

Mitigating Visual Hallucinations via Evidence-Guided Multi-Image Reasoning in Visual Retrieval-Augmented Generation

Anonymous ACL submission

Abstract

Visual Retrieval-Augmented Generation (VRAG) enhances Vision–Language Models (VLMs) by incorporating retrieved images as contextual evidence to support reasoning. However, existing VRAG systems often struggle to reliably perceive and integrate evidence across multiple images, leading to erroneous reasoning and visual hallucinations. In this paper, we propose EVisRAG, an end-to-end framework for evidence-guided multi-image reasoning that mitigates these issues by explicitly observing images, recording per-image evidence, and reasoning over aggregated evidence to derive the final answer. To train EVisRAG effectively, we introduce Reward-Scoped Group Relative Policy Optimization (RS-GRPO), which assigns fine-grained rewards to scope-specific tokens to jointly optimize visual perception and reasoning. Experiments on multiple visual question answering benchmarks show that EVisRAG consistently outperforms backbone VLMs, achieving an average improvement of 27%, while substantially reducing visual hallucinations through accurate evidence localization and grounded reasoning.

1 Introduction

Retrieval-Augmented Generation (RAG) equips Large Language Models (LLMs) with external knowledge retrieval to provide task-relevant context and mitigate hallucinations caused by limited parametric knowledge (Lewis et al., 2020; Liu et al., 2025b). However, much real-world knowledge is inherently non-textual, residing in modalities such as images, tables, and complex document layouts. Text-centric preprocessing pipelines that rely on image captioning or OCR to linearize these signals discard essential visual and spatial structure, limiting the model’s ability to access and reason over information originally present in images or document pages (Zhang et al., 2024b).

To address this limitation, Visual RAG (VRAG) (Yu et al., 2025; Faysse et al., 2024a) retrieves document page snapshots as units, preserving visual and spatial cues so VLMs can read evidence directly from images. Recent variants couple retrieval with reinforcement learning, inserting retrieved images into intermediate reasoning steps so the model can derive the correct answer from pixels rather than text alone (Peng et al., 2025; Wu et al., 2025). Despite these gains, many methods still transplant text-based RAG practices into the visual domain and neglect modality-specific capabilities such as perceiving information relevant to the question from images. Such a deficiency in visual perception often exacerbates visual hallucinations. In particular, models may fail to correctly attend to and utilize relevant visual evidence present in the images, leading to incorrect answers, or may erroneously perceive and reason about visual content that does not actually exist, resulting in spurious or hallucinated conclusions. Some works introduce perception-oriented actions or auxiliary agents to mitigate visual hallucinations (Wang et al., 2025b,a), which improves attention to visual detail but increases architectural complexity and computational cost, complicating end-to-end training and later reconfiguration.

To bypass these complexities, recent advances in vision-language reasoning models (VLRMs) have introduced promising strategies for enhancing visual perception on a single image during the reasoning process (Shen et al., 2025; Xu et al., 2025) by incorporating auxiliary rewards related to visual perception. Although these VLRMs perform well on single-image inputs, VRAG often retrieves multiple images, requiring cross-image localization and integration of visual evidence. Current methods lack a built-in per-image evidence collection and instead rely on external tools or agents, which increases complexity and instability. In addition, current VLRM training strategies typically opti-

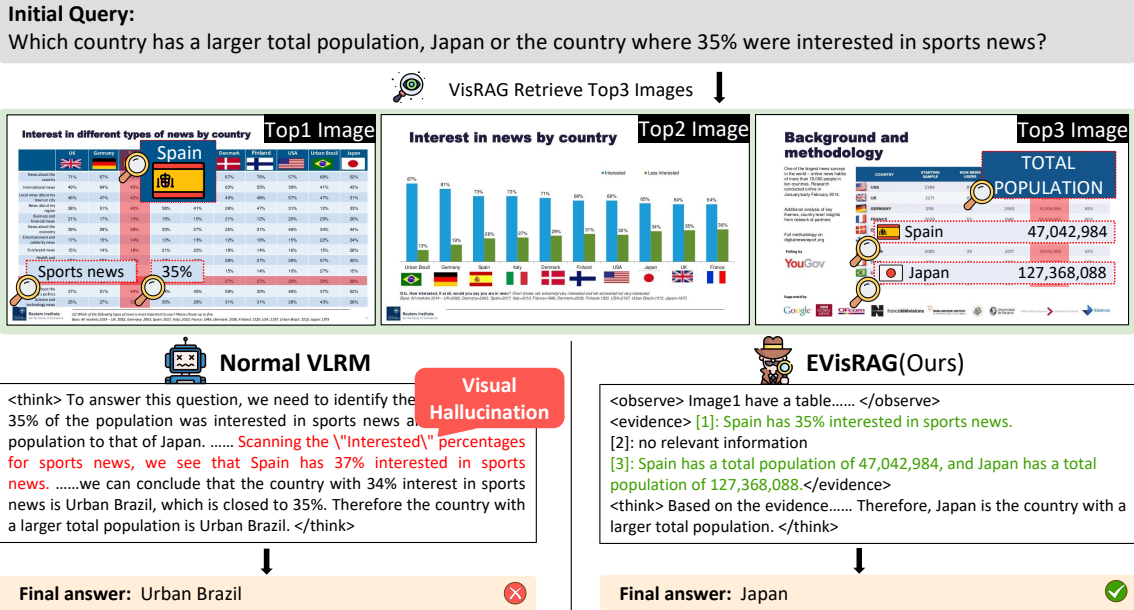


Figure 1: Comparison of normal vision-language reasoning model (VLRM) and EVisRAG

084 mize perception and reasoning with mixed rewards,
 085 overlooking the effective scope and objective differ-
 086 ences of each signal, which blurs credit assignment
 087 and causes interference.

088 Motivated by these challenges, we pro-
 089 pose Evidence-guided Vision Retrieval-augmented
 090 Generation (EVisRAG) to equip VLMs with pre-
 091 cise visual perception in multi-image scenarios. As
 092 illustrated in Figure 1, EVisRAG conducts a lin-
 093 guistic observation phase that sequentially gath-
 094 ers evidence from retrieved images, maintaining
 095 focus on them, and then performs reasoning on
 096 the collected evidence to derive the correct an-
 097 swer. To train EVisRAG effectively, we introduce
 098 Reward-Scoped Group Relative Policy Optimiza-
 099 tion (RS-GRPO), which applies fine-grained re-
 100 wards directly to in-scope tokens corresponding to
 101 visual evidence collection or reasoning. Unlike con-
 102 ventional joint optimization with mixed rewards,
 103 where perception and reasoning signals interfere
 104 and lead to unstable or suboptimal training, RS-
 105 GRPO explicitly scopes rewards to their effective
 106 regions, enabling precise credit assignment and sta-
 107 ble optimization. Without such reward scoping,
 108 jointly optimizing visual perception and reasoning
 109 in multi-image settings becomes unstable, as gra-
 110 dients from heterogeneous objectives compete across
 111 irrelevant tokens. Experiments on different VQA
 112 tasks demonstrate the effectiveness of EVisRAG,
 113 showing substantial improvements over different
 114 VRAGs. Powered by RS-GRPO, EVisRAG can

115 precisely identify question-relevant evidence im-
 116 age by image and reason over the recorded cues
 117 to produce grounded answers, much like a detec-
 118 tive assembling evidence. Moreover, EVisRAG
 119 demonstrates stronger visual perception and higher
 120 answer accuracy among other baselines, confirm-
 121 ing that richer visual perception improves the abil-
 122 ity of question understanding and alleviates visual
 123 hallucinations.

2 related work 124

125 Early research on retrieval-augmented generation
 126 (RAG) equips large language models (LLMs) with
 127 retrievers over curated corpora to provide task-
 128 relevant context and mitigate hallucinations (Lewis
 129 et al., 2020; Asai et al., 2024). However, a sub-
 130 stantial portion of real-world knowledge is non-
 131 textual, residing in images, tables, and documents
 132 with complex layouts. Pipelines that first lin-
 133 earize such signals via captioning or optical char-
 134 acter recognition and then provide only text to
 135 the model often discard critical visual and spa-
 136 tial cues, which degrades performance on down-
 137 stream reasoning tasks (Zhang et al., 2024b). To
 138 address this limitation, VisRAG (Yu et al., 2025)
 139 and Colpali (Faysse et al., 2024a) introduce Visual
 140 Retrieval-Augmented Generation (VRAG), which
 141 treats document page snapshots as retrieval units
 142 and enables vision-language models to directly
 143 read evidence from images.

144 Nevertheless, regardless of the modality of the

retrieved context, existing RAG and VRAG models remain prone to hallucinations and reasoning failures. In particular, models may confidently produce incorrect answers despite the presence of relevant evidence, become distracted by irrelevant or weakly related information, or fail to accurately extract and reason over key facts from long contextual inputs (Mishra et al., 2024; Cuconasu et al., 2024; Hsieh et al., 2024). Building on this observation, retrieval-augmented reasoning approaches (Shao et al., 2024; Rafailov et al., 2023; Schulman et al., 2017; Li et al., 2025; Song et al., 2025) acquire and exploit evidence at intermediate reasoning steps, using it to guide subsequent inference and reduce hallucinations arising from misused textual context. Despite these advances, many methods directly transplant text-centric RAG practices to the visual domain and insufficiently account for modality-specific challenges such as cross-image grounding, layout-aware reading, and region-level attention, resulting in unstable perception across multiple images (Wang et al., 2025a,b).

Recent advances in vision and language reasoning models (VLMs) have introduced effective strategies for strengthening visual perception during reasoning. Vision-R1 (Zhan et al., 2025), MM-Eureka (Meng et al., 2025), Ocean-R1 (Lingfeng et al., 2025), ThinkLite-VL (Wang et al., 2025c), and OpenVLThinker (Deng et al., 2025) show that directly applying GRPO, sometimes even without supervised fine-tuning, substantially promotes the emergence of chain of thought reasoning and can elicit “aha” moments. VLM-R1 (Shen et al., 2025) and Mixed-R1 (Xu et al., 2025) further improve perceptual grounding by augmenting answer correctness signals with auxiliary perception rewards, encouraging better use of image information. However, in VRAG settings that require reasoning over semantically rich content from multiple images, the remaining limitations in perceptual grounding often lead to misinterpretation of visual evidence, which in turn undermines the validity of the overall reasoning process.

3 Methodology

This section presents EVisRAG, a framework that enables vision–language models (VLMs) to reason over multiple images with rich visual evidence. We first outline its evidence-guided reasoning process, including observation, evidence recording, and answer reasoning (Section 3.1). We then introduce

Reward-Scoped Group Relative Policy Optimization (RS-GRPO), which enhances fine-grained perceptual grounding during reasoning (Section 3.2).

3.1 The Overview Framework of EVisRAG

Given a query q and a corpus of document-page images \mathcal{D} , EVisRAG performs evidence-guided multi-image reasoning to produce the final answer a . Instead of directly reasoning over raw images, EVisRAG explicitly decomposes the reasoning process into four structured stages: Observe, Record, Reason, and Answer, which together form an evidence-guided reasoning pipeline.

Information Retrieving. EVisRAG first retrieves a set of candidate document pages relevant to the query q from the corpus \mathcal{D} . Following prior Visual RAG frameworks, we obtain the top- k retrieved pages:

$$\mathcal{D}_r = \text{VisRAG_Ret}(q, \mathcal{D}), \quad (1)$$

where $\mathcal{D}_r \subset \mathcal{D}$ denotes the top- k retrieved document-page images relevant to the query q .

