# VisionThink: Smart and Efficient Vision Language Model via Reinforcement Learning

**Senqiao Yang**[* 1]    **Junyi Li**[* 2]    **Xin Lai**[* 1]    **Bei Yu**[1]    **Hengshuang Zhao**[2]    **Jiaya Jia**[1,3]

[1]CUHK        [2]HKU        [3]HKUST

Codes and models: `https://github.com/dvlab-research/VisionThink`

## Abstract

Recent advancements in vision-language models (VLMs) have improved performance by increasing the number of visual tokens, which are often significantly longer than text tokens. However, we observe that most real-world scenarios do not require such an extensive number of visual tokens. While the performance drops significantly in a small subset of OCR-related tasks, models still perform accurately in most other general VQA tasks with only 1/4 resolution. Therefore, we propose to dynamically process distinct samples with different resolutions, and present a new paradigm for visual token compression, namely, VisionThink. It starts with a downsampled image and smartly decides whether it is sufficient for problem solving. Otherwise, the model could output a special token to request the higher-resolution image. Compared to existing Efficient VLM methods that compress tokens using fixed pruning ratios or thresholds, VisionThink autonomously decides whether to compress tokens case by case. As a result, it demonstrates strong fine-grained visual understanding capability on OCR-related tasks, and meanwhile saves substantial visual tokens on simpler tasks. We adopt reinforcement learning and propose the LLM-as-Judge strategy to successfully apply RL to general VQA tasks. Moreover, we carefully design a reward function and penalty mechanism to achieve a stable and reasonable image resize call ratio. Extensive experiments demonstrate the superiority, efficiency, and effectiveness of our method.

## 1 Introduction

Recently, Vision-Language Models (VLMs) [31, 30, 33, 9, 3] have achieved remarkable performance in general visual question answering (General VQA) and various real-world scenarios by projecting and adapting visual tokens into the LLM space [61, 1, 93, 4]. However, as the performance of VLMs continues to improve, the consumption of visual tokens increases exponentially. For example, a 2048×1024 photo taken with a smartphone required 576 visual tokens in LLaVA 1.5 [33], but now requires 2,678 visual tokens in Qwen2.5-VL [5]. Therefore, it is imperative to avoid the excessive use of visual tokens.

Numerous works on visual token compression have been proposed [65, 55, 18, 22, 8, 91]. Most approaches prune or merge a fixed number of visual tokens using predetermined thresholds. However, redundancy levels vary across different questions and images, leading to a natural question: *Should we really apply a uniform token compression ratio across all scenarios?*

To answer this question, we simply reduced the image resolution to decrease the number of visual tokens and evaluated Qwen2.5-VL's[5] performance on several benchmarks. As shown in the left of Fig. 1, we found that for most real-world scenarios, such as MME and RealWorldQA, even reducing the image resolution by a factor of four, which significantly cuts visual tokens by 75%, has minimal impact on the model's performance. However, as shown in the right of Fig. 1, for benchmarks such
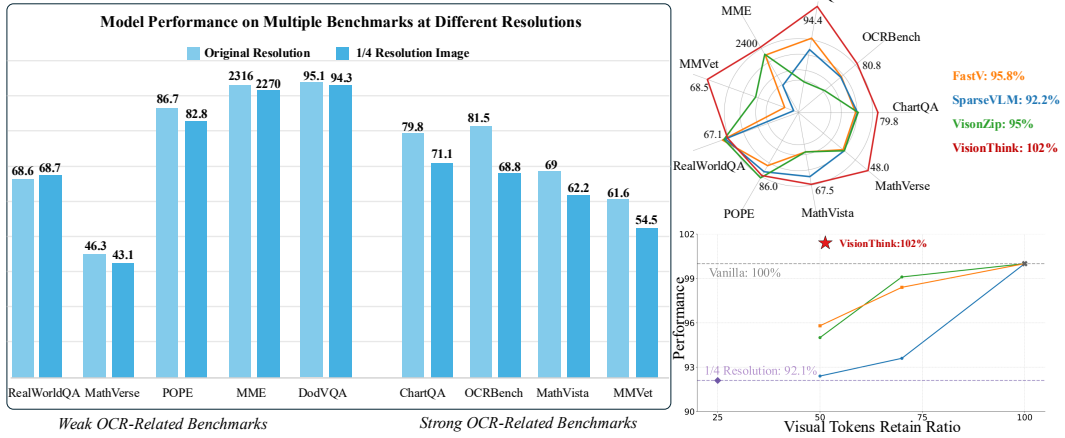
---

[*]Equal Contribution

Figure 1: **Our key observations and VisionThink performance and efficiency**. **Left**: We find that in most general scenarios, even reducing visual tokens by a factor of four results in only minimal performance drop. However, token compression leads to a significant performance drop on strong OCR-related benchmarks. **Right**: Our VisionThink significantly outperforms previous work in both performance and efficiency.

as ChartQA and OCRBench, which require detailed understanding and OCR-related capabilities, reducing the number of visual tokens leads to a significant drop in performance. Based on these observations, we find that most real-world questions do not require high-resolution images with long visual tokens, while a small subset of OCR-related tasks demand such detailed input much. Therefore, there is significant potential for efficiency optimization if we can dynamically distinguish between samples that require high-resolution processing and those that do not.

In this paper, we propose VisionThink, a new EfficientVLM paradigm that leverages the model's reasoning capabilities. Unlike prior methods that process full images and later discard redundant tokens, VisionThink directly inputs compressed visual tokens and allows the model to request the original high-resolution image when needed. This enables more efficient inference in most real-world scenarios, and meanwhile preserving performance on OCR-related tasks.

Although VisionThink offers a promising way to handle samples with varying levels of visual redundancy smartly, it still faces two key challenges:

**Effective Reinforcement Learning for General VQA.** Conventional rule-based reinforcement learning algorithms, typically used to optimize reasoning process, struggle with the diversity and complexity of general VQA. To overcome this issue, we propose the LLM-as-Judge approach, enabling semantic matching. Experiments show performance improvement across several general VQA benchmarks, highlighting the potential to extend vision-based reinforcement learning beyond visual math reasoning to broader VQA tasks.

**Determine When High Resolution is Worth.** To improve efficiency without compromising performance, the model must accurately determine when high-resolution input is necessary. We achieve this by carefully designing a balanced reward function to prevent the model from collapsing into always requiring high-resolution images or always using low-resolution images. With this mechanism, VisionThink maintains strong performance on OCR benchmarks while delivering significant speed-ups on non-OCR benchmarks, achieving up to 100% for DocVQA.

Overall, we present a simple yet effective pipeline—VisionThink. It introduces a new approach to visual token compression by dynamically determining compression based on the content of each sample, thereby achieving efficiency gains at the sample level. Consequently, it is compatible with other advanced spatial-level methods. We hope our work sheds new light on this area.

## 2 Preliminary

### 2.1 Large Language Models and Reinforcement Learning

Recent progress in improving the reasoning ability of large language models (LLMs)[16, 21] has shown that Reinforcement Learning (RL) is an effective training approach. In this work, we use Group Relative Policy Optimization (GRPO)[52] as our training method. GRPO removes the need for a separate critic model by using group scores to estimate baselines. This reduces computation cost, improves training stability, and leads to faster and more reliable performance gains.

During the training process, GRPO samples a group of outputs $\{o_1, o_2, \cdots, o_G\}$ based on the given question $q$ from the old policy $\pi_{\theta_{old}}$ and then optimizes the policy model $\pi_\theta$ by maximizing the following objective:

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}_{[q \sim \mathcal{D}, \{o_i\}_{i=1}^G \sim \pi_{\theta_{old}}(\cdot|q)]}$$

$$\frac{1}{G} \sum_{i=1}^G \left( \min \left( \frac{\pi_\theta(o_i|q)}{\pi_{\theta_{old}}(o_i|q)} A_i, \mathrm{clip} \left( \frac{\pi_\theta(o_i|q)}{\pi_{\theta_{old}}(o_i|q)}, 1 - \epsilon, 1 + \epsilon \right) A_i \right) - \beta \mathbb{D}_{KL} \left( \pi_\theta || \pi_{ref} \right) \right) \tag{1}$$

$$\mathbb{D}_{KL} \left( \pi_\theta || \pi_{ref} \right) = \frac{\pi_{ref}(o_i|q)}{\pi_\theta(o_i|q)} - \log \frac{\pi_{ref}(o_i|q)}{\pi_\theta(o_i|q)} - 1, \tag{2}$$

where $q$ represents the input questions drawn from the dataset $\mathcal{D}$, and $o$ denotes the generated text response. $\mathbb{D}_{KL}$ is the KL-divergence measure, while $\epsilon$ and $\beta$ are hyper-parameters. $A_i$ indicates the advantage, computed using a group of rewards $\{r_1, r_2, \ldots, r_G\}$ corresponding to the outputs within each group.

### 2.2 Computation Complexity

To evaluate the computational complexity of VLMs, we analyze key components, including the self-attention mechanism and the feedforward network (FFN). The total floating-point operations (FLOPs) are given by:

$$\text{Total FLOPs} = T \times (4nd^2 + 2n^2d + 2ndm)$$

where $T$ denotes the number of transformer layers, $n$ is the sequence length, $d$ is the size of the hidden dimension, and $m$ is the intermediate size of the FFN.

This equation shows that computational complexity is heavily driven by the sequence length $n$. In general VLM tasks, the sequence length $n$ is composed of $n_{sys} + n_{img} + n_{question}$, where $n_{img}$ is often significantly larger than the other two components. Hence, controlling the number of image tokens is key to achieving VLM efficiency.

## 3 Methodology

### 3.1 Overview

Our objective is to develop a smart and efficient VLM, capable of autonomously determining whether the information in the given image is sufficient to answer the question accurately. As shown in Fig. 2, the pipeline first processes a low-resolution image to minimize the computataion cost. It then smartly requests original high-resolution inputs when the information in the downsampled image is insufficient to answer the question. Ideally, this strategy maintains high performance while sharply reducing computational load. To achieve this goal, we must address two key challenges:

**Effective RL on General VQA.** Due to the diversity and complexity of general VQA, traditional rule-based RL algorithms are not directly applicable. To address this, we propose an LLM-as-Judge strategy, in which a large language model guides and evaluates the RL training process (Sec. 3.2). We further extend the Multi-Turn GRPO algorithm to suit our setting (Sec. 3.3).

**Enabling the model to decide when high resolution is necessary.** The model must learn to assess whether a downsampled image contains sufficient information to answer the question if the original
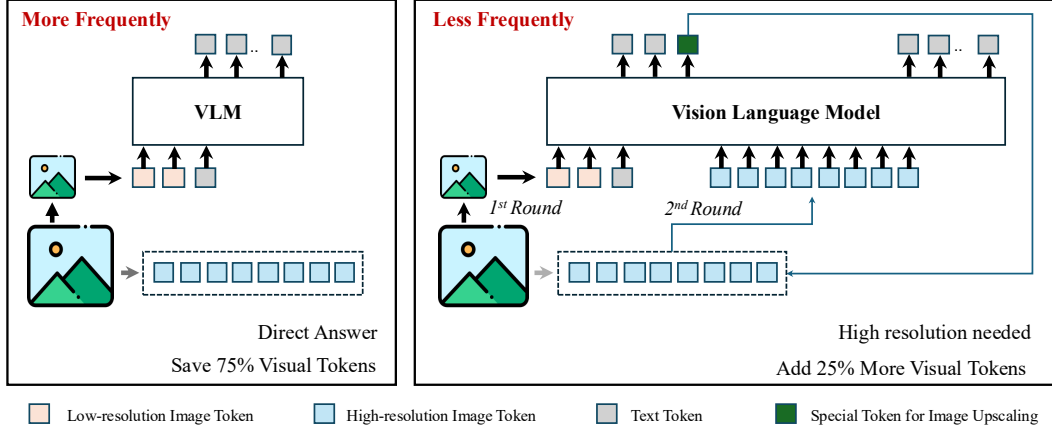
Figure 2: **Framework of VisionThink.** (a) The left image illustrates VisionThink processing an image with resolution reduced by a factor of four, where the VLM directly provides an answer. (b) The right image shows a case where the model detects insufficient information and requests a high-resolution image to answer the question.

high-resolution image is required. So that the model could balance the efficiency and performance. To this end, we design a reward function that encourages optimal resolution decisions (Sec. 3.4) and collect training data across multiple resolutions to support effective learning (Sec. 3.5).

