

GeometryZero: Advancing Geometry Solving via Group Contrastive Policy Optimization

Yikun Wang^{1,2} Yibin Wang^{1,2} Dianyi Wang^{1,2} Zimian Peng^{2,3}
Qipeng Guo^{2,4} Dacheng Tao⁵ Jiaqi Wang^{2,†}

¹Fudan University ²Shanghai Innovation Institute ³Zhejiang University
⁴Shanghai AI Laboratory ⁵Nanyang Technological University

Abstract

Recent progress in large language models (LLMs) has boosted mathematical reasoning, yet geometry remains challenging where auxiliary construction is often essential. Prior methods either underperform or depend on very large models (e.g., GPT-4o), making them costly. We argue that reinforcement learning with verifiable rewards (e.g., GRPO) can train smaller models to couple auxiliary construction with solid geometric reasoning. However, naively applying GRPO yields unconditional rewards, encouraging indiscriminate and sometimes harmful constructions. We propose Group Contrastive Policy Optimization (GCPO), an RL framework with two components: (1) *Group Contrastive Masking*, which assigns positive/negative construction rewards based on contextual utility, and (2) a *Length Reward* that encourages longer reasoning chains. On top of GCPO, we build GeometryZero, an affordable family of geometry reasoning models that selectively use auxiliary construction. Experiments on Geometry3K and MathVista show GeometryZero consistently outperforms RL baselines (e.g., GRPO, ToRL). The code has been available at <https://github.com/ekonwang/GeometryZero>.

1 Introduction

Recent advancements in large language models (LLMs) have demonstrated remarkable performance across domains (Ouyang et al., 2022; Team, 2024; DeepSeek-AI et al., 2025) including mathematics (Lu et al., 2022; Yue et al., 2024; Masry et al., 2022; Lu et al., 2023). Among them geometry problem solving is deemed as a challenging task, which requires both perception of visual contexts (i.e., geometric diagrams) and complex reasoning (Lu et al., 2021; Kazemi et al., 2023). Existing

training methods either utilize massive annotated data for supervised learning (Gao et al., 2023) or focus on algebraic-level formal deviation (Brehmer et al., 2023; Trinh et al., 2024). This makes current models show unsatisfying performance on this domain and lack self-correction capabilities in their reasoning chains due to their reliance on annotation quality (Lu et al., 2024).

Another sequence of works focuses on auxiliary lines, which are valuable either when diagrams are inherently complex or when the problem’s intrinsic properties benefit from such constructions, significantly reducing problem solving difficulty (Chervonyi et al., 2025). Several works including (Hu et al., 2024; Wang et al., 2025c) have attempted to enhance visual language models’ (Bi et al., 2025, 2024) utilization of contexts through modifying formal languages (e.g., code) for auxiliary construction, thereby improving their reasoning capabilities on complex geometry problems. Existing works validate that transforming visual contexts into formal languages and leveraging LLM yields better reasoning (Yang et al., 2025). AlphaGeometry2 (Chervonyi et al., 2025) also employs LLMs for auxiliary construction. However, these approaches rely on prompting or training colossal models (e.g., Gemini (Comanici et al., 2025), GPT-4o (OpenAI, 2024)), which incur expensive computational costs that limit their real-world deployment.

After the success of Deepseek-R1-Zero (DeepSeek-AI et al., 2025), GRPO has emerged as a generalizable and effective paradigm for both reasoning tasks and tool learning (Peng et al., 2025; Liu et al., 2025; Meng et al., 2025; Li et al., 2025; Wang et al., 2025b). This makes it particularly suitable for training moderate-sized models capable of auxiliary construction while achieving a strong geometric reasoning performance. However, directly applying the GRPO framework to geometric reasoning with auxiliary construction presents challenges: in

† Corresponding Author

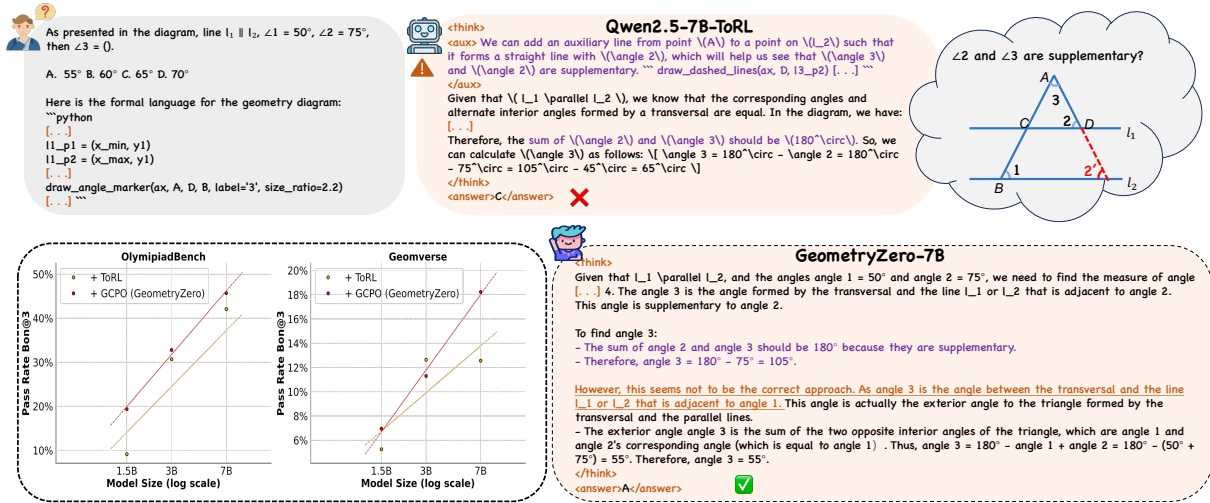


Figure 1: A Comparative Study between ToRL and our GCPO. (a) Two cases comparing GeometryZero-7B with Qwen2.5-7B-GRPO, revealing GeometryZero-7B judiciously determines to directly reason, while a ToRL-trained model indiscriminately conducts auxiliary construction. (b) Purple texts emphasize the erroneous reasoning steps both models undergo. The orange underlined texts amid reasoning process illustrate the critical reflection steps, which we identify as the model’s “aha moments” (DeepSeek-AI et al., 2025) in geometric problem solving, from which GeometryZero-7B benefits in geometric problem solving scenarios. (c) GeometryZero-7B showcases superior overall performance and better scaling effect across different model sizes compared to ToRL.

certain cases, forced or redundant auxiliary constructions prove unnecessary and potentially detrimental (Liu et al., 2023; Yoran et al., 2024; Du et al., 2025). Specifically, some problems can be solved through direct reasoning without auxiliary lines, where their forced inclusion may actually lead to incorrect solutions (Wang et al., 2025a; Shi et al., 2023). Current RL approaches for tool use typically rely on *unconditional rewards* (consistently positive signals across all examples) to encourage indiscriminate tool invocation (Li et al., 2025; Zhang et al., 2025; Hao et al., 2025). This approach lacks the flexibility to autonomously determine when auxiliary constructions are appropriate, thereby limiting RL’s effectiveness in geometric problem solving.

We posit that a flexible mechanism is needed, allowing models to learn through RL when to use auxiliary construction and when to abstain. To this end, we propose *Group Contrastive Policy Optimization*¹ (GCPO), a novel reinforcement learning approach that avoids the drawbacks of unconditional rewards. Specifically, GCPO differs crucially from traditional GRPO: it quantitatively estimates the benefits of auxiliary construction through two contrastive groups of rollouts, then provides flexi-

¹The term “group contrastive” refers to comparing outcomes between two groups of rollouts: one with auxiliary construction and one without, enabling the model to learn when construction is beneficial versus detrimental.

ble signals (*conditional reward*) including encouragements or penalties through Group Contrastive Masking. This mechanism enables GCPO to flexibly encourage auxiliary construction in clearly beneficial scenarios while punishing it in clearly detrimental situations. Inspired by LCPO (Aggarwal and Welleck, 2025), our work introduces a length reward to encourage multidimensional and more in-depth reasoning.

Building upon GCPO, we develop GeometryZero models, a series of lightweight (from 1.5B to 7B) LLMs specialized for geometric reasoning. Extensive experiments demonstrate that GeometryZero outperforms GRPO-trained models across multiple geometry problem-solving benchmarks, like Geometry3K (Lu et al., 2021) and MathVista (Lu et al., 2024). As shown in Figure 1, by judiciously selecting scenarios for auxiliary construction rather than applying it indiscriminately, our GeometryZero showcases remarkable geometric reasoning and reflection ability, while achieving superior overall performance and better scaling across different model sizes compared to RL method with unconditional reward methods like ToRL (Li et al., 2025).