Observe and Record. Conditioned on the query q and retrieved pages \mathcal{D}_r , the model sequentially observes each image to identify regions potentially relevant to the question. Based on these observations, EVisRAG records per-image evidence \mathcal{E}_i in a structured textual form. Evidence from each image is explicitly indexed, while images containing no relevant information are recorded as negative evidence. This explicit evidence recording step encourages faithful visual grounding and mitigates perceptual hallucinations.

Reason and Answer. After collecting evidence across images, EVisRAG performs reasoning \mathcal{R} over the aggregated evidence $\mathcal{E} = \{\mathcal{E}_i\}_{i=1}^k$ to derive the final answer. By separating evidence perception from reasoning, the model avoids directly hallucinating answers from raw visual input and instead bases its predictions on recorded evidence. The final answer is generated solely from the aggregated evidence and the query. Formally, the overall process can be abstracted as:

$$a = \text{EVisRAG}(q, \mathcal{D}_r, \mathcal{E}, \mathcal{R}), \quad (2)$$

where a is the final answer produced by reasoning over the retrieved images \mathcal{D}_r .

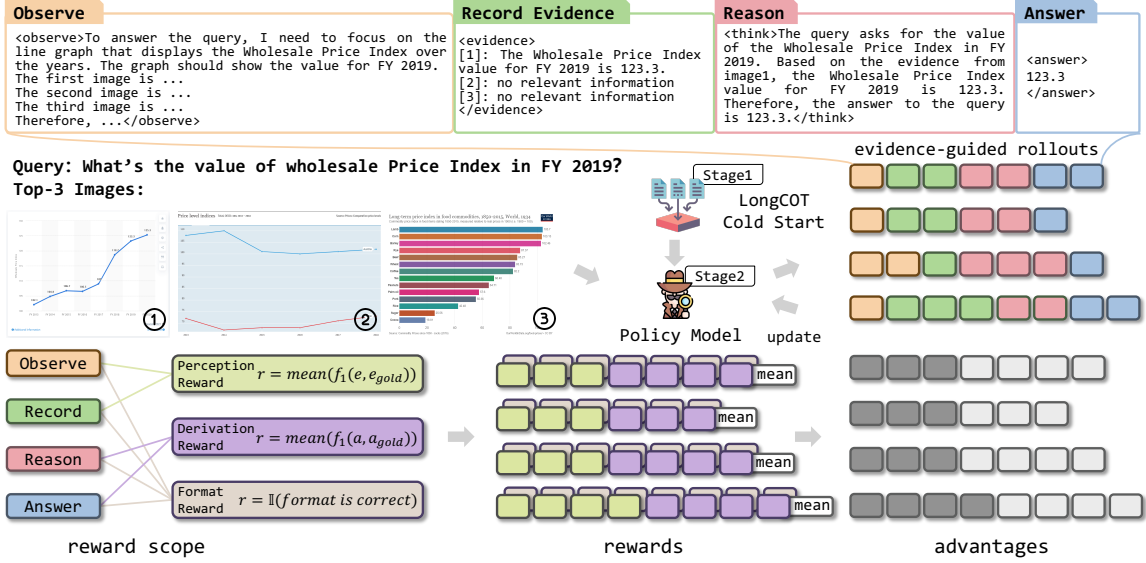


Figure 2: Overall framework of EVISRAG. Followed by the query and top-3 retrieved document pages, EVISRAG outputs four token scopes: observe, record evidence, reason, and answer. RS-GRPO assigns three fine-grained rewards to scope-specific tokens. In-scope rewards are then averaged and group-normalized to obtain token advantages for policy updates.

3.2 Optimizing VLMs to Evidence-guided Reason Using RS-GRPO

To enhance EVISRAG’s ability to accurately record evidence from multiple images and reason based on that evidence, we employ a two-stage training as shown in Figure 2. In the first stage, we apply supervised fine-tuning (SFT) as a cold start. In the second stage, we introduce Reward-Scoped Group Relative Policy Optimization (RS-GRPO), which extends GRPO to jointly optimize perception and reasoning ability of VLMs, with fine-grained rewards applied to their corresponding reward scopes.

Reward Scopes. To evaluate model outputs while encouraging the evidence-guided reasoning paradigm, we adopt three fine-grained rewards in a coordinated scheme. The format reward R_{format} enforces adherence to an evidence-guided reasoning paradigm by requiring the model to observe, record evidence, reason, and answer in a disciplined order, making intermediate steps explicit and supervision stable. The perception reward $R_{\text{perception}}$ checks whether question-relevant regions are correctly localized and summarized for each image based on the ground truth evidence generated by a larger VLM and allows an explicit no relevant information when evidence is absent. The derivation reward $R_{\text{derivation}}$ evaluates whether the model derives the correct final answer from its visual perception, ensuring the reasoning is grounded in the observed and recorded evidence. More details of the reward

design are shown in Appendix A.3

To jointly train the perception and reasoning ability of VLMs, we introduce Reward Scopes, which route supervision to scope-specific tokens to sharpen credit assignment, reduce interference, and stabilize training. Let $\mathcal{M}(t)$ denote the set of reward channels applicable to the token at position t . The output sequence is segmented by special tokens into four scopes, the observe scope \mathcal{T}_o , the record evidence scope \mathcal{T}_e , the reason scope \mathcal{T}_r , and the answer scope \mathcal{T}_a . Rewards act only where they are meaningful. $R_{\text{perception}}$ supervises tokens in \mathcal{T}_o and \mathcal{T}_e , guiding the model to summarize the right visual regions. $R_{\text{derivation}}$ supervises tokens in \mathcal{T}_r and \mathcal{T}_a , encouraging the model to derive the correct final answer from what was perceived. R_{format} applies to all tokens and keeps the evidence-guided workflow explicit and stable. Formally, we define the reward–scope mapping as:

$$\mathcal{M}(t) = \begin{cases} \{ R_{\text{perception}}, R_{\text{format}} \} & t \in \mathcal{T}_o \cup \mathcal{T}_e \\ \{ R_{\text{derivation}}, R_{\text{format}} \} & t \in \mathcal{T}_r \cup \mathcal{T}_a \end{cases} \quad (3)$$

For the i -th sampled output and its token at position t , let $R_t^{(m),i}$ denote the score from reward channel $m \in \mathcal{M}(t)$. The scope-aggregated token reward is the mean over its in-scope channels:

$$\bar{R}_t^i = \frac{1}{|\mathcal{M}(t)|} \sum_{m \in \mathcal{M}(t)} R_t^{(m),i} \quad (4)$$

RS-GRPO objective. To train both visual perception and reasoning, EVisRAG adopts an RS-GRPO objective that explicitly computes token advantages under reward scopes. Given a group of G sampled outputs, the token-level advantage is

$$\hat{A}_t^i = \frac{\bar{R}_t^i - \text{mean}(\{\bar{R}_t^1, \bar{R}_t^2, \dots, \bar{R}_t^G\})}{\text{std}(\{\bar{R}_t^1, \bar{R}_t^2, \dots, \bar{R}_t^G\})}, \quad (5)$$

where i indexes the i -th sample in the group, and G is the group size. We incorporate the resulting token-level advantages into DAPO to enhance exploration diversity and training stability, and optimize the model by minimizing the following objective:

$$\mathcal{L}_{\text{RS-GRPO}}(\theta) = -\frac{1}{\sum_{i=1}^G |o^i|} \sum_{i=1}^G \sum_{t=1}^{|o^i|} \min \left(r_t^i(\theta) \hat{A}_t^i, \text{clip}(r_t^i(\theta), 1 - \epsilon_{\text{low}}, 1 + \epsilon_{\text{high}}) \hat{A}_t^i \right), \quad (6)$$

where o^i is the i -th sampled output sequence, $r_t^i(\theta)$ is the importance ratio, and $\epsilon_{\text{low}}, \epsilon_{\text{high}}$ are the lower and upper clipping thresholds.

4 Experimental Methodology

This section describes the datasets, baselines, evaluation metrics, and implementation details.

Datasets. We first introduce the datasets used in our experiments, followed by the data statistics for golden reasoning trajectory construction.

We evaluate our EVisRAG on five visual question answering (VQA) tasks encompassing diverse document types, including ChartQA (Masry et al., 2022) and InfographicsVQA (Mathew et al., 2022) for various types of figures, MP-DocVQA (Tito et al., 2023) for industrial documents, SlideVQA (Tanaka et al., 2023) for presentation slides, and ViDoseek (Wang et al., 2025a) for multi-document scenarios. For each query, we utilize VisRAG-Ret (Yu et al., 2025) to retrieve the top-3 relevant images as context. Subsequently, each question is categorized according to whether the retrieved context provides sufficient information to answer the question with sufficient context or with insufficient context. More details of the test datasets are provided in Appendix A.1.

We collect 30,000 samples from the training sets of ChartQA and InfographicsVQA and split them into supervised fine-tuning (SFT) and GRPO subsets with an 8:2 ratio. For each query, a VisRAG-based retriever provides multiple candidate images

as visual context. To construct high-quality reasoning trajectories, we generate candidate chains of thought using Qwen2.5-VL (Bai et al., 2025) models and retain only those that lead to correct answers. Following (An et al., 2025), we further filter out samples that can be easily solved to focus training on more challenging cases. This process yields 60,000 samples for SFT and 4,000 samples for GRPO training. Detailed data construction procedures are provided in Appendix A.2.

Baselines. All baselines use VisRAG-Ret for retrieval. For each query, we fetch the top-3 documents, then the model answers using the retrieved images and the original question.

We compare three groups. General VLMs include Qwen2.5-VL-7B, Qwen2.5-VL-32B (Bai et al., 2025) and MiMo-VL-7B-RL (Xiaomi, 2025). VLRMs trained on Qwen2.5-VL-7B-Instruct include Vision-R1-7B (Zhan et al., 2025), MM-Eureka-7B (Meng et al., 2025), Ocean-R1-7B (Lingfeng et al., 2025), ThinkLite-VL-7B (Wang et al., 2025c) and OpenVLThinker-7B (Deng et al., 2025). VRAG methods with the same backbone include R1-Router (Peng et al., 2025), MMSearch-R1 (Wu et al., 2025), and VRAG-RL (Wang et al., 2025b). More implementation details of the baseline methods are provided in Appendix A.6

Evaluation Metrics. Due to inherent limitations in retrieval, the selected context may or may not provide sufficient information to answer the query. To rigorously assess both the perceptual and reasoning capabilities of the model while mitigating the visual hallucination, we categorize each query-context pair into two types: sufficient context and insufficient context (Joren et al., 2025).

For queries where the retrieved images provide sufficient evidence, we adopt the original reference answer as the ground truth. When the context is inadequate to support a correct answer, the model is required to output “insufficient to answer.” To evaluate overall performance under realistic VRAG settings, we report global *Accuracy* and *F1 Score* over all queries as comprehensive, dataset-level metrics.

Additional implementation details for the baseline methods are provided in Appendix A.6. We also compare three CoT approaches with our Evidence-guided prompt approach. Since these methods were not trained, we present them separately in Appendix A.4.