## 3.2 LLM-as-Judge for General VQA

**Challenges.** One of the central challenges in applying reinforcement learning to General VQA lies in evaluating model responses, especially when answers are open-ended or context-dependent. Most existing multi-modal RL efforts remain limited to structured tasks such as visual math, where ground-truth answers can be easily defined and verified via rules or exact matching. However, this approach breaks down in General VQA settings, where the diversity and ambiguity of valid answers make rule-based verification infeasible.

**Pure Text Accuracy Judgement.** To address this, we employ an external LLM as a judgment evaluator. Leveraging its broad knowledge and language understanding, the LLM assesses the correctness of model outputs in a human-aligned and flexible manner. Importantly, the evaluation is conducted purely in text by comparing the model's answer with the ground-truth. This design avoids biases from visual content and the limitations of VLM performance. Furthermore, to minimize potential misjudgment by the evaluator, the reward is discrete (either 0 or 1) rather than continuous. The detailed judgment prompt is shown in Appendix B.1.

**Effectiveness.** The LLM-as-Judge is flexible, one advantage is that most of the SFT data could be used. To verify the effectiveness of our proposed LLM-as-Judge, we collected 130K samples, which can be directly used to train the model with GRPO, without requiring any cold-start process. The results show significant improvement compared to the base model, Qwen2.5VL-Instruct. Further details are provided in Appendix B.5.

## 3.3 Mutli-Turn Training Algorithm

**Multi-Turn GRPO.** In our VisionThink framework, we first input the question and the downsampled image into the VLM. If the information is insufficient to answer the current question, the model will autonomously request a higher-resolution image and generate a new response. This process is essentially a multi-turn interaction. Therefore, we extend the original GRPO (Eq. 1) to a multi-turn GRPO, as shown in Eq. 3:
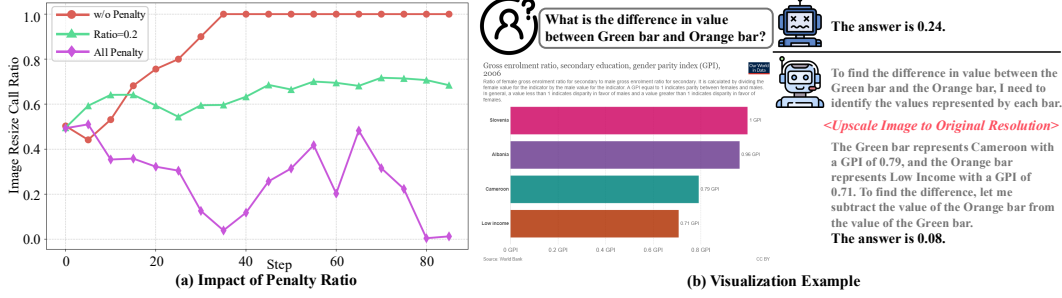
Figure 3: (a) Impact of the Penalty Ratio. Applying a penalty to all resize image requests or removing the penalty entirely will both lead to model collapse. (b) VisionThink correctly solves OCR-related problems by autonomously requesting high-resolution images.

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}_{q \sim \mathcal{D}, \{o_i\}_{i=1}^{G} \sim \pi_{\text{old}}(\cdot|q;\mathcal{I})} \left[ \frac{1}{G} \sum_{i=1}^{G} \frac{1}{\sum_{t=1}^{|o_i|} \mathbb{I}(o_{i,t})} \sum_{t=1}^{|o_i|} \mathbb{I}(o_{i,t}) \right.$$

$$\left. \cdot \min\left( p_{i,t}\hat{A}_{i,t}, \text{clip}\left( p_{i,t}, 1-\epsilon, 1+\epsilon \right)\hat{A}_{i,t} \right) - \beta \mathbb{D}_{KL}\left[ \pi_\theta || \pi_{\text{ref}} \right] \right], \tag{3}$$

where $p_{i,t} = \frac{\pi_\theta(o_{i,t}|q, o_{i,<t};\mathcal{I})}{\pi_{\text{old}}(o_{i,t}|q, o_{i,<t};\mathcal{I})}$, and $\mathbb{I}(o_t)$ is the token loss masking operation such that $\mathbb{I}(o_t) = 1$ if $o_t$ is the generated token from LLM and $\mathbb{I}(o_t) = 0$ if $o_t$ is the response token from the called tools. Intuitively, we masked all text and image tokens from the user and performed optimization solely based on the multi-turn output tokens generated by the VLM.

**How does the model signal the need for a high-resolution image?** To determine when the model requires a high-resolution image, we modify the prompt to instruct the model to output specific special tokens. Notably, this is a non-trivial process because our training does not introduce any cold-start phase, which leads to a performance drop in general VQA (Appendix C.2). Therefore, selecting an appropriate and effective prompt at the early stage of training is crucial. The prompt must ensure that the model is capable of outputting the required special tokens during multi-turn rollouts in a zero-shot setting. Otherwise, GRPO will fail to optimize correctly due to the absence of gradients. We conduct comparative ablation studies in Appendix C.3 and find that the Agent Prompt recommended by Qwen-2.5VL [5] suits VisionThink best. The prompt details are provided in Appendix B.1.

### 3.4 Reward Design

Different reward functions can lead the model toward different optimization directions and final performance outcomes. The reward function in our VisionThink framework consists of three components:

$$\mathcal{R}_{\text{overall}} = \mathcal{R}_{\text{accuracy}} + \mathcal{R}_{\text{format}} - \mathcal{P}_{\text{control}}, \tag{4}$$

where $\mathcal{R}$ represents the reward and $\mathcal{P}$ represents the penalty.

**Accuracy Reward.** The accuracy reward is designed following Sec.3.2, using LLM-as-Judge, where 0 indicates an incorrect answer and 1 indicates a correct answer.

**Format Reward.** To maintain the model's instruction-following capability and ensure that the trained model can more accurately call the image resize function, we apply a format reward. Specifically, we require the reasoning process to be enclosed in "<think></think>" tags, the final answer in "<answer></answer>" tags, and the function call to conform to the JSON format specified in Appendix B. If any of these formats are incorrect, the format score is 0. Only when all formats are correct can the model achieve the full format score of 0.5.

**Penalty Control.** The design of the penalty is a key component of the reward function. As shown in Fig. 3(a), since using high-resolution images generally improves performance, without any penalty, the model tends to collapse into always requesting high-resolution images. To prevent this, we initially

Table 1: **Effective Performance Compared to the Sota Model.** Our model is based on Qwen2.5-VL-7B-Instruct. VisionThink‡ represents a model trained on general VQA tasks using full image resolution with the LLM-as-Judge strategy, which does not contain efficiency capabilities. Qwen2.5-VL-7B* reports the results evaluate by lmms-eval[86].

| Method | MMMU | MMMU-Pro | MMBench | RealWorldQA | POPE | MME | MathVista | MathVerse | MMVet |
|---|---|---|---|---|---|---|---|---|---|
| | val | test | en_test | test | test | test | testmini | testmini | test |
| *Closed-Source Model* | | | | | | | | | |
| GPT-4o [48] | 69.1 | 54.0 | 83.4 | 58.6 | 85.6 | 2329 | 63.8 | 50.2 | 69.1 |
| Claude-3.5 Sonnet [2] | 68.3 | 55.0 | 82.6 | 59.9 | - | 1920 | 67.7 | 41.2 | 70.1 |
| Gemini-1.5-Pro [57] | 62.2 | 49.4 | 73.9 | 70.4 | 88.2 | - | 63.9 | - | 64.0 |
| *Open-Source General Model* | | | | | | | | | |
| Cambrain-1-8B [60] | 42.7 | - | 75.9 | 60.0 | 86.4 | 1803 | 49.0 | - | - |
| InternVL2-8B [12] | 49.3 | 32.5 | 81.7 | 64.4 | 84.2 | 2210 | 58.3 | - | 60.0 |
| LLaVA-OneVision-7B [28] | 48.8 | - | - | 66.3 | 88.4 | 1998 | 63.2 | - | 57.5 |
| MiniCPM-Llama-V-2.5-8B [81] | 45.8 | 19.6 | 77.2 | 63.0 | 86.7 | 2025 | 54.3 | - | - |
| MiniCPM-V-2.6-8B [81] | 49.8 | 27.2 | 78.0 | 65.0 | 83.2 | 2348 | 60.6 | - | - |
| IXC-2.5 [87] | 42.9 | - | 82.2 | 67.8 | - | 2229 | 63.8 | - | 51.7 |
| InternVL2.5-8B [11] | 56.0 | 38.2 | 84.6 | 70.1 | 90.6 | 2344 | 64.4 | 39.5 | 62.8 |
| *Reasoning Model* | | | | | | | | | |
| LLaVA-CoT-11B [71] | - | - | 75.0 | - | - | - | 54.8 | - | 60.3 |
| LLaVA-Reasoner-8B [89] | - | - | - | - | - | - | 50.6 | - | - |
| Insight-V-8B [14] | 50.2 | 24.9 | 82.3 | - | - | 2312 | 59.9 | - | - |
| Mulberry-7B [78] | 55.0 | - | - | - | - | 2396 | 63.1 | - | - |
| Vision-R1-LlamaV-CI-11B [19] | - | - | - | - | - | 2190 | 62.7 | 27.1 | - |
| *VisionThink* | | | | | | | | | |
| Qwen2.5-VL-7B* [5] | 50.3 | 37.7 | 82.6 | 68.6 | 86.7 | 2316 | 68.2 | 46.3 | 61.6 |
| VisionThink ‡ | 51.0 | 40.1 | 82.9 | 68.6 | 87.9 | 2307 | 71.2 | 48.8 | 67.5 |
| VisionThink | 51.2 | 38.9 | 80.0 | 68.5 | 86.0 | 2400 | 67.5 | 48.0 | 67.1 |

followed Search-R1 [23] and applied a 0.1 penalty for correct answers that relied on high-resolution images. However, this approach causes the model to favor direct answers, leading to a collapse where the model relies solely on direct answers, as indicated by the purple line in Fig. 3. The reason is that even blurry, low-resolution images sometimes allow the model to guess the correct answer, and the 0.1 penalty unintentionally reinforced this preference for direct answering.

To address this, we introduce a threshold to control the phenomenon of "lucky guesses". When the probability of correctly answering with a low-resolution image is low, we apply a 0.1 penalty to direct answers to encourage high-resolution requests; conversely, when the probability is high, we penalize high-resolution requests with a 0.1 penalty. In summary, the penalty is designed as below:

$$\mathcal{P}_{control} = 0.1 \cdot \left[ \mathbf{1}_{\text{direct}} \mathbb{I}(r < \theta) + \mathbf{1}_{\text{high}} \mathbb{I}(r \geq \theta) \right], \qquad r = \frac{C_{\text{direct}}}{C_{\text{direct}} + C_{\text{high}}}, \qquad (5)$$

where $C_{\text{direct}}$ and $C_{\text{high}}$ are the correct-answer counts for low- and high-resolution inputs, respectively, and $\mathbf{1}_{\text{action}}$ is the indicator of the chosen action, and we set $\theta$ as 0.2 here. We will discuss the impact of the threshold in Appendix C.3.

## 3.5 Data Preparation

To enable our model can decide when high resolution is necessary, we collect corresponding VQA samples, including both cases requiring high-resolution images and cases adequately answered using downsampled images. To achieve this, we use our base policy model, Qwen2.5VL-Instruct, to perform multiple rollouts on the training dataset and classify the samples based on accuracy. Specifically, we set the temperature to 1 and roll out each sample 8 times. If both the high-resolution and downsampled images yield correct answers in all 8 rollouts, we classify the sample as solvable

using low resolution. Conversely, if the number of correct answers using the high-resolution image exceeds that of the downsampled image by 6 or more, we classify the sample as requiring high resolution. By using the above method, we selected 10K samples that require high-resolution images and 10K samples that do not, to train our model.