In summary, our contributions can be summarized as follows:

- We validate that through auxiliary construction during their reasoning process, LLMs can

better solve complex tasks in geometric problem solving scenarios, where they utilize both contextual and altered formal languages for auxiliary construction.

- A novel reinforcement learning method called GCPO is proposed in our work, which flexibly provides either encouraging or punishing signals for auxiliary construction across different samples during reinforcement learning, avoiding models from indiscriminately applying auxiliary constructions while maintaining their benefits when strategically justified.
- We train GeometryZero, a family of lightweight geometric reasoning models that selectively use auxiliary constructions. We evaluate them against strong baselines, perform ablations on GCPO, and provide detailed analyses to uncover key insights.

2 Related Work

2.1 Geometry Problem Solving

With the development of large language models (LLMs), researchers have begun to apply LLMs to geometric problem solving (Trinh et al., 2024). However, some early work such as (Brehmer et al., 2023) primarily focuses on algebraic-level formal derivation, which has limited effectiveness in solving practical problems with numeric solutions. Other studies address the lack of geometry problem-solving data by proposing targeted benchmarks and datasets (Lu et al., 2021; Kazemi et al., 2023; Lu et al., 2024). Some recent work employs large-scale annotated data to perform supervised fine-tuning of models, aiming to enhance the performance of multimodal LLMs (Bi et al., 2024) on geometric problems (Gao et al., 2023).

Several approaches, such as GeoCoder (Sharma et al., 2024), attempt to utilize formal languages as context (e.g., code) to assist models in geometric reasoning. Other work explores the use of symbolic tools to strengthen models' geometric reasoning capabilities (Ning et al., 2025). Recent studies such as (Hu et al., 2024; Wang et al., 2025c; Chervonyi et al., 2025) propose encouraging models to construct auxiliary lines by modifying formal languages, thereby better leveraging the intrinsic properties of geometric context to solve the problems.

2.2 Reinforcement Learning with Verifiable Reward

Reinforcement learning has long been a significant focus in the LLM research community (Schulman et al., 2017b; Ouyang et al., 2022; Rafailov et al., 2024). Following the emergence of Deepseek-R1 (DeepSeek-AI et al., 2025), the research community has begun to focus on the application of reinforcement learning with verifiable reward, particularly GRPO (Shao et al., 2024), across various AI domains. Some studies attempt to reproduce GRPO's effectiveness in incentivizing reasoning capabilities on smaller LLMs (Peng et al., 2025). Others apply RLVR methods to multimodal LLMs to enhance their understanding of visual contexts (Meng et al., 2025; Liu et al., 2025). Additional work explores converting visual contexts into formal languages and utilizing reasoning LLMs for inference, aiming to surpass the capabilities of multimodal LLMs (Yang et al., 2025).

The GRPO algorithm was initially proposed in (Shao et al., 2024) and applied to mathematical reasoning. Compared to PPO (Schulman et al., 2017b,a), it simplifies the reinforcement learning pipeline and eliminates the need for a critic model. CPPO (Lin et al., 2025) attempts to optimize the efficiency of the GRPO algorithm through pruning, reducing training costs while maintaining accuracy. DAPO (Yu et al., 2025) introduces a clipping mechanism and dynamic sampling to improve training diversity and stability. (Liu et al., 2025) adapts GRPO's verifiable reward to visual perception tasks, enhancing model performance in visual reasoning. ToRL (Li et al., 2025) attempts to integrate tool-use rewards into GRPO, enhancing the model's tool invocation capability to improve its performance on mathematical reasoning. A separate work proposes a temporal reward coupled with accuracy reward to improve model grounding performance in video contexts (Feng et al., 2025).

3 Preliminary

Group Relative Policy Optimization (GRPO) is a novel algorithm that leverages objectively verifiable supervision signals to enhance model performance on tasks requiring strong reasoning, such as mathematical and code-related problems. Compared with previous approaches, e.g., Reinforcement Learning from Human Feedback, which rely on trained reward models, GRPO utilizes direct verification functions to provide reliable reward

feedback. This method simplifies the reward learning mechanism while enabling efficient alignment with the task’s intrinsics.

Specifically, given a question \mathbf{q} , the policy model π_θ generates a set of N sampled outputs $\mathbf{O} = \{\mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_N\}$, where each output \mathbf{o}_i receives a reward signal \mathbf{r}_i through predefined verifiable reward functions. GRPO then optimizes the following clipped objective:

$$\begin{aligned} \max_{\pi_\theta} \mathbb{E}_{\mathbf{o} \sim \pi_\theta(\mathbf{q})} & \left[\frac{1}{N} \sum_{i=1}^N \min \left(\frac{\pi_\theta(\mathbf{o}_i | \mathbf{q})}{\pi_{\theta_{\text{old}}}(\mathbf{o}_i | \mathbf{q})} \mathbf{A}_i, \right. \right. \\ & \left. \left. \text{clip} \left(\frac{\pi_\theta(\mathbf{o}_i | \mathbf{q})}{\pi_{\theta_{\text{old}}}(\mathbf{o}_i | \mathbf{q})}, 1 - \epsilon, 1 + \epsilon \right) \mathbf{A}_i \right) \right. \\ & \left. - \beta \mathbb{D}_{\text{KL}} [\pi_\theta(\mathbf{o} | \mathbf{q}) \| \pi_{\text{ref}}(\mathbf{o} | \mathbf{q})] \right]. \end{aligned} \quad (1)$$

Here, π_{old} denotes the policy before the current update, and π_{ref} is a fixed reference policy (e.g., the initial model). \mathbf{A}_i is the advantage estimate for output \mathbf{o}_i based on its reward signal \mathbf{r}_i , ϵ is the clipping threshold, and β is a hyperparameter controlling KL regularization to prevent excessive policy deviation.

Existing works, such as the DeepSeek R1-Zero (DeepSeek-AI et al., 2025) algorithm, abandon reliance on supervised fine-tuning and instead train entirely via reinforcement learning, particularly within the Group Relative Policy Optimization (GRPO) framework. In contrast to traditional reinforcement learning methods like PPO (Schulman et al., 2017b), GRPO does not require a critic model to evaluate the policy’s outputs. Given a question q , GRPO first generates G distinct responses $\mathbf{O} = \{\mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_G\}$ using the current policy $\pi_{\theta_{\text{old}}}$. Then, the reward function is applied to obtain a set of verifiable rewards $\{\mathbb{R}(\mathbf{o}_i)\}$. By computing the mean and standard deviation of these rewards, GRPO normalizes them and estimates the advantage value for each response \mathbf{o}_i as follows:

$$\mathbf{A}_i = \frac{\mathbb{R}(\mathbf{o}_i) - \mathbb{E}_{\mathbf{o}_i \sim \mathbf{O}}[\mathbb{R}(\mathbf{o}_i)]}{\text{std}(\{\mathbb{R}(\mathbf{o}_i)\})}, \quad (2)$$

where, \mathbf{A}_i is the advantage value corresponding to the i -th response, representing its relative quality, \mathbb{R} is the sum of verifiable rewards. The GRPO framework encourages the model to generate responses with higher verifiable rewards, thereby improving both reliability and correctness in reasoning-intensive tasks.

4 Method

4.1 Overview

Group Contrastive Policy Optimization (GCPO) introduces one key novel modification: it incorporates a crucial mechanism called group contrastive masking, which provides a positive mask for auxiliary reward in scenarios where auxiliary construction is beneficial, while applying a negative mask (i.e., penalty) in others. To achieve this objective, we propose the Group Contrastive Masking. GCPO also introduces an additional length reward optimized for longer completion, due to the requirement for auxiliary reasoning.

