

000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 GENERALIZABLE GEOMETRIC IMAGE CAPTION SYN- THESIS

Anonymous authors

Paper under double-blind review

ABSTRACT

Multimodal large language models (MLLMs) have various practical applications that demand strong reasoning abilities. Despite recent advancements, these models still struggle to solve complex geometric problems. A key challenge stems from the lack of high-quality image-text pair datasets for understanding geometric images. Furthermore, most symbolic data synthesis pipelines typically fail to generalize to questions beyond their predefined templates. In this paper, we bridge this gap by introducing a complementary process of Reinforcement Learning with Verifiable Rewards (RLVR) into the data generation pipeline. By adopting accuracy-guided RLVR to refine captions for symbolically synthesized geometric images, our pipeline successfully captures the key features of geometry problem-solving. This enables better task generalization and yields non-trivial improvements. Furthermore, even in out-of-distribution scenarios, the generated dataset GeoReasoning-10K achieves non-trivial performance gains, yielding accuracy improvements of 2.8%–4.8% in non-geometric subtasks of MathVista and MathVerse. This generalization ability is further validated in MMMU, where significant improvements of 2.4%–3.9% in Art & Design and Tech & Engineering tasks are observed.

1 INTRODUCTION

Multimodal Large Language Models have exhibited impressive capabilities across a variety of vision-related tasks, including Visual Question Answering (VQA), visual grounding, and image captioning. Recent MLLMs, such as Qwen2.5-VL, Intern2.5-VL, and LLaVA-Next (Bai et al., 2025; Chen et al., 2024; Liu et al., 2024), have shown superior performance compared to specialized vision models across a wide range of visual tasks. As the field advances, there has been increasing interest in enhancing the reasoning capabilities of MLLMs (Jaech et al., 2024; Shao et al., 2024), which is regarded as a crucial factor in extending the performance boundaries of these models. Among various reasoning tasks, mathematical reasoning (Zhang et al., 2024a) has attributed particular attention due to its structured problem-solving nature, offering a clear pathway for MLLMs to develop and improve their reasoning skills.

Research from MathVerse (Zhang et al., 2024a) indicates that MLLMs perform best when the input is purely textual, while their performance declines significantly with visual-only inputs. This highlights the urgent need for MLLMs to develop strong cross-modal reasoning capabilities, which involves accurately and comprehensively transferring information from the image to the text. Although numerous geometry and math datasets have been introduced (Lu et al., 2023; Zhang et al., 2024b; Wang et al., 2024) to boost various aspects of model performance, high-quality datasets explicitly designed for cross-modal reasoning remain scarce. That is primarily because in existing datasets, the alignment between images and captions is often asymmetrical. For instance, in geometric problems, two lines of equal length can be easily described textually but may not be correspondingly annotated or visually distinct in the image. Such discrepancies hinder the model’s ability to learn robust cross-modal reasoning.

Meanwhile, Reinforcement Learning (RL) has been demonstrated to significantly improve model reasoning and generalization capabilities (Guo et al., 2025). Its reward-driven framework is particularly effective for cross-modal reasoning, allowing models to optimize decision-making through interactive feedback (Deng et al., 2025; Peng et al., 2025; Huang et al., 2025a). Building on these insights, we

