

# Ground-R1: Thinking with Images via Scale Relative Policy Optimization

Anonymous ACL submission

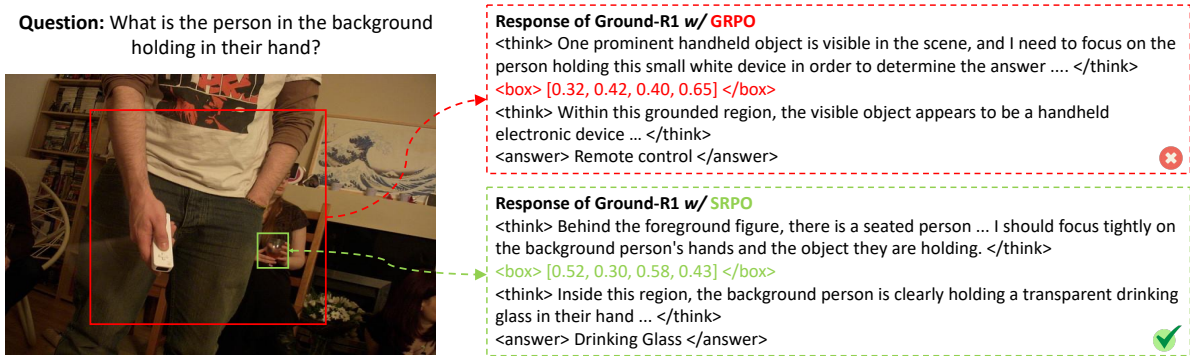


Figure 1: **Comparison between standard GRPO and our proposed SRPO.** GRPO grounds on a large, salient foreground region and produces an incorrect answer by relying on spurious context. In contrast, SRPO encourages more precise evidence localization, allowing the model to attend to the small yet critical background region.

## Abstract

Large Vision-Language Models (LVLMs) have become powerful general-purpose assistants, yet their predictions often lack reliability and interpretability due to insufficient grounding in visual evidence. The emerging thinking-with-images paradigm seeks to address this issue by explicitly anchoring reasoning to image regions. However, we empirically find that most existing methods suffer from a systematic *scale-driven bias* in optimization, where training rewards are dominated by large visual regions, suppressing learning from small but semantically critical evidence and leading to spurious grounding at inference time. To address this limitation, we propose **Ground-R1**, a de-biased thinking-with-images framework trained via a novel **Scale Relative Policy Optimization (SRPO)** objective that replaces standard GRPO. Specifically, our SRPO recalibrates reward learning across evidence regions of different sizes through scale-aware binning and intra-/inter-bin comparisons, enabling balanced credit assignment during training. Experimental results on general LVLM, high-resolution, and visual grounding benchmarks validate the effectiveness of Ground-R1 and show that SRPO yields consistent gains over standard GRPO in both response accuracy and evidence grounding.

## 1 Introduction

Large Vision-Language Models (LVLMs) (OpenAI, 2024; Anthropic, 2024; Bai et al., 2025) have emerged as versatile general-purpose assistants, capable of interpreting and executing a wide range of real-world tasks through the unified processing of visual and linguistic signals. Despite their impressive capabilities, LVLMs still suffer from limited answer reliability and poor interpretability. These issues often stem from the models' tendency to rely on spurious correlations in pre-training data distributions, rather than grounding their predictions in salient visual evidence from the input (Gupta et al., 2023; Xu et al., 2024b; Zhou et al., 2024).

To mitigate these limitations, recent studies advocate *thinking with images* (Shao et al., 2024a; Su et al., 2025; Lai et al., 2025; Zheng et al., 2025), *i.e.*, a paradigm that requires models to explicitly anchor reasoning processes to semantically grounded evidence regions. Despite this progress, most existing approaches build upon vanilla GRPO (Guo et al., 2025), which relies on outcome-level rewards without step-wise supervision, making accurate evidence grounding difficult to ensure. Empirically, we further identify a systematic *scale-driven bias*, where the training rewards are dominated by large



Figure 2: **Comparison of training reward trajectories** across evidence regions of different scales. **Top: vanilla GRPO**, where large regions consistently receive higher rewards. **Bottom: our proposed SRPO**, which alleviates scale-driven bias and yields more balanced reward signals across regions of different sizes. The large-, medium-, and small-scale regions are defined by relative area thresholds of less than 10%, between 10% and 30%, and greater than 30%, respectively.

visual regions. To substantiate this observation, we conduct diagnostic experiments to analyze reward distributions across evidence regions of different scales during training. Specifically, we track reward trajectories for large-, medium-, and small-scale regions, defined by relative area thresholds of less than 10%, between 10% and 30%, and greater than 30%, respectively. As shown in Figure 2, the large evidence regions consistently receive higher rewards than medium and small ones, indicating that they dominate the training signal. This leads to consistently negative advantages for small regions with lower mean rewards, causing their gradients to be suppressed or even clipped during optimization. Consequently, at inference time, models tend to favor visually prominent objects over small yet semantically critical evidence (*c.f.* Figure 1).

To address this issue, we propose **Ground-R1**, a de-biased thinking-with-images framework equipped with **Scale Relative Policy Optimization (SRPO)**, a novel training paradigm that recalibrates reward signals across evidence regions of different scales. Specifically, SRPO first discretizes evidence regions into scale-aware bins based on their relative areas. It then decomposes reward learning into complementary *intra-bin* and *inter-bin* components, capturing scale-calibrated relative comparisons within each scale while encouraging discriminative comparisons across scales. Furthermore, SRPO adopts a scale-aware advantage computation

strategy, where intra-bin advantages are computed using per-bin statistics and inter-bin advantages are estimated globally. This scale-relative design effectively mitigates the dominance of large regions in training, leading to more balanced reward propagation and more faithful visual evidence grounding.

We evaluate Ground-R1 across a comprehensive set of benchmarks covering general LVLM tasks, high-resolution scenarios, and visual grounding evaluations. Our method consistently delivers strong performance over supervised fine-tuning (SFT) counterparts, grounded visual reasoning approaches (Liu et al., 2024b; Shao et al., 2024a; Qi et al., 2024; Liu et al., 2024b), R1-series works (Yang et al., 2025; Peng et al., 2025; Huang et al., 2025) and concurrent thinking-with-image methods (Zheng et al., 2025; Lai et al., 2025). For example, built upon Qwen2.5-VL (Bai et al., 2025), Ground-R1 achieves a 11.9% absolute improvement on V\* (Wu and Xie, 2024). In addition, compared to standard GRPO, our proposed SRPO consistently delivers further performance gains while enabling more precise and faithful grounding of evidence regions. Specifically, Ground-R1 trained via SRPO achieves an 1.8% absolute improvement over the GRPO variant on the HR-8K benchmark. Moreover, SRPO demonstrates improved evidence region grounding over the course of training, resulting in higher grounding IoU values than GRPO (*c.f.* Figure 4).

In summary, our contributions are in three-folds:

- We introduce Ground-R1, a thinking-with-images paradigm that decomposes reasoning into grounding and answering stages to explicitly anchor predictions in visual evidence.
- We identify and systematically diagnose a scale-driven bias in GRPO-based grounded reasoning, and propose scale relative policy optimization to recalibrate rewards across evidence regions of varying sizes.
- Extensive experiments across general LVLMM, high-resolution, and visual grounding benchmarks demonstrate that Ground-R1 consistently outperforms prior methods, with SRPO yielding consistent gains over GRPO in both answering accuracy and evidence grounding.

## 2 Related Work

**Visual Reasoning in LVLMMs.** Recent advances in LVLMMs’ visual reasoning capabilities have driven sustained research efforts toward implementing cognitive processes through explicit problem-solving trajectories (Wang et al., 2025e; Huang and Chang, 2023; Xu et al., 2025a). The early attempts (Chen et al., 2023b; Yao et al., 2023; Besta et al., 2024; Wei et al., 2022) employ carefully designed prompts to guide the rationale generation. Another line of work enhances LVLMM reasoning capabilities through supervised fine-tuning on rationale annotations (Zhang et al., 2023; Wang et al., 2024a; Dong et al., 2024; Xu et al., 2024a; Thawakar et al., 2025; Yao et al., 2024). LLaVA-CoT (Xu et al., 2024a) explicitly defines the structured reasoning stages to enhance the process interpretability. Mulberry (Yao et al., 2024) augments Monte Carlo Tree Search by integrating collective knowledge from multiple LVLMMs, where negative exploration paths are leveraged to synthesize self-reflective data for iterative self-correction. Recent breakthroughs from DeepSeek-R1 (Guo et al., 2025) have catalyzed a wave of R1-inspired methodologies in vision-language research (Yang et al., 2025; Peng et al., 2025; Tan et al., 2025; Liu et al., 2025; Meng et al., 2025; Yu et al., 2025; Huang et al., 2025). Vision-R1 (Huang et al., 2025) introduces a progressive thinking suppression training strategy to address overthinking issues. LMM-R1 (Peng et al., 2025) focuses on unleashing general reasoning capabilities in 3B-parameter models, demonstrating progress in visual geometry and agent domains. Visual-RFT (Liu et al., 2025) further extends this

paradigm to broader vision tasks, including fine-grained classification and object detection. Despite recent progress, most LVLMMs still adopt a text-centric reasoning paradigm, whereas human cognition relies on interleaved visual–text reasoning grounded in visual evidence, motivating the need for more grounded reasoning frameworks.