Implementation Details. We use Qwen2.5-VL-

Methods	In Distribution				Out of Distribution						Average	
	ChartQA		InfoVQA		DocVQA		SlideVQA		ViDoSeek			
	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1
Qwen2.5-VL-7B	59.20	52.80	60.86	54.61	63.28	56.03	51.62	46.11	42.56	42.48	55.50	50.41
MiMo-VL-7B-RL	54.96	40.59	68.11	45.93	74.11	47.67	77.88	47.45	<u>48.34</u>	38.30	64.68	43.99
Qwen2.5-VL-32B	69.12	60.58	<u>78.13</u>	66.06	<u>83.93</u>	73.78	<u>78.42</u>	58.65	<u>47.55</u>	52.78	<u>71.43</u>	62.37
Vision-R1	56.16	50.84	29.53	27.73	32.49	30.04	52.34	47.51	39.05	37.82	41.91	38.79
Ocean-R1-7B	47.68	47.58	53.20	53.49	56.35	57.02	60.07	57.86	40.63	46.75	51.59	52.54
MM-Eureka	64.32	58.28	40.53	40.32	56.68	54.75	63.49	58.68	44.40	47.25	53.88	51.86
ThinkLite-VL-7B	57.60	53.53	61.70	61.62	62.61	62.37	65.29	63.30	45.18	48.40	58.48	57.84
OpenVLThinker	67.60	62.72	70.47	70.51	71.74	72.51	73.02	72.63	43.52	57.27	65.27	67.13
MMSearch-R1	63.28	59.89	57.94	57.71	61.59	60.82	65.29	60.97	44.40	54.34	58.50	58.75
VRAG-RL	47.00	10.03	64.86	12.21	73.22	22.39	73.85	15.37	43.82	18.14	60.55	15.63
R1-Router	60.72	15.53	60.58	15.17	75.97	25.25	75.36	17.21	44.66	12.53	63.46	17.14
EVisRAG(ours)												
EVisRAG-3B	<u>72.64</u>	<u>72.54</u>	71.03	<u>71.83</u>	78.17	<u>79.30</u>	75.84	<u>75.49</u>	45.71	<u>60.13</u>	68.68	<u>71.86</u>
EVisRAG-7B	76.80	76.60	79.39	79.80	85.45	86.82	81.29	80.28	52.10	65.78	75.01	77.86

Table 1: Overall Performance of EVisRAG and Baselines. **Bold** denotes the highest value, and underline denotes the second highest value.

7B (Bai et al., 2025) as the backbone for our proposed EVisRAG. We use LLaMA-Factory (Zheng et al., 2024) and Easy-R1 (Yaowei et al., 2025) for open-sourcing the training framework that we used for SFT and GRPO. All experiments were executed on GPU clusters with computational capabilities comparable to NVIDIA A100 80GB GPUs. Further details on the hyperparameters that we used for SFT, GRPO are provided in Appendix A.5.

5 Results and Analysis

5.1 Overall Performance

Table 1 reports the overall results for EVisRAG and all baselines. EVisRAG-7B consistently outperforms every comparator across all benchmarks, with substantial gains over the Qwen2.5-VL-7B backbone, averaging +19% in accuracy and +27% in F1 score. These improvements indicate that an evidence-guided reasoning paradigm, coupled with RS-GRPO, strengthens perceptual grounding and enables reasoning that is explicitly conditioned on grounded evidence. Compared with RL-trained VLRMs, EVisRAG’s explicit visual perception yields a clear advantage. Within the VLRM group, models emphasizing logical reasoning (e.g., OpenVLThinker (Deng et al., 2025)) do improve question answering performance, underscoring the value of stronger reasoning. Nevertheless, EVisRAG’s added perceptual grounding closes a further gap. Moreover, the three VRAG models improve

the extraction of key evidence from retrieved context and, owing to their generalization, outperform the backbone when reasoning over multiple images. Yet they remain more than ten percentage points below EVisRAG-7B, since they neglect the need for strong perceptual grounding over multiple images with rich visual information. EVisRAG further allows a 7B parameter model to exceed the performance of considerably larger 32B parameter models. Furthermore, thrived on our RS-GRPO algorithm, EVisRAG jointly improves both perception and reasoning capabilities of VLMs, leading to a more effective and adaptable RAG framework.

5.2 Ablation Study

This section reports ablation studies that isolate three training strategies, including the evidence-guided reasoning paradigm, perception reward, and Reward-Scoped Group Relative Policy Optimization (RS-GRPO), to assess the effectiveness of EVisRAG.

As shown in Table 2, EVisRAG achieves the best results across datasets, demonstrating that the evidence-guided reasoning paradigm combined with RS-GRPO enables VLMs to precisely localize question-relevant evidence in each image and reason over the recorded cues to produce grounded answers. Training under a think-then-answer paradigm (w/o Perception) on the same data yields only modest gains, reflecting the absence of ex-

Methods	In Distribution		Out of Distribution			Avg. Acc
	ChartQA	InfoVQA	DocVQA	SlideVQA	ViDoSeek	
EVisRAG (Ours)	76.8 \pm 0.6	79.2 \pm 0.7	85.5 \pm 1.2	81.3 \pm 1.0	51.8 \pm 0.7	74.9 \pm 0.8
w/o Perception	67.2 \pm 1.3	73.3 \pm 1.6	75.7 \pm 2.1	77.3 \pm 2.0	41.8 \pm 1.2	67.1 \pm 1.6
w/o Perception Reward	69.8 \pm 1.1	74.2 \pm 1.2	79.9 \pm 3.5	77.5 \pm 2.2	48.1 \pm 1.7	69.9 \pm 1.9
w/o RS-GRPO	72.0 \pm 2.2	75.7 \pm 2.1	80.0 \pm 2.2	77.9 \pm 1.7	48.7 \pm 1.9	70.9 \pm 2.0

Table 2: Ablation study on accuracy (%) averaged over 5 runs with different random seeds. We report mean \pm standard deviation: “w/o Perception” trains the model with a standard think-then-answer approach on the same data. “w/o Perception Reward” uses only answer correctness as the reward, omitting the additional Perception Reward. “w/o RS-GRPO” sums all rewards and applies them to every token, corresponding to the standard GRPO algorithm.

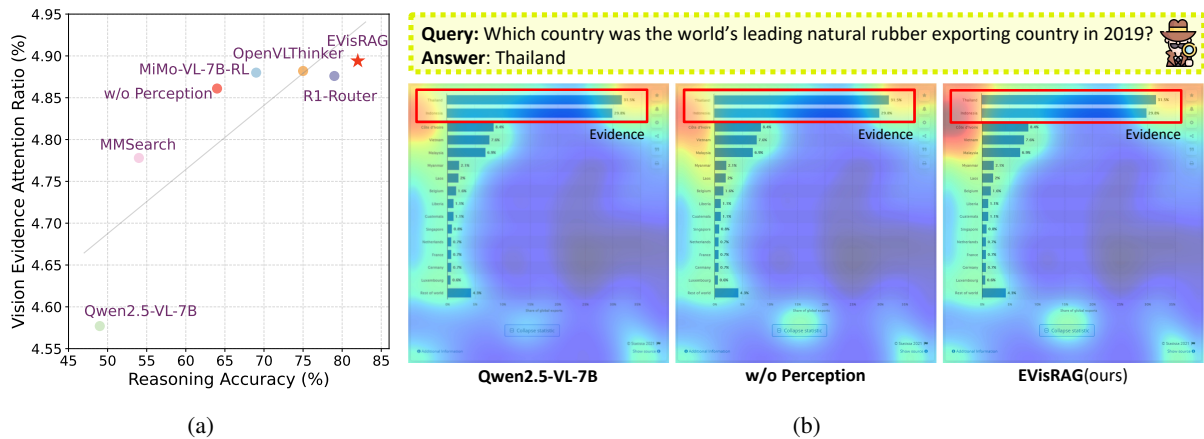


Figure 3: Comparison of models’ attention to question-relevant visual evidence. (a) Accuracy vs. attention ratio within human-annotated boxes; EVisRAG achieves the highest. (b) Qualitative maps: Compared with the baseline, EVisRAG better focuses on the top bar encoding the evidence.

447 plicit mechanisms for perceptual grounding. Introducing evidence-guided reasoning paradigm while
 448 rewarding only the final answer (w/o Perception Reward) leads to additional improvements, while
 449 underscores EVisRAG demonstrating that introducing perception reward can further enhance the per-
 450 ception ability of VLMs. Augmenting GRPO with a perception reward but without reward-scoped
 451 (w/o RS-GRPO) provides a further increment, yet the uniformly aggregated rewards dilute guid-
 452 ance across tokens. In contrast, RS-GRPO applies rewards directly within their designated re-
 453 ward scopes, sharpening credit assignment, stabilizing optimization, and ultimately delivering the
 454 strongest overall performance.

462 5.3 Evaluating Perceptual Ability through 463 Visual Attention

464 In this section, we evaluate EVisRAG’s visual perception using qualitative and quantitative evidence.
 465 Figure 3b shows a representative case for the question “Which country was the world’s leading nat-
 466 ural rubber exporting country in 2019?”. Training

467 the backbone under a think-then-answer paradigm (w/o Perception) modestly improves attention of
 470 VLMs to the legend and lower caption, helping the model read that the bars indicate shares of global
 471 exports. EVisRAG, with evidence-guided reasoning, further concentrates attention on the top bar
 472 region and correctly identifies Thailand as the leading exporter.

473 To quantify perception, we manually annotate evidence regions in more than 100 cases and compute
 474 the visual evidence attention ratio, defined as the percentage of attention mass falling inside the an-
 475 notated evidence box. As shown in Figure 3a, EVisRAG achieves the highest reasoning accuracy and
 476 the highest visual evidence attention ratio among all baselines. The scatter also reveals a clear posi-
 477 tive trend: higher attention to the evidence region is associated with higher answer accuracy.

487 5.4 Visual Evidence Density Comparison

488 We conduct a visual evidence density analysis to evaluate robustness under varying numbers of re-
 489 trieved images. As shown in Figure 4, we retrieve
 490

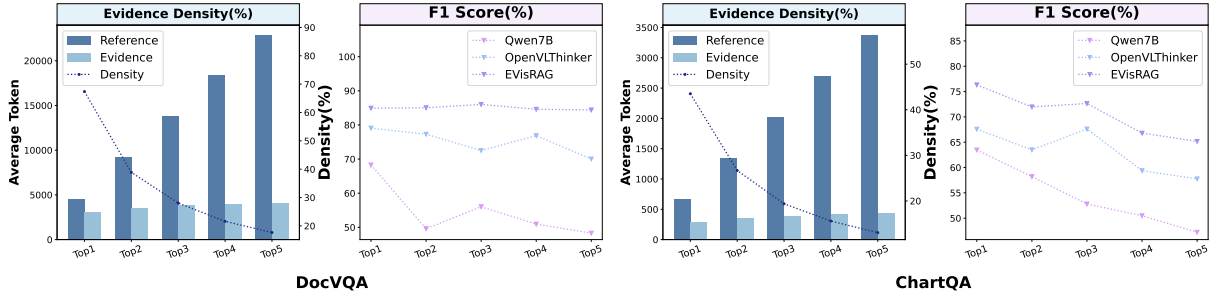


Figure 4: Performance comparison on different visual evidence density. Despite increasing noise with more retrieved images, EVisRAG maintains stable.