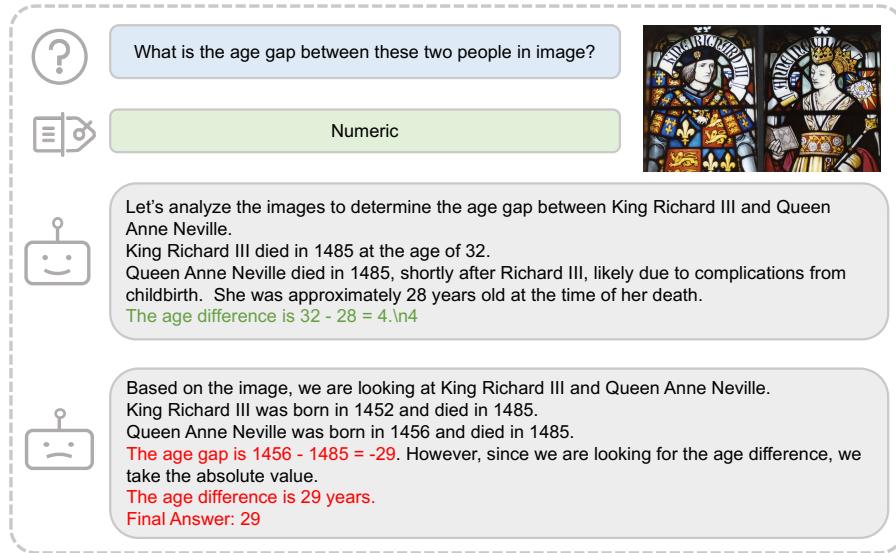


Figure 1: **Examples of generalization.** MLLMs learn from our synthetic geometric mathematical problems and generalize to algebraic cases with even non-geometric input images.

employ the RAFT method (Dong et al., 2023) and design a reward function that incorporates both reasoning and caption rewards. This facilitates the alternating optimization of dataset quality and model reasoning abilities, leading to improved results on complex multimodal tasks.

To bridge the gap between visual and linguistic modalities, we propose an RL-based data refinement engine that iteratively enhances data quality. Utilizing this pipeline, we introduce a novel geometry dataset comprising 10,000 image-caption pairs. To the best of our knowledge, this is the first high-quality dataset in which visual and textual information are fully aligned and generalize well to out-of-distribution tasks, making it a valuable resource for improving cross-modal reasoning. Experimental results demonstrate that our dataset significantly enhances models’ cross-modal reasoning abilities and their performance on geometric image textualization tasks. Furthermore, models trained on our dataset exhibit strong generalization capabilities on other mathematics-focused benchmarks, including MathVerse and MathVista, as shown in Figure 1. In summary, our main contributions are summarized as follows:

- We introduce **GeoReasoning-10K**, a new dataset containing 10,000 carefully constructed image-caption pairs where visual and textual information are fully equivalent. This dataset serves as a high-quality resource for training cross-modal reasoning models, outperforming previous geometric datasets like AutoGeo in both MathVerse and MathVista.
- We propose **Geo-Image-Textualization**, a scalable RL-based framework for generating and refining high-quality synthetic image-caption pairs in geometry. Our iterative RL-driven optimization significantly enhances data alignment and semantic accuracy, and demonstrates generalization to out-of-domain geometric tasks.
- Extensive experiments show that the improvements facilitated by GeoReasoning extend beyond geometric tasks and even generalize well to non-mathematical domains. In particular, GeoReasoning-10K brings an accuracy improvement of **+2.8%–4.8%** in non-geometric subtasks of MathVista and MathVerse, and **+2.4%–3.9%** in Art & Design and Tech & Engineering tasks in MMMU.

2 RELATED WORKS

2.1 DATA GENERATION

Recent studies have highlighted the scarcity of high-quality geometry image–caption datasets, which limits fine-grained cross-modal reasoning in geometric tasks. Following AlphaGeometry (Trinh

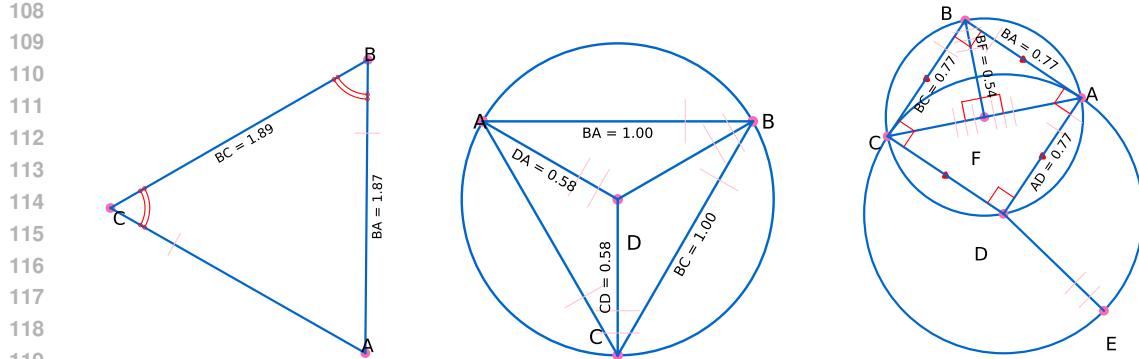


Figure 2: **Symbolically synthesized geometric images.** These geometry problems are symbolically composed from our relation library, corresponding to easy, medium, and hard difficulty levels, respectively, where the pink ticks and red arcs indicate equal-length segments and equal angles. The symbolic engine can generate images with infinite types and difficulties. For visual clarity, this figure has a fixed set of colors, font sizes, and line thicknesses compared to the original images in our constructed dataset. Please refer to the original dataset for precise details.

