathematical reasoning with math-vision dataset. *Advances in*  
600 *Neural Information Processing Systems*, 37:95095–95169, 2024.
- 601 Peiyu Wang, Yichen Wei, Yi Peng, Xiaokun Wang, Weijie Qiu, Wei Shen, Tianyidan Xie, Jiangbo  
602 Pei, Jianhao Zhang, Yunzhuo Hao, et al. Skywork r1v2: Multimodal hybrid reinforcement learn-  
603 ing for reasoning. *arXiv preprint arXiv:2504.16656*, 2025b.
- 604  
605 Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xiong-Hui Chen,  
606 Jianxin Yang, Zhenru Zhang, Yuqiong Liu, An Yang, Andrew Zhao, Yang Yue, Shiji Song, Bowen  
607 Yu, Gao Huang, and Junyang Lin. Beyond the 80/20 rule: High-entropy minority tokens drive  
608 effective reinforcement learning for LLM reasoning. In *The Thirty-ninth Annual Conference on*  
609 *Neural Information Processing Systems*, 2025c. URL [https://openreview.net/forum?](https://openreview.net/forum?id=yfcpdY4gMP)  
610 [id=yfcpdY4gMP](https://openreview.net/forum?id=yfcpdY4gMP).
- 611 Xiyao Wang, Zhengyuan Yang, Chao Feng, Hongjin Lu, Linjie Li, Chung-Ching Lin, Kevin Lin,  
612 Furong Huang, and Lijuan Wang. SoTA with less: MCTS-guided sample selection for data-  
613 efficient visual reasoning self-improvement. In *The Thirty-ninth Annual Conference on Neural*  
614 *Information Processing Systems*, 2025d. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=PHu9xJeAum)  
615 [PHu9xJeAum](https://openreview.net/forum?id=PHu9xJeAum).
- 616 Zhenhailong Wang, Xuehang Guo, Sofia Stoica, Haiyang Xu, Hongru Wang, Hyeonjeong Ha, Xiusi  
617 Chen, Yangyi Chen, Ming Yan, Fei Huang, et al. Perception-aware policy optimization for mul-  
618 timodal reasoning. *arXiv preprint arXiv:2507.06448*, 2025e.
- 619  
620 Lai Wei, Yuting Li, Kaipeng Zheng, Chen Wang, Yue Wang, Linghe Kong, Lichao Sun, and Weiran  
621 Huang. Advancing multimodal reasoning via reinforcement learning with cold start. *arXiv*  
622 *preprint arXiv:2505.22334*, 2025.
- 623 Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan Gao, Kun Yan, Sheng-ming Yin, Shuai  
624 Bai, Xiao Xu, Yilei Chen, et al. Qwen-image technical report. *arXiv preprint arXiv:2508.02324*,  
625 2025.
- 626 Tong Xiao, Xin Xu, Zhenya Huang, Hongyu Gao, Quan Liu, Qi Liu, and Enhong Chen. Advancing  
627 multimodal reasoning capabilities of multimodal large language models via visual perception  
628 reward. *arXiv preprint arXiv:2506.07218*, 2025.
- 629  
630 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,  
631 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint*  
632 *arXiv:2505.09388*, 2025a.
- 633 Yi Yang, Xiaoxuan He, Hongkun Pan, Xiyan Jiang, Yan Deng, Xingtao Yang, Haoyu Lu, Dacheng  
634 Yin, Fengyun Rao, Minfeng Zhu, Bo Zhang, and Wei Chen. R1-onevision: Advancing general-  
635 ized multimodal reasoning through cross-modal formalization. In *Proceedings of the IEEE/CVF*  
636 *International Conference on Computer Vision (ICCV)*, pp. 2376–2385, October 2025b.
- 637 Zheng Yaowei, Lu Junting, Wang Shenzhi, Feng Zhangchi, Kuang Dongdong, and Xiong Yuwen.  
638 Easyr1: An efficient, scalable, multi-modality rl training framework. [https://github.com/](https://github.com/hiyoga/EasyR1)  
639 [hiyoga/EasyR1](https://github.com/hiyoga/EasyR1), 2025.
- 640  
641 En Yu, Kangheng Lin, Liang Zhao, jisheng yin, Yana Wei, Yuang Peng, Haoran Wei, Jian-  
642 jian Sun, Chunrui Han, Zheng Ge, Xiangyu Zhang, Daxin Jiang, Jingyu Wang, and Wenbing  
643 Tao. Perception-r1: Pioneering perception policy with reinforcement learning. In *The Thirty-*  
644 *ninth Annual Conference on Neural Information Processing Systems*, 2025a. URL [https://](https://openreview.net/forum?id=BeXcXrXetA)  
645 [openreview.net/forum?id=BeXcXrXetA](https://openreview.net/forum?id=BeXcXrXetA).
- 646 Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian  
647 Fan, Gaohong Liu, Lingjun Liu, et al. Dapo: An open-source llm reinforcement learning system  
at scale. *arXiv preprint arXiv:2503.14476*, 2025b.