**Thinking with Images.** Recent advances in multimodal reasoning have moved beyond text-centric paradigms toward thinking with images (Shao et al., 2024a; Su et al., 2025; Lai et al., 2025; Zheng et al., 2025), where visual information is actively involved as intermediate representations during reasoning rather than serving as static context. Existing methods can be broadly categorized into three types according to how visual information is manipulated: 1) *Tool-driven reasoning* (Wu et al., 2024; Liu et al., 2023; Qi et al., 2024; Shao et al., 2024a; Hong et al., 2025; Zheng et al., 2025; Lai et al., 2025), where models act as high-level planners that orchestrate predefined visual tools, e.g., visual, numerical, and information retrieval tools, to optimize grounded reasoning policies; 2) *Programmatic reasoning* (Gupta and Kembhavi, 2023; Surís et al., 2023; Liu et al., 2025), which treats visual reasoning as a code generation process, allowing models to compose customized and verifiable visual operations, particularly effective for structured tasks such as geometry (Mallis et al., 2025) and mathematical reasoning (Wang et al., 2025b); 3) *Intrinsic reasoning* (Team, 2024; Sun et al., 2024b; Chen et al., 2025; Xu et al., 2025b), where models internally synthesize visual representations within a closed reasoning loop, unifies visual generation and reasoning by making visual imagination a native component of the reasoning process. Our Ground-R1 follows the tool-driven reasoning paradigm, avoiding the computational overhead of programmatic execution and the uncertainty of intrinsic visual imagination (Hu et al., 2024). Moreover, we identify a systematic bias in reward learning and propose SRPO as an alternative to GRPO that recalibrates reward signals in a scale-aware manner.

## 3 Ground-R1

A schematic illustration of Ground-R1 is shown in Figure 3. The overall pipeline is described in Section 3.1, while the proposed scale-relative policy optimization is detailed in Section 3.2.

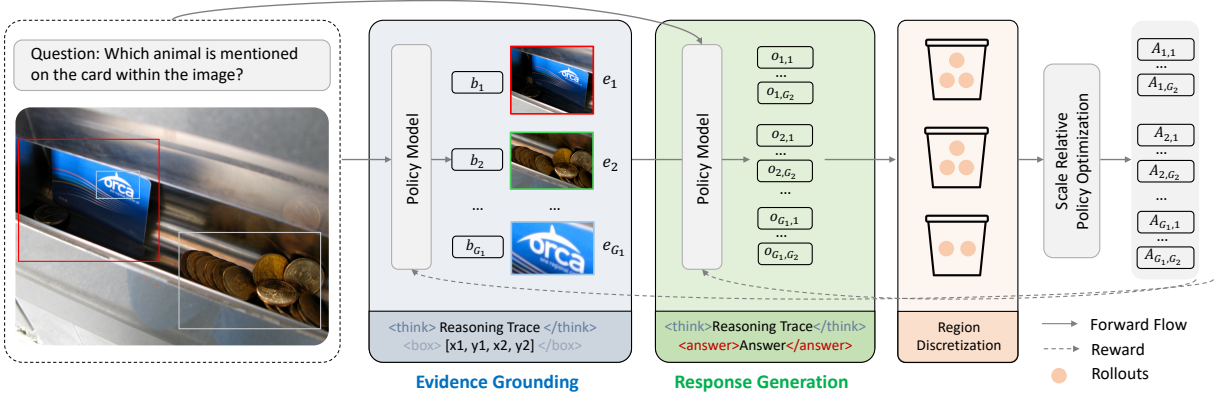


Figure 3: **Schematic illustrations of Ground-R1.** The **grounding phase** analyzes input instructions and generates evidence region rollouts.  $\mathbf{b}_i \in \mathbb{R}^4$  denotes the axis-aligned bounding box coordinates and  $e_i$  is the corresponding cropped evidence region,  $i \in [1, G_1]$ . The **answering phase** takes the input image, question, and the generated evidence regions as input and delivers final answers.  $o_{i,j}$  denote the  $j$ -th rollout answers based on the  $i$ -th evidence region  $e_i$ .  $j \in [1, G_2]$ . **Region discretization** partitions rollouts into  $K$  equal-sized bins based on their relative areas.  $A_{i,j}$  is the computed advantages (*c.f.* Eq. (5)) and the proposed SGPO is optimized via Eq. (6).

### 3.1 Overall Pipeline

**Grounding Rollout.** Given the input question  $q$  and image  $v$ , the grounding rollout process synthesizes critical evidence regions through iterative reasoning. We formally structure cognitive processes through dedicated markup tags: analytical reasoning flows are encapsulated within `<think>` and `</think>` elements, while visual grounding outputs are delimited by `<box>` and `</box>` tags.

Following GRPO (Shao et al., 2024b), a set of grounding rollouts is generated by an existing policy  $\pi_{\theta_{\text{old}}}$  as follows.

$$\mathbf{b} = \{\mathbf{b}_i\}_{i=1}^{G_1} \sim \pi_{\theta_{\text{old}}}(\cdot \mid \mathbf{q}, \mathbf{v}), \quad (1)$$

where  $\mathbf{b}_i \in \mathbb{R}^4$  denotes the axis-aligned bounding box represented by its top-left and bottom-right corner coordinates.  $G_1$  is the number of sampled grounding rollouts.

**Answer Rollout.** Leveraging the evidence region coordinates generated from the grounding rollout process, we extract corresponding image regions through spatial cropping. During the answer rollout phase, the policy model integrates multiple inputs including the original image, textual query, and evidence-aligned image regions to synthesize final responses. Consider a generated bounding box  $\mathbf{b}_i$  as an illustrative case, the corresponding evidence region  $e_i$  is derived through zoom-in and cropping operations. The answer rollout process is formally expressed as follows:

$$\mathbf{o}_i = \{\mathbf{o}_{i,j}\}_{j=1}^{G_2} \sim \pi_{\theta_{\text{old}}}(\cdot \mid \mathbf{q}, \mathbf{v}, e_i), \quad (2)$$

where  $o_{i,j}$  denote the  $j$ -th rollout answers based on the  $i$ -th evidence region  $e_i$ .  $i \in [1, G_1]$ ,  $j \in [1, G_2]$ .  $G_2$  denotes the number of answer rollouts generated per evidence region, yielding a total of  $G_1 \cdot G_2$  reasoning trajectories.

### 3.2 Scale Relative Policy Optimization

Our Ground-R1 is trained via the proposed SRPO, which extends the conventional GRPO (Shao et al., 2024b) by conditioning advantage normalization on the region scale.

**Evidence Region Discretization.** We first perform the *area-based discretization* on the evidence regions  $e_i$  generated in the grounding rollout process, partitioning them into  $K$  equal-sized bins based on their relative areas, as follows:

$$s(i) = \min(\lfloor K a_i \rfloor + 1, K) \in \{1, \dots, K\}. \quad (3)$$

where  $s(i)$  denotes the bucket index of the  $i$ -th evidence region and  $a_i$  is the relative area of  $e_i$ .

**Reward Scoring.** Since the rewards vary significantly across different area bins, we apply binary rewards including both *intra-* and *inter-bin* rewards.

- **Intra-bin reward:** It is designed to capture the relative magnitudes within each bin. Specifically, it incorporates two complementary components including a *format reward* which ensure that the responses conform to the desired format and an *answer reward* which evaluates the correctness of the final prediction. Formally, the intra-bin reward is denoted as  $r_{i,j}^{\text{intra}}$ ,  $i \in [1, G_1]$ ,  $j \in [1, G_2]$ . For each bin  $\mathcal{S}_k = \{(i, j) \mid s(i) = k\}$ ,  $k \in [1, K]$ , we compute the mean and standard devi-

ation of the intra-bin rewards, denoted as  $\mu_k$  and  $\sigma_k$ , respectively.