## 4 Experiments

### 4.1 Evaluation Setup

**Benchmarks.** We evaluate VisionThink on several general VQA benchmarks, including ChartQA [43], OCRBench [37], MathVista [40], MMVet [83], RealWorldQA [68], and Math-Verse [88], etc. Notably, benchmarks such as ChartQA, OCRBench, and MathVista are strongly OCR-related, requiring the model to possess a high level of detail comprehension. The detailed descriptions of these benchmarks are shown in Appendix B.4.

**Implementation Details.** We conduct experiments based on Qwen2.5-VL-7B-Instruct[5]. For training, we employ veRL[54] framework and use a total batch size of $512$, with a mini-batch size of 32, we set the policy LLM learning rate to $1e-6$ and sample 16 responses per prompt, ensuring a stable and effective training process. For inference, we use the vLLM framework and set the temperature to 0. Further details are shown in Appendix B.3.

### 4.2 Reinforcement Learning Enables VLM to Be More Effective

**Main Results.** To demonstrate the effectiveness of our VisionThink, we compare our VisionThink with the current open-source and closed-source state-of-the-art (sota) method. As shown in Table 6, VisionThink ‡ is used to demonstrate the effectiveness of the LLM-as-Judge strategy on general VQA tasks. It represents a model trained with full image resolution using only accuracy and format rewards, and thus does not incorporate efficiency capabilities. The results show that our VisionThink achieves comparable or even superior performance on general VQA tasks while being more efficient. Specifically, MathVerse and MMVet achieve scores of 48.0 and 67.1, representing improvements of 3.7% and 8.9%, respectively, over the base model. Furthermore, our model performs comparably to closed-source models on several benchmarks such as MathVista and MMBench, and even surpasses all closed-source models on MME, achieving a score of 2400. Besides, as shown in Fig. 3(b), by introducing the LLM-as-Judge for test-time scaling, VisionThink's answer outperforms the vanilla model's short direct answer. Moreover, we scale up the data size to 130K, and further demonstrate the effectiveness of LLM-as-Judge on General VQA Tasks. The results are shown in Appendix B.5.

### 4.3 Reinforcement Learning Enables VLM to Be More Efficient

**Comparison with the Reasoning Model.** To demonstrate the efficiency of our model, we first compare our VisionThink with QwenRL and QwenRL 1/4, both of which are reasoning models trained using the LLM-as-Judge strategy based on Qwen2.5-VL-7B Instruct. QwenRL and QwenRL 1/4 represent inference using the full-resolution image and the 1/4-resolution image, respectively. As shown in Fig. 4, we compare the inference time costs of the three models. Notably, the reported inference times reflect the actual time consumed during vLLM inference, which we believe best represents efficiency in real-world applications. The results show that on most benchmarks, our model's inference time is close to that of QwenRL 1/4, which uses 1/4 of the image tokens, and significantly better than the QwenRL model that processes all image tokens. Specifically, on the DocVQA benchmark, our VisionThink model is more than twice as fast as QwenRL. It also outperforms the baseline by approximately one-third in terms of inference time on benchmarks such as MME and POPE. It is worth noting that on strongly OCR-dependent benchmarks like ChartQA, our model consumes more time than the baseline QwenRL. This is because VisionThink identifies that most questions cannot be answered correctly at low resolution and thus autonomously requests high-resolution images. As a result, the total number of image tokens used by VisionThink exceeds that of the baseline, which we consider reasonable. However, such strongly OCR-dependent benchmarks are relatively rare, so the overall efficiency of VisionThink remains high.

**Comparison with the Previous Efficient VLM.** To further show the effectiveness of our Vision-Think, we compare it with the previous Efficient VLM method FastV and SparseVLM. Notably, all
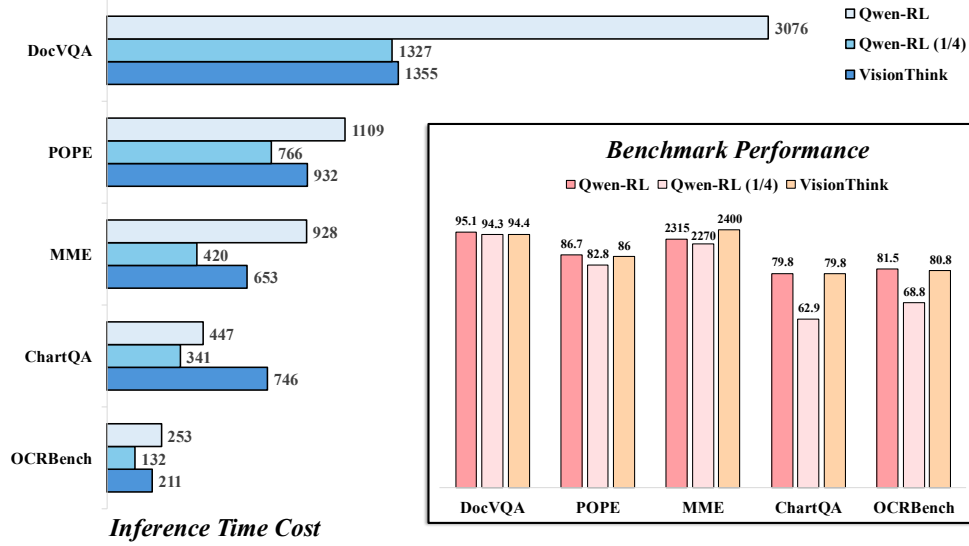
Figure 4: **Inference Time Cost and Benchmark Performance Comparison for Reasoning Model**. Qwen-RL and Qwen-RL (1/4) represent leveraging the LLM-as-Judge on the Qwen2.5-VL-Instruct Model and inference on full resolution image and 1/4 resolution image, respectively.

these methods require computing attention scores to prune visual tokens, which makes them difficult to optimize with FlashAttention2 and may lead to increased memory usage. Furthermore, they are not directly compatible with the efficient inference framework vLLM. Therefore, to ensure a fair comparison, we evaluate model performance while keeping visual token consumption as consistent as possible. As shown in Table 2, our VisionThink outperforms previous methods on average across nine benchmarks. Furthermore, previous approaches require a predefined pruning ratio threshold, whereas our method can autonomously decide whether to reduce tokens based on the question and image content. As a result, on OCR-Related benchmarks such as ChartQA and OCR Bench, our method significantly surpasses FastV and SparseVLM by 9.0% and 8.3%, respectively.

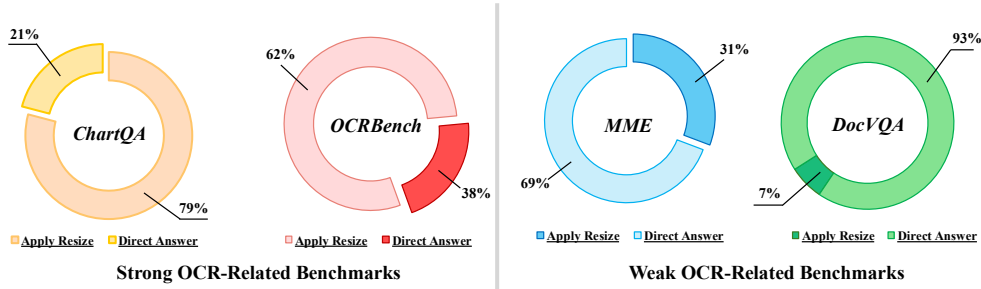## 4.4 Reinforcement Learning Enables VLM to Be Smarter



Figure 5: **VisionThink smartly determine the high-resolution image ratio.** Apply Resize indicates that the model autonomously requests to view the original high-resolution image, while Direct Answer indicates that the model is able to answer the question using only the 1/4-sized image.

In this section, we present the proportion of samples across different benchmarks for which our VisionThink gives direct answers versus those for which it requests high-resolution images. This illustrates the model's ability to smartly determine whether the information in the downsampled image is sufficient. As shown in Fig. 5, we observe that on benchmarks such as ChartQA and OCRBench, which require detailed visual understanding, our model shows a higher ratio of requests for high-resolution images. In contrast, for benchmarks like MME and DocVQA, at least 70% of the samples can be answered directly using low-resolution images at 1/4 of the original resolution.

8

Table 2: **Comparison with Traditional Efficient VLM Methods.** Vanilla represents the Qwen2.5-VL-7B-Instrcut. The retained ratio of the baseline methods is a predefined hyperparameter, while for VisionThink, the ratio is determined autonomously by the model and reported as a statistical value. Note that *Down-Sample* refers to the model's performance when directly fed images with their resolution reduced by half. Additional baseline comparison results are shown in Table. 7

| Method | ChartQA[†] | OCRBench | DocVQA | MME | MMVet | RealWorldQA | POPE | MathVista | MathVerse | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|
| | test | test | val | test | test | test | test | testmini | testmini | |
| *Retain 100% Visual Tokens Across All Benchmarks* | | | | | | | | | | |
| Vanilla | 79.8 | 81.5 | 95.1 | 2316 | 61.6 | 68.6 | 86.7 | 68.2 | 46.3 | 100% |
| | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | |
| *Retain 25% Visual Tokens Across All Benchmarks* | | | | | | | | | | |
| Down-Sample | 62.9 | 68.8 | 94.3 | 2270 | 54.5 | 68.8 | 82.8 | 62.2 | 43.1 | 92.1% |
| | 78.8% | 84.4% | 99.1% | 98.0% | 88.5% | 100.3% | 95.5% | 91.2% | 93.1% | |
| *Retain 50% Visual Tokens Across All Benchmarks* | | | | | | | | | | |
| SparseVLM (ICML 2025) | 73.2 | 75.6 | 66.8 | 2282 | 51.5 | 68.4 | 85.5 | 66.6 | 45.1 | **92.2%** |
| | 91.7% | 92.7% | 70.2% | 98.5% | 83.6% | 99.7% | 98.6% | 97.6% | 97.4% | |
| FastV (ECCV 2024) | 72.6 | 75.8 | 93.6 | 2308 | 52.8 | 68.8 | 84.7 | 63.7 | 45.0 | **95.8%** |
| | 91.0% | 93.0% | 98.4% | 99.6% | 85.7% | 100.3% | 97.7% | 93.4% | 97.2% | |
| *Retain 70% Visual Tokens Across All Benchmarks* | | | | | | | | | | |
| SparseVLM (ICML 2025) | 75.8 | 79.3 | 68.7 | 2276 | 53.7 | 68.5 | 85.4 | 66.3 | 45.1 | **93.6%** |
| | 94.9% | 97.3% | 72.2% | 98.3% | 87.2% | 99.8% | 98.5% | 97.2% | 97.4% | |
| FastV (ECCV 2024) | 77.2 | 82.2 | 94.4 | 2342 | 56.0 | 68.6 | 85.9 | 65.9 | 46.9 | **98.4%** |
| | 96.7% | 100.8% | 99.3% | 101.1% | 90.9% | 100% | 99.1% | 96.6% | 101.3% | |
| *Retain Approximately 51.3% Visual Tokens Across All Benchmarks* | | | | | | | | | | |
| VisionThink | 79.8 | 80.8 | 94.4 | 2400 | 68.5 | 67.1 | 86.0 | 67.5 | 48.0 | <span style="color:red">101.4%</span> |
| | 100% | 99.1% | 99.3% | 103.6% | 111.2% | 97.8% | 99.2% | 99.0% | 103.7% | |

These results align with human intuition: most daily questions do not require high-resolution images, and only OCR-related tasks truly depend on them. Furthermore, to better demonstrate the 'smart' capabilities of VisionThink, we conduct case studies in Appendix D.

## 4.5 Relationship of the EfficientVLM methods and VisionThink.

**Key Differences.** Traditional EfficientVLM methods take a redundant image as input and attempt to remove the redundancy during inference. However, this process typically relies on fixed thresholds, which may yield acceptable performance on standard VQA tasks but result in poor performance on OCR-related or detail-sensitive scenarios, limiting their practical applicability. In contrast, Vision-Think inputs compressed visual tokens and enables the model to autonomously determine whether a higher-resolution image is needed. Ideally, this approach avoids any performance degradation.