648 Tianyu Yu, Bo Ji, Shouli Wang, Shu Yao, Zefan Wang, Ganqu Cui, Lifan Yuan, Ning Ding, Yuan  
649 Yao, Zhiyuan Liu, et al. Rlpr: Extrapolating rlvr to general domains without verifiers. *arXiv*  
650 *preprint arXiv:2506.18254*, 2025c.

651  
652 Jingyi Zhang, Jiaxing Huang, Huanjin Yao, Shunyu Liu, Xikun Zhang, Shijian Lu, and Dacheng Tao.  
653 R1-vl: Learning to reason with multimodal large language models via step-wise group relative  
654 policy optimization. In *Proceedings of the IEEE/CVF International Conference on Computer*  
655 *Vision (ICCV)*, pp. 1859–1869, October 2025a.

656 Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou,  
657 Pan Lu, Kai-Wei Chang, Yu Qiao, et al. Mathverse: Does your multi-modal llm truly see the  
658 diagrams in visual math problems? In *European Conference on Computer Vision*, pp. 169–186.  
659 Springer, 2024.

660 Yi-Fan Zhang, Xingyu Lu, Xiao Hu, Chaoyou Fu, Bin Wen, Tianke Zhang, Changyi Liu, Kaiyu  
661 Jiang, Kaibing Chen, Kaiyu Tang, et al. R1-reward: Training multimodal reward model through  
662 stable reinforcement learning. *arXiv preprint arXiv:2505.02835*, 2025b.

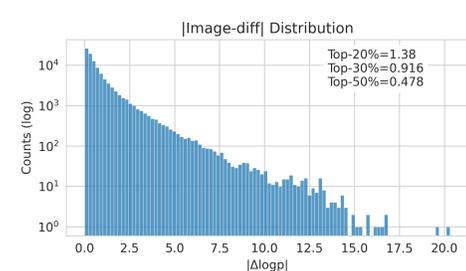
663 Yi-Fan Zhang, Xingyu Lu, Shukang Yin, Chaoyou Fu, Wei Chen, Xiao Hu, Bin Wen, Kaiyu  
664 Jiang, Changyi Liu, Tianke Zhang, et al. Thyme: Think beyond images. *arXiv preprint*  
665 *arXiv:2508.11630*, 2025c.

666  
667 Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen  
668 Duan, Weijie Su, Jie Shao, et al. Internvl3: Exploring advanced training and test-time recipes for  
669 open-source multimodal models. *arXiv preprint arXiv:2504.10479*, 2025.

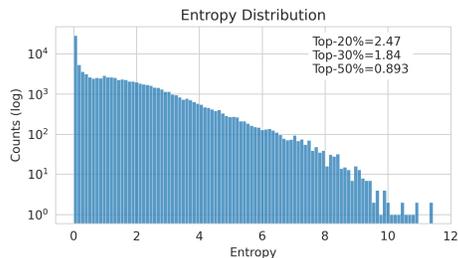
## 671 A USAGE OF LLMs

672  
673 We employed large language models (LLMs) solely for polishing the language of this paper. All  
674 content was originally drafted by the authors, and the use of LLMs was limited to refining pre-  
675 organized text and paragraphs. All suggested modifications were carefully reviewed by the authors  
676 to ensure accuracy and consistency with the intended meaning.

## 678 B TOKEN DISTRIBUTION AND CHARACTERISTIC ANALYSIS



691 Figure 8: Distribution of log-probability differ-  
692 ences for Qwen-2.5-VL-7B on Hallusion-  
693 Bench, used to identify perception-related to-  
694 kens (Guan et al., 2024).



691 Figure 9: Distribution of entropy values for  
692 Qwen-2.5-VL-7B on HallusionBench, used  
693 to identify reasoning-related tokens (Guan  
694 et al., 2024).