- *Inter-bin reward*: To encourage discriminative comparisons across different bins, we introduce an additional inter-bin reward that favors rollouts associated with the bin exhibiting the highest mean reward. Specifically, the inter-bin reward  $r_{i,j}^{\text{inter}}$  is computed as follows.

$$r_{i,j}^{\text{inter}} = \begin{cases} 1, & \text{if } i = \arg \max_{k \in \{1, \dots, K\}} \mu_k \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

**Scale-aware Advantage Estimation.** For intra-bin rewards, advantages are estimated using *per-bin* statistics instead of global normalization to ensure scale-consistent within each bin. By contrast, inter-bin rewards capture cross-bin discrimination and thus use standard *global* advantage estimation.

$$A_{i,j} = \frac{r_{i,j}^{\text{intra}} - \mu_{s(i)}}{\sigma_{s(i)}} + \frac{r_{i,j}^{\text{inter}} - \text{mean}(r_{i,j}^{\text{inter}})}{\text{std}(r_{i,j}^{\text{inter}})}, \quad (5)$$

where  $\mu_{s(i)}$  and  $\sigma_{s(i)}$  denote the mean and standard deviation of the intra-bin rewards for the bin to which the  $i$ -th evidence region belongs.

**Optimization Objective.** Our proposed SRPO is optimized by maximizing the following objective:

$$\begin{aligned} \mathcal{J}(\theta) = \mathbb{E}_{\mathbf{q}, \mathbf{v}, \{\mathbf{o}_{i,j}\}_{i=1, j=1}^{G_1, G_2} \sim \pi_{\theta_{\text{old}}}} & \left[ \frac{1}{G_1 G_2} \sum_{i=1}^{G_1} \sum_{j=1}^{G_2} \right. \\ \min & \left( \rho_{i,j}(\theta) A_{i,j}, \text{clip}(\rho_{i,j}(\theta), 1 - \varepsilon, 1 + \varepsilon) A_{i,j} \right) \\ & \left. - \beta \mathbb{D}_{\text{KL}}(\pi_{\theta} \parallel \pi_{\text{ref}}) \right], \quad (6) \end{aligned}$$

where  $\varepsilon$  is clipping-related hyper-parameters for stabilizing training and the importance sampling ratio  $\rho_{i,j}(\theta)$  is defined as follow:

$$\rho_{i,j}(\theta) = \frac{\pi_{\theta}(\mathbf{o}_{i,j} \mid \mathbf{q}, \mathbf{o}_{i, < j})}{\pi_{\theta_{\text{old}}}(\mathbf{o}_{i,j} \mid \mathbf{q}, \mathbf{o}_{i, < j})}. \quad (7)$$

## 4 Experiments

### 4.1 Experimental Settings

**Training Dataset.** We adopt the dataset introduced by DeepEyes (Zheng et al., 2025) including fine-grained visual search (Wu and Xie, 2024), arXivQA (Li et al., 2024b), ThinkLite-VL (Wang et al., 2025d) for RL training. Note that we do not employ curated SFT dataset for cold-start training, as the distribution of SFT data can substantially alter the grounding behavior of LVLMS (Lai et al., 2025),

which would hinder our investigation of the base model’s inherent behavior under the thinking-with-images paradigm.

**Implementation Details.** We conducted experiments using Qwen-2.5-VL-7B-Instruct (Bai et al., 2025). The training was performed for 1,000 steps with a batch size of 8 and a learning rate of  $1 \times e^{-6}$ . The rollout numbers for the grounding ( $G_1$ ) and answering ( $G_2$ ) stages were set to 4 and 2, respectively. The number of equal-sized bins  $K$  was set to 3. We used a default sampling temperature of 1 and a maximum response length of 512 tokens. All experiments were run on  $8 \times \text{H100}$  GPUs, with a total training time of approximately 12 hours.

### 4.2 Evaluation on LVLMS Benchmarks.

**Results on General Benchmarks.** Table 1 reports the quantitative comparisons on six widely used benchmarks, including MME (Fu et al., 2024), MM-Vet (Yu et al., 2023), SEED-Bench (Li et al., 2023a), MME-RealWorld-Lite (Zhang et al., 2024), RealworldQA (xAI, 2024), and POPE (Li et al., 2023b). These benchmarks jointly evaluate general multi-modal reasoning and hallucination resistance. We include strong proprietary models (e.g., GPT-4V and GPT-4o), state-of-the-art open-source LVLMS (e.g., InternVL2.5, Qwen2-VL), and recent thinking-with-images (e.g., DeepEyes, Mini-o3) and R1-style approaches (e.g., Vision-R1, LMM-R1) for comprehensive comparison.

As shown in Table 1, Ground-R1 consistently outperforms its base model Qwen2.5-VL-7B across all evaluated benchmarks, e.g., achieving notable absolute improvements of +83.4 on MME and +17.1% on MME-RWL. These results indicate that Ground-R1 substantially enhances both general multimodal reasoning and grounding-sensitive capabilities without sacrificing robustness.

**Results on High-Resolution Benchmarks.** We provide the comparison results on high-resolution benchmarks, including V (Wu and Xie, 2024) and HR-Bench (Wang et al., 2025c), which pose more challenging evaluations due to high-resolution images (4K/8K) and the presence of small-scale objects. As shown, Ground-R1 consistently outperforms Qwen2.5-VL-7B as well as prior methods across all high-resolution settings, achieving notable absolute improvements of +11.9% on V\*, +6.8% on HR-4K, and +8.4% on HR-8K. These gains indicate that Ground-R1 is particularly effective over high-resolution scenarios where critical evidence often occupies a small spatial extent.

Table 1: **Evaluation results on LVLB benchmarks** including general and high-resolution scenarios. HR-4K/8K denotes HR-Bench (Wang et al., 2025c) with the high resolution of 4K/8K. MME-RWL denotes the MME-RealWorld lite version (Zhang et al., 2024). † denotes results reproduced by us under the same experimental setting.

| Method                              | MME           | MM-Vet      | SEED        | MME-RWL     | RWQA        | POPE        | V*          | HR-4K       | HR-8K       |
|-------------------------------------|---------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| GPT-4V (OpenAI, 2023)               | 1926.6        | 67.5        | 53.8        | –           | 61.4        | –           | –           | –           | –           |
| GPT-4o (OpenAI, 2024)               | –             | 69.1        | 72.0        | –           | 75.4        | 86.9        | 65.2        | 62.0        | 58.3        |
| InternVL2-8B (Chen et al., 2024b)   | 2210.3        | 54.2        | –           | –           | 64.4        | 86.9        | –           | –           | –           |
| InternVL2.5-8B (Chen et al., 2024a) | 2344.1        | 62.8        | –           | –           | 70.1        | 90.6        | –           | –           | –           |
| Qwen2-VL-7B (Wang et al., 2024b)    | 2326.8        | 62.0        | 75.1        | –           | 70.1        | 88.1        | –           | –           | –           |
| CogCoM (Qi et al., 2024)            | –             | 46.1        | –           | –           | –           | –           | –           | –           | –           |
| Chain of Spot (Liu et al., 2024b)   | 1501.1        | 30.8        | 59.7        | –           | –           | 86.4        | –           | –           | –           |
| FaST (Sun et al., 2024a)            | 1517.4        | 31.0        | 60.1        | –           | –           | 86.4        | –           | –           | –           |
| Vision-R1-7B (Huang et al., 2025)   | 2306.2        | 58.3        | 73.0        | 41.0        | 50.2        | 88.7        | –           | –           | –           |
| LMM-R1 (Peng et al., 2025)          | 2210.6        | 64.5        | 74.4        | 35.4        | 52.2        | 86.5        | –           | –           | –           |
| R1-Onevision (Yang et al., 2025)    | 2192.2        | 67.5        | 66.5        | 35.1        | 46.5        | 84.9        | –           | –           | –           |
| LLaVA-OneVision (Li et al., 2024a)  | –             | –           | –           | –           | –           | –           | 75.4        | 63.0        | 59.8        |
| DeepEyes† (Zheng et al., 2025)      | –             | –           | –           | –           | –           | 87.7        | 83.3        | 73.2        | 69.5        |
| Mini-o3† (Lai et al., 2025)         | –             | –           | –           | –           | –           | –           | 86.9        | 74.9        | 70.6        |
| Pixel Reasoner (Wang et al., 2025a) | –             | –           | –           | –           | –           | –           | 86.3        | 74.0        | 66.9        |
| Qwen2.5-VL-7B (Bai et al., 2025)    | 2327.0        | 67.1        | 76.5        | 33.8        | 65.9        | 83.6        | 75.5        | 68.2        | 62.7        |
| Ground-R1 (Ours)                    | <b>2410.4</b> | <b>69.8</b> | <b>77.8</b> | <b>50.9</b> | <b>68.8</b> | <b>89.8</b> | <b>87.4</b> | <b>75.0</b> | <b>71.1</b> |
| Δ (vs. Qwen2.5-VL-7B)               | +83.4         | +2.7        | +1.3        | +17.1       | +2.9        | +6.2        | +11.9       | +6.8        | +8.4        |

### 4.3 Visual Grounding

To ensure the broad coverage of multi-modal competencies, we also evaluate the visual grounding capabilities of Ground-R1 on RefCOCO (Kazemzadeh et al., 2014), RefCOCO+ (Kazemzadeh et al., 2014) and RefCOCOg (Mao et al., 2016). We use the following prompts as follows, *i.e.*, “Locate <ref> in this image and output the bbox coordinates in JSON format.”, where <ref> refers to the specific expression.