**Integration Potential.** Our proposed VisionThink essentially introduces a new paradigm for reading images, which can be integrated with existing Efficient VLMs. In this paper, to provide a straightforward validation of VisionThink, we chose to use image resizing perform token compression. We believe that adopting more advanced token compression techniques could further improve the model's direct answering accuracy, consequently, enhance its overall efficiency. Further discussions are shown in Appendix C.

## 5 Related Works

**Vision Language Model Reasoning.** With the advancement of LLM reasoning capabilities [17], many studies have aimed to improve the reasoning abilities of VLMs [90, 46, 41]. One common approach is using Chain-of-Thought (CoT) prompting to construct SFT datasets. However, the CoTs generated often lack natural human cognitive processes, limiting their effectiveness and generalization. Furthermore, inspired by DeepSeek-R1 [17], some studies have attempted to transfer this reasoning paradigm to vision tasks [79, 59, 20, 77]. However, most current approaches remain limited to the visual math and fail to generalize to general VQA tasks. In contrast, VisionThink successfully applies effective reinforcement learning to general VQA by leveraging the LLM-as-Judge strategy. Due to

space limitations, additional related work on efficient VLMs and LLM-based reasoning is presented in Appendix A.

# 6 Concluding Remarks

## 6.1 Summary

In this work, we introduce VisionThink, a novel paradigm for General VQA that enhances efficiency and performance. By initially processing a downsampled image and using reinforcement learning to selectively upscale to higher resolution when needed, VisionThink optimizes computational resources while preserving accuracy. Leveraging the LLM-as-Judge strategy and a tailored reward function, our approach outperforms prior state-of-the-art models across diverse VQA benchmarks, particularly in tasks requiring fine-grained details like OCR. We believe VisionThink demonstrates the potential of reinforcement learning in vision-language models and encourages the development of more effective and efficient AI systems.

## 6.2 Limitations and Future work

In this work, we focus on the setting of 2x resolution upscaling and at most two turns of conversations and yield promising results. However, it has not been extended to the setting of flexible resolution upscaling. Besides, incorporating more visual tools such as cropping would further bring benefits in both efficiency and performance. Furthermore, multi-turn (for example, more than 5 turns) image tool calls could gain more in solving complex visual problems.

Additionally, our paper utilizes image resizing to reduce the number of visual tokens. This simple method achieves a good balance between performance and efficiency via reinforcement learning. We hope this work inspires further research in the field of efficient reasoning vision language models, especially on making models smarter and more human-like. We will continue to explore the path toward building more general, powerful, and efficient vision-language models.

# References

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv:2303.08774*, 2023.

[2] Anthropic. Claude 3.5 sonnet, 2024.

[3] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-VL: A frontier large vision-language model with versatile abilities. *arXiv:2308.12966*, 2023.

[4] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023.

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

[6] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.

[7] Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17682–17690, 2024.

[8] Liang Chen, Haozhe Zhao, Tianyu Liu, Shuai Bai, Junyang Lin, Chang Zhou, and Baobao Chang. An image is worth 1/2 tokens after layer 2: Plug-and-play inference acceleration for large vision-language models, 2024.

[9] Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv:2311.12793*, 2023.

[10] Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks, 2023.

[11] 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.

[12] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. Internvl2: Better than the best—expanding performance boundaries of open-source multimodal models with the progressive scaling strategy, 2024.

[13] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. *See https://vicuna. lmsys. org (accessed 14 April 2023)*, 2(3):6, 2023.

[14] Yuhao Dong, Zuyan Liu, Hai-Long Sun, Jingkang Yang, Winston Hu, Yongming Rao, and Ziwei Liu. Insight-v: Exploring long-chain visual reasoning with multimodal large language models. *arXiv preprint arXiv:2411.14432*, 2024.

[15] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Zhenyu Qiu, Wei Lin, Jinrui Yang, Xiawu Zheng, et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv:2306.13394*, 2023.

[16] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

[17] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

[18] Yefei He, Feng Chen, Jing Liu, Wenqi Shao, Hong Zhou, Kaipeng Zhang, and Bohan Zhuang. Zipvl: Efficient large vision-language models with dynamic token sparsification and kv cache compression. *arXiv preprint arXiv:2410.08584*, 2024.

[19] 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.

[20] 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.

[21] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.

[22] Yiren Jian, Tingkai Liu, Yunzhe Tao, Chunhui Zhang, Soroush Vosoughi, and Hongxia Yang. Expedited training of visual conditioned language generation via redundancy reduction. *arXiv preprint arXiv:2310.03291*, 2023.

[23] Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *arXiv preprint arXiv:2503.09516*, 2025.

[24] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, Quan Vuong, Thomas Kollar, Benjamin Burchfiel, Russ Tedrake, Dorsa Sadigh, Sergey Levine, Percy Liang, and Chelsea Finn. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024.

[25] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.

[26] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. *arXiv preprint arXiv:2308.00692*, 2023.

[27] Xin Lai, Zhuotao Tian, Yukang Chen, Senqiao Yang, Xiangru Peng, and Jiaya Jia. Step-dpo: Step-wise preference optimization for long-chain reasoning of llms, 2024.

[28] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024.

[29] Jingyao Li, Senqiao Yang, Sitong Wu, Han Shi, Chuanyang Zheng, Hong Xu, and Jiaya Jia. Logits-based finetuning. *arXiv preprint arXiv:2505.24461*, 2025.

[30] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, 2023.

[31] Yanwei Li, Yuechen Zhang, Chengyao Wang, Zhisheng Zhong, Yixin Chen, Ruihang Chu, Shaoteng Liu, and Jiaya Jia. Mini-gemini: Mining the potential of multi-modality vision language models. *arXiv:2403.18814*, 2024.

[32] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. *arXiv:2305.10355*, 2023.

[33] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv:2310.03744*, 2023.

[34] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024.

[35] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 2024.

[36] Jiaming Liu, Mengzhen Liu, Zhenyu Wang, Lily Lee, Kaichen Zhou, Pengju An, Senqiao Yang, Renrui Zhang, Yandong Guo, and Shanghang Zhang. Robomamba: Multimodal state space model for efficient robot reasoning and manipulation. *arXiv preprint arXiv:2406.04339*, 2024.

[37] Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xucheng Yin, Cheng lin Liu, Lianwen Jin, and Xiang Bai. Ocrbench: On the hidden mystery of ocr in large multimodal models. *arXiv:2305.07895*, 2023.

[38] Yuqi Liu, Bohao Peng, Zhisheng Zhong, Zihao Yue, Fanbin Lu, Bei Yu, and Jiaya Jia. Seg-zero: Reasoning-chain guided segmentation via cognitive reinforcement. *arXiv preprint arXiv:2503.06520*, 2025.

[39] Ziyu Liu, Zeyi Sun, Yuhang Zang, Xiaoyi Dong, Yuhang Cao, Haodong Duan, Dahua Lin, and Jiaqi Wang. Visual-rft: Visual reinforcement fine-tuning. *arXiv preprint arXiv:2503.01785*, 2025.

[40] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *ICLR*, 2024.

[41] Bozhi Luan, Hao Feng, Hong Chen, Yonghui Wang, Wengang Zhou, and Houqiang Li. Textcot: Zoom in for enhanced multimodal text-rich image understanding. *arXiv preprint arXiv:2404.09797*, 2024.

[42] Trung Quoc Luong, Xinbo Zhang, Zhanming Jie, Peng Sun, Xiaoran Jin, and Hang Li. Reft: Reasoning with reinforced fine-tuning, 2024.

[43] Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*, 2022.

[44] Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. Docvqa: A dataset for vqa on document images. In *WACV*, 2021.

[45] Fanqing Meng, Lingxiao Du, Zongkai Liu, Zhixiang Zhou, Quanfeng Lu, Daocheng Fu, Tiancheng Han, Botian Shi, Wenhai Wang, Junjun He, et al. Mm-eureka: Exploring the frontiers of multimodal reasoning with rule-based reinforcement learning. *arXiv preprint arXiv:2503.07365*, 2025.

[46] Chancharik Mitra, Brandon Huang, Trevor Darrell, and Roei Herzig. Compositional chain-of-thought prompting for large multimodal models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14420–14431, 2024.

[47] OpenAI, :, Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko, Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally Bennett, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrew Duberstein, Andrew Kondrich, Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Barret Zoph, Behrooz Ghorbani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew, Borys Minaiev, Botao Hao, Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon Eastman, Camillo Lugaresi, Cary Bassin, Cary Hudson, Chak Ming Li, Charles de Bourcy, Chelsea

Voss, Chen Shen, Chong Zhang, Chris Koch, Chris Orsinger, Christopher Hesse, Claudia Fischer, Clive Chan, Dan Roberts, Daniel Kappler, Daniel Levy, Daniel Selsam, David Dohan, David Farhi, David Mely, David Robinson, Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Freeman, Eddie Zhang, Edmund Wong, Elizabeth Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang, Felipe Petroski Such, Filippo Raso, Florencia Leoni, Foivos Tsimpourlas, Francis Song, Fred von Lohmann, Freddie Sulit, Geoff Salmon, Giambattista Parascandolo, Gildas Chabot, Grace Zhao, Greg Brockman, Guillaume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart Andrin, Hessam Bagherinezhad, Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichan, Ian O'Connell, Ian Osband, Ignasi Clavera Gilaberte, Ilge Akkaya, Ilya Kostrikov, Ilya Sutskever, Irina Kofman, Jakub Pachocki, James Lennon, Jason Wei, Jean Harb, Jerry Twore, Jiacheng Feng, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joaquin Quiñonero Candela, Joe Palermo, Joel Parish, Johannes Heidecke, John Hallman, John Rizzo, Jonathan Gordon, Jonathan Uesato, Jonathan Ward, Joost Huizinga, Julie Wang, Kai Chen, Kai Xiao, Karan Singhal, Karina Nguyen, Karl Cobbe, Katy Shi, Kayla Wood, Kendra Rimbach, Keren Gu-Lemberg, Kevin Liu, Kevin Lu, Kevin Stone, Kevin Yu, Lama Ahmad, Lauren Yang, Leo Liu, Leon Maksin, Leyton Ho, Liam Fedus, Lilian Weng, Linden Li, Lindsay McCallum, Lindsey Held, Lorenz Kuhn, Lukas Kondraciuk, Lukasz Kaiser, Luke Metz, Madelaine Boyd, Maja Trebacz, Manas Joglekar, Mark Chen, Marko Tintor, Mason Meyer, Matt Jones, Matt Kaufer, Max Schwarzer, Meghan Shah, Mehmet Yatbaz, Melody Y. Guan, Mengyuan Xu, Mengyuan Yan, Mia Glaese, Mianna Chen, Michael Lampe, Michael Malek, Michele Wang, Michelle Fradin, Mike McClay, Mikhail Pavlov, Miles Wang, Mingxuan Wang, Mira Murati, Mo Bavarian, Mostafa Rohaninejad, Nat McAleese, Neil Chowdhury, Neil Chowdhury, Nick Ryder, Nikolas Tezak, Noam Brown, Ofir Nachum, Oleg Boiko, Oleg Murk, Olivia Watkins, Patrick Chao, Paul Ashbourne, Pavel Izmailov, Peter Zhokhov, Rachel Dias, Rahul Arora, Randall Lin, Rapha Gontijo Lopes, Raz Gaon, Reah Miyara, Reimar Leike, Renny Hwang, Rhythm Garg, Robin Brown, Roshan James, Rui Shu, Ryan Cheu, Ryan Greene, Saachi Jain, Sam Altman, Sam Toizer, Sam Toyer, Samuel Miserendino, Sandhini Agarwal, Santiago Hernandez, Sasha Baker, Scott McKinney, Scottie Yan, Shengjia Zhao, Shengli Hu, Shibani Santurkar, Shraman Ray Chaudhuri, Shuyuan Zhang, Siyuan Fu, Spencer Papay, Steph Lin, Suchir Balaji, Suvansh Sanjeev, Szymon Sidor, Tal Broda, Aidan Clark, Tao Wang, Taylor Gordon, Ted Sanders, Tejal Patwardhan, Thibault Sottiaux, Thomas Degry, Thomas Dimson, Tianhao Zheng, Timur Garipov, Tom Stasi, Trapit Bansal, Trevor Creech, Troy Peterson, Tyna Eloundou, Valerie Qi, Vineet Kosaraju, Vinnie Monaco, Vitchyr Pong, Vlad Fomenko, Weiyi Zheng, Wenda Zhou, Wes McCabe, Wojciech Zaremba, Yann Dubois, Yinghai Lu, Yining Chen, Young Cha, Yu Bai, Yuchen He, Yuchen Zhang, Yunyun Wang, Zheng Shao, and Zhuohan Li. Openai o1 system card, 2024.