695 To further investigate the characteristics of perception- and reasoning-related tokens, we analyze  
696 the token distribution of Qwen-2.5-VL-7B on both perception and reasoning benchmarks. Specif-  
697 ically, we report the entropy distribution (Figures 9, 11) and the image log-probability difference  
698 distribution (Figures 8, 10) on *HallusionBench* (Guan et al., 2024) and *WeMath* (Qiao et al., 2024),  
699 respectively. For each case, we further highlight the values corresponding to the top 20%, 30%, and  
700 50% of tokens, ranked by their entropy or log-probability differences.

701 Our observations are as follows:

702  
703  
704  
705  
706  
707  
708  
709  
710  
711

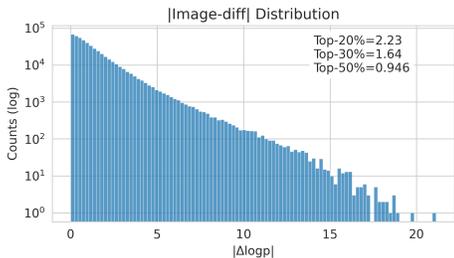


Figure 10: Distribution of log-probability differences for Qwen-2.5-VL-7B on WeMath, used to identify perception-related tokens (Qiao et al., 2024).

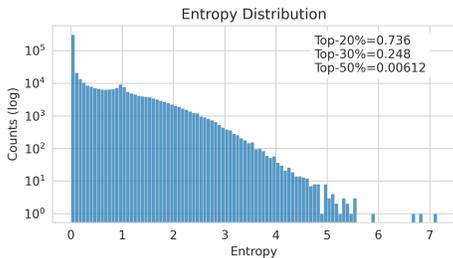


Figure 11: Distribution of entropy values for Qwen-2.5-VL-7B on WeMath, used to identify reasoning-related tokens (Qiao et al., 2024).

712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729

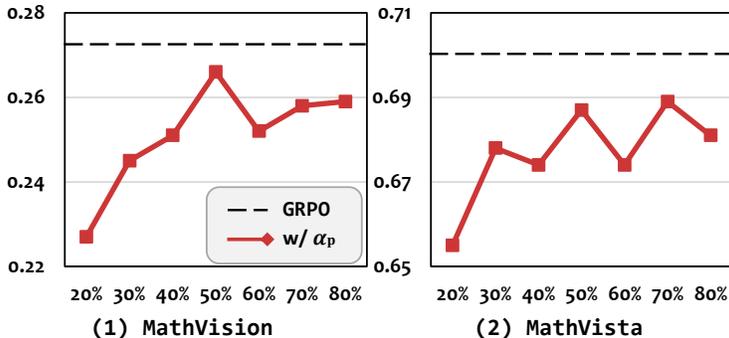


Figure 12: Additional experimental results with different ratios of perception-related tokens ( $\alpha_p$ ) on *MathVision* and *MathVista*.

730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750

① Perception-related benchmark (Guan et al., 2024). High-entropy tokens exhibit large variations, with more than 50% of tokens exceeding 0.89. In contrast, the image log-probability difference of perception-related tokens is relatively stable: more than 70% of tokens are less than 0.916. This indicates that a small set of perception-related tokens is consistently stable and plays a critical role in capturing visual cues, suggesting that focusing on these tokens effectively enhances the model’s perceptual ability.

② Reasoning-related benchmark (Qiao et al., 2024). In reasoning tasks, the distribution shows the opposite trend. Perception-related tokens present large variations, with more than 50% exceeding 0.946. Meanwhile, high-entropy tokens remain relatively stable, with over 70% less than 0.248. This implies that a small subset of reasoning-related tokens is more robust and can effectively capture the key reasoning process, highlighting their importance in enhancing the model’s reasoning capability.

These results demonstrate that perception and reasoning rely on different types of stable tokens: perception emphasizes stability in a small number of visually sensitive tokens, while reasoning relies on the robustness of high-entropy tokens. This contrast validates our token re-weighting strategy that explicitly leverages both perception- and reasoning-related tokens.

### C ADDITIONAL EXPERIMENTAL RESULTS

751  
752  
753  
754  
755

In this section, we provide additional ablation results on two reasoning-focused benchmarks, *MathVision* (Wang et al., 2024) and *MathVista* (Lu et al., 2024). Specifically, we evaluate the impact of different ratios of perception-related tokens ( $\alpha_p$ ), reasoning-related tokens ( $\alpha_r$ ), and their combinations. The results are illustrated in Figure 12, Figure 13, and Figure 14, respectively.