As shown in Table 2, Ground-R1 demonstrates superior visual grounding capabilities, achieving 93.1% accuracy on RefCOCO val and surpassing most LVLBs, while narrowing the performance gap between generalist and specialist models (*e.g.*, Grounding DINO (Liu et al., 2024a)). Although Ground-R1 slightly underperforms CogCoM (Qi et al., 2024) on specific benchmarks (*e.g.*, test-A of RefCOCO), this discrepancy is explainable: CogCoM leverages a specialized grounded VQA dataset (Wang et al., 2024c) for additional training to enhance visual grounding, whereas our Ground-R1 maintains dual competence in both visual reasoning and grounding without requiring dedicated grounding-specific training.

### 4.4 Ablation Studies

To justify the design of our approach, we compare Ground-R1 against the following model variants:

- **Vanilla-SFT**: This variant replaces the RL-based training in Ground-R1 with SFT.
- **Vanilla-R1**: This variant omits the grounding

phase of Ground-R1 and leverages RL training.

- **Ground-R1-GRPO**: This variant replaces the proposed SRPO with the standard GRPO.
- **Ground-R1-Intra/Inter**: These two variants use only the intra-bin or inter-bin rewards in SRPO, respectively (*c.f.* Section 3.2).

The comparison results are demonstrated in Table 3. Building upon these model variants, we seek to address four critical inquiries:

**Q1**: Does RL-based methodology outperform SFT-based approaches in visual reasoning?

By comparing Ground-R1 with Vanilla-SFT, we observe that RL-based methods consistently outperform SFT-based counterparts across all the benchmarks. For example, Ground-R1 achieves an absolute improvement of 14.1% over Vanilla-SFT on V\* Bench. These results indicate that RL-based optimization more effectively aligns model behavior with task objectives by providing direct feedback on reasoning outcomes.

**Q2**: Does the “grounding-then-answering” paradigm outperform direct answer generation?

To answer this question, we compare Ground-R1 with Vanilla-R1, where the latter removes the explicit grounding stage and directly generates answers. As shown in Table 3, Ground-R1 consistently outperforms Vanilla-R1 across all the evaluated benchmarks. These results indicate that explicitly separating grounding from answering enables the model to localize task-relevant visual evidence more effectively, leading to more accurate reasoning than direct answer generation.

Table 2: Evaluation results on visual grounding benchmarks including RefCOCO, RefCOCO+ and RefCOCOg.

| Method                              | RefCOCO     |             |             | RefCOCO+    |             |             | RefCOCOg    |             |
|-------------------------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
|                                     | val         | test-A      | test-B      | val         | test-A      | test-B      | val         | test        |
| Gemini 1.5 Pro (Reid et al., 2024)  | 73.2        | 72.9        | 74.6        | 62.5        | 63.9        | 65.0        | 75.2        | 76.2        |
| Grounding DINO (Liu et al., 2024a)  | 90.6        | 93.2        | 88.2        | 88.2        | 89.0        | 75.9        | 86.1        | 87.0        |
| Shikra-7B (Chen et al., 2023a)      | 87.0        | 90.6        | 80.2        | 81.6        | 87.4        | 72.1        | 82.3        | 82.2        |
| Shikra-13B (Chen et al., 2023a)     | 87.8        | 91.1        | 81.8        | 82.9        | 87.8        | 74.4        | 82.6        | 83.2        |
| Qwen2-VL (Bai et al., 2023)         | 89.4        | 92.3        | 85.3        | 83.1        | 88.3        | 77.2        | 85.6        | 85.5        |
| CogVLM (Wang et al., 2024c)         | 92.5        | 94.0        | 88.7        | 87.5        | 91.8        | 81.4        | 89.5        | 90.1        |
| CogCoM (Qi et al., 2024)            | 92.3        | 94.6        | 89.2        | 88.2        | 92.8        | 82.1        | 89.3        | 90.5        |
| Qwen2.5-VL-7B (Bai et al., 2025)    | 90.0        | 92.5        | 85.4        | 84.2        | 89.1        | 76.9        | 87.2        | 87.2        |
| InternVL2.5-8B (Chen et al., 2024a) | 90.3        | 94.5        | 85.9        | 85.2        | 91.5        | 78.8        | 86.7        | 87.6        |
| Vision-R1-7B (Huang et al., 2025)   | 58.3        | 64.2        | 48.5        | 50.9        | 61.1        | 38.5        | 51.9        | 52.6        |
| LMM-R1 (Peng et al., 2025)          | 87.3        | 91.4        | 83.6        | 79.2        | 87.1        | 72.9        | 84.1        | 84.1        |
| R1-Onevision (Yang et al., 2025)    | 45.0        | 51.6        | 35.3        | 47.2        | 57.0        | 37.1        | 53.0        | 48.2        |
| DeepEyes (Zheng et al., 2025)       | 89.8        | –           | –           | 83.6        | –           | –           | 86.7        | –           |
| Qwen2.5-VL-7B (Bai et al., 2025)    | 90.0        | 92.5        | 85.4        | 84.2        | 89.1        | 76.9        | 87.2        | 87.2        |
| Ground R1 (Ours)                    | <b>93.1</b> | <b>94.0</b> | <b>88.2</b> | <b>86.9</b> | <b>91.2</b> | <b>78.8</b> | <b>89.8</b> | <b>90.2</b> |
| $\Delta$ (vs. Qwen2.5-VL-7B)        | +3.1        | +1.5        | +2.8        | +2.7        | +2.1        | +1.9        | +2.6        | +3.0        |

Table 3: Ablation studies of Ground-R1. Refer to Section 4.4 for the configuration of each model variants.

| Method          | MME           | MM-Vet      | SEED        | MME-RWL     | RWQA        | POPE        | V*          | HR-4K       | HR-8K       |
|-----------------|---------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| Ground-R1       | <b>2410.4</b> | <b>69.8</b> | <b>77.8</b> | <b>50.9</b> | <b>68.8</b> | <b>89.8</b> | <b>87.4</b> | <b>75.0</b> | <b>71.1</b> |
| Vanilla-SFT     | 2337.0        | 68.0        | 76.7        | 34.2        | 67.1        | 87.6        | 73.3        | 66.6        | 60.0        |
| Vanilla-R1      | 2338.8        | 53.2        | 75.4        | 35.1        | 66.5        | 88.4        | 80.1        | 71.3        | 67.9        |
| Ground-R1-GRPO  | 2395.1        | 68.8        | 77.6        | 49.4        | 67.7        | 89.7        | 85.3        | 73.8        | 69.3        |
| Ground-R1-Intra | 2358.7        | 59.2        | 76.7        | 43.6        | 63.9        | 88.8        | 81.7        | 71.8        | 68.5        |
| Ground-R1-Inter | 2390.5        | 67.9        | 76.9        | 44.1        | 65.5        | 89.3        | 81.2        | 72.6        | 68.8        |

**Q3:** To what extent does the proposed SRPO outperform standard GRPO?

To quantify the effectiveness of SRPO, we compare Ground-R1 with Ground-R1-GRPO, which shares the same architecture but replaces SRPO with standard GRPO. As shown in Table 3, Ground-R1 consistently outperforms Ground-R1-GRPO across all evaluated benchmarks. Notably, SRPO brings substantial gains on high-resolution benchmarks, achieving an improvement of +2.1% on V\*, +1.2% on HR-4K, and +1.8% on HR-8K. These results demonstrate that SRPO provides additional benefits over standard GRPO by more effectively recalibrating reward signals across evidence regions of different scales.

**Q4:** How do the intra- and inter-bin reward components of SRPO contribute to the performance?

To assess the contribution of each reward component in SRPO, we compare Ground-R1 with its two ablated variants, Ground-R1-Intra and Ground-R1-Inter, which retain only the intra-bin or inter-bin reward, respectively. As shown in Table 3, both variants underperform the full Ground-R1 model across all benchmarks, indicating that neither component alone is sufficient. In particular, Ground-R1-Intra exhibits notable degradation, suggesting that

intra-bin normalization alone cannot adequately distinguish evidence across different scales.

**Q5:** Compared to GRPO, does SRPO lead to more accurate grounded evidence region?

To investigate this, we analyze the training dynamics of Ground-R1 trained with SRPO and standard GRPO. As shown in Figure 4, we report the normalized average reward, grounding IoU, and tool call count over the course of training. Following (Zheng et al., 2025), grounding IoU is evaluated on the fine-grained visual search split of the training dataset, which provides ground-truth bounding boxes closely aligned with the target answers.