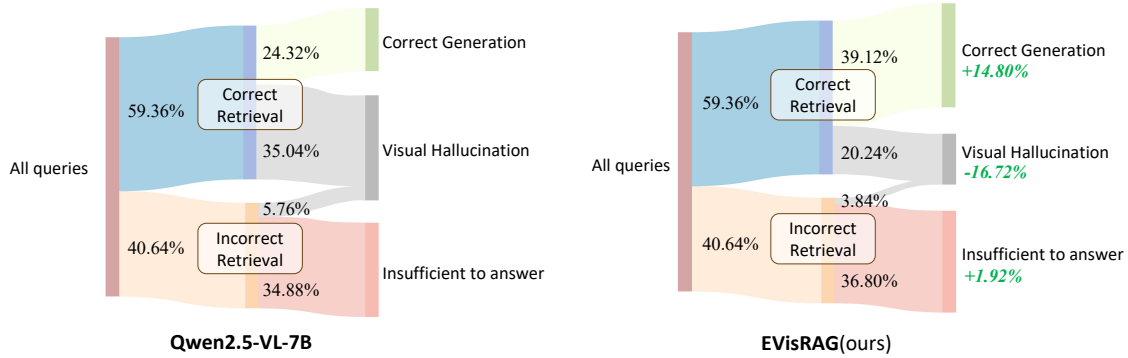


Figure 5: Model performance comparisons in different retrieval scenarios on ChartQA. Compared with the backbone, EVisRAG remains more faithful to the retrieved content in both correct and incorrect retrieval scenarios.

the top-1 to top-5 images per question and define answer-supporting image tokens as Evidence. As more images are retrieved, the total amount of Evidence increases while the evidence density decreases. We compare our method with Qwen-7B and OpenVLThinker across different evidence densities, and observe that EVisRAG consistently outperforms both baselines in terms of F1 score. Notably, on DocVQA, our method maintains stable performance even at low evidence densities, demonstrating strong resistance to visual hallucinations in multi-image settings.

5.5 Impact of Training on Model Performance

As shown in Figure 5, we examine the model’s reasoning under varying degrees of contextual sufficiency to evaluate its balance between informativeness and hallucination. Before training, the backbone model displays a strong hallucination tendency even with correct retrieval—only 24.32% of queries yield correct generations, whereas 35.04% produce incorrect responses. Under incorrect retrieval, it predominantly abstains, reflecting a conservative strategy typical of smaller models. After training with our method, EVisRAG achieves a markedly better performance: the

correct-generation rate in correctly retrieved contexts rises substantially, while a modest increase in abstention under incorrect retrieval is accompanied by a controlled reduction in incorrect generations. Overall, EVisRAG exhibits strengthened evidence-sensitive reasoning and reduced visual hallucination in underdetermined scenarios.

6 Conclusion

In this paper, we propose EVisRAG, a framework that enables Vision–Language Models (VLMs) to observe and localize evidence across multiple images during reasoning, improving visual grounding and mitigating visual hallucinations in complex multi-image scenarios. To this end, we introduce Reward-Scoped Group Relative Policy Optimization (RS-GRPO), which assigns reward signals to scope-specific token spans, stabilizing long chain-of-thought training and encouraging accurate evidence grounding. Extensive experiments show that EVisRAG effectively aggregates key visual evidence, leading to substantial gains in reasoning accuracy while reducing hallucinated responses. Overall, EVisRAG represents a promising step toward more reliable and hallucination-resistant visual retrieval-augmented generation systems.

541 Limitations

542 This work focuses on visual retrieval-augmented
543 generation and evaluates EVisRAG exclusively in
544 image-based settings. While the proposed frame-
545 work is designed to improve evidence-guided rea-
546 soning over multiple images, it does not currently
547 support hybrid multimodal contexts that jointly in-
548 volve text and images. Extending EVisRAG to
549 handle mixed text–image evidence remains an im-
550 portant direction for future work.

551 In addition, our framework is formulated and
552 evaluated under a single-turn reasoning setup,
553 where each query is answered independently based
554 on the retrieved visual context. We do not consider
555 multi-turn or interactive scenarios in which the
556 model must maintain and update evidence across
557 turns. Adapting EVisRAG to multi-turn retrieval-
558 augmented dialogue and interactive reasoning sys-
559 tems is a promising avenue for future research.

560 References

561 Chenxin An, Zhihui Xie, Xiaonan Li, Lei Li, Jun Zhang,
562 Shansan Gong, Ming Zhong, Jingjing Xu, Xipeng
563 Qiu, Mingxuan Wang, and Lingpeng Kong. 2025. *Pol-
564 aris: A post-training recipe for scaling reinforcement
565 learning on advanced reasoning models.*

566 Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and
567 Hannaneh Hajishirzi. 2024. Self-rag: Learning to
568 retrieve, generate, and critique through self-reflection.
569 In *The Twelfth International Conference on Learning
570 Representations, ICLR 2024, Vienna, Austria, May
571 7-11, 2024*. OpenReview.net.

572 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wen-
573 bin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie
574 Wang, Jun Tang, and 1 others. 2025. Qwen2.5-vl
575 technical report. *CoRR*, abs/2502.13923.

576 Florin Cuconasu, Giovanni Trappolini, Federico Sicil-
577 iano, Simone Filice, Cesare Campagnano, Yoelle
578 Maarek, Nicola Tonello, and Fabrizio Silvestri.
579 2024. The power of noise: Redefining retrieval for
580 rag systems. In *Proceedings of the 47th International
581 ACM SIGIR Conference on Research and Develop-
582 ment in Information Retrieval*, pages 719–729.

583 Yihe Deng, Hritik Bansal, Fan Yin, Nanyun Peng, Wei
584 Wang, and Kai-Wei Chang. 2025. Openvlthinker:
585 An early exploration to complex vision-language
586 reasoning via iterative self-improvement. *CoRR*,
587 abs/2503.17352.

588 Manuel Faysse, Hugues Sibille, Tony Wu, Bilel Omrani,
589 Gautier Viaud, Céline Hudelot, and Pierre Colombo.
590 2024a. Colpali: Efficient document retrieval with
591 vision language models. *CoRR*, abs/2407.01449.

Manuel Faysse, Hugues Sibille, Tony Wu, Bilel Omrani,
Gautier Viaud, Céline Hudelot, and Pierre Colombo.
2024b. Colpali: Efficient document retrieval with vi-
sion language models. *Preprint*, arXiv:2407.01449. 592
593
594
595

596 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song,
597 Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma,
598 Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-
599 r1: Incentivizing reasoning capability in llms via
600 reinforcement learning. *CoRR*, abs/2501.12948.

601 Michael Günther, Saba Sturua, Mohammad Kalim
602 Akram, Isabelle Mohr, Andrei Ungureanu, Sedigheh
603 Eslami, Scott Martens, Bo Wang, Nan Wang, and
604 Han Xiao. 2025. jina-embeddings-v4: Universal
605 embeddings for multimodal multilingual retrieval.
606 *Preprint*, arXiv:2506.18902.

607 Cheng-Yu Hsieh, Yung-Sung Chuang, Chun-Liang Li,
608 Zifeng Wang, Long Le, Abhishek Kumar, James
609 Glass, Alexander Ratner, Chen-Yu Lee, Ranjay Kr-
610 ishna, and 1 others. 2024. Found in the middle: Cali-
611 brating positional attention bias improves long con-
612 text utilization. In *Findings of the Association for
613 Computational Linguistics: ACL 2024*, pages 14982–
614 14995.

615 Hailey Joren, Jianyi Zhang, Chun-Sung Ferng, Da-
616 Cheng Juan, Ankur Taly, and Cyrus Rashtchian. 2025.
617 Sufficient context: A new lens on retrieval augmented
618 generation systems. In *The Thirteenth International
619 Conference on Learning Representations, ICLR 2025,
620 Singapore, April 24-28, 2025*. OpenReview.net.

621 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio
622 Petroni, Vladimir Karpukhin, Naman Goyal, Hein-
623 rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-
624 täschel, and 1 others. 2020. Retrieval-augmented
625 generation for knowledge-intensive NLP tasks. In
626 *Advances in Neural Information Processing Systems
627 33: Annual Conference on Neural Information Pro-
628 cessing Systems 2020, NeurIPS 2020, December 6-
629 12, 2020, virtual*.

630 Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang,
631 Yujia Zhou, Yutao Zhu, Peitian Zhang, and Zhicheng
632 Dou. 2025. Search-o1: Agentic search-enhanced
633 large reasoning models. *CoRR*, abs/2501.05366.

634 Ming Lingfeng, Li Yadong, Chen Song, Xu Jian-
635 hua, Zhou Zenan, and Chen Weipeng. 2025.
636 Ocean-r1: An open and generalizable large vision-
637 language model enhanced by reinforcement learn-
638 ing. <https://github.com/VLM-RL/Ocean-R1>.
639 Accessed: 2025-04-03.

640 Chengzhi Liu, Zhongxing Xu, Qingyue Wei, Juncheng
641 Wu, James Zou, Xin Eric Wang, Yuyin Zhou, and
642 Sheng Liu. 2025a. More thinking, less seeing? as-
643 sessing amplified hallucination in multimodal reason-
644 ing models. *CoRR*, abs/2505.21523.

645 Zhenghao Liu, Pengcheng Huang, Zhipeng Xu, Xinze
646 Li, Shuliang Liu, Chunyi Peng, Haidong Xin, Yukun
647 Yan, Shuo Wang, Xu Han, and 1 others. 2025b.
648 Knowledge intensive agents. *Available at SSRN
649 5459034*.

762	Tsung-Han Wu, Giscard Biamby, Jerome Quenum,	Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan	817
763	Ritwik Gupta, Joseph E Gonzalez, Trevor Dar-	Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma.	818
764	rell, and David M Chan. 2024. Visual haystacks:	2024. Llamafactory: Unified efficient fine-tuning of	819
765	A vision-centric needle-in-a-haystack benchmark.	100+ language models. <i>CoRR</i> , abs/2403.13372.	820
766	<i>arXiv preprint arXiv:2407.13766</i> .		
767	LLM-Core-Team Xiaomi. 2025. Mimo-vl technical		
768	report . <i>Preprint</i> , arXiv:2506.03569.		
769	Shilin Xu, Yanwei Li, Rui Yang, Tao Zhang, Yueyi Sun,		
770	Wei Chow, Linfeng Li, Hang Song, Qi Xu, Yunhai		
771	Tong, and 1 others. 2025. Mixed-r1: Unified reward		
772	perspective for reasoning capability in multimodal		
773	large language models. <i>CoRR</i> , abs/2505.24164.		
774	Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo		
775	Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao,		
776	Zhihui He, Qianyu Chen, Huarong Zhou, Zhensheng		
777	Zou, Haoye Zhang, Shengding Hu, Zhi Zheng, Jie		
778	Zhou, Jie Cai, Xu Han, and 4 others. 2024. Minicpm-		
779	v: A gpt-4v level mllm on your phone . <i>arXiv preprint</i>		
780	<i>arXiv:2408.01800</i> .		
781	Zheng Yaowei, Lu Junting, Wang Shenzhi, Feng		
782	Zhangchi, Kuang Dongdong, and Xiong Yuwen.		
783	2025. Easyr1: An efficient, scalable, multi-		
784	modality rl training framework. https://github.		
785	com/hiyouga/EasyR1 .		
786	Shi Yu, Chaoyue Tang, Bokai Xu, Junbo Cui, Jun-		
787	hao Ran, Yukun Yan, Zhenghao Liu, Shuo Wang,		
788	Xu Han, Zhiyuan Liu, and 1 others. 2025. Vis-		
789	rag: Vision-based retrieval-augmented generation on		
790	multi-modality documents. In <i>The Thirteenth Inter-</i>		
791	<i>national Conference on Learning Representations,</i>		
792	<i>ICLR 2025, Singapore, April 24-28, 2025</i> . OpenRe-		
793	view.net.		
794	Yufei Zhan, Yousong Zhu, Shurong Zheng, Hongyin		
795	Zhao, Fan Yang, Ming Tang, and Jinqiao Wang. 2025.		
796	Vision-r1: Evolving human-free alignment in large		
797	vision-language models via vision-guided reinforce-		
798	ment learning. <i>CoRR</i> , abs/2503.18013.		
799	Daoan Zhang, Junming Yang, Hanjia Lyu, Zijian Jin,		
800	Yuan Yao, Mingkai Chen, and Jiebo Luo. 2024a.		
801	Cocot: Contrastive chain-of-thought prompting for		
802	large multimodal models with multiple image inputs.		
803	<i>CoRR</i> , abs/2401.02582.		
804	Ge Zhang, Scott Qu, Jiaheng Liu, Chenchen Zhang,		
805	Chenghua Lin, Chou Leuang Yu, Danny Pan, Esther		
806	Cheng, Jie Liu, Qunshu Lin, and 1 others. 2024b.		
807	Map-neo: Highly capable and transparent bilingual		
808	large language model series. <i>CoRR</i> , abs/2405.19327.		
809	Ge Zheng, Bin Yang, Jiajin Tang, Hong-Yu Zhou, and		
810	Sibe Yang. 2023. Ddcot: Duty-distinct chain-of-		
811	thought prompting for multimodal reasoning in lan-		
812	guage models. In <i>Advances in Neural Information</i>		
813	<i>Processing Systems 36: Annual Conference on Neu-</i>		
814	<i>ral Information Processing Systems 2023, NeurIPS</i>		
815	<i>2023, New Orleans, LA, USA, December 10 - 16,</i>		
816	<i>2023</i> .		