et al., 2024), Autoge (Huang et al., 2025b) proposed an automatic generation engine to produce image-caption pairs, constructing a 100K dataset named AutoGeo-100k. MATHGLANCE (Sun et al., 2025) introduced GeoPeP, a perception-oriented dataset of 200K structured geometry image-text pairs explicitly annotated with geometric primitives and spatial relationships. MagicGeo (Wang et al., 2025) formulates diagram synthesis as a coordinate optimization problem, ensuring formal geometric correctness via solvers before coordinate-aware text generation.

Despite the advances, existing pipelines still struggle to guarantee full modality alignment, i.e., captions frequently omit visual details, while images lack exhaustively aligned textual descriptions.

2.2 IMAGE CAPTIONING

Image captioning aims to generate comprehensive descriptions of visual content and has been widely studied for natural images. While general-purpose MLLMs such as mPLUG-Owl2 (Ye et al., 2024), MiniGPT-4 (Zhu et al., 2023), and BLIP-3 (Xue et al., 2024) can perform captioning to some extent, their effectiveness is often limited by insufficient fine-grained cross-modal alignment. Image-Textualization (Pi et al., 2024) mitigates this issue by integrating multiple vision experts to produce more detailed and accurate captions.

However, the potential of utilizing image captioning to enhance geometric reasoning capacity remains underexplored. OmniCaptioner (Lu et al., 2025) proposes a unified visual captioning framework that converts diverse images into fine-grained textual descriptions. Nonetheless, its geometric annotations are derived from AutoGeo and MAVIS, largely relying on synthetic or loosely aligned pairs rather than fully equivalent visual-textual representations. Moreover, the scarcity of high-quality geometric image-caption pairs makes it difficult to accurately extract and align geometric information. As a result, current models underperform on geometric image textualization compared to natural image captioning and general visual reasoning benchmarks.

3 METHODS

In this section, we introduce our Geo-Image-Textualization data generation pipeline first, followed by our RAFT method used for data refinement.

3.1 GEO-IMAGE-TEXTUALIZATION DATA GENERATION ENGINE

The proposed data generation pipeline mainly contains three parts: the relation sampling, image-caption pair generation, and question-answer generation procedure, as shown in Figure 3.