As shown in Figure 4, our SRPO consistently achieves higher average rewards and substantially improves grounding IoU compared to GRPO, indicating more accurate and faithful evidence region grounding. Notably, this performance gap becomes particularly pronounced after approximately 600 training steps. We attribute this to the inherent optimization bias of GRPO, which primarily favors large and visually salient bounding boxes (*c.f.* Figure 2). As training progresses, the policy increasingly exploits previously high-reward large regions, leading to an *exploration collapse* that prevents revisiting small but potentially critical evidence.

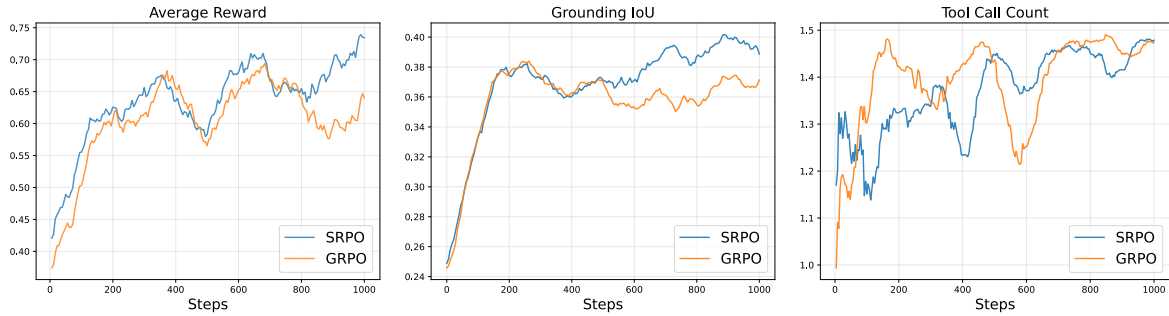




Figure 4: **Training dynamics of Ground-R1** including normalized average reward, grounding IoU, and tool call count under standard GRPO and our proposed SRPO.



**Ground-R1 w/ GRPO:**

```
<think>
Let me think through this carefully. The image mainly depicts a residential neighborhood with houses, greenery, and an open field in the foreground. Residential areas often contain roads and parked cars, so it is reasonable to assume that vehicles might be present somewhere within the scene. However, without being able to clearly isolate or confirm the exact locations of individual vehicles, it is difficult to determine whether a black vehicle or a silver vehicle is visible.
</think>
<box> [0.25, 0.39, 1.00, 0.69] </box>
<think> Given this uncertainty and the inability to clearly identify distinct vehicles in the image, I conclude that the vehicles cannot be reliably seen, and therefore their relative distance to the camera cannot be determined. </think>
<answer> The vehicles are not visible in the image, it cannot be determined which one is closer. </answer>
```



**Question:**  
Which one is closer to the camera, the black vehicle or the silver vehicle?



**Ground-Truth  
Grounded Region**

**Ground-R1 w/ SRPO:**

```
<think>
Let me reason through this step by step. The image shows a grassy field in the foreground and a residential area behind it. Between the houses, there is a road where vehicles may be parked. Both a black vehicle and a silver vehicle are not obvious. The most relevant evidence lies along the road area in front of the houses where a silver vehicle might appear.
</think>
<box> [0.60, 0.35, 0.75, 0.50] </box>
<think>
Focusing on this cropped region, the silver vehicle can be seen parked on the road, with clear spatial cues such as its relative size and vertical placement. The silver vehicle appears slightly larger and lower in scale and sits a bit lower in the frame compared with the black vehicle farther up the road, so it is nearer to the camera than the black vehicle located farther back.
</think>
<answer> The silver vehicle is closer to the camera. </answer>
```




Figure 5: **Comparisons of reasoning trajectories** of Ground-R1 under standard GRPO and our proposed SRPO.

## 4.5 Visualizations

Figure 5 provides a qualitative comparison between Ground-R1 trained with standard GRPO and our proposed SRPO. Under GRPO, Ground-R1 fails to localize the relevant evidence and instead grounds its reasoning on a large, visually salient region containing multiple houses and background context. As a result, the model defaults to an uncertain or incorrect response. In contrast, Ground-R1 trained with SRPO exhibits more precise evidence grounding by focusing on a small, task-relevant region where the vehicles are visible. By accurately localizing the silver vehicle and leveraging relative size cues, SRPO enables the model to correctly determine the relative distance to the camera.

## 5 Conclusion

In this work, we investigate a fundamental limitation of existing thinking-with-images approaches, namely a systematic scale-driven bias that favors large visual regions during optimization. To address this issue, we propose Ground-R1, equipped with scale relative policy optimization, which recalibrates reward signals across evidence regions of different sizes. Extensive experiments demonstrate that Ground-R1 not only improves response accuracy but also yields more precise and faithful evidence grounding, offering a more robust pathway toward trustworthy and interpretable vision-language reasoning.

## Ethical Considerations

Our work contributes to enhancing the reliability and interpretability of LVLMs through a de-biased thinking-with-images pattern. By explicitly mitigating the scale-driven bias, Ground-R1 with SRPO reduces the model’s tendency to rely on visually prominent yet semantically irrelevant regions, thereby lowering the risk of generating predictions based on spurious correlations. This advancement in evidence-based reasoning is particularly relevant for high-stakes applications, such as medical imaging analysis or autonomous systems, where overlooking subtle but critical details could have serious consequences. All data used in our experiments are from publicly available benchmark datasets intended for research purposes, and no personally identifiable information is involved. However, as with any foundation model enhancement, the improved capability could be deployed in dual-use scenarios. We advocate for responsible development and application of such technologies, including thorough testing in domain-specific contexts before real-world deployment.

## Limitation

While Ground-R1 significantly improves evidence grounding and response accuracy, our approach has the following two limitations that warrant further investigation. First, Ground-R1 remains dependent on the quality of region proposals, *i.e.*, suboptimal proposals or missing candidate regions may still lead to incorrect grounding, particularly for highly occluded or visually ambiguous objects. Second, our evaluation focuses on established benchmarks while the performance and grounding fidelity of Ground-R1 in more open-world, long-tail, or domain-shifted scenarios remain to be explored. Future work should aim to integrate more robust region generation strategies and extend evaluations to broader real-world settings with greater visual and semantic diversity.

## References

Anthropic. 2024. [Claude 3.5 sonnet](#).

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

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, and 1 others. 2025. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*.

Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and 1 others. 2024. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17682–17690.

Jiuhai Chen, Zhiyang Xu, Xichen Pan, Yushi Hu, Can Qin, Tom Goldstein, Lifu Huang, Tianyi Zhou, Saining Xie, Silvio Savarese, and 1 others. 2025. Blip3-o: A family of fully open unified multimodal models-architecture, training and dataset. *arXiv preprint arXiv:2505.09568*.

Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. 2023a. Shikra: Unleashing multimodal llm’s referential dialogue magic. *arXiv preprint arXiv:2306.15195*.

Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, and 1 others. 2024a. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*.

Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, and 1 others. 2024b. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*.

Zhenfang Chen, Qinhong Zhou, Yikang Shen, Yining Hong, Hao Zhang, and Chuang Gan. 2023b. See, think, confirm: Interactive prompting between vision and language models for knowledge-based visual reasoning. *arXiv preprint arXiv:2301.05226*.

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

Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, and 1 others. 2024. Mme: A comprehensive evaluation benchmark for multimodal large language models, 2024. URL <https://arxiv.org/abs/2306.13394>, 2.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shiron Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.