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767

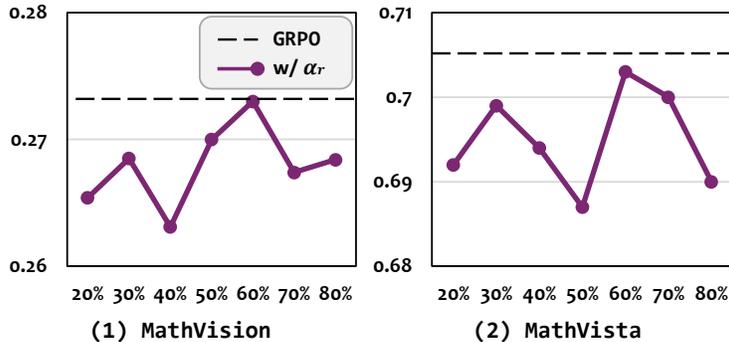


Figure 13: Additional experimental results with different ratios of reasoning-related tokens ( $\alpha_r$ ) on *MathVision* and *MathVista*.

771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783

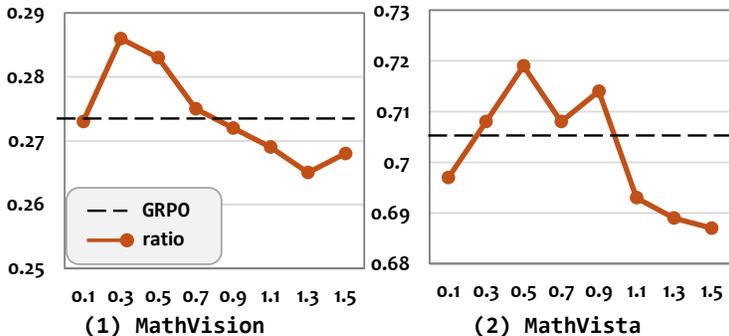


Figure 14: Additional experimental results with combined ratios of perception- and reasoning-related tokens ( $\gamma_p = 1, \gamma_r$  from 0.1 to 1.5) on *MathVision* and *MathVista*.

From Figure 12 and Figure 13, we can observe that relying solely on perception- or reasoning-related tokens underperforms the full GRPO optimization, while focusing on reasoning tokens achieves relatively better results than perception tokens. Moreover, from Figure 14, the re-weighting strategy consistently improves performance, and the effectiveness of using 0.5 perception tokens across both benchmarks further validates the robustness of our re-weighting strategy.

#### D PERCEPTION-REASONING INTERDEPENDENCE

In this section, we demonstrate the interdependence between perception and reasoning tokens. Specifically, we visualize the relationship between reasoning uncertainty and perception strength over the training and validation set of Geo3K (Figure 15 & Figure 16), where we can observe a strong push-pull dynamic relationship between the reasoning and perception. Moreover, for clear illustration, we illustrate the relationship in Figure 17, and the learning dynamics as in Figure 18.

#### E PERCEPTION TOKEN IDENTIFICATION WITH VARIOUS CRITERIA

In this section, we compare different perception tokens selection criteria, including: “logp-diff”, “probs-diff”, “entropy-diff”, and “attention scores” to the image, Experimental results across various datasets are listed in Figure 19 → Figure 27, where we can observe that “logp-diff” better trades how much the image matter and how meaningful is the change.

803  
804  
805  
806  
807  
808  
809

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

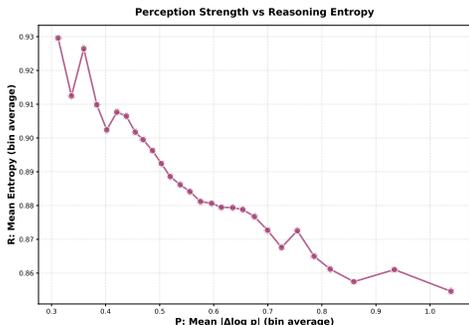


Figure 15: Relationship between (1) reasoning uncertainty: the entropy value of reasoning tokens and (2) perceptron strength: the logp diff between perception tokens for each response over Geo3K training set, response sampled from Qwen-2.5-VL 7B.