4.2 Reinforcement Learning for Auxiliary Construction

To enable models to incorporate auxiliary construction reasoning—a form of tool utilization (i.e. attempt to construct auxiliary lines in thinking process with formal language like tikz code)—we introduce an auxiliary reward that teaches the “*how-to*” capability, as in ToRL (Li et al., 2025), where an additional tool related reward is introduced for using coding for mathematical reasoning. During training, the textual context contains either TikZ code or logic form as detailed in Appendix A.1, which strictly corresponds to geometric diagrams, and the model is prompted to autonomously decide whether to include auxiliary line construction in its reasoning process. For executable TikZ code, the auxiliary reward is positive if the altered tikz code in thinking process can execute successfully and render a diagram; for logic forms, we detect the presence of special tokens `<aux>` and `</aux>` indicating attempts to modify the logic form for auxiliary lines construction. Thus, the auxiliary reward for a given response \mathbf{o}_i is defined as follow:

$$\mathbf{R}_{aux}(\mathbf{o}_i) = \begin{cases} 1 & \text{if model constructs auxiliary lines,} \\ 0 & \text{otherwise.} \end{cases} \quad (3)$$

4.3 Conditional Reward for Auxiliary Construction

Inspired by existing works (Chervonyi et al., 2025; Hu et al., 2024), we aim to endow models with the capability of auxiliary construction with formal language in geometric problem solving. However, it is crucial to note that while teaching models “*how to*” is important (Li et al., 2025), we must also teach them “*when to do it*”. Although GRPO has

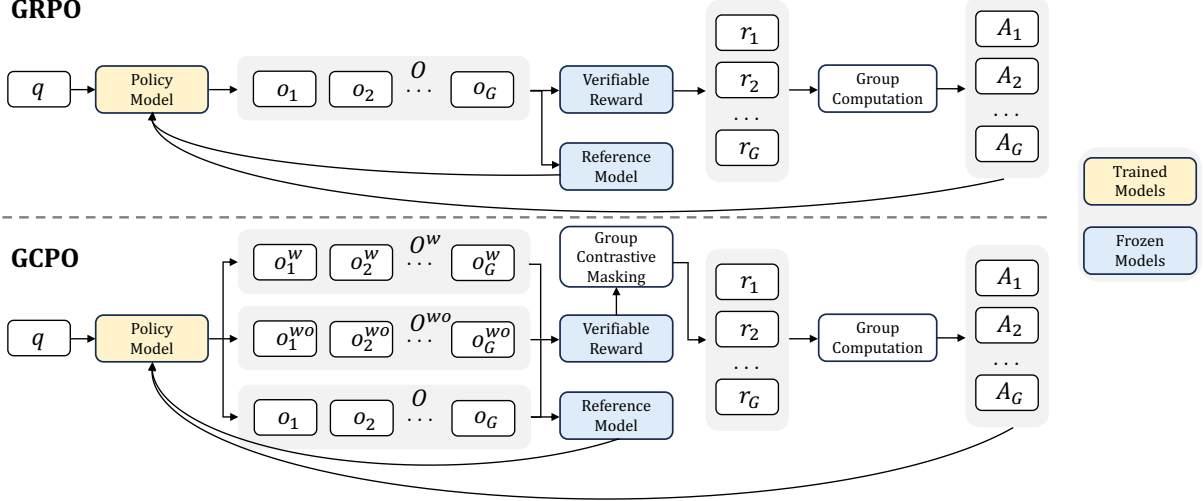


Figure 2: **The Illustration of GCPO.** One key difference between our GCPO and GRPO is Group Contrastive Masking: (1) GCPO samples two additional rollout groups O^w and O^{wo} for evaluating the quantitative benefits via *accuracy reward*. (2) The auxiliary reward signals of GCPO are dynamically masked to positive, negative, and zero during training as (Eq. 4). Another difference is that a novel length reward mechanism is incorporated.

proven effective in enhancing reasoning capabilities, it lacks conditional rewards for tool usage and relies solely on unconditional rewards to encourage desired behaviors, which may cause indiscriminate use of the tool in certain scenarios. To address this limitation, we propose Group Contrastive Policy Optimization (**GCPO**), which introduces a group contrastive reward mechanism for a conditional reward signal that flexibly provides encouragements or penalties during training, allowing the flexibility of the trained models of whether to apply the tool (i.e. auxiliary construction) or not.

The key insight of GCPO stems from the observation that in geometric problem solving, while auxiliary lines can enhance reasoning in many cases, they may be unnecessary or even detrimental in some cases as well. Unconditional encouragement of auxiliary construction could lead to sub-optimal performance for complex scenarios like geometry problems, thus, we need to incorporate a conditional reward during reinforcement learning.

4.4 Components of Group Contrastive Policy Optimization

Group Contrastive Masking The core idea of GCPO is that during training, models should become aware of the quantifiable benefits of using code to draw auxiliary lines - employing this capability when beneficial and avoiding it when counterproductive. As shown in Figure 2, during response sampling, the model generates G rollouts $\mathbf{O} = \{\mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_G\}$, along with another two

contrastive rollout groups: one requiring auxiliary line thought process group $\mathbf{O}^w = \{\mathbf{o}_i^w\}$ and one group prohibiting it $\mathbf{O}^{wo} = \{\mathbf{o}_i^{wo}\}$. The group contrastive masking function is defined as:

$$\text{Mask}(\mathbf{R}_{aux}(\mathbf{O})) = \begin{cases} \mathbf{R}_{aux}(\mathbf{O}) & \text{if } E(\mathbf{R}_{acc}(\mathbf{O}^w)) > E(\mathbf{R}_{acc}(\mathbf{O}^{wo})) + \epsilon, \\ -\mathbf{R}_{aux}(\mathbf{O}) & \text{if } E(\mathbf{R}_{acc}(\mathbf{O}^w)) + \epsilon < E(\mathbf{R}_{acc}(\mathbf{O}^{wo})), \\ 0 & \text{otherwise,} \end{cases} \quad (4)$$

where $\mathbf{R}_{acc}(\mathbf{o}_i)$ represents the accuracy reward function indicating whether response \mathbf{o}_i contains the correct final answer, $\mathbf{R}_{aux}(\mathbf{o}_i)$ refers to the auxiliary reward in (Eq. 3) and ϵ (set to 0.05 in experiments) is a threshold hyperparameter controlling reward masking. Following standard mathematical conventions, the function \mathbf{R}_{aux} can naturally extend from a single response to a set of responses: $\mathbf{R}_{aux}(\mathbf{O}) = \{\mathbf{R}_{aux}(\mathbf{o}_1), \dots, \mathbf{R}_{aux}(\mathbf{o}_G)\}$, which also holds for $\mathbf{R}_{acc}(\mathbf{O}^w)$ and $\mathbf{R}_{acc}(\mathbf{O}^{wo})$.

Length Reward Auxiliary construction thinking requires deeper, multi-dimensional analysis, necessitating longer reasoning processes. Inspired by Length Controlled Policy Optimization (LCPO) (Aggarwal and Welleck, 2025), we adapt by introducing a simplified length reward, where $\text{len}(\mathbf{o}_i)$ counts tokens in response \mathbf{o}_i and l_{\max} is the maximum allowed completion length.

$$\mathbf{R}_{length}(\mathbf{o}_i) = \min\left\{1, \frac{\text{len}(\mathbf{o}_i)}{l_{\max}}\right\} \quad (5)$$

Verifiable Reward Combination As a variant of RLVR, the verifiable reward $\mathbb{R}(\mathbf{o}_i)$ of GCPO

Table 1: **The main empirical results.** The BoN@3 accuracy rate across in-domain benchmarks including Geomverse, Geometry3k and out-of-domain results on MathVista and OlympiadBench, where the best results are **bold** and the second best are underlined. Results from our GeometryZero (w.r.t., + GCPO) models are in gray.

Method	Geomverse	Geometry3k	MathVista	OlympiadBench	Avg.
<i>1.5B models</i>					
Qwen2.5-1.5B-Instruct	4.20	41.76	47.70	13.44	26.78
+ SFT	4.80	44.25	43.11	<u>14.51</u>	26.67
+ GRPO	<u>5.76</u>	53.35	<u>57.79</u>	<u>14.51</u>	<u>32.85</u>
+ ToRL	5.26	<u>57.01</u>	<u>57.79</u>	11.29	32.84
GeometryZero-1.5B (ours)	6.95	60.24	60.65	16.47	36.08
<i>3B models</i>					
Qwen2.5-3B-Instruct	10.53	65.83	67.88	32.25	44.12
+ SFT	10.20	71.65	73.08	30.64	46.39
+ GRPO	<u>12.13</u>	<u>75.87</u>	82.87	31.72	<u>50.65</u>
+ ToRL	11.63	76.31	80.34	<u>32.87</u>	50.29
GeometryZero-3B (ours)	11.30	79.25	<u>82.56</u>	35.48	52.15
<i>7B models</i>					
G-Llava-7B	6.23	49.31	46.92	27.82	32.57
GNS-Llava-1.5-7B	5.21	62.00	51.40	33.54	38.04
Qwen2.5-7B-Instruct	14.76	70.99	68.19	39.24	48.30
+ SFT	15.36	75.98	76.14	41.93	52.35
+ GRPO	<u>16.93</u>	79.03	<u>86.23</u>	40.32	<u>55.63</u>
+ ToRL	12.56	78.75	83.48	<u>44.08</u>	54.72
GeometryZero-7B (ours)	18.23	<u>78.81</u>	87.15	45.69	57.47

combines multiple weighted components as below, where $\mathbf{R}_{\text{GRPO}}(\mathbf{o}_i)$ includes a accuracy reward and a format reward ensuring proper output structure, while hyperparameter λ representing auxiliary reward weight and hyperparameter β representing length reward weight are both set to 0.5.

$$\mathbb{R}(\mathbf{o}_i) = \mathbf{R}_{\text{GRPO}}(\mathbf{o}_i) + \lambda \cdot \text{Mask}(\mathbf{R}_{\text{aux}}(\mathbf{o}_i)) + \beta \cdot \mathbf{R}_{\text{length}}(\mathbf{o}_i) \quad (6)$$

In essence, GCPO uses masked auxiliary rewards to teach appropriate tool usage contexts, while length rewards ensure sufficient reasoning capacity. The framework largely inherits GRPO’s verifiable reward framework, employing outcome-based reinforcement learning for model training.