|     |   |     |
|-----|---|-----|
| 610 | Shashank Gupta, Vaishnavi Shrivastava, Ameet Deshpande, Ashwin Kalyan, Peter Clark, Ashish Sabharwal, and Tushar Khot. 2023. Bias runs deep: Implicit reasoning biases in persona-assigned llms. <i>arXiv preprint arXiv:2311.04892</i> .   | 664 |
| 611 |   | 665 |
| 612 |   | 666 |
| 613 |   | 667 |
| 614 |   |     |
| 615 | Tanmay Gupta and Aniruddha Kembhavi. 2023. Visual programming: Compositional visual reasoning without training. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 14953–14962.  | 668 |
| 616 |   | 669 |
| 617 |   | 670 |
| 618 |   | 671 |
| 619 |   | 672 |
| 620 | Jack Hong, Chenxiao Zhao, ChengLin Zhu, Weiheng Lu, Guohai Xu, and Xing Yu. 2025. Deepeyesv2: Toward agentic multimodal model. <i>arXiv preprint arXiv:2511.05271</i> .   | 673 |
| 621 |   | 674 |
| 622 |   | 675 |
| 623 |   | 676 |
| 624 | Yushi Hu, Otilia Stretcu, Chun-Ta Lu, Krishnamurthy Viswanathan, Kenji Hata, Enming Luo, Ranjay Krishna, and Ariel Fuxman. 2024. Visual program distillation: Distilling tools and programmatic reasoning into vision-language models. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 9590–9601. | 677 |
| 625 |   | 678 |
| 626 |   | 679 |
| 627 |   | 680 |
| 628 |   | 681 |
| 629 |   | 682 |
| 630 |   |     |
| 631 | Jie Huang and Kevin Chen-Chuan Chang. 2023. Towards reasoning in large language models: A survey. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 1049–1065.  | 683 |
| 632 |   | 684 |
| 633 |   | 685 |
| 634 |   | 686 |
| 635 | Wenxuan Huang, Bohan Jia, Zijie Zhai, Shaosheng Cao, Zheyu Ye, Fei Zhao, Zhe Xu, Yao Hu, and Shaohui Lin. 2025. Vision-rl: Incentivizing reasoning capability in multimodal large language models. <i>arXiv preprint arXiv:2503.06749</i> .   | 687 |
| 636 |   | 688 |
| 637 |   | 689 |
| 638 |   | 690 |
| 639 |   | 691 |
| 640 | Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. 2014. Referitgame: Referring to objects in photographs of natural scenes. In <i>Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)</i> , pages 787–798.  | 692 |
| 641 |   | 693 |
| 642 |   | 694 |
| 643 |   | 695 |
| 644 |   | 696 |
| 645 |   | 697 |
| 646 | Xin Lai, Junyi Li, Wei Li, Tao Liu, Tianjian Li, and Hengshuang Zhao. 2025. Mini-o3: Scaling up reasoning patterns and interaction turns for visual search. <i>arXiv preprint arXiv:2509.07969</i> .  | 698 |
| 647 |   | 699 |
| 648 |   |     |
| 649 |   |     |
| 650 | Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024a. Llava-onevision: Easy visual task transfer. <i>arXiv preprint arXiv:2408.03326</i> .   | 700 |
| 651 |   | 701 |
| 652 |   | 702 |
| 653 |   | 703 |
| 654 |   | 704 |
| 655 | Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. 2023a. Seed-bench: Benchmarking multimodal llms with generative comprehension. <i>arXiv preprint arXiv:2307.16125</i> .   | 705 |
| 656 |   | 706 |
| 657 |   | 707 |
| 658 |   | 708 |
| 659 | Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong Feng, Lingpeng Kong, and Qi Liu. 2024b. Multimodal arxiv: A dataset for improving scientific comprehension of large vision-language models. <i>arXiv preprint arXiv:2403.00231</i> .   | 709 |
| 660 |   | 710 |
| 661 |   | 711 |
| 662 |   | 712 |
| 663 |   | 713 |
|     |   | 714 |
|     |   | 715 |
|     |   | 716 |
|     |   | 717 |
|     |   | 718 |
|     | Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023b. Evaluating object hallucination in large vision-language models. <i>arXiv preprint arXiv:2305.10355</i> .   |     |
|     |   |     |
|     | Shilong Liu, Hao Cheng, Haotian Liu, Hao Zhang, Feng Li, Tianhe Ren, Xueyan Zou, Jianwei Yang, Hang Su, Jun Zhu, and 1 others. 2023. Llava-plus: Learning to use tools for creating multimodal agents. <i>arXiv preprint arXiv:2311.05437</i> .   |     |
|     |   |     |
|     | Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, and 1 others. 2024a. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In <i>European Conference on Computer Vision</i> , pages 38–55. Springer.  |     |
|     |   |     |
|     | Ziyu Liu, Zeyi Sun, Yuhang Zang, Xiaoyi Dong, Yuhang Cao, Haodong Duan, Dahua Lin, and Jiaqi Wang. 2025. Visual-rft: Visual reinforcement fine-tuning. <i>arXiv preprint arXiv:2503.01785</i> .   |     |
|     |   |     |
|     | Zuyan Liu, Yuhao Dong, Yongming Rao, Jie Zhou, and Jiwen Lu. 2024b. Chain-of-spot: Interactive reasoning improves large vision-language models. <i>arXiv preprint arXiv:2403.12966</i> .  |     |
|     |   |     |
|     | Dimitrios Mallis, Ahmet Serda Karadeniz, Sebastian Cavada, Danila Rukhovich, Niki Foteinopoulou, Kseniya Cherenkova, Anis Kacem, and Djamila Aouada. 2025. Cad-assistant: tool-augmented vllms as generic cad task solvers. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 7284–7294.                      |     |
|     |   |     |
|     | Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L Yuille, and Kevin Murphy. 2016. Generation and comprehension of unambiguous object descriptions. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 11–20.  |     |
|     |   |     |
|     | Fanqing Meng, Lingxiao Du, Zongkai Liu, Zhixiang Zhou, Quanfeng Lu, Daocheng Fu, Tiancheng Han, Botian Shi, Wenhai Wang, Junjun He, and 1 others. 2025. Mm-eureka: Exploring the frontiers of multimodal reasoning with rule-based reinforcement learning. <i>arXiv preprint arXiv:2503.07365</i> .   |     |
|     |   |     |
|     | OpenAI. 2023. Gpt-4v(ision) system card. <a href="https://api.semanticscholar.org/CorpusID:263218031">https://api.semanticscholar.org/CorpusID:263218031</a> .  |     |
|     |   |     |
|     | OpenAI. 2024. Hello gpt-4o. <i>OpenAI Blog</i> .  |     |
|     |   |     |
|     | Yingzhe Peng, Gongrui Zhang, Miaosen Zhang, Zhiyuan You, Jie Liu, Qipeng Zhu, Kai Yang, Xingzhong Xu, Xin Geng, and Xu Yang. 2025. Lmm-rl: Empowering 3b llms with strong reasoning abilities through two-stage rule-based rl. <i>arXiv preprint arXiv:2503.07536</i> .   |     |
|     |   |     |
|     | Ji Qi, Ming Ding, Weihang Wang, Yushi Bai, Qingsong Lv, Wenyi Hong, Bin Xu, Lei Hou, Juanzi Li, Yuxiao Dong, and 1 others. 2024. Cogcom:  |     |