[48] OpenAI. Hello gpt-4o, 2024.

[49] Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with gpt-4. *arXiv:2304.03277*, 2023.

[50] Guanqiao Qu, Qiyuan Chen, Wei Wei, Zheng Lin, Xianhao Chen, and Kaibin Huang. Mobile edge intelligence for large language models: A contemporary survey. *arXiv preprint arXiv:2407.18921*, 2024.

[51] Tianyuan Qu, Longxiang Tang, Bohao Peng, Senqiao Yang, Bei Yu, and Jiaya Jia. Does your vision-language model get lost in the long video sampling dilemma? *arXiv preprint arXiv:2503.12496*, 2025.

[52] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.

[53] Haozhan Shen, Peng Liu, Jingcheng Li, Chunxin Fang, Yibo Ma, Jiajia Liao, Qiaoli Shen, Zilun Zhang, Kangjia Zhao, Qianqian Zhang, et al. Vlm-r1: A stable and generalizable r1-style large vision-language model. *arXiv preprint arXiv:2504.07615*, 2025.

[54] Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv: 2409.19256*, 2024.

[55] Dachuan Shi, Chaofan Tao, Ying Jin, Zhendong Yang, Chun Yuan, and Jiaqi Wang. Upop: Unified and progressive pruning for compressing vision-language transformers. In *International Conference on Machine Learning*, pages 31292–31311. PMLR, 2023.

[56] Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters, 2024.

[57] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

[58] Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, Chuning Tang, Congcong Wang, Dehao Zhang, Enming Yuan, Enzhe Lu, Fengxiang Tang, Flood Sung, Guangda Wei, Guokun Lai, Haiqing Guo, Han Zhu, Hao Ding, Hao Hu, Hao Yang, Hao Zhang, Haotian Yao, Haotian Zhao, Haoyu Lu, Haoze Li, Haozhen Yu, Hongcheng Gao, Huabin Zheng, Huan Yuan, Jia Chen, Jianhang Guo, Jianlin Su, Jianzhou Wang, Jie Zhao, Jin Zhang, Jingyuan Liu, Junjie Yan, Junyan Wu, Lidong Shi, Ling Ye, Longhui Yu, Mengnan Dong, Neo Zhang, Ningchen Ma, Qiwei Pan, Qucheng Gong, Shaowei Liu, Shengling Ma, Shupeng Wei, Sihan Cao, Siying Huang, Tao Jiang, Weihao Gao, Weimin Xiong, Weiran He, Weixiao Huang, Wenhao Wu, Wenyang He, Xianghui Wei, Xianqing Jia, Xingzhe Wu, Xinran Xu, Xinxing Zu, Xinyu Zhou, Xuehai Pan, Y. Charles, Yang Li, Yangyang Hu, Yangyang Liu, Yanru Chen, Yejie Wang, Yibo Liu, Yidao Qin, Yifeng Liu, Ying Yang, Yiping Bao, Yulun Du, Yuxin Wu, Yuzhi Wang, Zaida Zhou, Zhaoji Wang, Zhaowei Li, Zhen Zhu, Zheng Zhang, Zhexu Wang, Zhilin Yang, Zhiqi Huang, Zihao Huang, Ziyao Xu, and Zonghan Yang. Kimi k1.5: Scaling reinforcement learning with llms, 2025.

[59] Omkar Thawakar, Dinura Dissanayake, Ketan More, Ritesh Thawkar, Ahmed Heakl, Noor Ahsan, Yuhao Li, Mohammed Zumri, Jean Lahoud, Rao Muhammad Anwer, et al. Llamav-o1: Rethinking step-by-step visual reasoning in llms. *arXiv preprint arXiv:2501.06186*, 2025.

[60] Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula, Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, et al. Cambrian-1: A fully open, vision-centric exploration of multimodal llms. *arXiv preprint arXiv:2406.16860*, 2024.

[61] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timo-thée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv:2302.13971*, 2023.

[62] Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. Cogvlm: Visual expert for pretrained language models. *arXiv:2311.03079*, 2023.

[63] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023.

[64] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837. Curran Associates, Inc., 2022.

[65] Yuxin Wen, Qingqing Cao, Qichen Fu, Sachin Mehta, and Mahyar Najibi. Efficient vision-language models by summarizing visual tokens into compact registers. *arXiv preprint arXiv:2410.14072*, 2024.

[66] Jinming Wu, Zihao Deng, Wei Li, Yiding Liu, Bo You, Bo Li, Zejun Ma, and Ziwei Liu. Mmsearch-r1: Incentivizing lmms to search. *arXiv preprint arXiv:2506.20670*, 2025.

[67] xAI. Grok. `https://x.ai/`, 2023. Large language model.

[68] X.AI. Grok-1.5 vision preview. `https://x.ai/blog/grok-1.5v`, 2024.

15

[69] Yuxi Xie, Anirudh Goyal, Wenyue Zheng, Min-Yen Kan, Timothy P. Lillicrap, Kenji Kawaguchi, and Michael Shieh. Monte carlo tree search boosts reasoning via iterative preference learning, 2024.

[70] Long Xing, Qidong Huang, Xiaoyi Dong, Jiajie Lu, Pan Zhang, Yuhang Zang, Yuhang Cao, Conghui He, Jiaqi Wang, Feng Wu, et al. Pyramiddrop: Accelerating your large vision-language models via pyramid visual redundancy reduction. *arXiv preprint arXiv:2410.17247*, 2024.

[71] Guowei Xu, Peng Jin, Li Hao, Yibing Song, Lichao Sun, and Li Yuan. Llava-o1: Let vision language models reason step-by-step. *arXiv e-prints*, pages arXiv–2411, 2024.

[72] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.

[73] Senqiao Yang, Yukang Chen, Zhuotao Tian, Chengyao Wang, Jingyao Li, Bei Yu, and Jiaya Jia. Visionzip: Longer is better but not necessary in vision language models. *arXiv preprint arXiv:2412.04467*, 2024.

[74] Senqiao Yang, Jiaming Liu, Ray Zhang, Mingjie Pan, Zoey Guo, Xiaoqi Li, Zehui Chen, Peng Gao, Yandong Guo, and Shanghang Zhang. Lidar-llm: Exploring the potential of large language models for 3d lidar understanding. *arXiv preprint arXiv:2312.14074*, 2023.

[75] Senqiao Yang, Tianyuan Qu, Xin Lai, Zhuotao Tian, Bohao Peng, Shu Liu, and Jiaya Jia. An improved baseline for reasoning segmentation with large language model. *arXiv preprint arXiv:2312.17240*, 2023.

[76] Senqiao Yang, Zhuotao Tian, Li Jiang, and Jiaya Jia. Unified language-driven zero-shot domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 23407–23415, June 2024.

[77] Yi Yang, Xiaoxuan He, Hongkun Pan, Xiyan Jiang, Yan Deng, Xingtao Yang, Haoyu Lu, Dacheng Yin, Fengyun Rao, Minfeng Zhu, et al. R1-onevision: Advancing generalized multimodal reasoning through cross-modal formalization. *arXiv preprint arXiv:2503.10615*, 2025.

[78] Huanjin Yao, Jiaxing Huang, Wenhao Wu, Jingyi Zhang, Yibo Wang, Shunyu Liu, Yingjie Wang, Yuxin Song, Haocheng Feng, Li Shen, et al. Mulberry: Empowering mllm with o1-like reasoning and reflection via collective monte carlo tree search. *arXiv preprint arXiv:2412.18319*, 2024.

[79] Huanjin Yao, Jiaxing Huang, Wenhao Wu, Jingyi Zhang, Yibo Wang, Shunyu Liu, Yingjie Wang, Yuxin Song, Haocheng Feng, Li Shen, et al. Mulberry: Empowering mllm with o1-like reasoning and reflection via collective monte carlo tree search. *arXiv preprint arXiv:2412.18319*, 2024.

[80] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822, 2023.

[81] Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv:2408.01800*, 2024.

[82] Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, et al. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*, 2025.

[83] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. In *ICML*, 2024.

[84] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. In *International conference on machine learning*. PMLR, 2024.

[85] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of CVPR*, 2024.

[86] Kaichen Zhang, Bo Li, Peiyuan Zhang, Fanyi Pu, Joshua Adrian Cahyono, Kairui Hu, Shuai Liu, Yuanhan Zhang, Jingkang Yang, Chunyuan Li, et al. Lmms-eval: Reality check on the evaluation of large multimodal models. *arXiv preprint arXiv:2407.12772*, 2024.

[87] Pan Zhang, Xiaoyi Dong, Yuhang Zang, Yuhang Cao, Rui Qian, Lin Chen, Qipeng Guo, Haodong Duan, Bin Wang, Linke Ouyang, et al. Internlm-xcomposer-2.5: A versatile large vision language model supporting long-contextual input and output. *arXiv preprint arXiv:2407.03320*, 2024.

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

[89] Ruohong Zhang, Bowen Zhang, Yanghao Li, Haotian Zhang, Zhiqing Sun, Zhe Gan, Yinfei Yang, Ruoming Pang, and Yiming Yang. Improve vision language model chain-of-thought reasoning. *arXiv preprint arXiv:2410.16198*, 2024.

[90] Ruohong Zhang, Bowen Zhang, Yanghao Li, Haotian Zhang, Zhiqing Sun, Zhe Gan, Yinfei Yang, Ruoming Pang, and Yiming Yang. Improve vision language model chain-of-thought reasoning. *arXiv preprint arXiv:2410.16198*, 2024.

[91] Yuan Zhang, Chun-Kai Fan, Junpeng Ma, Wenzhao Zheng, Tao Huang, Kuan Cheng, Denis Gudovskiy, Tomoyuki Okuno, Yohei Nakata, Kurt Keutzer, et al. Sparsevlm: Visual token sparsification for efficient vision-language model inference. *arXiv preprint arXiv:2410.04417*, 2024.

[92] Zhisheng Zhong, Chengyao Wang, Yuqi Liu, Senqiao Yang, Longxiang Tang, Yuechen Zhang, Jingyao Li, Tianyuan Qu, Yanwei Li, Yukang Chen, et al. Lyra: An efficient and speech-centric framework for omni-cognition. *arXiv preprint arXiv:2412.09501*, 2024.

[93] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.

# Contents

## A  Related Works

### A.1  Efficient Vision Language Models

Large Language Models (LLMs) have demonstrated remarkable progress in language understanding and generation [1, 61, 4, 49, 13, 29]. Building on their success, VLMs have rapidly advanced by integrating visual information into LLM architectures [33, 34, 31, 60, 62, 35, 93, 26, 75]. Prominent models such as LLaVA [33] utilize visual encoders followed by the projection layers to convert images into token sequences compatible with LLMs. However, as the performance of vision-language models continues to improve, the number of visual tokens grows rapidly, leading to increased computational costs. This trend limits the practical deployment of such models in scenarios like edge computing, autonomous driving, and robotics [24, 36, 50, 74, 76, 81, 51]. Therefore, it is imperative to avoid the excessive use of visual tokens.

Recently, some studies [8, 91, 70, 65, 55, 18, 73, 92] have also recognized the redundancy in visual tokens and proposed various methods to address it. Most of these works input images containing redundancy and use the attention scores assigned by the model to prune or merge tokens for token compression. Furthermore, they typically apply a fixed threshold to compress the same proportion of redundant tokens across all data samples. Although these methods maintain good performance on general VQA tasks, they perform poorly on OCR-related benchmarks. In contrast to previous works, our proposed VisionThink initially inputs compressed tokens and allows the model to autonomously determine whether token compression is sufficient or if a high-resolution image is required. Through this approach, our method achieves efficiency while maintaining strong performance on OCR-related benchmarks. Additionally, VisionThink is not a specific token-level compression strategy but represents a new paradigm that can be integrated with existing EfficientVLM methods.