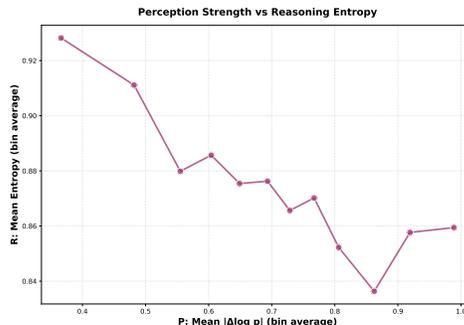


Figure 16: Relationship between (1) reasoning uncertainty: the entropy value of reasoning tokens and (2) perceptron strength: the logp diff between perception tokens for each response over Geo3K validation set, response sampled from Qwen-2.5-VL 7B.

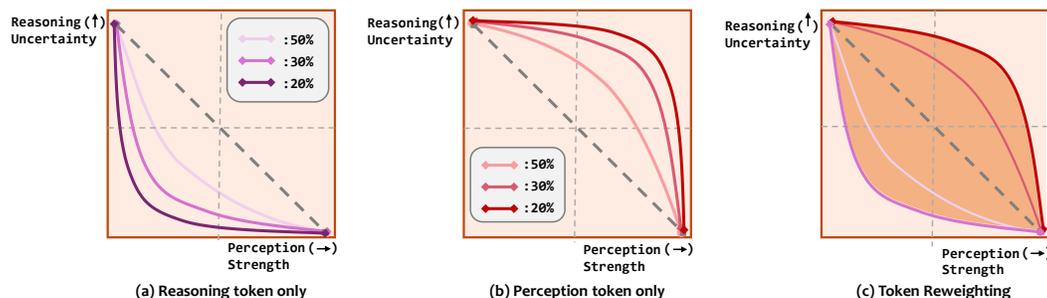


Figure 17: The comparison of perception strengths and reasoning uncertainty with different token ratios for GRPO optimization.

## F DISTRIBUTION OF SELECTED TOKENS OVER VARIOUS ROLLOUTS.

In this section, we show the distribution of selected tokens over a batch of rollouts as in Figure 28. We can find that different rollouts receive a comparable overall amount of optimization, but with different mixtures of tokens: harder groups are optimized more on reasoning, easier groups more on perception, while both token types remain well represented across the rollout batch.

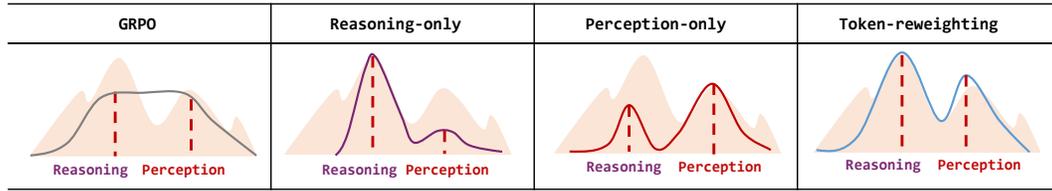


Figure 18: Comparison between vanilla GRPO, GRPO with reasoning-tokens only, GRPO with perception-tokens only, and GRPO with Token-reweighting.

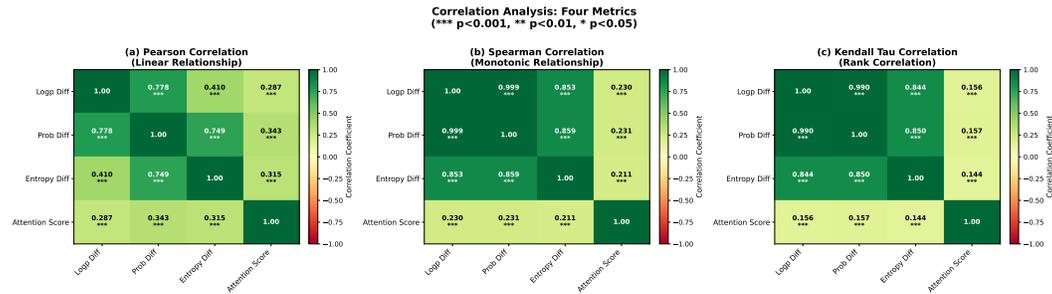


Figure 19: Correlation heatmaps between different image token selection strategies over the validation set.