|     |   |  |   |
|-----|---|--|---|
| 719 | Train large vision-language models diving into details through chain of manipulations. <i>arXiv preprint arXiv:2402.04236</i> .   |  |   |
| 720 |   |  |   |
| 721 |   |  |   |
| 722 | Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, and 1 others. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. <i>arXiv preprint arXiv:2403.05530</i> . |  |   |
| 723 |   |  |   |
| 724 |   |  |   |
| 725 |   |  |   |
| 726 |   |  |   |
| 727 |   |  |   |
| 728 |   |  |   |
| 729 | Hao Shao, Shengju Qian, Han Xiao, Guanglu Song, Zhuofan Zong, Letian Wang, Yu Liu, and Hongsheng Li. 2024a. Visual cot: Advancing multimodal language models with a comprehensive dataset and benchmark for chain-of-thought reasoning. <i>Advances in Neural Information Processing Systems</i> , 37:8612–8642.                  |  |   |
| 730 |   |  |   |
| 731 |   |  |   |
| 732 |   |  |   |
| 733 |   |  |   |
| 734 |   |  |   |
| 735 |   |  |   |
| 736 | Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, and 1 others. 2024b. Deepseek-math: Pushing the limits of mathematical reasoning in open language models. <i>arXiv preprint arXiv:2402.03300</i> .   |  |   |
| 737 |   |  |   |
| 738 |   |  |   |
| 739 |   |  |   |
| 740 |   |  |   |
| 741 |   |  |   |
| 742 | Zhaochen Su, Peng Xia, Hangyu Guo, Zhenhua Liu, Yan Ma, Xiaoye Qu, Jiaqi Liu, Yanshu Li, Kaide Zeng, Zhengyuan Yang, and 1 others. 2025. Thinking with images for multimodal reasoning: Foundations, methods, and future frontiers. <i>arXiv preprint arXiv:2506.23918</i> .  |  |   |
| 743 |   |  |   |
| 744 |   |  |   |
| 745 |   |  |   |
| 746 |   |  |   |
| 747 |   |  |   |
| 748 | Guangyan Sun, Mingyu Jin, Zhenting Wang, Chenglong Wang, Siqi Ma, Qifan Wang, Tong Geng, Ying Nian Wu, Yongfeng Zhang, and Dongfang Liu. 2024a. Visual agents as fast and slow thinkers. <i>arXiv preprint arXiv:2408.08862</i> .   |  |   |
| 749 |   |  |   |
| 750 |   |  |   |
| 751 |   |  |   |
| 752 |   |  |   |
| 753 | Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiyang Yu, Yueze Wang, Yongming Rao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. 2024b. Generative multimodal models are in-context learners. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 14398–14409.                 |  |   |
| 754 |   |  |   |
| 755 |   |  |   |
| 756 |   |  |   |
| 757 |   |  |   |
| 758 |   |  |   |
| 759 |   |  |   |
| 760 | Dídac Surís, Sachit Menon, and Carl Vondrick. 2023. Vipergpt: Visual inference via python execution for reasoning. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pages 11888–11898.   |  |   |
| 761 |   |  |   |
| 762 |   |  |   |
| 763 |   |  |   |
| 764 |   |  |   |
| 765 | Huajie Tan, Yuheng Ji, Xiaoshuai Hao, Minglan Lin, Pengwei Wang, Zhongyuan Wang, and Shanghang Zhang. 2025. Reason-rft: Reinforcement fine-tuning for visual reasoning. <i>arXiv preprint arXiv:2503.20752</i> .  |  |   |
| 766 |   |  |   |
| 767 |   |  |   |
| 768 |   |  |   |
| 769 |   |  |   |
| 770 | Chameleon Team. 2024. Chameleon: Mixed-modal early-fusion foundation models. <i>arXiv preprint arXiv:2405.09818</i> .   |  |   |
| 771 |   |  |   |
| 772 |   |  |   |
| 773 | Omkar Thawakar, Dinura Dissanayake, Ketan More, Ritesh Thawkar, Ahmed Heakl, Noor Ahsan, Yuhao  |  |   |
| 774 |   |  |   |
|     |   | Li, Mohammed Zumri, Jean Lahoud, Rao Muhammad Anwer, and 1 others. 2025. Llamav-o1: Rethinking step-by-step visual reasoning in llms. <i>arXiv preprint arXiv:2501.06186</i> .   | 775<br>776<br>777<br>778                      |
|     |   | Haozhe Wang, Alex Su, Weiming Ren, Fangzhen Lin, and Wenhui Chen. 2025a. Pixel reasoner: Incentivizing pixel-space reasoning with curiosity-driven reinforcement learning. <i>arXiv preprint arXiv:2505.15966</i> .  | 779<br>780<br>781<br>782<br>783               |
|     |   | Ke Wang, Junting Pan, Linda Wei, Aojun Zhou, Weikang Shi, Zimu Lu, Han Xiao, Yunqiao Yang, Houxing Ren, Mingjie Zhan, and 1 others. 2025b. Mathcoder-vl: Bridging vision and code for enhanced multimodal mathematical reasoning. <i>arXiv preprint arXiv:2505.10557</i> .   | 784<br>785<br>786<br>787<br>788<br>789        |
|     |   | Lei Wang, Yi Hu, Jiabang He, Xing Xu, Ning Liu, Hui Liu, and Heng Tao Shen. 2024a. T-sciq: Teaching multimodal chain-of-thought reasoning via large language model signals for science question answering. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 19162–19170.                           | 790<br>791<br>792<br>793<br>794<br>795        |
|     |   | Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, and 1 others. 2024b. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. <i>arXiv preprint arXiv:2409.12191</i> .   | 796<br>797<br>798<br>799<br>800<br>801        |
|     |   | Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Song XiXuan, and 1 others. 2024c. Cogvlm: Visual expert for pretrained language models. <i>Advances in Neural Information Processing Systems</i> , 37:121475–121499.  | 802<br>803<br>804<br>805<br>806<br>807        |
|     |   | Wenbin Wang, Liang Ding, Minyan Zeng, Xiabin Zhou, Li Shen, Yong Luo, Wei Yu, and Dacheng Tao. 2025c. Divide, conquer and combine: A training-free framework for high-resolution image perception in multimodal large language models. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 7907–7915. | 808<br>809<br>810<br>811<br>812<br>813<br>814 |
|     |   | Xiyao Wang, Zhengyuan Yang, Chao Feng, Hongjin Lu, Linjie Li, Chung-Ching Lin, Kevin Lin, Furong Huang, and Lijuan Wang. 2025d. Sota with less: Mcts-guided sample selection for data-efficient visual reasoning self-improvement. <i>arXiv preprint arXiv:2504.07934</i> .  | 815<br>816<br>817<br>818<br>819<br>820        |
|     |   | Yaoting Wang, Shengqiong Wu, Yuecheng Zhang, William Wang, Ziwei Liu, Jiebo Luo, and Hao Fei. 2025e. Multimodal chain-of-thought reasoning: A comprehensive survey. <i>arXiv preprint arXiv:2503.12605</i> .   | 821<br>822<br>823<br>824<br>825               |
|     |   | Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.  | 826<br>827<br>828<br>829<br>830<br>831        |

|     |  |  |   |
|-----|--|--|---|
| 832 | Penghao Wu and Saining Xie. 2024. V?: Guided visual search as a core mechanism in multimodal llms. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 13084–13094.  | models for integrated capabilities. <i>arXiv preprint arXiv:2308.02490</i> .   | 887<br>888                                    |
| 837 | Wenshan Wu, Shaoguang Mao, Yadong Zhang, Yan Xia, Li Dong, Lei Cui, and Furu Wei. 2024. Mind’s eye of llms: visualization-of-thought elicits spatial reasoning in large language models. <i>Advances in Neural Information Processing Systems</i> , 37:90277–90317.                                | Yi-Fan Zhang, Huanyu Zhang, Haochen Tian, Chaoyou Fu, Shuangqing Zhang, Junfei Wu, Feng Li, Kun Wang, Qingsong Wen, Zhang Zhang, and 1 others. 2024. Mme-realworld: Could your multimodal llm challenge high-resolution real-world scenarios that are difficult for humans? <i>arXiv preprint arXiv:2408.13257</i> . | 889<br>890<br>891<br>892<br>893<br>894<br>895 |
| 843 | xAI. 2024. <a href="#">Grok-1.5v</a> . Accessed: 2024-05-16.   | Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. 2023. Multimodal chain-of-thought reasoning in language models. <i>arXiv preprint arXiv:2302.00923</i> .  | 896<br>897<br>898<br>899                      |
| 844 | Fengli Xu, Qian Yue Hao, Zefang Zong, Jingwei Wang, Yunke Zhang, Jingyi Wang, Xiaochong Lan, Jiahui Gong, Tianjian Ouyang, Fanjin Meng, and 1 others. 2025a. Towards large reasoning models: A survey of reinforced reasoning with large language models. <i>arXiv preprint arXiv:2501.09686</i> . | Ziwei Zheng, Michael Yang, Jack Hong, Chenxiao Zhao, Guohai Xu, Le Yang, Chao Shen, and Xing Yu. 2025. Deepeyes: Incentivizing "thinking with images" via reinforcement learning. <i>arXiv preprint arXiv:2505.14362</i> .   | 900<br>901<br>902<br>903<br>904               |
| 850 | Guowei Xu, Peng Jin, Li Hao, Yibing Song, Lichao Sun, and Li Yuan. 2024a. Llava-o1: Let vision language models reason step-by-step. <i>arXiv preprint arXiv:2411.10440</i> .   | Hanzhang Zhou, Zijian Feng, Zixiao Zhu, Junlang Qian, and Kezhi Mao. 2024. Unibias: Unveiling and mitigating llm bias through internal attention and fn manipulation. <i>arXiv preprint arXiv:2405.20612</i> .   | 905<br>906<br>907<br>908                      |
| 854 | Wenda Xu, Guanglei Zhu, Xuandong Zhao, Liangming Pan, Lei Li, and William Yang Wang. 2024b. Pride and prejudice: Llm amplifies self-bias in self-refinement. <i>arXiv preprint arXiv:2402.11436</i> .  |  |   |
| 858 | Yi Xu, Chengzu Li, Han Zhou, Xingchen Wan, Caiqi Zhang, Anna Korhonen, and Ivan Vulić. 2025b. Visual planning: Let’s think only with images. <i>arXiv preprint arXiv:2505.11409</i> .  |  |   |
| 862 | Yi Yang, Xiaoxuan He, Hongkun Pan, Xiyan Jiang, Yan Deng, Xingtao Yang, Haoyu Lu, Dacheng Yin, Fengyun Rao, Minfeng Zhu, and 1 others. 2025. R1-onevision: Advancing generalized multimodal reasoning through cross-modal formalization. <i>arXiv preprint arXiv:2503.10615</i> .                  |  |   |
| 868 | Huanjin Yao, Jiaying Huang, Wenhao Wu, Jingyi Zhang, Yibo Wang, Shunyu Liu, Yingjie Wang, Yuxin Song, Haocheng Feng, Li Shen, and 1 others. 2024. Mulberry: Empowering mllm with o1-like reasoning and reflection via collective monte carlo tree search. <i>arXiv preprint arXiv:2412.18319</i> . |  |   |
| 874 | Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. <i>Advances in neural information processing systems</i> , 36:11809–11822.                                       |  |   |
| 879 | En Yu, Kangheng Lin, Liang Zhao, Jisheng Yin, Yana Wei, Yuang Peng, Haoran Wei, Jianjian Sun, Chunrui Han, Zheng Ge, and 1 others. 2025. Perception-r1: Pioneering perception policy with reinforcement learning. <i>arXiv preprint arXiv:2504.07954</i> .   |  |   |
| 884 | Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2023. Mm-vet: Evaluating large multimodal  |  |   |

## A Appendix

This appendix provides additional experimental analyses and qualitative visualizations to complement the main paper:

- Hyper-parameter Analysis.
- Prompts for Ground-R1.
- Illustrations of the intra-bin reward design.
- Visualizations of completion length.
- Qualitative comparisons.