5 Experiments

5.1 Experimental Setting

For dataset, we apply a synthesized dataset with data sampled from two popular geometry problem solving (GPS) dataset including Geomverse (Kazemi et al., 2023) and Geometry3k (Lu et al., 2021), for training SFT and RL models. The training dataset recipe is detailed in Appendix A.1. For training details, we utilize the vLLM inference

framework (Kwon et al., 2023) during training and evaluation. More details are presented in Appendix A.2.

As for models, inspired by (Yang et al., 2025), we turn the visual diagrams into formal language contexts for better reasoning performance. Thus, we use Qwen2.5 (Qwen et al., 2025) series language models for training. For benchmarks, besides the in-domain benchmarks of Geometry3k and Geomverse, the OOD geometry benchmarks comprise MathVista (Lu et al., 2024) and OlympiadBench (He et al., 2024). The details are provided in Appendix A.3.

Baselines To fully demonstrate the effectiveness of GCPO, we compare against the following baselines: **(1)** SFT. The model undergoes supervised fine-tuning using prompt-response pairs, where responses are either human-annotated or distilled from capable LLMs. **(2)** GRPO (Shao et al., 2024). A reinforcement learning via verifiable outcome algorithm where the model generates multiple responses for baseline advantage estimation to encourage reasoning and improve problem-solving, which eliminates the usage of a critic model or reward model. **(3)** ToRL (Li et al., 2025). An RL

Table 2: **Results on Qwen2.5-VL-7B-Instruct.** GCPO generalizes to vision-language models that directly process visual geometric diagrams, where the best results are **bold** and the second best are underlined. Results from our method are in `gray`.

Method	Geomverse	Geometry3k	MathVista	OlympiadBench	Avg.
Qwen2.5-VL-7B-Instruct	11.34	65.38	63.29	30.83	42.71
+ SFT	12.78	69.54	71.95	32.36	46.66
+ GRPO	<u>13.58</u>	<u>73.46</u>	<u>81.23</u>	31.82	<u>50.02</u>
+ ToRL	9.35	72.39	78.46	<u>34.59</u>	48.70
GeometryZero-VL-7B (ours)	15.84	74.68	83.64	37.19	52.84

algorithm building on GRPO that appends an additional reward function (Eq. 3) to unconditionally encourage tool usage (i.e., auxiliary construction).

5.2 Results

As shown in the table 1, our experimental results across four geometry problem-solving benchmarks demonstrate GCPO’s effectiveness in enhancing model capabilities on geometric problems. Key findings are summarized below:

SFT Memorizes while RL Generalizes. We observe that SFT models (Qwen2.5-1.5B-SFT and Qwen2.5-3B-SFT) show consistent improvements over original Instruct models on in-domain benchmarks like Geomverse and Geometry3k. For instance, Qwen2.5-1.5B-SFT and Qwen2.5-3B-SFT gains an improvement of 2.49% and 5.83% on Geometry3k. However, these SFT models exhibit either performance drops or smaller gains compared to RL methods on OOD benchmarks like MathVista and OlympiadBench. For instance, while Qwen2.5-1.5B-SFT exhibits a performance decline of 4.59% compared to the base model on the OOD benchmark MathVista, Qwen2.5-1.5B-GRPO demonstrates a notable improvement of 10.09%. Overall, RL approaches including GRPO, ToRL, and GCPO achieve more consistent improvements across both in-domain and OOD benchmarks, surpassing SFT and proving the effectiveness of reinforcement learning.

Group Contrastive Policy Optimization Works. Compared to GRPO, ToRL models unconditionally encourage auxiliary construction during reasoning process across all examples with an unconditional reward design (Eq. 3). The empirical results demonstrate that ToRL models has no clear advantage over GRPO across various model scales, indicating that this coarse-grained policy fails to provide significant benefits for auxiliary construction

in geometric problem-solving scenarios. For instance, while ToRL demonstrates a marginal 0.64% advantage over GRPO on 3B models, it exhibits a 0.91% performance reduction on 7B models. In contrast, GCPO improves model performance on both in-domain and OOD benchmarks, achieving consistently better average performance on most benchmarks across model sizes. This indicates discerning when to incorporate auxiliary reasoning ultimately improving problem-solving capabilities. As shown in Figure 3, GCPO enhances geometric problem-solving by generating auxiliary constructions during the reasoning process, we also provide more case studies in Appendix D.

5.3 Generalization to Vision-Language Models

To validate that GCPO generalizes beyond text-only models, we apply it to the vision-language model Qwen2.5-VL-7B-Instruct (Bai et al., 2025), which directly processes visual diagrams. In this setting, auxiliary constructions generated by the model are rendered via TikZ/Python code and appended as updated diagrams to the visual input, allowing the VL model to directly *see* the auxiliary lines. As shown in Table 2, GeometryZero-VL-7B consistently outperforms all baselines across benchmarks, demonstrating that GCPO effectively transfers to the vision-language setting.

Interestingly, despite having access to visually rendered auxiliary constructions, GeometryZero-VL-7B (52.84 avg.) still underperforms the text-only GeometryZero-7B (57.47 avg.) by a notable margin. This suggests that when geometric diagrams are fully described in formal language, reasoning in pure text space holds an advantage over multimodal reasoning — likely because formal language provides a more precise and unambiguous representation that LLMs can manipulate more effectively during chain-of-thought reasoning.

Table 3: **The ablation study with Qwen2.5-3B-Instruct.** The components includes auxiliary reward (AR), group contrastive (GC) masking and length reward (LR). We report mean and standard deviation across 3 runs.

Method	AR	GC	LR	Geomverse	Geometry3k	MathVista	OlympiadBench	Avg.
GRPO	✗	✗	✗	11.85 (± 0.52)	75.12 (± 0.32)	82.53 (± 1.56)	31.75 (± 1.45)	50.31 (± 0.37)
GCPO (/wo AR)	✗	✗	✓	12.45 (± 0.56)	75.36 (± 0.43)	81.54 (± 2.04)	33.87 (± 1.50)	50.81 (± 0.62)
GCPO (/wo LR, /wo GC)	✓	✗	✗	12.58 (± 0.90)	76.37 (± 0.58)	80.74 (± 1.32)	33.77 (± 1.23)	50.61 (± 0.54)
GCPO (/wo LR)	✓	✓	✗	12.80 (± 1.10)	77.61 (± 0.45)	82.11 (± 1.68)	31.94 (± 1.28)	51.12 (± 0.43)
GeometryZero	✓	✓	✓	12.66 (± 1.58)	79.36 (± 0.39)	82.78 (± 1.77)	35.68 (± 0.97)	52.62 (± 0.53)