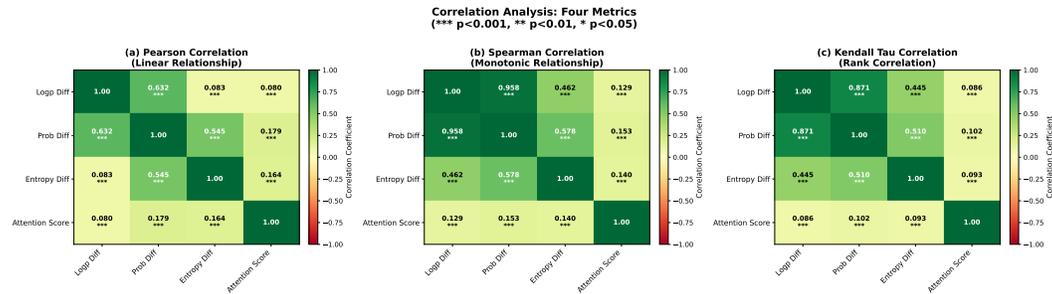


Figure 20: Correlation heatmaps between different image token selection strategies over the halubench benchmark.

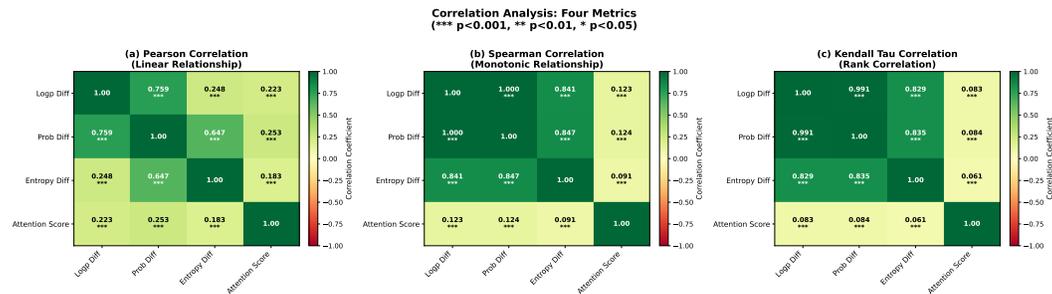


Figure 21: Correlation heatmaps between different image token selection strategies over the wemath benchmark.

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

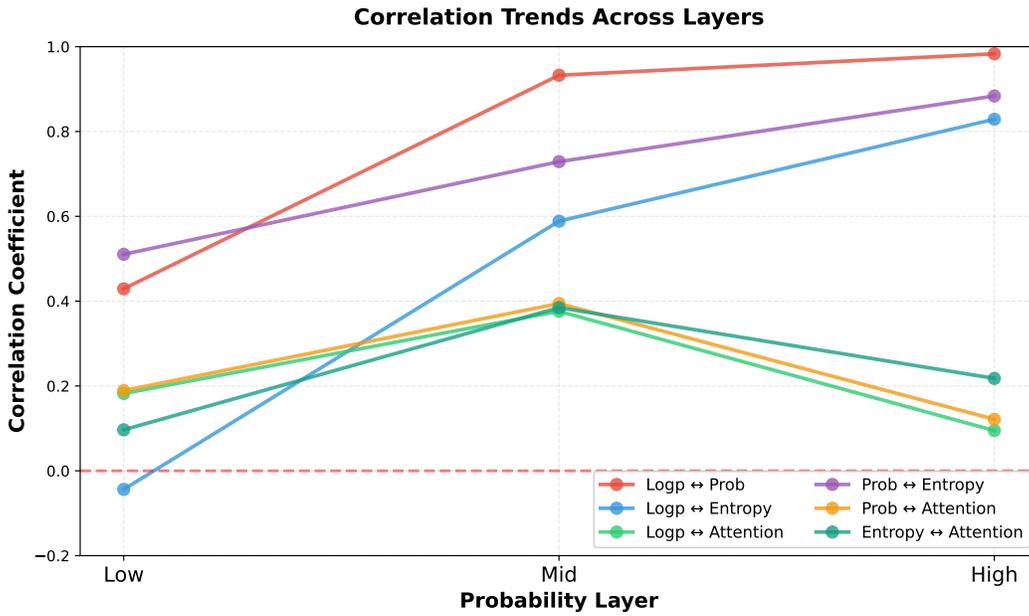


Figure 22: Correlation comparison across layers between different image token selection strategies over the validation set.