### A.2 Large Language Model Reasoning

Recent advances in large language models (LLMs) [64, 42, 47, 58, 17, 67, 57, 72] have significantly improved their reasoning capabilities through methods that simulate human-like stepwise thinking. One foundational technique, Chain-of-Thought (CoT) prompting [63], encourages models to decompose complex tasks into intermediate steps, enhancing performance on a variety of reasoning benchmarks. Furthermore, researchers have explored more structured and dynamic reasoning paradigms, such as Tree-of-Thought and Graph-of-Thought [80, 7], which organize reasoning as branching or interconnected processes. Complementary approaches like Program-of-Thought (PoT) [10] further improve reasoning fidelity by integrating external computational tools to verify or simplify logic steps.

Besides, recent work has also shifted attention from model architecture design and train-time scaling to test-time scaling [56], such as Monte Carlo Tree Search (MCTS) [69], stepwise preference optimization [27], and reinforcement learning [42] are used to refine outputs during inference. Models such as DeepSeek-R1 [17], OpenAI-O1 [47] demonstrate the effectiveness of combining large-scale RL with reward functions that prioritize both correctness and reasoning quality. Although LLMs have shown remarkable progress in structured reasoning, extending these abilities to Vision Language Models remains an open challenge.

### A.3 Vision Language Model Reasoning

With the advancement of LLM reasoning capabilities [17], many studies have aimed to improve the reasoning abilities of VLMs [90, 46, 41]. One common approach is using Chain-of-Thought (CoT) prompting to construct SFT datasets. However, the CoTs generated often lack natural human cognitive processes, limiting their effectiveness and generalization. Furthermore, inspired by DeepSeek-R1 [17], several studies have attempted to transfer this reasoning paradigm to vision tasks [79, 59, 20, 77, 38, 53, 45, 39, 66]. Most of these efforts, by collecting CoT data to perform a cold start and then training the model using a reinforcement learning strategy such as GRPO. While this approach achieves performance improvements on specific tasks, it significantly degrades the model's general performance. Moreover, current methods remain limited to visual math or segmentation tasks, failing to generalize to broader general VQA tasks. In this paper, we propose VisionThink, which effectively applies reinforcement learning to general VQA tasks by leveraging the LLM-as-Judge strategy.

## B Additional Experiments

### B.1 Prompt Details

#### B.1.1 LLM-as-Judge Prompt Design

In this section, we detail the prompt design for our LLM-as-Judge strategy. As shown in Table 3, the placeholders Ground Truth and Prediction are dynamically replaced with the corresponding question, ground truth answer, and model prediction during evaluation. Specifically, the judgment process is conducted entirely in text. Our findings indicate that, compared to VLMs, current LLMs achieve higher judgment accuracy and exhibit fewer hallucinations. Moreover, by eliminating the need for visual token inputs, it significantly reduce the overall evaluation cost.

Furthermore, we require the LLM to return a discrete value, with 1 indicating a correct prediction and 0 indicating an incorrect one, rather than a continuous score representing the degree of correctness. This binary format further reduces the likelihood of misjudgment. In a user study of 1,000 cases, no misclassifications were observed.

#### B.1.2 VisionThink Image Resize Prompt

As shown in Table 4, we present the detailed system and user prompts used in our proposed Vision-Think. Specifically, we integrate image resizing as a tool-call function. Following the Qwen2.5-VL cookbook [6], we employ an Agent Prompt that enables the model to output special tokens to trigger image resizing. This prompt design allows the model to exhibit distinct behaviors such as requesting image resizing or directly answering the question. These behaviors introduce differentiable gradients,

Table 3: **Judgment Prompt Template.** Question, Ground Truth and Prediction are dynamically replaced with the specific question, ground truth and model prediction during evaluation.

---

***SYSTEM PROMPT:***
You are an intelligent chatbot designed for evaluating the correctness of generative outputs for question-answer pairs.
Your task is to compare the predicted answer with the correct answer and determine if they match meaningfully. Here's how you can accomplish the task:
INSTRUCTIONS:
- Focus on the meaningful match between the predicted answer and the correct answer.
- Consider synonyms or paraphrases as valid matches.
- Evaluate the correctness of the prediction compared to the answer.

---

***USER PROMPT:***
I will give you a question related to an image and the following text as inputs:
1. **Question Related to the Image**: Question
2. **Ground Truth Answer**: Ground Truth
3. **Model Predicted Answer**: Prediction
Your task is to evaluate the model's predicted answer against the ground truth answer, based on the context provided by the question related to the image. Consider the following criteria for evaluation:
- **Relevance**: Does the predicted answer directly address the question posed, considering the information provided by the given question?
- **Accuracy**: Compare the predicted answer to the ground truth answer. You need to evaluate from the following two perspectives:
(1) If the ground truth answer is open-ended, consider whether the prediction accurately reflects the information given in the ground truth without introducing factual inaccuracies. If it does, the prediction should be considered correct.
(2) If the ground truth answer is a definitive answer, strictly compare the model's prediction to the actual answer. Pay attention to unit conversions such as length and angle, etc. As long as the results are consistent, the model's prediction should be deemed correct.
**Output Format**:
Your response should include an integer score indicating the correctness of the prediction: 1 for correct and 0 for incorrect. Note that 1 means the model's prediction strictly aligns with the ground truth, while 0 means it does not.
The format should be Score: 0 or 1

---

which make it feasible to apply the GRPO algorithm. Furthermore, we analyze the impact of different prompts in Sec. C.3.

## B.2   Details of the Format Reward

The format reward has a total score of 0.5, which is awarded only when all formatting requirements are fully satisfied. Specifically, as shown in the VisionThink Prompt (Table. 4), the first requirement is that the model's output must include both the the <answer></answer> and <think></think> tags, which denote the final answer and the reasoning process, respectively. The second requirement states that for responses involving an image resize operation, the model must output a correctly formatted <tool_call></tool_call> tag containing a valid JSON content.

## B.3   Implementation Details

**Training Details.** In this paper, we conduct experiments using Qwen2.5-VL-7B-Instruct [5] as the base model, trained with the veRL framework [54]. We use a total batch size of 512 with mixed-precision (FP16) training. The mini-batch size is set to 32, and the KL divergence coefficient is 0.001. The policy model is optimized using an initial learning rate of $1 \times 10^{-6}$. For each prompt, we generate 16 candidate responses using a temperature of 1.0, and apply duplicate and empty response filtering, similar to DAPO [82].

Table 4: **VisionThink Image Resize Prompt Template.** <span style="color:red">Question</span> will be replaced with the specific question during training and inference.

---

*SYSTEM PROMPT:*
You are a helpful assistant.
# Tools
You may call the function tool shown below to assist with the user query.
You are provided with the function signature within <tools></tools> XML tags:
<tools>
{
   "type": "function",
   "function":{
     "name_for_human": "resize_image",
     "name": "resize_image",
     "description": "Resize the image resolution.",
       "parameters": {
      "properties": {
        "action": {
          "description": "The action to perform. The available actions are:
               **resize**: Double the resolution of the current image. You should only use this tool if you are unable to obtain the critical information needed to answer the question from the current resolution.",
          "enum": ["resize"],
          "type": "string"
          }
        }
       "required": ["action"],
       "type": "object",
     },
     "args_format": "Format the arguments as a JSON object."
    }
}
</tools>
For each function call, return a json object with the function name and the corresponding argument within <tool_call></tool_call> XML tags:
<tool_call> {"name": <function-name>, "arguments": <args-json-object>} </tool_call>

---

*USER PROMPT:*
Answer the question based on the image provided. You must conduct reasoning within `<think>` and `</think>` first in each of your reasoning steps. You may call ONE function tool per step to help you better solve the problem. Place the function tool within `<tool_call>` and `</tool_call>` at the end of each step to perform a function call. You should continue your reasoning process based on the content returned by the function tool. Once you confirm your final answer, place the final answer inside `<answer>` and `</answer>`. For mathematical or multiple-choice problem, wrap the answer value or choice with `\boxed{}`. Here is the image and question: <span style="color:red">Question</span>.

---

**Inference Details.** In this paper, we use the lmms-eval [86] to evaluate the model's performance. Besides, in order to save the GPU memory and improve the inference speed, we utilize the vLLM[25] framework and set the temperature to zero for inference.

## B.4 Benchmark Datasets and Evaluation Metrics

We conduct experiments on these widely used visual understanding benchmarks.

**ChartQA.** ChartQA [43] is a benchmark designed to evaluate how well multimodal models answer questions about charts, emphasizing both visual understanding and logical reasoning. It includes various chart types, such as bar charts and line graphs, with a mix of human-written and automatically generated questions to assess complex reasoning abilities. Notably, ChartQA is a strongly OCR-

**Republicans and Democrats say they can't agree on 'basic facts'**

*On important issues facing the country, most Republican voters and Democratic voters ... (%)*

Not only disagree over plans and policies, but also cannot agree on basic facts

• 26%

73%

Can agree on basic facts, even if they often disagree over plans and policies

Source: Survey of U.S. adults conducted Sept. 3-15, 2019.

PEW RESEARCH CENTER

***Question***: What's the gap percent between the bad and the better scenarios being studied?

***GT***: 47

***Pred1***: 47 ✅

***Pred2***: 46 ✅

***Pred3***: 47% ❌

***Pred3***: 0.47 ❌

Figure 6: An example illustrating the original evaluation method used in ChartQA.

dependent benchmark that requires fine-grained visual understanding, as models must extract textual information from charts and reason over it.

```python
def _to_float(text: str):
    try:
        if text.endswith("%"):
            return float(text.rstrip("%")) / 100.0
        else:
            return float(text)
    except ValueError:
        return None


prediction_float = _to_float(prediction)
target_float = _to_float(target)

if prediction_float is not None and target_float is not None:
    relative_change = abs(prediction_float - target_float) / abs(
        target_float)
    return relative_change <= max_relative_change  # 0.05
else:
    return prediction.lower() == target.lower()
```

Listing 1: Core evaluation code from the original ChartQA assessment method..

Furthermore, we observe that the evaluation process of ChartQA in lmms-eval [86] relies on a float-value comparison method, which presents several limitations in practical evaluation scenarios. The corresponding implementation is shown in Listing 1, and an illustrative example is provided in Fig. 6 for further analysis.

As shown in Fig. 6, for the question "What's the gap percent between the bad and the better scenarios being studied?", the intuitive answer derived from the image is 47%. And the `_to_float()` function (Line 1 in Listing 1) converts both `0.47` and `47%` to `0.47`, while converting the ground truth value 47 to `47.0`. Hence, the comparison at Line 14 treats both `0.47` and `47%` as incorrect predictions, leading to an erroneous evaluation result. Moreover, when the model incorrectly predicts 46, the current evaluation method still considers it correct, as the relative error compared to the ground truth 47 is:

$$\frac{|46 - 47|}{47} = 0.02 < 0.05.$$

which is also a wrong judgment result of the evaluation.

Based on this observation, all ChartQA evaluations in this paper are conducted using a combination of GPT-4o-Judge and human verification, denoted as ChartQA$^{\dagger}$.

**MME.** The MME benchmark [15] assesses multimodal models on 14 subtasks that reflect both perceptual processing and cognitive reasoning abilities. By utilizing carefully crafted instruction-response pairs, MME aims to minimize the risk of training data contamination, ensuring a fair and rigorous evaluation process.

**OCRBench.** OCRBench [37] is a comprehensive benchmark for evaluating the OCR capabilities of vision language models. It covers five key tasks: text recognition, scene text-centric VQA, document-oriented VQA, key information extraction, and handwritten mathematical expression recognition. With 29 datasets and 10,000 human-verified QA pairs across 31 scenarios. Its scenarios span street scenes, receipts, and formulas, testing models on multilingual, handwritten, non-semantic, and mathematical text.

**DocVQA.** DocVQA [44] is a dataset for VQA on document images, comprising 50,000 questions defined on over 12,000 document images. It covers various document types, including forms, receipts, and scientific papers, testing models' ability to understand and reason about document content, such as textual information, tables, and visual elements.