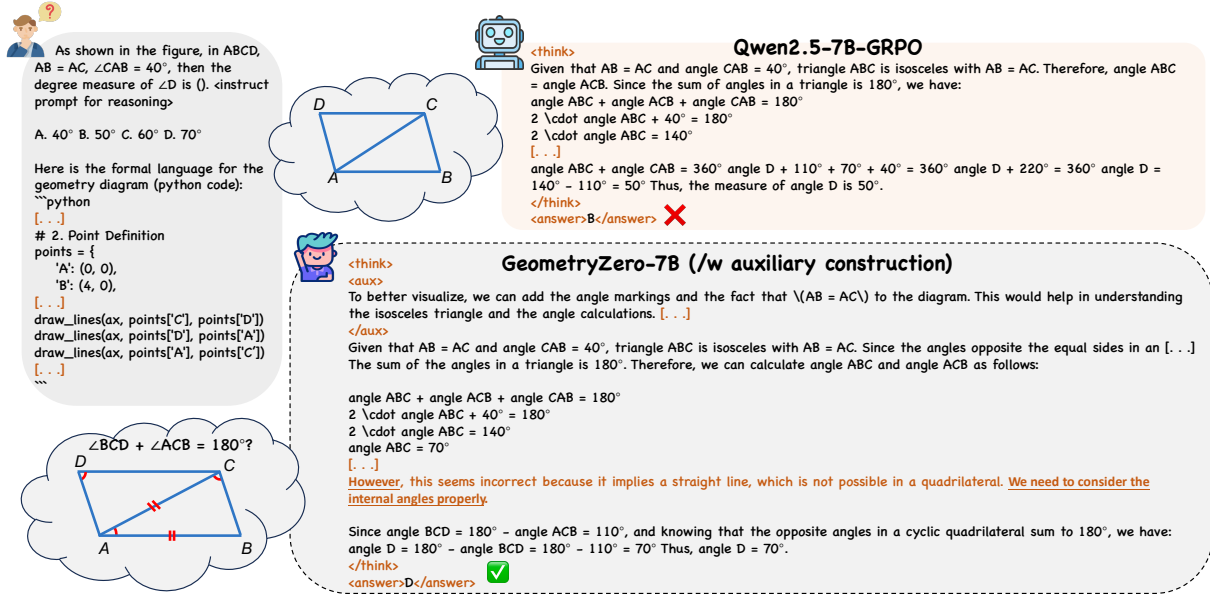


Figure 3: **A Case Study between GRPO and GCPO.** Two reponses comparing Qwen2.5-7B-GRPO with our GeometryZero-7B for a MathVista problem, revealing how GeometryZero-7B effectively constructs auxiliary elements during its reasoning process. The orange underlined texts during reasoning process are reflection process in geometric problem solving.

5.4 Ablation Study

To better understand the contributions of components in GCPO, we conduct an ablation study to evaluate three variants of GeometryZero and compare them with the GRPO model and GeometryZero, where the descriptions of the variants are further detailed in Appendix C.1.

Our findings show that GeometryZero (/wo LR) achieves higher performance than GeometryZero (/wo LR, /wo GC). Both GeometryZero (/wo AR) and GeometryZero (/wo LR) demonstrate better average performance across benchmarks compared to Qwen-2.5-3B-GRPO respectively, while these two variants show lower average performance than GeometryZero. We also provide more ablation studies on different-sized models in Appendix C.2.

The experimental results indicate that removing either the auxiliary reward or its corresponding group contrastive masking leads to performance

degradation across benchmarks. Similarly, eliminating the length reward in GCPO also poses negative effects. These results validate the effectiveness of our proposed method.

5.5 Completion Length of Models

Response length serves as a crucial metric for observing training dynamics in RL (Meng et al., 2025). We monitor the variation in response length during the training process of GeometryZero and GRPO models as shown in Figure 4. For 7B models, we observe the following trends in response length: During the initial few steps, the model’s response length increases rapidly, subsequently it decreases, after reaching the lowest point, it then begins to rise again.

The phenomenon aligns with observations in llm-r1 (Peng et al., 2025). We hypothesize that in the first phase, the model is encouraged by format rewards to learn reasoning patterns that gener-

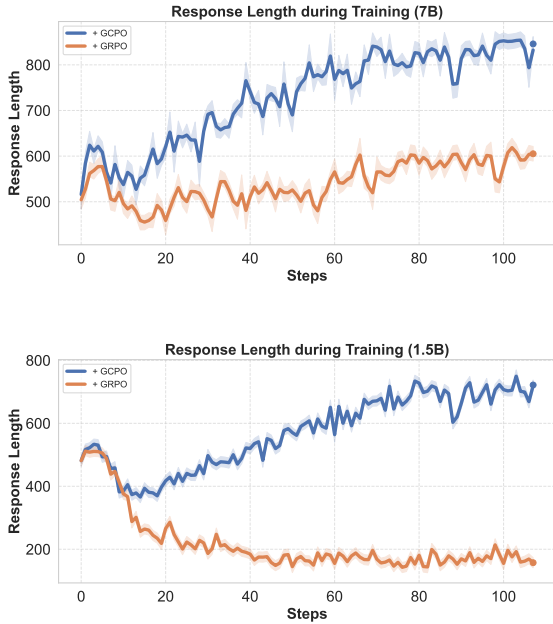


Figure 4: (TOP) The trend of **completion length** during reinforcement learning of 7B models. (BELOW) The trend of **completion length** during reinforcement learning of 1.5B models. We observe that the completion length of GCPO models follows a distinct pattern during training: initially increasing, then decreasing, before rising again, same observation for 7B GRPO models.

ate thoughts before answers, leading to increased output length. In the second phase, as training progresses, the model begins to optimize the reward function, particularly the accuracy reward, causing it to reduce redundant outputs while maintaining the required format, resulting in decreased response length. For the third phase, we speculate that in later training stages, the model learns more sophisticated reasoning patterns and attempts to generate more complex reasoning steps, leading to the length recovery.

The observation differs for 1.5B models. GeometryZero-1.5B exhibits the rise-fall-rise pattern in response length, while the GRPO model shows no recovery in response length in the last stage. We attribute this to the model’s limited capacity due to smaller parameter size, which prevents it from learning more comprehensive and profound reasoning processes through GRPO alone in later training stages.

We also demonstrate more in-depth discussions for epsilon settings in Appendix B and the performance of GeometryZero models on geometry proving tasks in Appendix F.

6 Conclusion

In this paper, we propose **Group Contrastive Policy Optimization**, a novel reinforcement learning framework that incorporates verifiable rewards to optimize conditional reward particularly for auxiliary construction in geometric reasoning. GCPO dynamically adapts to different problem scenarios, supporting an autonomous strategy of tool-assisted and tool-free reasoning. Building upon this framework, we introduce GeometryZero, a series of geometric reasoning models that autonomously learn when and how to apply auxiliary constructions during the reasoning process. Extensive experiments demonstrate the effectiveness of our approach, while detailed analyses provide insights for future research directions.

Acknowledgments

We gratefully acknowledge the Shanghai Innovation Institute for their generous support of this research on GPUs and other computational resources.

Limitations

While GCPO demonstrates strong performance, several limitations warrant discussion. First, our method assumes access to verifiable reward signals, which may not be available for all geometry problem types. Additionally, due to compute constraints, we limited our experiments to moderate model sizes (under 7B parameters). These limitations point to valuable directions for future research in reasoning systems for geometric problems.

Ethics and Potential Risks

While GeometryZero demonstrates improved efficiency and stronger geometric reasoning performance, it is ultimately built on a large language model and thus inherits common reliability limitations of LLMs. In real-world deployments, the model may produce incorrect yet plausible-looking solutions, exhibit brittleness to distribution shifts (e.g., unseen problem styles or noisy contexts). We therefore position GeometryZero as a research system and recommend human verification and domain-specific safeguards before practical uses.

References

Pranjal Aggarwal and Sean Welleck. 2025. **L1: Controlling how long a reasoning model thinks with reinforcement learning**. *Preprint*, arXiv:2503.04697.

- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, and 8 others. 2025. [Qwen2.5-vl technical report](#). *Preprint*, arXiv:2502.13923.
- Jinhe Bi, Yifan Wang, Danqi Yan, Xun Xiao, Artur Hecker, Volker Tresp, and Yunpu Ma. 2025. Prism: Self-pruning intrinsic selection method for training-free multimodal data selection. *arXiv preprint arXiv:2502.12119*.
- Jinhe Bi, Yujun Wang, Haokun Chen, Xun Xiao, Artur Hecker, Volker Tresp, and Yunpu Ma. 2024. Visual instruction tuning with 500x fewer parameters through modality linear representation-steering. *arXiv preprint arXiv:2412.12359*.
- Johann Brehmer, Pim de Haan, Sönke Behrends, and Taco Cohen. 2023. [Geometric algebra transformer](#). *Preprint*, arXiv:2305.18415.
- Jiaqi Chen, Tong Li, Jinghui Qin, Pan Lu, Liang Lin, Chongyu Chen, and Xiaodan Liang. 2022. [Uni-geo: Unifying geometry logical reasoning via reformulating mathematical expression](#). *Preprint*, arXiv:2212.02746.
- Yuri Chervonyi, Trieu H. Trinh, Miroslav Olsák, Xiaomeng Yang, Hoang Nguyen, Marcelo Menegali, Junehyuk Jung, Vikas Verma, Quoc V. Le, and Thang Luong. 2025. [Gold-medalist performance in solving olympiad geometry with alphageometry2](#). *CoRR*, abs/2502.03544.
- Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, Luke Marris, Sam Petulla, Colin Gaffney, Asaf Aharoni, Nathan Lintz, Tiago Cardal Pais, Henrik Jacobsson, Idan Szpektor, Nan-Jiang Jiang, and 3416 others. 2025. [Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities](#). *Preprint*, arXiv:2507.06261.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, and 81 others. 2025. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#). *CoRR*, abs/2501.12948.
- Yufeng Du, Minyang Tian, Srikanth Ronanki, Subendhu Rongali, Sravan Bodapati, Aram Galstyan, Azton Wells, Roy Schwartz, Eliu A Huerta, and Hao Peng. 2025. [Context length alone hurts llm performance despite perfect retrieval](#). *Preprint*, arXiv:2510.05381.
- Kaituo Feng, Kaixiong Gong, Bohao Li, Zonghao Guo, Yibing Wang, Tianshuo Peng, Benyou Wang, and Xiangyu Yue. 2025. [Video-r1: Reinforcing video reasoning in mllms](#). *Preprint*, arXiv:2503.21776.
- Jiahui Gao, Renjie Pi, Jipeng Zhang, Jiacheng Ye, Wanjun Zhong, Yufei Wang, Lanqing Hong, Jianhua Han, Hang Xu, Zhenguo Li, and Lingpeng Kong. 2023. [Gllava: Solving geometric problem with multi-modal large language model](#). *CoRR*, abs/2312.11370.
- Bingguang Hao, Zengzhuang Xu, Maolin Wang, Yuntao Wen, Yicheng Chen, Cunyin Peng, Long Chen, Dong Wang, Xiangyu Zhao, Jinjie Gu, Chenyi Zhuang, and Ji Zhang. 2025. [Reasoning through exploration: A reinforcement learning framework for robust function calling](#). *Preprint*, arXiv:2508.05118.
- Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan Liu, and Maosong Sun. 2024. [Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems](#). *Preprint*, arXiv:2402.14008.
- Yushi Hu, Weijia Shi, Xingyu Fu, Dan Roth, Mari Ostendorf, Luke Zettlemoyer, Noah A. Smith, and Ranjay Krishna. 2024. [Visual sketchpad: Sketching as a visual chain of thought for multimodal language models](#). In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*.
- Mehran Kazemi, Hamidreza Alvari, Ankit Anand, Jialin Wu, Xi Chen, and Radu Soricut. 2023. [Geomverse: A systematic evaluation of large models for geometric reasoning](#). *Preprint*, arXiv:2312.12241.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. [Efficient memory management for large language model serving with pagedattention](#). *Preprint*, arXiv:2309.06180.
- Xuefeng Li, Haoyang Zou, and Pengfei Liu. 2025. [Torl: Scaling tool-integrated rl](#). *Preprint*, arXiv:2503.23383.
- Zhihang Lin, Mingbao Lin, Yuan Xie, and Rongrong Ji. 2025. [CPPO: accelerating the training of group relative policy optimization-based reasoning models](#). *CoRR*, abs/2503.22342.
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. [Lost in the middle: How language models use long contexts](#). *Preprint*, arXiv:2307.03172.
- Ziyu Liu, Zeyi Sun, Yuhang Zang, Xiaoyi Dong, Yuhang Cao, Haodong Duan, Dahua Lin, and Jiaqi Wang. 2025. [Visual-rft: Visual reinforcement fine-tuning](#). *Preprint*, arXiv:2503.01785.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2024. [Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts](#). *Preprint*, arXiv:2310.02255.

- Pan Lu, Ran Gong, Shibiao Jiang, Liang Qiu, Siyuan Huang, Xiaodan Liang, and Song-Chun Zhu. 2021. [Inter-gps: Interpretable geometry problem solving with formal language and symbolic reasoning](#). *Preprint*, arXiv:2105.04165.
- Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Taffjord, Peter Clark, and Ashwin Kalyan. 2022. [Learn to explain: Multimodal reasoning via thought chains for science question answering](#). *Preprint*, arXiv:2209.09513.
- Pan Lu, Liang Qiu, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, Tanmay Rajpurohit, Peter Clark, and Ashwin Kalyan. 2023. [Dynamic prompt learning via policy gradient for semi-structured mathematical reasoning](#). *Preprint*, arXiv:2209.14610.
- Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. [Chartqa: A benchmark for question answering about charts with visual and logical reasoning](#). *Preprint*, arXiv:2203.10244.
- Fanqing Meng, Lingxiao Du, Zongkai Liu, Zhixiang Zhou, Quanfeng Lu, Daocheng Fu, Tiancheng Han, Botian Shi, Wenhai Wang, Junjun He, Kaipeng Zhang, Ping Luo, Yu Qiao, Qiaosheng Zhang, and Wenqi Shao. 2025. [Mm-eureka: Exploring the frontiers of multimodal reasoning with rule-based reinforcement learning](#). *Preprint*, arXiv:2503.07365.
- Maizhen Ning, Zihao Zhou, Qiufeng Wang, Xiaowei Huang, and Kaizhu Huang. 2025. [GNS: solving plane geometry problems by neural-symbolic reasoning with multi-modal llms](#). In *AAAI-25, Sponsored by the Association for the Advancement of Artificial Intelligence, February 25 - March 4, 2025, Philadelphia, PA, USA*, pages 24957–24965. AAAI Press.
- OpenAI. 2024. [Gpt-4o system card](#). *Preprint*, arXiv:2410.21276.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. [Training language models to follow instructions with human feedback](#). *Preprint*, arXiv:2203.02155.
- Yingzhe Peng, Gongrui Zhang, Miaosen Zhang, Zhiyuan You, Jie Liu, Qipeng Zhu, Kai Yang, Xingzhong Xu, Xin Geng, and Xu Yang. 2025. [Lmm-r1: Empowering 3b llms with strong reasoning abilities through two-stage rule-based rl](#). *Preprint*, arXiv:2503.07536.
- Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, and 25 others. 2025. [Qwen2.5 technical report](#). *Preprint*, arXiv:2412.15115.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2024. [Direct preference optimization: Your language model is secretly a reward model](#). *Preprint*, arXiv:2305.18290.
- John Schulman, Sergey Levine, Philipp Moritz, Michael I. Jordan, and Pieter Abbeel. 2017a. [Trust region policy optimization](#). *Preprint*, arXiv:1502.05477.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017b. [Proximal policy optimization algorithms](#). *Preprint*, arXiv:1707.06347.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. [Deepseekmath: Pushing the limits of mathematical reasoning in open language models](#). *CoRR*, abs/2402.03300.
- Aditya Sharma, Aman Dalmia, Mehran Kazemi, Amal Zouaq, and Christopher J. Pal. 2024. [Geocoder: Solving geometry problems by generating modular code through vision-language models](#). *CoRR*, abs/2410.13510.
- Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed Chi, Nathanael Schärli, and Denny Zhou. 2023. [Large language models can be easily distracted by irrelevant context](#). *Preprint*, arXiv:2302.00093.
- Gemini Team. 2024. [Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context](#). *Preprint*, arXiv:2403.05530.
- Trieu Trinh, Yuhuai Wu, Quoc Le, He He, and Thang Luong. 2024. [Solving olympiad geometry without human demonstrations](#). *Nature*.
- Hongru Wang, Cheng Qian, Wanjun Zhong, Xiushi Chen, Jiahao Qiu, Shijue Huang, Bowen Jin, Mengdi Wang, Kam-Fai Wong, and Heng Ji. 2025a. [Acting less is reasoning more! teaching model to act efficiently](#). *Preprint*, arXiv:2504.14870.
- Yikun Wang, Zuyan Liu, Ziyi Wang, Han Hu, Pengfei Liu, and Yongming Rao. 2025b. [Geovista: Web-augmented agentic visual reasoning for geolocalization](#). *Preprint*, arXiv:2511.15705.
- Yikun Wang, Siyin Wang, Qinyuan Cheng, Zhaoye Fei, Liang Ding, Qipeng Guo, Dacheng Tao, and Xipeng Qiu. 2025c. [Visuothink: Empowering lvlm reasoning with multimodal tree search](#). *Preprint*, arXiv:2504.09130.
- Yi Yang, Xiaoxuan He, Hongkun Pan, Xiyan Jiang, Yan Deng, Xingtao Yang, Haoyu Lu, Dacheng Yin, Fengyun Rao, Minfeng Zhu, Bo Zhang, and Wei Chen. 2025. [R1-onevision: Advancing generalized multimodal reasoning through cross-modal formalization](#). *Preprint*, arXiv:2503.10615.

Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan Berant. 2024. [Making retrieval-augmented language models robust to irrelevant context](#). *Preprint*, arXiv:2310.01558.

Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, Hang Zhu, and 16 others. 2025. [Dapo: An open-source llm reinforcement learning system at scale](#). *Preprint*, arXiv:2503.14476.

Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, and 3 others. 2024. [Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi](#). *Preprint*, arXiv:2311.16502.