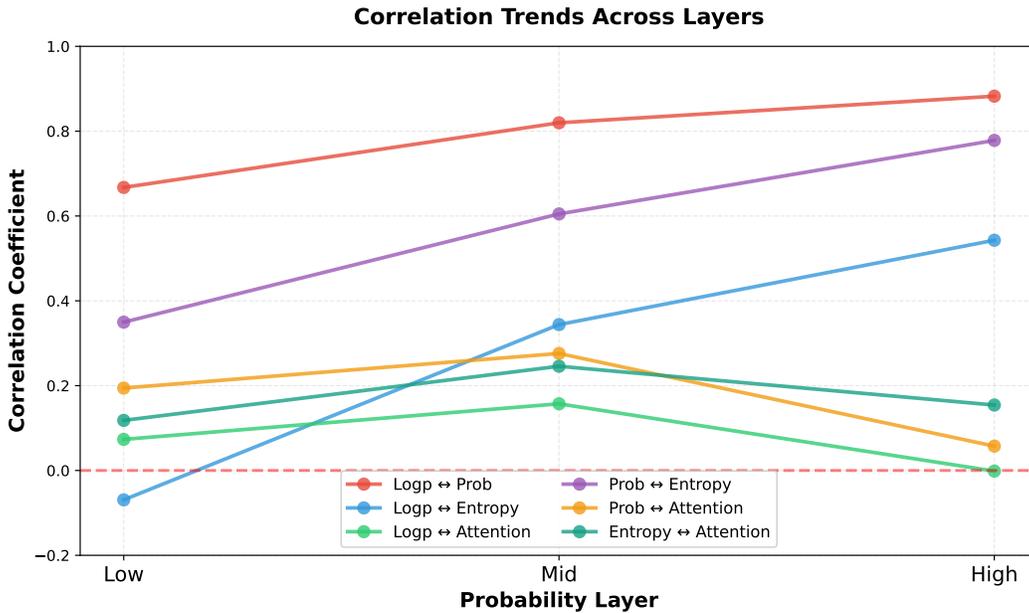


Figure 23: Correlation comparison across layers between different image token selection strategies over the hallubench benchmark.

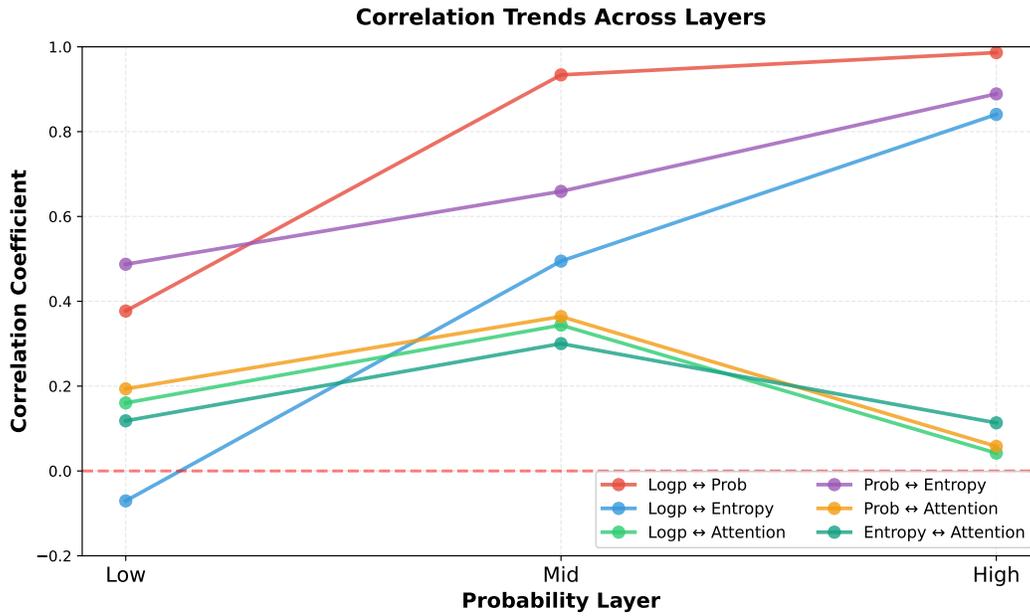


Figure 24: Correlation comparison across layers between different image token selection strategies over the wemath benchmark.

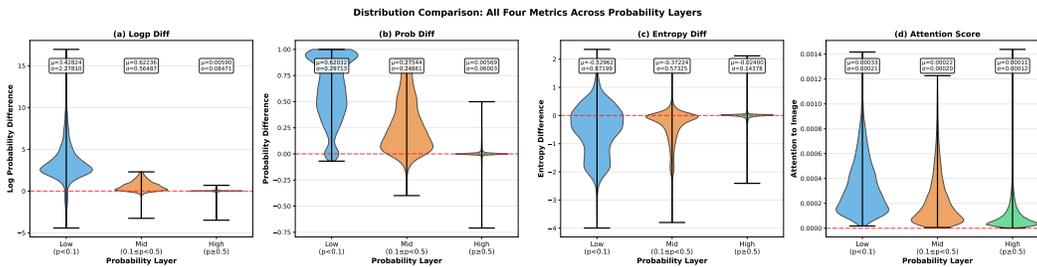


Figure 25: violin plot between different image token selection strategies over the validation set.