**RealWorldQA.** RealWorldQA [68] is a benchmark designed to evaluate the real-world spatial understanding capabilities of VLMs. It consists of over 700 images, each accompanied by a question and a verifiable answer, drawn from real-world scenarios, including those captured from vehicles. The benchmark assesses how well models comprehend physical environments and spatial relationships, which are crucial for applications in navigation, robotics, and general AI assistance.

**MMVet.** MMVet [84] introduces a structured framework to assess six foundational vision-language skills: recognition, OCR, knowledge, language generation, spatial awareness, and math. These capabilities are combined in 16 evaluation configurations to test how well multimodal systems can integrate them for solving complex tasks, offering a detailed and quantitative performance analysis.

**POPE.** POPE [32] is designed to measure object hallucination in vision-language models using binary-choice questions that verify whether specific objects are present in given images. It employs metrics such as Accuracy, Recall, Precision, and F1 Score across three distinct sampling strategies, delivering a robust and fine-grained evaluation of hallucination tendencies. In our paper, the result of POPE is F1-score.

**MMMU.** MMMU [85] serves as a benchmark for assessing multimodal models on intricate, college-level tasks that demand both extensive knowledge and reasoning capabilities. It comprises 11.5K carefully selected questions sourced from exams, quizzes, and academic textbooks, spanning six broad fields: Art & Design, Business, Science, Health & Medicine, Humanities & Social Sciences, and Technology & Engineering. These questions encompass 30 academic subjects and 183 specialized areas, incorporating a wide variety of visual formats such as diagrams, graphs, and chemical formulas. MMMU is designed to push models toward expert-level performance by testing their ability to understand and reason across disciplines and modalities.

**MathVista.** MathVista [40] is a benchmark for evaluating the mathematical reasoning capabilities of foundation models within visual contexts. It includes 6,141 examples, derived from 28 existing multimodal datasets involving mathematics and three newly created datasets: IQTest, FunctionQA, and PaperQA. These tasks require fine-grained visual understanding and compositional reasoning, often involving the interpretation of graphs, equations, and other mathematical visuals. MathVista aims to systematically study the ability of VLMs to solve mathematical problems presented in visual formats, highlighting the need for models that can seamlessly integrate visual perception with mathematical reasoning.

**MathVerse.** MathVerse [88] is a benchmark for rigorously evaluating the capabilities of VLMs in interpreting and reasoning with visual information in mathematical problems. MathVerse consists of 2,612 high-quality, multi-subject math problems with diagrams, each transformed into six distinct versions with varying degrees of information content in multi-modality, resulting in 15,000 test samples.

### B.5 Scaling-up Reinforcement Learning on General VQA Tasks

Due to the diversity and complexity inherent in general VQA tasks, traditional rule-based reinforcement learning algorithms are not directly applicable. To overcome this limitation, we introduce an

LLM-as-Judge strategy, which enables our model to be trained via reinforcement learning on the General VQA task. To further demonstrate the effectiveness of our method, we scale up the dataset size to 130K to validate its effectiveness.

**Dataset.** Since the LLM-as-Judge approach is flexible, one advantage is that most of the SFT data can be utilized. Therefore, we only filter out subjective open-ended questions whose answers are not unique and can be correctly addressed from different perspectives, such as image descriptions, essay writing, and similar tasks. Based on this, we ultimately filtered 130K QA pairs to train the VLM via reinforcement learning, without requiring any cold-start phase. All the data will be open-sourced.

**Prompt.** To verify the effectiveness of the LLM-as-Judge strategy on general VQA tasks, we conduct experiments with minimal modifications to both the system and user prompts. The detailed prompts are provided in Table 5.

**Reward.** Since the entire training process in this setting does not involve any decision-making regarding the need for high-resolution images, the total reward function in reinforcement learning is designed to focus solely on answer quality and response formatting. Specifically, the reward comprises two components: The first component is an accuracy reward, evaluated by the LLM-as-Judge. This component assesses the correctness of the model's answer against the ground truth, with a maximum of 1 point awarded for a fully correct response. The second component is a formatting reward, worth 0.5 points. This is granted when the model correctly wraps its response using both the <answer></answer> and <think></think> tags. These tags are critical for maintaining consistent output formatting and enabling downstream interpretability.

Table 5: **Prompt Template for VisionThink♠.** VisionThink♠ refers to a model trained on general VQA tasks using full image resolution and the LLM-as-Judge strategy. The Question placeholder is replaced with the actual question during training and inference.

---
***SYSTEM PROMPT:***
You FIRST think about the reasoning process as an internal monologue and then provide the final answer.
The reasoning process MUST BE enclosed within `<think>` `</think>` tags. The final answer MUST BE put within `<answer>` `</answer>` tags. For mathematical or multiple-choice problems, wrap the answer value or choice with `\boxed{}`.

---
***USER PROMPT:***
Question.

---

**Experimental Results.** As shown in Table 6, we compare our model against state-of-the-art open-source and closed-source vision-language models across several general VQA benchmarks. In this evaluation, VisionThink ♠ denotes our model variant trained using the proposed LLM-as-Judge strategy with the above reward function and 130K QA pairs.

The experimental results demonstrate that our method outperforms the baseline model, Qwen2.5VL-7B-Instruct, across multiple benchmarks. The improvement is particularly notable on the MMVet, where our model achieves a significant performance gain of 7.9% over the baseline. This highlights the model's superior capability in handling general VQA tasks. Furthermore, on the recently popular benchmark MathVista [40], which is designed to assess both mathematical reasoning and general visual question answering reasoning abilities, our model achieves a score of 71.2. This result not only surpasses all existing open-source models but also outperforms several closed-source models. These findings provide strong empirical evidence for the effectiveness and generalizability of our LLM-as-Judge strategy in enhancing the reasoning capabilities of VLMs across general VQA tasks.

## B.6 Comparison with Previous Efficient VLM

To further demonstrate the effectiveness of our proposed VisionThink, we conduct a comparative analysis against an additional efficient Vision-Language Model (VLM), VisionZip [73]. While previous methods such as FastV [8] and SparseVLM [91] perform token compression within the language model component based on attention scores, VisionZip applies compression directly within the vision encoder using a similar attention-based mechanism.

Table 6: **Effectiveness of LLM-as-Judge Accuracy Reward Design.** VisionThink♠ is a model we developed by training with the LLM-as-Judge on 130K filtered General VQA datasets and leverages Qwen2.5-VL-7B-Instruct as the base model. Qwen2.5-VL-7B* reports the results evaluate by lmms-eval[86].

| Method | MMMU | MMMU-Pro | MMBench | RealWorldQA | POPE | MME | MathVista | MathVerse | MMVet |
|---|---|---|---|---|---|---|---|---|---|
| | val | test | en_test | test | test | test | testmini | testmini | test |
| *Closed-Source Model* | | | | | | | | | |
| GPT-4o [48] | 69.1 | 54.0 | 83.4 | 58.6 | 85.6 | 2329 | 63.8 | 50.2 | 69.1 |
| Claude-3.5 Sonnet [2] | 68.3 | 55.0 | 82.6 | 59.9 | - | 1920 | 67.7 | 41.2 | 70.1 |
| Gemini-1.5-Pro [57] | 62.2 | 49.4 | 73.9 | 70.4 | 88.2 | - | 63.9 | - | 64.0 |
| *Open-Source General Model* | | | | | | | | | |
| Cambrain-1-8B [60] | 42.7 | - | 75.9 | 60.0 | 86.4 | 1803 | 49.0 | - | - |
| InternVL2-8B [12] | 49.3 | 32.5 | 81.7 | 64.4 | 84.2 | 2210 | 58.3 | - | 60.0 |
| LLaVA-OneVision-7B [28] | 48.8 | - | - | 66.3 | 88.4 | 1998 | 63.2 | - | 57.5 |
| MiniCPM-Llama-V-2.5-8B [81] | 45.8 | 19.6 | 77.2 | 63.0 | 86.7 | 2025 | 54.3 | - | - |
| MiniCPM-V-2.6-8B [81] | 49.8 | 27.2 | 78.0 | 65.0 | 83.2 | 2348 | 60.6 | - | - |
| IXC-2.5 [87] | 42.9 | - | 82.2 | 67.8 | - | 2229 | 63.8 | - | 51.7 |
| InternVL2.5-8B [11] | 56.0 | 38.2 | 84.6 | 70.1 | 90.6 | 2344 | 64.4 | 39.5 | 62.8 |
| *Reasoning Model* | | | | | | | | | |
| LLaVA-CoT-11B [71] | - | - | 75.0 | - | - | - | 54.8 | - | 60.3 |
| LLaVA-Reasoner-8B [89] | - | - | - | - | - | - | 50.6 | - | - |
| Insight-V-8B [14] | 50.2 | 24.9 | 82.3 | - | - | 2312 | 59.9 | - | - |
| Mulberry-7B [78] | 55.0 | - | - | - | - | 2396 | 63.1 | - | - |
| Vision-R1-LlamaV-CI-11B [19] | - | - | - | - | - | 2190 | 62.7 | 27.1 | - |
| *VisionThink* | | | | | | | | | |
| Qwen2.5-VL-7B* [5] | 50.3 | 37.7 | 82.6 | 68.6 | 86.7 | 2316 | 68.2 | 46.3 | 61.6 |
| VisionThink ♠ | **52.7** | **41.1** | **83.4** | 66.5 | **88.6** | 2314 | **71.2** | **48.3** | **69.5** |

As shown in Table 7, although previous efficient VLM methods achieve competitive performance on general VQA benchmarks, their accuracy drops significantly on OCR-related tasks. This degradation is particularly evident even when a substantial portion of the visual token is retained (70%), as demonstrated on the ChartQA dataset.

In contrast, our proposed model, VisionThink, can smartly decide whether to request the original high-resolution image based on the complexity and demands of each sample. This adaptive strategy enables the model to maintain high accuracy on general VQA tasks while substantially improving performance on benchmarks requiring detailed textual recognition. Through this capability, VisionThink demonstrates stronger fine-grained visual understanding and addresses a key limitation of previous efficient VLMs—namely, their poor performance on OCR-related tasks, which has constrained their applicability in real-world scenarios.

Notably, VisionZip‡ refers to the variant fine-tuned on the 130K dataset [73]. However, compared to the training-free version of VisionZip, this fine-tuned model does not show any performance improvement. We attribute this to the limited coverage and diversity of the fine-tuning dataset, which falls short of the supervised fine-tuning data used by the official Qwen team. This observation indirectly suggests that, compared to supervised fine-tuning, reinforcement learning provides better generalization, which we further discuss in Sec. C.1.