**Hyper-parameter analysis.** Table 4 presents an ablation study on the key hyper-parameters of Ground-R1, including the number of grounding rollouts  $G_1$ , answer rollouts  $G_2$ , and the number of scale bins  $K$ . We observe that a balanced allocation between grounding and answering stages is crucial for optimal performance. In particular, the configuration  $G_1=4$  and  $G_2=2$  achieves the best overall results across all high-resolution benchmarks. Moreover, varying the number of scale bins shows that  $K=3$  consistently outperforms smaller or larger values, suggesting that moderate scale discretization provides a favorable balance between granularity and training stability.

**Prompts for Ground-R1.** The prompts for Ground-R1 is illustrated in Table 5.

**Illustrations of the intra-bin reward design.** The intra-bin reward in SRPO integrates two complementary reward mechanisms:

- *The format reward* ensures coherence by structuring reasoning traces within `<think>` and `</think>` tags, spatial coordinates within `<box>` and `</box>` tags, and final answers within `<answer>` and `</answer>` tags.
- *The answer reward* is defined according to the question type: for multiple-choice questions, it is a binary score based on exact answer matching, while for free-form questions, it is computed as lexical alignment using the average of ROUGE-1, ROUGE-2, and ROUGE-L scores against the ground-truth answers.

**Visualizations of completion length.** Figure 6 depicts the training dynamics of completion length for Ground-R1 trained with SRPO and standard GRPO. As shown, both methods exhibit comparable completion length trajectories throughout training, with similar magnitudes and fluctuations across steps. This observation indicates that the performance improvements achieved by SRPO are not attributable to longer or more verbose generations. Instead, SRPO enhances reasoning quality and evidence

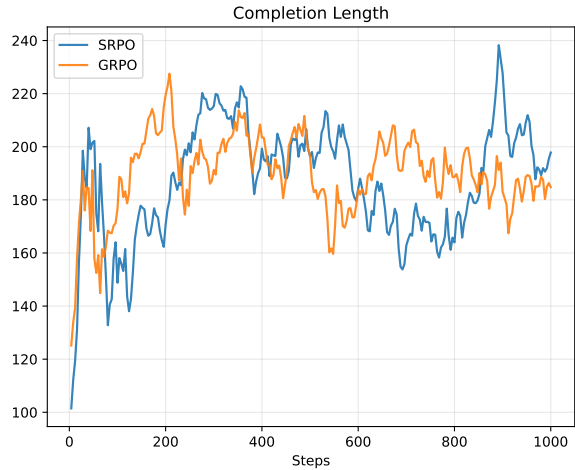


Figure 6: Training dynamics of completion length for Ground-R1 under SRPO and standard GRPO.

Table 4: Ablations of hyper-parameters.

| $G_1$    | $G_2$    | $K$      | $V^*$       | HR-4K       | HR-8K       |
|----------|----------|----------|-------------|-------------|-------------|
| 1        | 8        | 3        | 77.0        | 67.0        | 59.1        |
| 2        | 4        | 3        | 83.8        | 73.8        | 69.5        |
| <b>4</b> | <b>2</b> | <b>3</b> | <b>87.4</b> | <b>75.0</b> | <b>71.1</b> |
| 8        | 1        | 3        | 86.4        | 74.7        | 70.6        |
| 4        | 2        | 2        | 86.9        | 73.5        | 70.0        |
| 4        | 2        | 4        | 86.9        | 74.1        | 70.8        |

grounding without increasing output length or introducing additional inference overhead.

**Qualitative comparisons.** Figures 7–9 present qualitative comparisons of reasoning trajectories produced by Ground-R1 trained with standard GRPO and our proposed SRPO. Figure 7 focuses on a fine-grained recognition scenario that requires identifying a small numerical marker on a bridge. Under standard GRPO, the model grounds its reasoning on a large, visually salient region and incorrectly associates an unrelated numeric sign with the bridge. In contrast, SRPO enables more precise localization of the relevant bridge region, allowing the model to correctly identify the target number. Figure 8 illustrates a foreground object identification task involving multiple people in the scene. While GRPO grounds on a prominent background figure and produces an incorrect answer, SRPO successfully focuses on the closest person to the camera and correctly recognizes the object being held. Figure 9 further examines a counting task involving small-scale objects, where GRPO overgeneralizes from a large region and undercounts the solar panels, whereas SRPO isolates the small but relevant rooftop region and accurately counts

Table 5: **Prompts for Ground-R1.** `{input}` will be replaced with the specific question and image during training and inference.

---

Question: `{input}`

Please think about this question as if you were a human pondering deeply. Provide detailed reasoning between the `<think>` and `</think>` tags. Reason whether there exists a region in the image that could help you answer the question better. If such a region exists, provide one bounding box coordinate in the format `[x1,y1,x2,y2]` inside the `<box>` and `</box>` tags.

Then, you will receive a cropped image based on the bounding box. Use all the images to continue reasoning. You may conduct multiple rounds of grounding to refine your region as you want. The bounding box you provide should always be selected based on the original image. If at any point you determine no further visual information is needed, you may directly provide the final answer inside the `<answer>` and `</answer>` tags.

Format Example:

`<think>` Reasoning Process `</think>` `<box>` `[x1,y1,x2,y2]` `</box>` or,  
`<think>` Reasoning Process `</think>``<answer>` Final Answer `</answer>`

---


984 all instances.

985 Overall, these examples consistently demon-  
986 strate that SRPO mitigates the tendency of standard  
987 GRPO to favor large, visually salient regions and  
988 instead encourages grounding on small yet seman-  
989 tically critical evidence. As a result, Ground-R1  
990 trained with SRPO exhibits more faithful evidence  
991 localization and more reliable reasoning outcomes  
992 across diverse visual reasoning scenarios.



Figure 7: Comparisons of reasoning trajectories of Ground-R1 under standard GRPO and our proposed SRPO.

**Question:**  
What is the person closest to the camera holding in their hand?



**Ground-R1 w/ GRPO:**  
 <think> The scene contains multiple people. One person is standing in the background near the couch, holding an object in their hand. This person is clearly holding something, I focus on this individual to determine what object is being held. </think>  
 <box> [0.45, 0.12, 0.62, 0.48] </box>  
 <think> Within this area, the standing person is holding a bottle-shaped object with a narrow neck, which resembles a beverage container. </think>  
 <answer> A beer bottle. </answer>

**Ground-R1 w/ SRPO:**  
 <think> The question asks about the person closest to the camera. In this image, the closest person appears in the foreground, seated and leaning forward. To answer correctly, I need to focus on the hands of this foreground person, as proximity to the camera and interaction with the object are the key cues. </think>  
 <box> [0.42, 0.56, 0.56, 0.79] </box>  
 <think> In this region, the object in the foreground person's hand is clearly visible. It has a rounded head and a handle, and it is being held close to the mouth, which strongly indicates it is a microphone used for singing or speaking. </think>  
 <answer> A microphone. </answer>

Figure 8: Comparisons of reasoning trajectories of Ground-R1 under standard GRPO and our proposed SRPO.

**Question:** How many solar panels are there?



**Ground-R1 w/ GRPO:**

*<think>* The question asks for the number of solar panels in the image. Since solar installations are often placed on rooftops or visible structures, I will focus on the lower portion of the image where rooftops and mounted objects are present. This area contains several man-made elements that could potentially be solar panels. *</think>*

*<box>* [0.10, 0.65, 1.00, 1.00] *</box>*

*<think>* Within this region, one rectangular panel is installed on a rooftop. And one round panel is visible on the lower rooftops. I conclude that there are two different solar panels present in the image *</think>*

*<answer>* 2 *</answer>*

**Ground-R1 w/ SRPO:**

*<think>* To answer how many solar panels are present, I need to identify flat, rectangular modules with a grid-like appearance that are mounted together. In the lower-right area there is a small cluster of rectangular panels installed on a rooftop. These panels are visually distinct from walls and decorative elements because of their uniform shape. I should focus on that specific rooftop area to count the panels precisely. *</think>*

*<box>* [0.67, 0.68, 0.83, 0.79] *</box>*

*<think>* Within this cropped region, the solar panels are clearly visible and arranged in a neat grid. By counting each distinct rectangular panel, I can identify a total of eight individual solar panels. *</think>*

*<answer>* 8 *</answer>*

Figure 9: Comparisons of reasoning trajectories of Ground-R1 under standard GRPO and our proposed SRPO.