A Appendix

A.1 Datasets

We evaluated five VQA benchmarks: InfoVQA, DocVQA, and SlideVQA were obtained from the VisRAG release (Yu et al., 2025), ChartQA from its test split (Masry et al., 2022), and ViDoSeek from ViDoRAG (Wang et al., 2025a). Each dataset provides ground-truth answer image IDs. For each question, we retrieved the top-3 images using VisRAG-Ret as contexts. Following Joren et al. (2025), we labeled it sufficient if all ground-truth images were included, otherwise insufficient. The number of questions and sufficient context ratio in the dataset are shown in Table 3.

A.2 Data Construction of Golden Reasoning Trajectories

For model training, we collected 30,000 samples from the ChartQA (Masry et al., 2022) and InfographicsVQA (Mathew et al., 2022) datasets, which were randomly divided into two subsets for SFT and GRPO in an 8:2 ratio. During the retrieval stage, VisRAG-Ret retrieves the top five candidate images for each query. While evaluation uses only the top three images as context, training leverages a variable number of retrieved images (top-1 to top-5) for data augmentation. Reasoning trajectories are constructed by generating candidate chains of thought with Qwen2.5-VL-72B and Qwen2.5-VL-7B (Bai et al., 2025), followed by a filtering process that retains only those trajectories yielding correct answers. This procedure generated 60,000 high-quality samples for SFT training, from which we extracted evidence as Ground Truth Evidence.

In the GRPO phase, we adopt a curriculum learning strategy following (An et al., 2025). Specifically, the SFT-trained model generates eight candidate completions for each sample, which are ranked according to their scores. Completions with perfect scores are excluded to mitigate overfitting. In addition, we incorporate 400 more challenging multi-hop examples from MMLongBench (Ma et al., 2024). The final GRPO training set consists of 4,000 carefully curated samples, organized to ensure a smooth progression from simple to complex instances, with a deliberate emphasis on more difficult cases to strengthen the model’s reasoning robustness. The distributions of data difficulty before and after filtering are illustrated in Figure 6.

A.3 More Details on the Fine-grained Reward

In addition to stabilizing training to ensure accurate perceptual grounding and evidence-guided reasoning in VLMs, we further introduce five reward components, namely:

Perception reward. For text-only language models, using only answer accuracy as the reward signal together with GRPO training can elicit emergent “aha moments” and strengthen reasoning abilities (Guo et al., 2025). For VLMs, however, directly optimizing answer accuracy may improve reasoning while failing to improve perceptual accuracy (Liu et al., 2025a). To optimize perceptual grounding and reasoning at the same time during training, we introduce a fine-grained perception reward:

$$R_{\text{perception}} = \frac{\sum_{i=1}^n r_i}{\sum_{i=1}^n (y_i \cdot k_{\text{pos}} + (1 - y_i) \cdot 1)},$$
$$r_i = \begin{cases} k_{\text{pos}} \cdot f_1(e_i^{\text{pred}}, e_i^{\text{gold}}), & \text{if } y_i = 1, \\ \mathbb{I}(e_i^{\text{pred}} = \text{“no relevant information”}), & \text{if } y_i = 0. \end{cases} \quad (7)$$

The perception reward assesses whether the model extracts useful visual information. For each image, the evidence recorded by the model is compared with the corresponding gold evidence. For images that contain information relevant to the question, the reward is the F1 score between the predicted and gold evidence. For images that are irrelevant, the reward equals 1 if the model correctly indicates the absence of evidence and 0 otherwise. The final perception reward is the normalized sum of the image level rewards.

Derivation reward. We employ the F1-score between the predicted answer and the gold truth as the reasoning reward, where the gold truth is set to the fixed response “insufficient to answer” when the context is incomplete.

$$R_{\text{derivation}} = f_1(a^{\text{pred}}, a^{\text{gold}}),$$
$$a^{\text{gold}} = \begin{cases} a^{\text{gold}}, & \text{if sufficient context} \\ \text{“insufficient to answer”}, & \text{if insufficient context} \end{cases} \quad (8)$$

where a^{pred} denotes the model’s predicted answer, a^{gold} denotes the ground-truth answer, and Acc_{evi} indicates whether the model’s evidence predictions for all images are correct (assigned 1 if all are correct, and 0 otherwise).

Format reward. Beyond the accuracy-based reward, we also incorporate a format reward model

Name	#Questions	Description	Sufficient Context Ratio
ChartQA	1250	Visual and Logical Reasoning about Charts	59.36%
InfoVQA	718	Question Answering on Infographic Images	92.90%
DocVQA	591	Document Visual Question Answering	83.59%
SlideVQA	556	Question Answering based on Multiple Slides	89.93%
ViDoSeek	1142	Retrieval and Reasoning on Visually Rich Documents	84.24%

Table 3: Datasets used in our experiments.

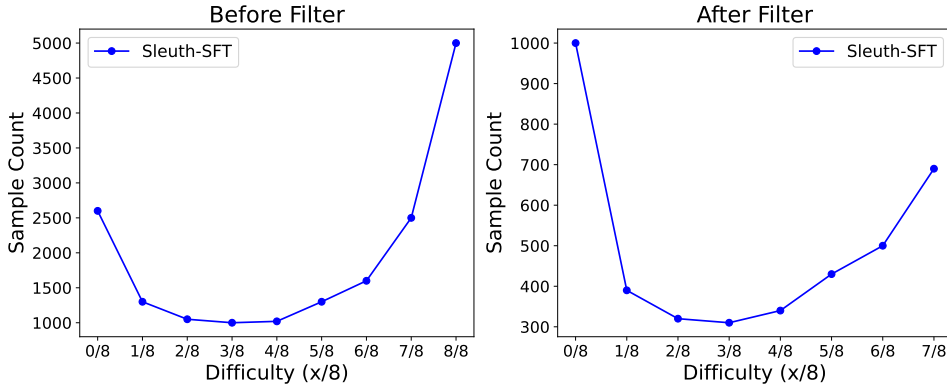


Figure 6: Data Difficulty Distribution of Before-Filtering and After-Filtering.

that compels the model to follow our CoT design by sequentially performing observation, evidence recording, reasoning, and answering, with each stage encapsulated by its corresponding special tag (<observe>, <evidence>, <think>, <answer>).

$$R_{\text{format}}(a_i) = \begin{cases} 1, & \text{if the format of } a_i \text{ is correct} \\ 0, & \text{otherwise} \end{cases} \quad (9)$$

A.4 Impact of Evidence-Guided Reasoning

To evaluate the effectiveness of Evidence-Guided Reasoning, which explicitly encourages the VLM to first observe and record visual evidence before reasoning, we conducted two additional experiments. First, we compared our reasoning paradigm against three MCOT baselines, which also avoid additional training but attempt to enhance perception and reasoning by enforcing fixed prompting patterns. As shown in Table 4, our approach consistently outperforms the baselines across five datasets. Although these MCOT strategies also prompt the model to improve perception by extensively describing image details, they tend to neglect the actual question. This often amplifies hallucinations by encouraging excessive descriptions. In contrast, our method records only question-relevant visual evidence, ensuring conciseness and enabling a more coherent and effective reasoning process. The three MCOT baselines are summarized as follows:

DDCOT (Zheng et al., 2023). A prompting strategy that decomposes complex questions into sub-questions and explicitly distinguishes between those requiring visual information and those that do not, thereby mitigating hallucinations and enhancing multimodal reasoning.

CCOT (Mitra et al., 2024). A prompting approach that leverages scene graphs as compact linguistic representations to enrich both image and task prompts, enabling LMMs to handle a wider range of vision-language tasks.

COCOT (Zhang et al., 2024a). A prompting strategy that improves the model’s ability to capture fine-grained details in multi-image tasks by guiding it to explicitly identify similarities and differences between images.

Moreover, we further evaluated the generality of the Evidence-Guided Prompting approach across models of different scales. As illustrated in Figure 7, even without additional training, prompting the model to first record visual evidence and then reason upon it consistently improves both perception and reasoning across four different model sizes. This demonstrates the broad applicability and robustness of our proposed paradigm. prompt templates used by the EVISRAG are shown in Figure 11.

Methods	Single-hop						Multi-hop				Average	
	ChartQA		InfoVQA		DocVQA		SlideVQA		ViDoSeek		Acc	F1
	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1		
MCOT												
COCOT	49.52	46.71	31.34	26.96	40.95	32.66	35.25	29.52	33.80	33.93	38.17	33.96
CCOT	50.32	48.54	36.91	35.29	41.96	39.83	51.80	46.51	36.16	41.38	43.43	42.31
DDCOT	51.68	45.55	43.73	33.53	62.10	57.98	54.14	40.37	42.21	49.04	50.77	45.29
Evidence-Guided Prompt(Ours)	62.72	54.13	65.94	62.73	70.05	65.27	66.73	61.77	46.50	56.20	62.39	60.02

Table 4: Overall Performance of EvidenceCOT and Other MCOT.



Figure 7: Performance comparison of Evidence-Guided Prompt Approach across different model sizes on the SlideVQA dataset.

A.5 More Implementation Details

We acknowledge the contributions of LLaMA-Factory and EasyR1 (Yaowei et al., 2025) for releasing the training frameworks utilized in our SFT and GRPO experiments. We adopt Qwen2.5-VL-7B (Bai et al., 2025) as the backbone model for our proposed EVisRAG. EVisRAG is trained on 8x NVIDIA A100-80GB GPUs, with hyperparameters as shown in Tables 5 and 6.

A.6 More Implementation Details of the Baseline Methods

In this section, we provide comprehensive implementation details and prompt templates of the baseline methods evaluated in our study.

General VLMs. We assessed general vision-language models across different scales, namely Qwen2.5-VL-7B and Qwen2.5-VL-32B (Bai et al., 2025), as well as MiMo-VL-7B-RL (Xiaomi, 2025).

VisRAG-Gen. We additionally evaluate two generation strategies described in VisRAG (Yu et al., 2025).

Page Concatenation. Page Concatenation forms a single composite image by horizontally concatenating the top- k retrieved pages and feeds it to

a single-image VLM. In our implementation, we adopt Qwen2.5-VL-7B (Bai et al., 2025) as the backbone VLM to ensure a fair comparison with other strong VRAG systems.