Table 7: **Comparison with Previous Efficient VLM Methods.** Vanilla represents the Qwen2.5-VL-7B-Instrcut. The retained ratio of the baseline methods is a predefined hyperparameter, while for VisionThink, the ratio is determined autonomously by the model and reported as a statistical value. Note that *Down-Sample* refers to the model's performance when directly fed images with their resolution reduced by half. VisionZip‡ represents using the 130K data to finetuning the model.

| Method | ChartQA† | OCRBench | DocVQA | MME | MMVet | RealWorldQA | POPE | MathVista | MathVerse | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|
| | test | test | val | test | test | test | test | testmini | testmini | |
| *Retain 100% Visual Tokens Across All Benchmarks* | | | | | | | | | | |
| Vanilla | 79.8 | 81.5 | 95.1 | 2316 | 61.6 | 68.6 | 86.7 | 68.2 | 46.3 | 100% |
| | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | |
| *Retain 25% Visual Tokens Across All Benchmarks* | | | | | | | | | | |
| Down-Sample | 62.9 | 68.8 | 94.3 | 2270 | 54.5 | 68.8 | 82.8 | 62.2 | 43.1 | 92.1% |
| | 78.8% | 84.4% | 99.1% | 98.0% | 88.5% | 100.3% | 95.5% | 91.2% | 93.1% | |
| *Retain 50% Visual Tokens Across All Benchmarks* | | | | | | | | | | |
| FastV (ECCV 2024) | 72.6 | 75.8 | 93.6 | 2308 | 52.8 | 68.8 | 84.7 | 63.7 | 45.0 | **95.8%** |
| | 91.0% | 93.0% | 98.4% | 99.6% | 85.7% | 100.3% | 97.7% | 93.4% | 97.2% | |
| SparseVLM (ICML 2025) | 73.2 | 75.6 | 66.8 | 2282 | 51.5 | 68.4 | 85.5 | 66.6 | 45.1 | **92.2%** |
| | 91.7% | 92.7% | 70.2% | 98.5% | 83.6% | 99.7% | 98.6% | 97.6% | 97.4% | |
| VisionZip (CVPR 2025) | 73.4 | 70.5 | 93.8 | 2209 | 57.0 | 68.6 | 86.3 | 63.7 | 45.1 | **95.0%** |
| | 92.0% | 86.5% | 98.6% | 95.4% | 92.5% | 100% | 99.5% | 93.4% | 97.4% | |
| VisionZip‡ (CVPR 2025) | 77.3 | 77.9 | 93.8 | 2244 | 50.1 | 69.2 | 91.2 | 63.1 | 39.4 | **95.0%** |
| | 96.9% | 95.6% | 98.6% | 96.9% | 81.3% | 100.9% | 107.5% | 92.5% | 85.1% | |
| *Retain 70% Visual Tokens Across All Benchmarks* | | | | | | | | | | |
| FastV (ECCV 2024) | 77.2 | 82.2 | 94.4 | 2342 | 56.0 | 68.6 | 85.9 | 65.9 | 46.9 | **98.4%** |
| | 96.7% | 100.8% | 99.3% | 101.1% | 90.9% | 100.0% | 99.1% | 96.6% | 101.3% | |
| SparseVLM (ICML 2025) | 75.8 | 79.3 | 68.7 | 2276 | 53.7 | 68.5 | 85.4 | 66.3 | 45.1 | **93.6%** |
| | 94.9% | 97.3% | 72.2% | 98.3% | 87.2% | 99.8% | 98.5% | 97.2% | 97.4% | |
| VisionZip (CVPR 2025) | 76.8 | 80.9 | 94.5 | 2334 | 60.0 | 68.2 | 86.4 | 68.9 | 45.8 | **99.1%** |
| | 96.2% | 99.3% | 99.4% | 100.8% | 97.4% | 99.4% | 99.7% | 101.0% | 98.9% | |
| VisionZip‡ (CVPR 2025) | 78.2 | 81.3 | 94.1 | 2230 | 52.5 | 68.6 | 92.5 | 64.8 | 41.8 | **96.7%** |
| | 98.0% | 99.8% | 98.9% | 96.3% | 85.3% | 100% | 106.7% | 95.0% | 90.3% | |
| *Retain Approximately 51.3% Visual Tokens Across All Benchmarks* | | | | | | | | | | |
| VisionThink | 79.8 | 80.8 | 94.4 | 2400 | 68.5 | 67.1 | 86.0 | 67.5 | 48.0 | <span style="color:red">**101.4%**</span> |
| | 100.0% | 99.1% | 99.3% | 103.6% | 111.2% | 97.8% | 99.2% | 99.0% | 103.7% | |

# C   Further Discussions

## C.1   Why Use RL Instead of SFT?

In this paper, we train a smart and efficient vision-language model via reinforcement learning. A natural question arises: why use reinforcement learning instead of supervised fine-tuning to achieve this goal?

To answer this question, we conduct a comparative SFT experiment. Firstly, we construct the SFT training set. Specifically, compared to RL, which can directly utilize QA pairs and autonomously learn both the reasoning process and whether a high-resolution image is needed, SFT requires manually crafting both the reasoning steps and dialogue that involves high-resolution image requests. To overcome this limitation, we use GPT-4o to simulate both the high-resolution image requests and the corresponding reasoning process, enabling the SFT training data to closely approximate the behavior of the RL-trained model. Finally, we convert the original RL training data into a format compatible with SFT, maintaining a 1:1 ratio between high-resolution requests and direct answers.

As shown in Table 8, we compare the proportion of high-resolution image requests made by the SFT and RL models across evaluation benchmarks. Compared to the RL-trained model, which can smartly

Table 8: Comparison of image resize call ratios for RL and SFT trained models over multiple evaluation benchmarks.

| Method | ChartQA$^\dagger$ | OCRBench | DocVQA | MME | RealWorldQA | POPE |
|---|---|---|---|---|---|---|
| RL | 79.1% | 62.3% | 6.5% | 30.7% | 29.9% | 9.5% |
| SFT | 95.1% | 64.0% | 14.1% | 39.0% | 62.1% | 37.8% |

Table 9: **Performance comparison of the cold start model and no cold start model.** The without cold start model represents our VisionThink.

| Method | Type | ChartQA$^\dagger$ | OCRBench | DocVQA | MME | MMVet | RealWorldQA | POPE |
|---|---|---|---|---|---|---|---|---|
| w/o Cold Start | RL | 79.8 | 80.8 | 94.4 | 693/1707 | 68.5 | 67.1 | 86.0 |
| Cold Start (2K) | Base | 76.4 | 78.7 | 92.4 | 444/1354 | 58.3 | 47.2 | 86.6 |
| | RL | 77.7 | 80.2 | 93.0 | 622/1624 | 62.3 | 52.8 | 86.2 |
| Cold Start (8K) | Base | 76.8 | 78.2 | 90.5 | 525/1368 | 60.5 | 36.5 | 84.8 |
| | RL | 79.2 | 79.4 | 92.5 | 609/1637 | 66.2 | 55.8 | 85.6 |

decide when to answer directly and when to request a high-resolution image, the SFT model exhibits a significantly higher image resize calling ratio across all benchmarks. This behavior is especially evident in the RealWorldQA benchmark, where high-resolution images are generally unnecessary, yet the SFT model still issues requests 62.1% of the time.

Based on this observation, we find that SFT does not enable the model to become "smart" enough to accurately determine whether a high-resolution image is necessary and also requires constructing the training set using GPT-4o. In contrast, RL makes the VLM smarter and more generalizable, and can directly use the original QA pairs without additional formatting.

## C.2 Why not Cold-Start?

Currently, most explorations of RL in VLMs require a cold start stage. In this section, we explore why we do not use a cold start and instead train the model directly with RL.

To investigate this problem, we collect datasets of 2K and 8K samples to cold start our model. The cold-start data are constructed similarly to Sec. C.1, where GPT-4o is used to simulate both the requests for high-resolution images and the corresponding reasoning processes, enabling the SFT training data to closely mimic the behavior of the model trained via RL.

We first compare the performance of the cold-started models before RL training, as shown in the 'Base' lines of Table 9. Although performance improves with increasing data size, the cold-start models still fall short of the original Qwen2.5VL. We believe this is primarily due to the limited diversity and coverage of our data compared to the SFT data used by the Qwen team, resulting in the observed performance gap.

Furthermore, we compare the performance of models after RL training, using the same RL setup as VisionThink, as shown in the 'RL' lines of Table 9. While RL improves the performance of cold-start models, they still fall short compared to models trained from vanilla Qwen2.5VL with RL.

Based on this observation, we conclude that due to the lower diversity and coverage of the cold-start data compared to the original Qwen2.5VL SFT data, introducing cold-start training may improve performance in specific domains covered by the cold-start data but significantly reduces the model's general capability. This limitation restricts its broader applicability. Therefore, in this paper, we do not utilize the cold-start stage.

## C.3 Different Prompt Impact

Since we do not adopt a cold-start stage, designing an appropriate initial prompt becomes crucial to ensure the model begins from a good starting point. The base model Qwen2.5VL-Instruct, which has not undergone RL training, typically tends to answer questions directly and lacks the ability to smartly request high-resolution images when needed.



Figure 7: **Impact of Prompt Choice.** Prompts lead to substantial variation in image resize call ratios, with the Qwen official agent prompt demonstrating the most effective performance.

Therefore, it is essential for the base model, when conditioned on our prompt, to show some preference for calling image resizing. Otherwise, if it is overly biased toward direct answering, the GRPO training process will fail to optimize effectively and may collapse into the direct-answering mode. To address this, we compare three prompt settings. The first is the official agent prompt from Qwen's cookbook, shown in Table 4. The other two are our custom-designed prompts, detailed in Table 10.

Table 10: **Two custom prompts for analyzing the impact of different prompts.** The Question placeholder will be replaced with the specific question during training and inference.

| |
|---|
| ***Simple Resize System Prompt:*** |
| You are a helpful assistant. |
| ***Simple Resize User Prompt:*** |
| Answer the user's question based on the image provided. You can place `<resize></resize>` at the end of your response to call the image resize tool, it will return the resized image with its resolution doubled to help you better answer the question. Once you confirm your final answer, place the final answer inside `<answer>` and `</answer>`. Here is the image and question: Question. |
| ***Simple Think Resize System Prompt:*** |
| You are a helpful assistant. Answer the user's question based on the image provided. You can place `<resize>` at the end of your response to call the image resize tool, it will return the resized image with its resolution doubled to help you better answer the question. Once you confirm your final answer, place the final answer inside `<answer>` and `</answer>`. |
| ***Simple Think Resize User Prompt:*** |
| Enclose reasoning in `<think></think>` Question. |

28

Figure 8: **Ablation Study on Penalty Ratio Threshold.** As the threshold increases, the model progressively favors requesting image resizing instead of providing direct answers.

As shown in Fig. 7, our model using the official agent prompt demonstrates a higher image resize call ratio on the effective batch, which consists of an equal mix of high-resolution-required samples and direct-answer samples. In contrast, the other two custom-designed prompts perform poorly, both in their initial behavior and in the final trained model's ability to correctly request image resizing, and ultimately collapse into consistently producing direct answers. Based on these analyses, we find that the Qwen official agent prompt, likely optimized during the pretraining or supervised fine-tuning stages by the Qwen team, is more suitable for VisionThink.

## C.4 Ablation Study on Penalty Control Threshold

In the main paper Eq. 5, we design the penalty ratio as below:

$$\mathcal{P}_{control} = 0.1 \cdot \left[ \mathbf{1}_{\text{direct}} \mathbb{I}(r < \theta) + \mathbf{1}_{\text{high}} \mathbb{I}(r \geq \theta) \right], \qquad r = \frac{C_{\text{direct}}}{C_{\text{direct}} + C_{\text{high}}}, \tag{6}$$

where $\theta$ is the threshold.

Intuitively, the larger the value of $\theta$, the more likely the model is to penalize direct answers, thereby encouraging it to request high-resolution images. Conversely, a smaller $\theta$ leads the model to penalize responses that call for image resizing, thus promoting direct answers. Based on this intuition, we experimented with different threshold values and recorded the proportion of high-resolution image requests within the effective batch. The results are shown in Fig. 8. As indicated by the Eq. 6, increasing the threshold gradually shifts the model's behavior from favoring direct answers to favoring image resizing requests. Eventually, the model collapses into always requesting high-resolution images. However, within an appropriate range, the model's behavior is not highly sensitive to the exact threshold value.

# D   Qualitative Results

In this section, we present a case study comparing our proposed VisionThink with other efficient VLM methods: FastV [8], SparseVLM [91], and VisionZip [73]. As shown in the three cases below, for OCR-related or detail-intensive samples, our proposed VisionThink model can smartly determine when a high-resolution image is needed. In contrast to previous efficient VLMs, which suffer performance degradation due to fixed compression ratios, VisionThink avoids such issues by making adaptive decisions based on the input, thereby maintaining strong performance.

# E    Broader Impact Statement

This paper is conducted solely for academic research purposes. All data used in this work were collected through compliant and ethical channels, ensuring adherence to relevant data protection and usage guidelines. Furthermore, all models employed in this study comply with their respective license agreements. As such, this research upholds high standards of integrity and responsibility, with no foreseeable negative societal impact.

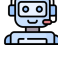Are there any fire hydrants here?

**FastV**: No

**SparseVLM**: No

**VisionZip**: No

**VisionThink**:

*[Think]*
The image shows a street scene with cars, buildings, and some urban infrastructure. To determine if there are fire hydrants, I'll need to carefully examine the image for any visible fire hydrants.

*<Upscale Image to Original Resolution>*

*[Think]*
Upon reviewing the resized image, I can see a fire hydrant on the right side of the image, near the curb. This is evident from the orange color and shape typical of fire hydrants.

*[Answer]*
**Yes**