Shaokun Zhang, Yi Dong, Jieyu Zhang, Jan Kautz, Bryan Catanzaro, Andrew Tao, Qingyun Wu, Zhiding Yu, and Guilin Liu. 2025. [Nemotron-research-tool-n1: Exploring tool-using language models with reinforced reasoning](#). *Preprint*, arXiv:2505.00024.

A Implementation Details

A.1 Training Dataset Construction

To ensure the model adequately learns geometric problem solving, we select two mainstream geometric problem solving (GPS) datasets. Our training data comes from Geometry3k and Geomverse.

- Geometry3k (Lu et al., 2021). We randomly select 1443 training samples from Geometry3k. For the SFT experiments, this dataset lacks supervised sequences, so we use Qwen2.5-72B-Instruct (Qwen et al., 2025) to generate CoT reasoning processes with known answers. These reasoning processes are concatenated with the solutions to form supervised responses. For RL-based methods like GRPO and GCPO, we only utilize the problems in the dataset and employ the final answers as supervision.
- Geomverse (Kazemi et al., 2023). We randomly choose $2k$ training samples from Geomverse. Since this dataset already contains human-annotated CoT processes, we directly use them for SFT experiments. We also only employ the problems in the dataset and the final answers as supervision for RL-based methods.

A.2 Training Details

We set train batch size to 32 and micro train batch size to 1, for response sampling we apply a rollout batch size of 64 and a micro rollout batch size of 2. We set max prompt length to 2048 and max completion length l_{max} to 1024. We use full parameter tuning rather than PEFT methods (Bi et al., 2025).

We set G to 8, with both the SFT learning rate and the GRPO learning rate at $3e-7$ and the format reward weight set to 0.5. Due to the limited training data and absence of significant policy shift concerns, we set the KL coefficient to 0 to achieve better tuning performance. As for compute hardware, we use 4 Nvidia H100 GPUs for training and later evaluation.

A.3 Evaluation Benchmarks

To comprehensively evaluate the model’s performance on geometric problem solving, we conduct evaluations on several mainstream geometric problem benchmarks. Besides using Geometry3k and the Geomverse D2 subset to test the model’s in-domain geometric capabilities, for out-of-distribution problems, we also evaluate the model’s performance on MathVista and OlympiadBench.

Besides the in-domain benchmarks, the OOD geometry benchmarks comprise:

- MathVista (Lu et al., 2024). A consolidated mathematical reasoning benchmark within visual contexts. To evaluate LLMs on geometric problems, GPT-4o converts visual contexts from MathVista testmini into textual Python code using ReACT and Self-Vote mechanisms. We then manually verify that the code-generated graphics match the original visual contexts.
- OlympiadBench (He et al., 2024). The benchmark is an Olympiad-level multimodal scientific benchmark. We extract all geometry problems and filter for those with only one solution to ensure single-solution supervision. Using the same pipeline as MathVista, we convert visual contexts into LLM-comprehensible Python code, obtaining an evaluation set to assess model performance on Olympiad-level geometry problems.

Dataset	Sample Size	Code Type	Code Executable
Geomverse	2k	Tikz Code	✓
Geometry3k	1443	Logic Form	✗

Table 4: **The training dataset construction details.** The training data are sampled from two popular geometry problem solving (GPS) dataset including Geomverse and Geometry3k.

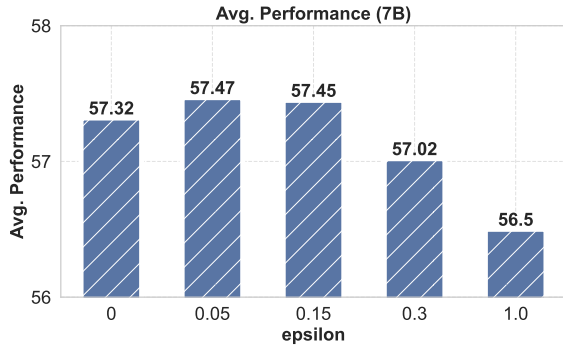


Figure 5: The average performance of GeometryZero with different hyperparameter epsilon settings in GCPO training.

B The Impact of Hyperparameter ϵ of Group Contrastive Masking

To provide more insightful analysis of our method, we conduct a comparative study with different epsilon hyperparameter settings. We set epsilon values at 0, 0.05, 0.15, 0.3, and 1.0 separately for training GeometryZero and evaluating their benchmark performance. As presented in Figure 5, we find that as epsilon increases from 0 to 1.0, the algorithm’s performance first improves slightly and then declines.

We speculate that when epsilon is too low, the algorithm applies positive or negative masks to cases where the benefit of auxiliary construction is uncertain, leading to unstable training in these cases and ultimately affecting model performance. When epsilon is too high, the threshold for group contrastive masking becomes excessively strict, causing auxiliary rewards to be zero in most cases, which effectively renders the auxiliary reward mechanism inoperative. We conclude that GCPO performs best in the epsilon range of 0.05 to 0.15, and thus we keep epsilon at 0.05 in our experiments.

C Ablation Study

C.1 Variant Models in Ablation Study

Here are the model variants used in ablation study, serving as a supplementary material for section 5.4:

- GeometryZero (/wo AR), which excludes the auxiliary construction reward (Eq. 3) and consequently removes the Group Contrastive Masking mechanism (Eq. 4), retaining only the length penalty term (Eq. 5);
- GeometryZero (/wo LR, /wo GC), which only retains the auxiliary reward (Eq. 3) encouraging auxiliary construction thinking during the reasoning phase but excludes the corresponding Group Contrastive Masking (Eq. 4), equivalent to ToRL using unconditional auxiliary reward;
- GeometryZero (/wo LR), which excludes the length reward (Eq. 5) in GCPO that encourages longer reasoning chains, retaining other components of GCPO.

C.2 Ablation study on 3B Model

We present the ablation study of the 7B GeometryZero models, please refer to tab.5.

D Case Study

Please refer to the Figure 6 for the case study, amid which the reasoning process the model outputs executable tikz code to construct auxiliary lines for geometric reasoning.

E Training Dynamics during Reinforcement Learning

E.1 Accuracy Reward

We provide the training dynamics of the accuracy reward of GeometryZero models during RL, please refer to Figure 7 and Figure 8 for inspection of GCPO and GRPO respectively.

E.2 Format Reward

We also present the training dynamics of the *format reward* of GeometryZero models during RL, please refer to Figure 9 and Figure 11 for inspection of GCPO and GRPO respectively.

Table 5: **The ablation study of GCPO components** on 7B models. The components includes auxiliary reward (AR), group contrastive (GC) masking and length reward (LR).

Method	AR	GC	LR	Geomverse	Geometry3k	MathVista	OlympiadBench	Avg.
GRPO	✗	✗	✗	16.93	79.03	86.23	40.32	55.63
GCPO (/wo AR)	✗	✗	✓	17.40	79.53	87.15	41.93	56.50
GCPO (/wo LR, /wo GC)	✓	✗	✗	12.56	78.75	83.48	44.08	54.72
GCPO (/wo LR)	✓	✓	✗	17.47	79.89	84.40	41.76	55.88
GeometryZero	✓	✓	✓	18.23	78.81	87.15	45.69	57.47

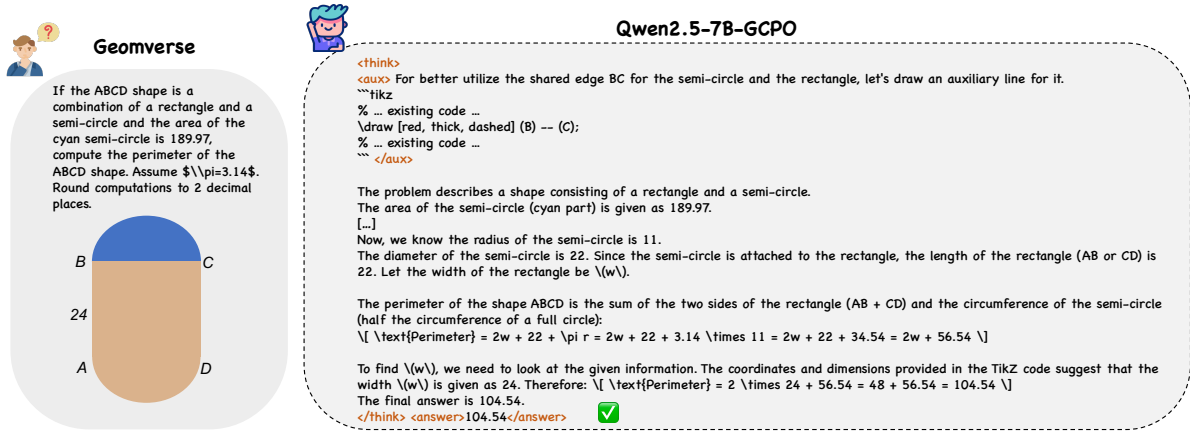


Figure 6: A case study example from Geomverse (Kazemi et al., 2023) of GeometryZero-7B (Qwen2.5-7B-GCPO).