Weighted Selection. Weighted Selection instead generates an answer for each retrieved page independently and selects the final output based on the highest confidence, where the confidence weight combines the generation likelihood and the normalized retrieval score. For this method, we use the official implementation and pretrained MiniCPM-V-2 (Yao et al., 2024) model released by the authors. Together, these two variants represent the canonical generation pipelines of VisRAG and serve as competitive baselines in our evaluation.

Vision-Language Reasoning Models (VLRMs). We compare five fine-tuned VLRMs, all initialized from Qwen2.5-VL-7B-Instruct, each employing distinct strategies to enhance reasoning capabilities:

Vision-R1-7B. Vision-R1-7B (Zhan et al., 2025) introduces a reinforcement learning-based fine-tuning approach that incentivizes reasoning through vision-guided feedback. It circumvents the need for human-curated preference data by adopting a criterion-driven reward function.

Epoch	1
Data type	bf16
Learning rate	5e-7
Global batch size	32
Scheduler	Cosine
Warmup ratio	0.1
Num train epochs	1
Image max pixels	3920000

Table 5: SFT hyperparameters.

Epoch	4
Rollout batch size	32
Global batch size	32
Max grad norm	1.0
Data type	bf16
Learning rate	1e-6
Weight decay	1e-2
Warmup ratio	0.0
Rollout temperature	1.2
epsilon	0.2
epsilon_high	0.28
Image max pixels	1568000

Table 6: GRPO hyperparameters.

Methods	In Distribution				Out of Distribution						Average	
	ChartQA		InfoVQA		DocVQA		SlideVQA		ViDoSeek		Acc	F1
	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1		
VisRAG-Gen												
Page Concatenation	59.20	52.80	52.92	46.42	60.58	47.84	64.57	48.45	45.01	41.37	56.63	47.63
Weighted Selection	32.24	32.32	25.07	27.36	33.67	37.11	33.81	36.44	21.98	31.64	29.35	32.97
EVisRAG(ours)												
EVisRAG-3B	72.64	72.54	71.03	71.83	78.17	79.30	75.84	75.49	45.71	60.13	68.68	71.86
EVisRAG-7B	76.80	76.60	79.39	79.80	85.45	86.82	81.29	80.28	52.10	65.78	75.01	77.86

Table 7: Overall Performance of EVisRAG and VisRAG-Gen. **Bold** denotes the highest value.

OpenVLThinker-7B. OpenVLThinker-7B (Deng et al., 2025) follows an iterative two-stage training scheme, alternating between supervised fine-tuning (SFT) and reinforcement learning (RL). Starting from distilled reasoning competencies in text-only domains, the model progressively refines its reasoning by generating its own training data through RL and then using that data to further supervise and fine-tune itself.

MM-Eureka-7B. MM-Eureka-7B (Meng et al., 2025) extends rule-based reinforcement learning (RL) to multimodal reasoning by incorporating new algorithms such as Online Filter, ADORA, and DAPO, which enhance reasoning efficiency and stability across multimodal tasks.

Ocean-R1-7B. Ocean-R1-7B (Lingfeng et al., 2025) builds upon a structured chain-of-thought evaluation framework that leverages knowledge graph exploration (e.g., OCEAN) to provide rich offline feedback, thereby aligning generated reasoning paths with factual knowledge.

ThinkLite-VL-7B. ThinkLite-VL-7B (Wang et al., 2025c) employs Monte Carlo Tree Search

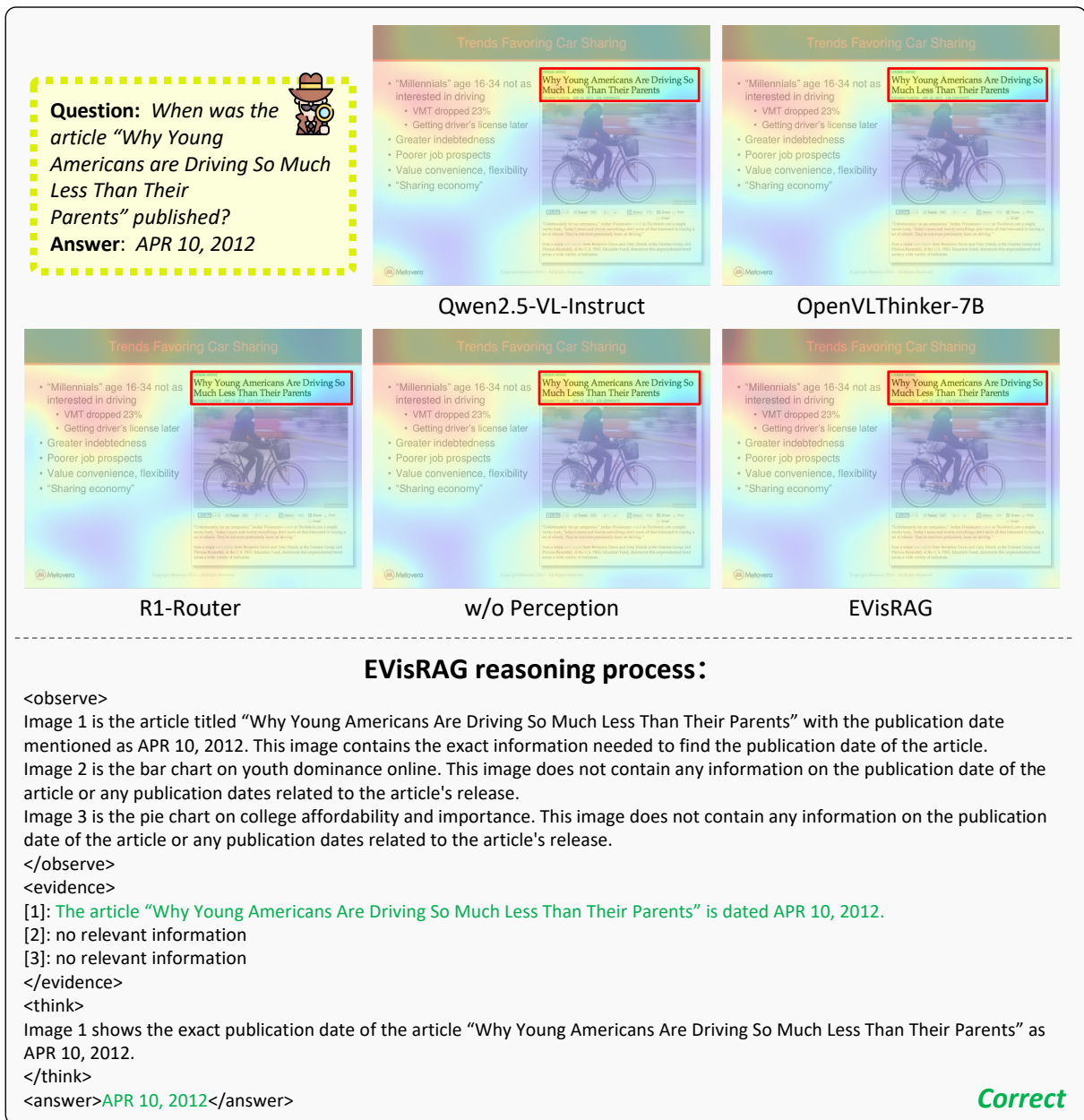
(MCTS)-guided sample selection to identify and train on genuinely challenging examples from a small dataset (11k samples), achieving state-of-the-art visual reasoning performance with high data efficiency.

VRAGs (Visual Retrieval-Augmented Generation). We further examine three advanced VRAG methods, all built upon the Qwen2.5-VL-7B-Instruct architecture:

R1-Router. R1-Router (Peng et al., 2025) employs a dynamic routing mechanism trained via Step-wise Group Relative Policy Optimization (Step-GRPO). R1-Router generates intermediate queries during the model’s reasoning process and directs them selectively to the most appropriate knowledge base (e.g., text, image, table KB), harnessing the evolving reasoning state. This fine-grained routing capability enhances retrieval efficiency and reasoning precision by minimizing unnecessary retrievals while adaptively integrating external evidence.

MMSearch-R1. MMSearch-R1 (Wu et al., 2025) integrates multimodal search into the reasoning

1062	loop, employing cross-modal retrieval mechanisms	demonstrating more stable optimization dynam-	1113
1063	to fetch contextually aligned information in both	ics. In contrast, removing any of the proposed	1114
1064	visual and textual forms.	components not only reduces accuracy but also	1115
1065	<i>VRAG-RL</i> . VRAG-RL (Wang et al., 2025b) in-	increases variance, highlighting the necessity and	1116
1066	corporates a reinforcement learning-based fine-	robustness of each part of our reward design and	1117
1067	tuning schema, enabling the model to progressively	training paradigm.	1118
1068	gather visual evidence from coarse to fine granu-		
1069	larity and support multi-turn reasoning via an opti-	A.8 More Visual Attention Cases of EVisRAG	1119
1070	mized retrieval-and-generation pipeline.	We present in Figure 8 a qualitative comparison of	1120
1071	The prompt templates employed for each base-	attention alignment with question-relevant visual	1121
1072	line are shown in Figure 12. For the three MCOT-	evidence. The query asks: When was the article	1122
1073	based comparisons (DDCOT, CCOT, and COCOT),	“Why Young Americans are Driving So Much Less	1123
1074	we adapted their original prompting strategies into	Than Their Parents” published? A human reader	1124
1075	an end-to-end chain-of-thought generation frame-	would first attend to the headline to verify topical	1125
1076	work compatible with our setup. Their correspond-	relevance, then shift gaze to the metadata directly	1126
1077	ing prompt templates are detailed in Figures 13, 14,	beneath it, where the publication date “APR 10,	1127
1078	and 15, respectively.	2012” appears. As shown in the figure, EVisRAG	1128
1079	VisRAG-Gen . We also compared two genera-	places greater attention mass on these evidence re-	1129
1080	tion methods in VisRAG (Yu et al., 2025), results	gions than the baselines, and in its reasoning trace,	1130
1081	shown in 7:	explicitly observes and records the date “APR 10,	1131
1082	<i>Page Concatenation</i> . Page Concatenation con-	2012,” yielding the correct answer. This case illus-	1132
1083	catenates all retrieved images into a single compos-	trates that EVisRAG enhances perception during	1133
1084	ite image, which is then fed to a vision–language	reasoning by aligning attention with task-critical	1134
1085	model that only supports single-image inputs; we	visual evidence.	1135
1086	adopt Qwen2.5-7B-VL as the backbone model for		
1087	this setting.	A.9 Evidence-Guided Reasoning	1136
1088	<i>Weighted Selection</i> . Following VisRAG, we	Optimization via RS-GRPO	1137
1089	prompts the model to generate an answer for each	To evaluate the effectiveness of our proposed	1138
1090	retrieved image independently and selects the final	Reward-Scoped Group Relative Policy Optimiza-	1139
1091	answer based on a confidence-weighted scoring	tion (RS-GRPO), we compare its training dynamics	1140
1092	scheme provided by the original implementation.	with the standard GRPO baseline. As shown in Fig-	1141
1093	A.7 Extended Ablation with Bootstrap	ure 10a, RS-GRPO consistently achieves higher	1142
1094	Confidence Intervals	answer rewards throughout training. While both	1143
1095	Figure 9 presents an extended visualization of the	methods exhibit fluctuations in early stages, RS-	1144
1096	ablation study in Table 2, where we report 95%	GRPO demonstrates a more stable upward trend	1145
1097	<i>bootstrap confidence intervals</i> computed from five	and converges to a substantially higher reward level.	1146
1098	runs using different random seeds. The intervals	This indicates that the fine-grained reward signals	1147
1099	are estimated via 10,000 bootstrap resamples for	applied to in-scope tokens allow RS-GRPO to bet-	1148
1100	each method–dataset pair, providing a more reli-	ter align visual perception with reasoning, leading	1149
1101	able characterization of uncertainty compared to	to more reliable improvements. Overall, these re-	1150
1102	reporting only the mean and standard deviation.	sults confirm that RS-GRPO provides more effec-	1151
1103	Across all five benchmarks, the full EVisRAG	tive optimization than GRPO, enabling EVisRAG	1152
1104	model consistently achieves the highest accuracy,	to achieve superior reasoning quality.	1153
1105	with its confidence intervals being well separated		
1106	from those of the ablated variants in nearly all cases.	A.10 Inference Efficiency of EVisRAG	1154
1107	This non–overlapping behavior indicates that the	The results in Figure 10b compare inference accu-	1155
1108	performance gains from Perception modeling, Per-	racy, latency, and output length on the ViDoSeek	1156
1109	ception Reward, and RS-GRPO are statistically	dataset across different approaches. Baseline mod-	1157
1110	significant rather than fluctuations due to random	els such as Qwen2.5-VL-7B-Instruct and Open-	1158
1111	initialization. Moreover, the bootstrap intervals	VLThinker exhibit relatively low inference time	1159
1112	of our complete method are noticeably narrower,	(around 90–95 seconds) and short outputs (approx-	1160
		imately 270–330 tokens), but their accuracy remains	1161