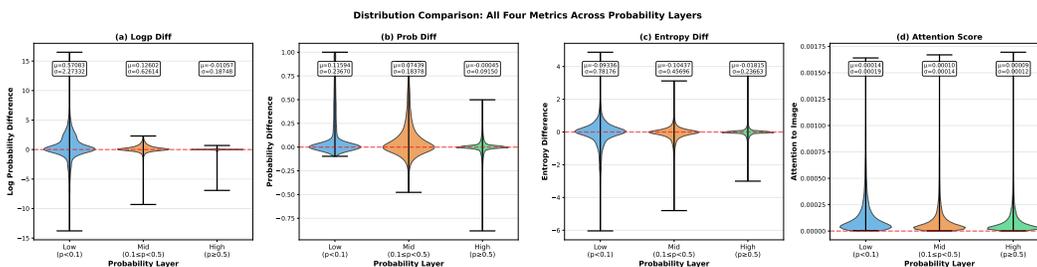


Figure 26: violin plot between different image token selection strategies over the hallubench benchmark.

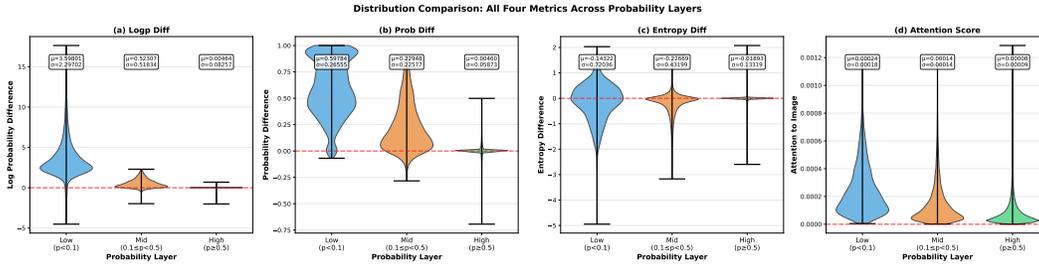


Figure 27: violin plot between different image token selection strategies over the wemath benchmark.

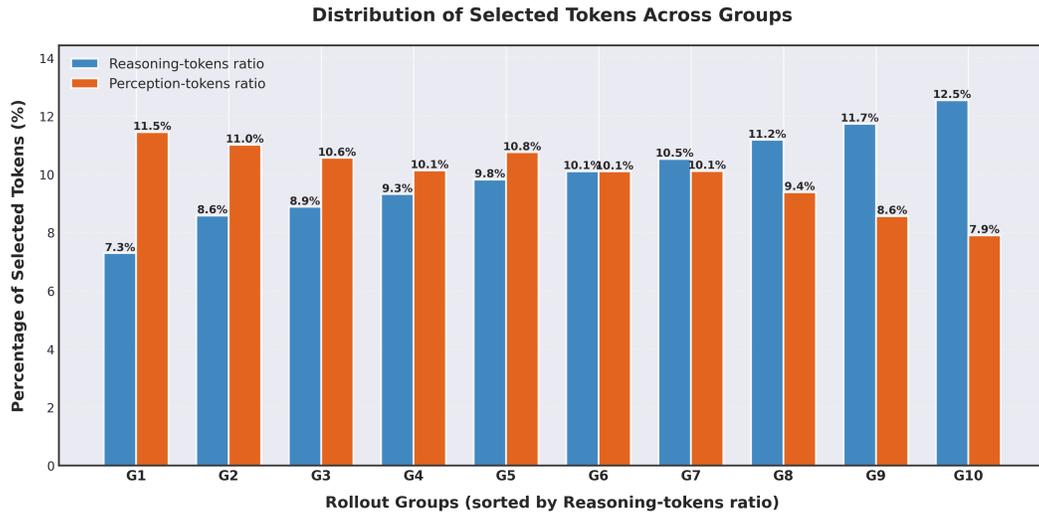


Figure 28: Distribution of selected tokens over the sampled rollouts, where we employ the Qwen-VL-2.5 7B model with a batch of 512 samples, each sample generates 16 rollout responses.