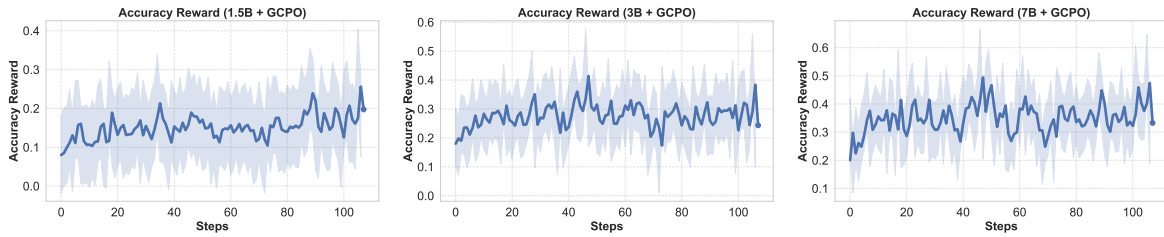


Figure 7: The trend of accuracy reward of **GeometryZero** (GCPO) models during training.

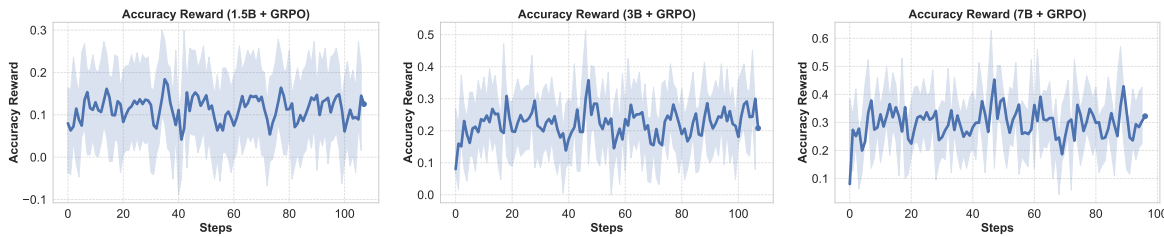


Figure 8: The trend of accuracy reward of GRPO models during training.

E.3 Mask Ratio

According to (Eq. 4), our method’s characteristic is that during Group Masking, it applies positive masks to auxiliary rewards for some cases, negative masks to others, while zero-masking cases where the mean accuracy reward gap does not exceed epsilon. As presented in Figure 10, we observe that while the overall proportions of positive and negative masks fluctuate, they remain generally stable

during training, with positive masks consistently outnumbering negative masks.

This phenomenon demonstrates that the roll-out group with auxiliary construction (i.e. O^w) achieves higher accuracy rewards than the group without auxiliary construction (i.e. O^{w^0}) in reward score computing, indicating that auxiliary construction generally contributes to obtaining correct solutions and thus validating its effectiveness. More

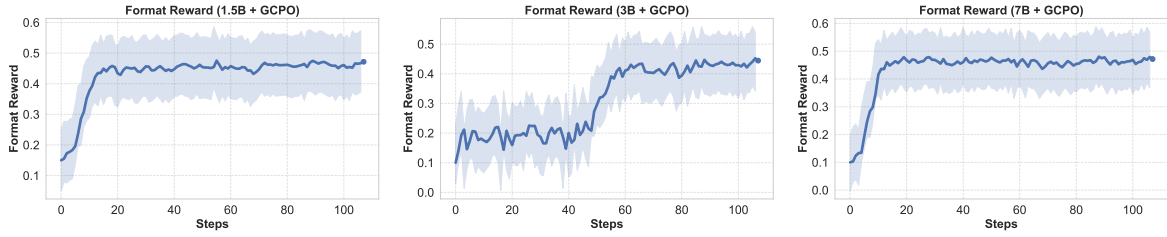


Figure 9: The trend of format reward of **GeometryZero** (GCPO) models during training.

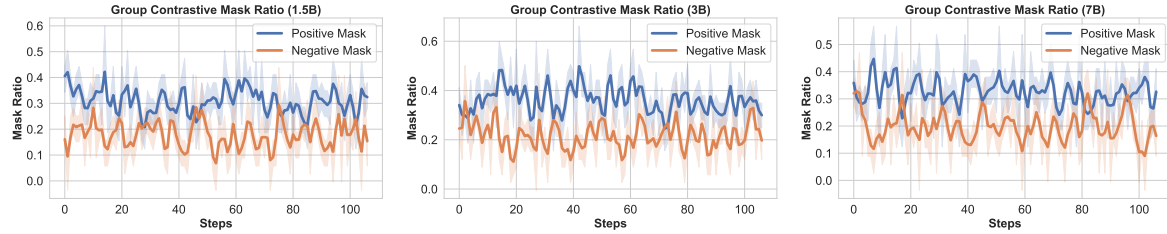


Figure 10: **The record of group mask ratio.** The positive group mask and negative group mask ratio in group contrastive masking for 1.5B, 3B and 7B models.

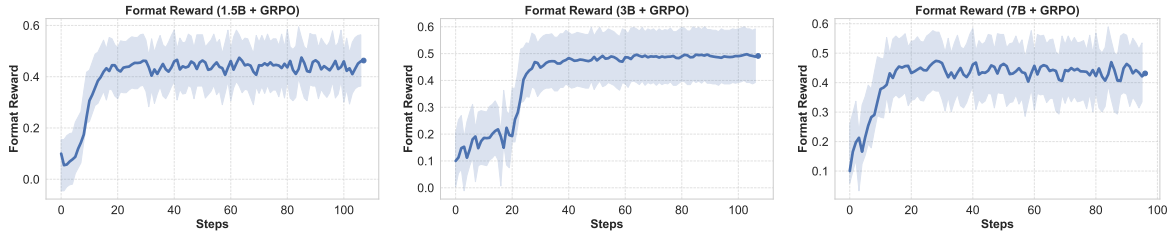


Figure 11: The trend of format reward of GRPO models during training.

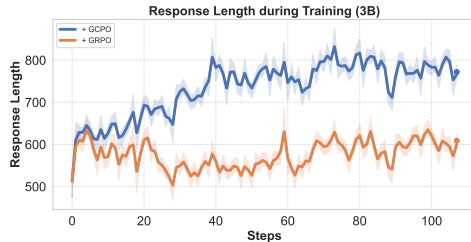


Figure 12: **The trend of response length of GCPO and GRPO during training on 3B models.** For 3B models, We also observe the completion length of follows a distinct pattern during training: initially increasing, then decreasing or stagnating, before rising again.

records of group mask ratio are presented in appendix E.3.

F GeometryZero on Geometric Proving Tasks

In widely used geometry benchmarks, UniGeo (Chen et al., 2022) contains a subset of geometric proof problems. For efficient comparison, we

Table 6: The performance of different models including AG (Trinh et al., 2024) and GeometryZero-14B on UniGeo (geometric proof part).

Model	UniGeo (proof part)
AlphaGeometry	94.4%
GPT-4o	74.1%
Qwen2.5-14B-Instruct	64.8%
GeometryZero-14B	72.2%

selected 108 problems of this subset for our additional experiments.

Since AlphaGeometry (Trinh et al., 2024) requires a strict geometric DSL (formal language describing points, lines, circles, relations), we first used GPT-4o to batch-formalize the 108 UniGeo problems into DSL. The correctness of the proofs was then verified using an automated validation script. For GPT-4o and GeometryZero, we generated complete proof sequences and compared them with golden sequences to measure accuracy on proof problems.

AG’s primary bottleneck lies in formalizing

problems into DSL, which accounts for the imperfection of its accuracy. Actually, the difficulty of UniGeo problems does not necessitate AG’s symbolic search process. GeometryZero-14B and GPT-4o achieve comparable performance, with GeometryZero-14B showing a 7.4% improvement over Qwen2.5-14B-Instruct, despite the absence of proof problems in its training data. This highlights the strong generalization capability of GCPO.

G Completion Length

As a supplement to the 1.5B and 7B length curves discussed in the main text, we report the response length dynamics of GCPO and GRPO on the 3B model in Figure 12. The 3B curves exhibit a pattern consistent with what we observe at other scales: completion length first increases in the early stage of training, then decreases or stagnates mid-training, and finally rises again as the policy refines its reasoning trajectories. Notably, GCPO maintains a longer reasoning chain than GRPO throughout most of training, which aligns with our length reward design and further supports the claim that GCPO encourages the model to produce more elaborate and verifiable reasoning.