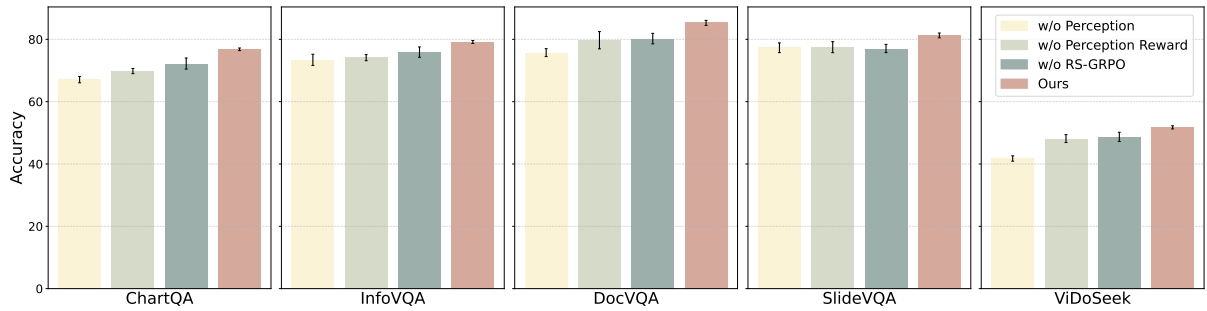


Figure 9: Ablation study(%): “w/o Perception” trains the model with a standard think-then-answer approach on the same data. “w/o Perception Reward” uses only answer correctness as the reward, omitting the additional Perception Reward. “w/o RS-GRPO” sums all rewards and applies them to every token, corresponding to the standard GRPO algorithm. Results are averaged over 5 runs with different random seeds, and error bars indicate 95% bootstrap confidence intervals.

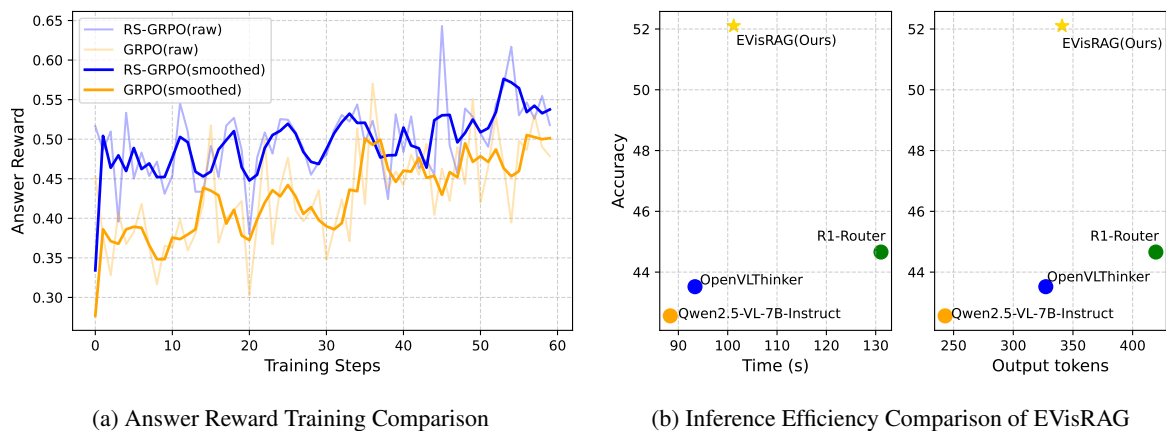


Figure 10: Training and Inference Efficiency Comparison of EVisRAG

1189 examples. For the 10- and 50-image settings, we
 1190 employ clip-vit-large-patch14-336(Radford
 1191 et al., 2021) to retrieve the top-3 most relevant im-
 1192 ages, which are then fed into the model.

1193 We compare our trained model against the origi-
 1194 nal *Qwen7B* model as the baseline. As shown in
 1195 Table 8, our approach achieves more than 20% ab-
 1196 solute accuracy improvement in the in-distribution
 1197 settings (2, 3, and 5 images). In the out-of-dis-
 1198 tribution settings (10 and 50 images), even when
 1199 relying on a relatively small CLIP model with im-
 1200 perfect retrieval quality, our method still yields
 1201 on average more than 10% absolute improvement.
 1202 These results demonstrate that our approach re-
 1203 mains highly effective on natural-image tasks and
 1204 is not limited to document-centric scenarios.

1205 A.12 Robustness to Different Retrievers

1206 To examine whether our approach depends on a
 1207 specific retrieval module, we further evaluate the
 1208 trained model under multiple independent retriev-
 1209 ers. Although our method is trained with VisRAG-

1210 Ret as the retrieval component, at inference time,
 1211 we replace the retriever with two alternative models
 1212 of different architectures and scales: Colpali-v1.3
 1213 and Jina-embeddings-v4. For each retriever, we
 1214 obtain the top-3 relevant images and feed them into
 1215 our QA model without any retraining or adapta-
 1216 tion. The results in Table 9 demonstrate that our
 1217 method yields consistent and substantial improve-
 1218 ments across all retrievers, which confirms that our
 1219 approach is retriever-agnostic.

1220 A.13 Case Studies of EVisRAG

1221 In this section, we present three case studies il-
 1222 lustrating EVisRAG’s effectiveness: (i) single-hop
 1223 QA, (ii) multi-hop QA, and (iii) alignment of atten-
 1224 tion with question-relevant visual evidence, each
 1225 compared against strong baselines.

1226 We begin with the single-hop case illustrated
 1227 in Figure 16, drawn from the DocVQA dataset.
 1228 In this example, the question asks for the name
 1229 of the chemist listed in the document. EVisRAG
 1230 correctly perceives and records the visual evidence,

	In Distribution (all images as input)			Out of Distribution (top-3 recall)	
	2 images	3 images	5 images	10 images	50 images
Qwen7b	64.67	64.44	62.00	57.8	56.2
EVisRAG(Ours)	86.22 <small>+21.55</small>	85.33 <small>+20.89</small>	83.11 <small>+21.11</small>	71.6 <small>+13.8</small>	64.5 <small>+8.3</small>

Table 8: Performance on natural-image QA.

	Sufficient Ratio	qwen7b-Acc	qwen7b-F1	evisrag-Acc	evisrag-F1
VisRAG-ret (Yu et al., 2025)	84.24	42.56	42.48	52.10 <small>+9.5</small>	65.78 <small>+23.3</small>
Colpali-v1.3 (Faysse et al., 2024b)	84.33	42.23	40.95	50.79 <small>+8.6</small>	63.82 <small>+22.9</small>
jina-embeddings-v4 (Günther et al., 2025)	85.11	40.72	53.02	49.37 <small>+8.7</small>	64.01 <small>+11.0</small>

Table 9: Performance of EVisRAG and Qwen7b with different retrievers on ViDoSeek.

1231 identifying that Richard W. Mann is annotated with
1232 the title chief chemist, and subsequently produces
1233 the correct answer, Richard W. Mann. In contrast,
1234 both OpenVLThinker and R1-Router misperceive
1235 the visual annotations during reasoning, mistakenly
1236 attributing the role of chemist to other individuals
1237 and thus generating incorrect answers.

1238 We then analyze the multi-hop case in Figure 17
1239 from the SlideVQA dataset. The question asks
1240 for the number of major languages in the country
1241 that governs mainland China and the largely self-
1242 governing territories of Hong Kong (since 1997)
1243 and Macau (since 1999). Answering requires in-
1244 tegrating evidence from two slides: one identi-
1245 fies the country as China. The other enumerates
1246 China’s major languages, including Mandarin, Yue
1247 (Cantonese), Wu (Shanghainese), Minbei (Fuzhou),
1248 Minnan (Hokkien–Taiwanese), Xiang, Gan, and
1249 Hakka, a total of eight. EVisRAG correctly records
1250 the provenance of each piece of evidence and pro-
1251 duces the correct answer, demonstrating both re-
1252 liable visual perception and cross-page reasoning.
1253 In contrast, OpenVLThinker and R1-Router fail:
1254 OpenVLThinker infers the correct subgoal but, hav-
1255 ing missed the second slide’s list, predicts that no
1256 answer exists. R1-Router locates both slides but
1257 misperceives the list and counts seven instead of
1258 eight.

1259 A.14 License

1260 We strictly comply with the original licenses and
1261 release terms of all datasets used in this work and
1262 do not redistribute any third-party raw images or
1263 proprietary data. All datasets are used solely for
1264 research and evaluation purposes. For datasets
1265 without explicitly stated open licenses, we follow

1266 their original release conditions, restrict usage to
1267 non-commercial academic research, and do not re-
1268 distribute any raw data or images. All released
1269 artifacts (including code, model checkpoints, and
1270 processed metadata) exclude any proprietary or re-
1271 stricted content.

You are an AI Visual QA assistant. I will provide you with a question and several images. Please follow the four steps below:

Step 1: Observe the Images

First, analyze the question and consider what types of images may contain relevant information. Then, examine each image one by one, paying special attention to aspects related to the question. Identify whether each image contains any potentially relevant information.

Wrap your observations within `<observe></observe>` tags.

Step 2: Record Evidences from Images

After reviewing all images, record the evidence you find for each image within `<evidence></evidence>` tags.

If you are certain that an image contains no relevant information, record it as: [i]: no relevant information (where i denotes the index of the image).

If an image contains relevant evidence, record it as: [j]: [the evidence you find for the question] (where j is the index of the image).

Step 3: Reason Based on the Question and Evidences

Based on the recorded evidences, reason about the answer to the question.

Include your step-by-step reasoning within `<think></think>` tags.

Step 4: Answer the Question

Provide your final answer based only on the evidences you found in the images.

Wrap your answer within `<answer></answer>` tags.

Avoid adding unnecessary contents in your final answer, like if the question is a yes/no question, simply answer "yes" or "no".

If none of the images contain sufficient information to answer the question, respond with `<answer>insufficient to answer</answer>`.

Formatting Requirements:

Use the exact tags `<observe>`, `<evidence>`, `<think>`, and `<answer>` for structured output.

It is possible that none, one, or several images contain relevant evidence.

If you find no evidence or few evidences, and insufficient to help you answer the question, follow the instruction above for insufficient information.

Question and images are provided below. Please follow the steps as instructed.

Question: {query}

Figure 11: The Prompt Template for EVisRAG(SFT&GRPO)

You are an AI assistant. I will provide a question and some images.

Put your reasoning process within <think></think>.

Please answer the questions based on the multiple pictures given to you, and put your your final answer in <answer></answer>.

Please try to remove irrelevant content in the final answer.

If you think there are no relevant information from the picture that can help you answer the question, answer <answer>insufficient to answer</answer> after your thinking.

Question: {query}

Figure 12: The Prompt Template for baselines.

You are an AI assistant. I will provide a query and some images. Follow these two steps:

In the first step:

Please think step-by-step about the preliminary knowledge to answer the question, deconstruct the question as completely as possible down to necessary sub-questions based on context, questions and options. Then with the aim of helping humans answer the original question, try to answer the sub-questions. The expected answering form is as follows:

Sub-questions:

1. <sub-question 1>
2. <sub-question 2>

...

Sub-answers:

1. <sub-answer 1> or 'Uncertain'
2. <sub-answer 2> or 'Uncertain'

...

For a question, assume that you do not have any information about the picture, but try to answer the sub-questions and prioritize whether your general knowledge can answer it, and then consider whether the context can help. If sub-questions can be answered, then answer in as short a sentence as possible. If sub-questions cannot be determined without information in images, please formulate corresponding sub-answer into "Uncertain".

In the second step:

Put your your final answer in <answer></answer> based on the scene graphs.

Please try to remove irrelevant content in the final answer. Like if the question is asking for yes or no, then only answer <answer>yes</answer> after your thinking.

If you think there are no relevant information from the picture that can help you answer the question, answer <answer>insufficient to answer</answer> after your thinking.

Question: {query}

Figure 13: The Prompt Template for DDCOT

You are an AI assistant. I will provide a query and some images. Follow these two steps:

In the first step:

For the provided images and its associated question, **generate a scene graph** for each images includes the following:

1. Objects that are relevant to answering the question
2. Object attributes that are relevant to answering the question
3. Object relationships that are relevant to answering the question

Then reason the answer of question based on scene graphs.

Put these process within <think></think>.

In the second step:

Put your your final answer in <answer></answer> based on the scene graphs.

Please try to remove irrelevant content in the final answer. Like if the question is asking for yes or no, then only answer <answer>yes</answer> after your thinking.

If you think there are no relevant information from the picture that can help you answer the question, answer <answer>insufficient to answer</answer> after your thinking.

Question: {query}

Figure 14: The Prompt Template for CCOT

You are an AI assistant. I will provide a query, and some images. Follow these two steps:

In the first step:

Find the **similarities and differences** of these images.

Output separately all the same points and all the differences you find.

Then reason the answer of question based on your findings.

Put these process within <think></think>.

In the second step:

Put your your final answer in <answer></answer>.

Please try to remove irrelevant content in the final answer. Like if the question is asking for yes or no, then only answer <answer>yes</answer> after your thinking.

If you think there are no relevant information from the picture that can help you answer the question, answer <answer>insufficient to answer</answer> after your thinking.

Question: {query}

Figure 15: The Prompt Template for COCOT

Question: *Who the only chemist is in the list?*
Ground Truth: *Richard W. Mann*



Evidence

EVISRAG reasoning process:

<observe>
 Image 1 is the cover of the book "Evaluating the Safety of Food Chemicals" by the National Academy of Sciences. It does not contain any information on any chemist's name, their profession, or any list of individuals and their occupations that could help identify the only chemist in the list. This image does not contain any relevant information.
 Image 2 is the list of guests with their titles and affiliations. It does not contain any information on the profession of the individuals, their specific roles as chemists, or any list that confirms who the only chemist is in the list. No direct information on the profession of each individual as a chemist or any other professional title that confirms their role in chemistry.
 Image 3 is the list of names with their titles and affiliations. It shows that Richard W. Mann is listed as the "Chief Chemist" from the Great Western Sugar Co. No other individual in the list is titled as a chemist or has a title that directly confirms their profession as a chemist. No other profession or title confirms any other individual as the only chemist in the list.
 </observe>
 <evidence>
 [1]: no relevant information
 [2]: no relevant information
 [3]: Richard W. Mann is listed as the "Chief Chemist" from the Great Western Sugar Co.
 </evidence>
 <think>
 Based on the evidence from Image 3, the only chemist in this list is Richard W. Mann.
 </think>
 <answer>Richard W. Mann</answer>

Correct

OpenVLThinker reasoning process:

<think>
 Okay, so I need to figure out who the only chemist is in the list provided. Let me start by looking at the images given. The first image is a book titled "Evaluating the Safety of Food Chemicals" by the National Academy of Sciences. That doesn't really help me with the list, but it might be related to the topic.
 Moving on to the second image, it's a list of guests with their titles and affiliations. I'll go through each one to see if anyone is a chemist. The list includes people like Mr. Joseph F. Abbott, Miss Charlotte Adams, Mrs. Edith M. Barber, and so on. None of these names have "chemist" in their titles.
 The third image is a list of names and their roles. It starts with Richard W. Mann, who is a MSG Plant in Johnstown, Colorado. Then there's Carl Peters, who is a Beet Grower Director. William Stewart is an Owner of KLMO. Ed Lehman is an Owner of Longmont Times Call. John Meyer is the President of Longmont National Bank. Robert Whyte is the President of St. Vrain Valley Bank. Reggie Batt is the President of First National Bank. Arthur G. Salberg is the President of First National Johnstown. Dr. R. K. Oldemeyer is a Director of Agricultural Research. Larry Flanders is a Company Attorney. Lloyd Dickens is a Farmer-Director of Production Credit. A. Frank Barnes is a Farmer. Harold Nesbitt is a Manager of A.S.C. Office. Art Stromquist is a Farmer.
 Looking through this list, the only person with a role that could be related to chemistry is Dr. R. K. Oldemeyer, who is a Director of Agricultural Research. However, the question is specifically about a chemist, not an agricultural researcher. So, I need to see if anyone else has a chemist title.
 After going through the list, I don't see any other chemists. The other roles are various business and agricultural positions. Therefore, the only chemist mentioned is Dr. R. K. Oldemeyer, but the question is asking for the only chemist in the list, which implies there's only one. Since the list doesn't mention any other chemists, the answer should be Dr. R. K. Oldemeyer.
 </think>
 The only chemist mentioned in the list is Dr. R. K. Oldemeyer, who is a Director of Agricultural Research. However, the question asks for the only chemist in the list, which implies there's only one. Since the list doesn't mention any other chemists, the answer should be Dr. R. K. Oldemeyer.
 <answer>Dr. R. K. Oldemeyer</answer>

Incorrect

R1-Router reasoning process:

<think>
 The question asks for the only chemist in a list. The information provided includes a list of individuals with their respective titles and affiliations. The list includes various professionals from different fields, including chemists. The relevant information can be found in the list of individuals, where titles such as "Chief Chemist" and "Chemist" are present.
 </think>
 <answer>
 The only chemist in the list is Dr. Walter H. Eddy, who is listed as an American Public Health Association member. This information is based on the list of individuals provided, where Dr. Walter H. Eddy is listed with the title "American Public Health Association" next to his name.
 </answer>

Incorrect

Figure 16: A Case Study on DocVQA

Question: How many major languages are there in the country that has control over mainland China and the largely self-governing territories of Hong Kong (since 1997) and Macau (since 1999)?

Ground Truth: 8



EVisRAG reasoning process :

<observe>

Image 1 is the Country Profile of the People's Republic of China (PRC). It mentions that the PRC has control over mainland China and the largely self-governing territories of Hong Kong (since 1997) and Macau (since 1999). This image contains the information needed to identify the country and its major languages.

Image 2 is the China Fact Sheet. It lists the major languages in China as Mandarin, Yue (Cantonese), Wu (Shanghainese), Minbei (Fuzhou), Minnan (Hokkien-Taiwanese), Xiang, Gan, and Hakka dialects. This image contains the exact information needed to find the number of major languages.

Image 3 is the map on Ethnolinguistic Groups in Afghanistan. This image does not contain any information on the major languages in the country that has control over mainland China and its territories, nor does it contain any information on the number of major languages in China or its control over the mentioned territories.

</observe>

<evidence>

[1]: The People's Republic of China (PRC) has control over mainland China and the largely self-governing territories of Hong Kong (since 1997) and Macau (since 1999).

[2]: The major languages in the People's Republic of China (PRC) are Mandarin, Yue (Cantonese), Wu (Shanghainese), Minbei (Fuzhou), Minnan (Hokkien-Taiwanese), Xiang, Gan, and Hakka dialects.

[3]: no relevant information

</evidence>

<think>

Image 2 lists the major languages in the People's Republic of China as Mandarin, Yue (Cantonese), Wu (Shanghainese), Minbei (Fuzhou), Minnan (Hokkien-Taiwanese), Xiang, Gan, and Hakka dialects. This image contains the exact information needed to find the number of major languages in the country that has control over mainland China and its territories.

</think>

<answer>8</answer>

Correct

OpenVLThinker reasoning process :

<think>

Okay, so I need to figure out how many major languages are spoken in the country that has control over mainland China and the self-governing territories of Hong Kong and Macau. Let me start by looking at the images provided.

First, I'll look at the Country Profile image. It mentions that the People's Republic of China (PRC) has control over mainland China and the self-governing territories of Hong Kong and Macau. That's the country we're looking for.

Next, I'll check the China Fact Sheet image. It lists the major languages spoken in China. The list includes Mandarin, Yue (Cantonese), Wu (Shanghainese), Minbei (Fuzhou), Minnan (Hokkien-Taiwanese), Xiang, Gan, and Hakka dialects. So, the question is asking about the number of major languages. The image mentions several dialects, but it doesn't specify how many are considered major.

However, the image doesn't provide a clear cut-off point for what constitutes a major language versus a dialect.

Given that the image lists multiple dialects, it's possible that each one is considered a major language. However, the image doesn't provide a definitive number. It's also worth noting that the image doesn't mention any other languages besides these Chinese dialects. Therefore, based on the information provided, the number of major languages is not explicitly stated. The image suggests multiple dialects, but it doesn't give a specific number. So, the answer would be that there are multiple major languages, but the exact number is not provided.

</think>

The image mentions several dialects, but it doesn't specify how many are considered major languages. Therefore, the number of major languages is not explicitly stated.

<answer>insufficient to answer</answer>

Incorrect

R1-Router reasoning process :

<think>

The question asks for the number of major languages in the country that has control over mainland China and the largely self-governing territories of Hong Kong (since 1997) and Macau (since 1999). The images provided show information about China, including its official name, capital city, area, population, literacy rate, GDP per capita, major languages, currency, and time zone. The images do not contain information about Afghanistan, which is unrelated to the question.

The relevant information from the images is: 1. The People's Republic of China (PRC) has control over mainland China and the largely self-governing territories of Hong Kong (since 1997) and Macau (since 1999). 2. The major languages in China are Mandarin, Yue (Cantonese), Wu (Shanghainese), Minbei (Fuzhou), Minnan (Hokkien-Taiwanese), Xiang, Gan, and Hakka dialects. Based on the information provided, the answer to the question is:

</think>

<answer>

There are 7 major languages in the country that has control over mainland China and the largely self-governing territories of Hong Kong (since 1997) and Macau (since 1999). These languages are Mandarin, Yue (Cantonese), Wu (Shanghainese), Minbei (Fuzhou), Minnan (Hokkien-Taiwanese), Xiang, and Gan. Hakka dialects are also considered major languages in China.

</answer>

Incorrect

Figure 17: A Case Study on SlideVQA