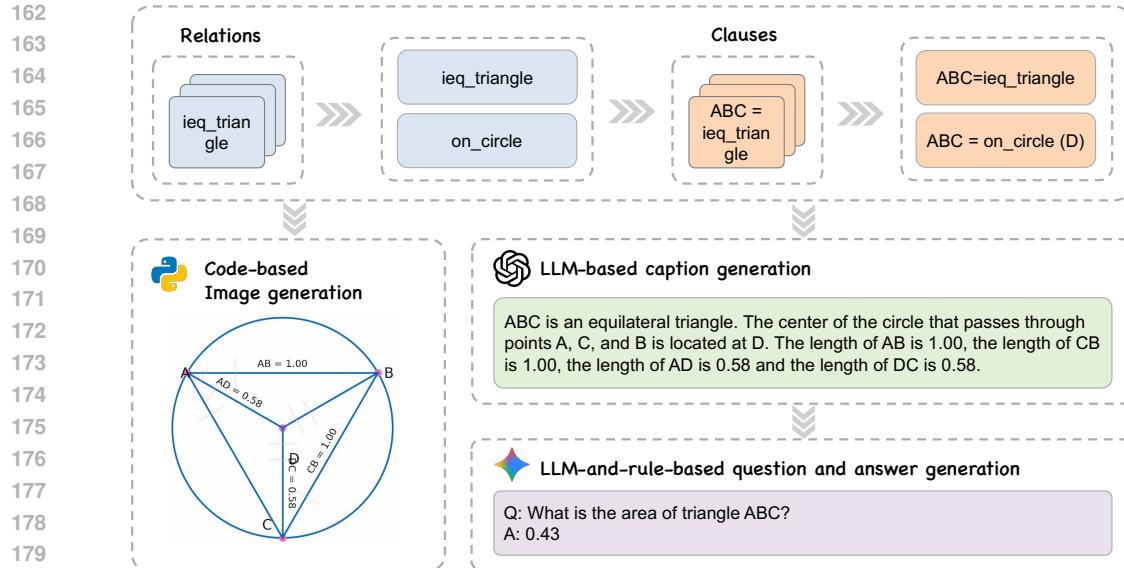


Figure 3: **The geometry data synthesis pipeline**, where a graph-based representation similar to AutoGeo (Huang et al., 2025b) is employed for generating the final geometry images. The relation library comprises over 50 basic geometric relationships that can be composed into complex ones, providing comprehensive coverage for geometric problems of various difficulties. The image-caption pair is utilized for the SFT stage, while the caption-QA pair for the RLVR stage.

3.1.1 RELATION SAMPLING

We develop the Geo-Image-Textualization pipeline upon the data generation procedure in Alpha-Geometry. In our framework, *Relations* act as fundamental construction operations that systematically generate diverse yet semantically coherent geometric premises for subsequent theorem synthesis. Each relation (e.g., `angle_mirror`, `circumcenter`, etc.) encodes a precise geometric procedure—such as reflecting a point across an angle bisector or locating the circumcenter of a triangle. In addition to the construction rule, each definition maintains dependency metadata, specifying which primitive objects (points, lines, circles) and prior constructions it depends on. This enables the symbolic engine to get the minimal set of premises required to derive a given theorem.

After sampling relations, each relation is converted into a clause, with associated point variables. For instance, the relation `angle_mirror` x a b c denotes: given points a , b , and c , construct point x as the reflection of c across angle $\angle ABC$. Finally, the system constructs a graph-based representation in AutoGeo (Huang et al., 2025b) to model geometric problems. Each clause, corresponding to either a geometric construction or a relational assertion, is incorporated into the graph by instantiating nodes for geometric entities (e.g., points, lines, circles) and establishing edges that encode their interdependencies. Before adding each clause, the system verifies the logical correctness of the selected set of predefined geometric rules.

3.1.2 IMAGE-CAPTION PAIR GENERATION

The geometric graph encodes all relevant entities, including points, lines, and circles, enabling the straightforward rendering of basic geometric elements, similar to AutoGeo. However, a fundamental limitation of AutoGeo is that the captions cannot be directly inferred from the rendered images because the visual content and the textual description are not semantically aligned. We argue that this misalignment constitutes a critical bottleneck in cross-modal reasoning.

To address this issue, we introduce a set of visual augmentation strategies that explicitly encode semantic relationships within the image, following most conventions in geometry problems (Dimmell & Herbst, 2015):

1. **Segment Equality Representation:** We use short perpendicular ticks to indicate equal-length line segments. When multiple pairs of equal segments exist, we distinguish them using a different number of ticks (e.g., one tick, two ticks)¹.
2. **Angle Annotations:** For angles that are integer multiples of 15° within the range [15°, 165°], we directly annotate the angle values within the image.
3. **Parallel and Perpendicular Indicators:** Parallel lines are marked using matching directional triangles, and right angles are indicated using a small square symbol at the vertex.
4. **Equal Angle Representation:** Equal angles are denoted by marking them with the same number of arcs, consistent with conventional geometric notation.
5. **Intersection and Collinearity:** Dashed lines are used to explicitly indicate intersections and collinearity relationships among points or segments.

For each clause in the symbolic representation, we apply a refined, rule-based template to convert it into natural language. These captions accurately describe the geometric diagram, including object relationships, special angle values, and other visual properties. Additionally, the captions explicitly state the lengths of specific line segments when such information is visually annotated in the image. By ensuring that all semantic content present in the image is mirrored in the caption, we achieve full cross-modal alignment.

3.2 QUESTION-ANSWER PAIR GENERATION

The most fundamental requirements for generating questions lie in three aspects. First, the generated question should be based on the caption, i.e., should not be irrelevant to the caption. Second, any information already present in the caption should be removed, as this would dilute the impact of the caption and make the evaluation of caption quality harder. Last, the question should be compatible with the given information, so that it can be logically answered based on what is provided.

Based on these requirements, we propose a rule-and-LLM-based pipeline to systematically generate the question and answer based on the pre-generated caption. Specifically, we prompt the large model (Gemini 2.5 Flash) with rubric-based instructions to generate initial questions based on the caption, while also letting the model flag those questions inconsistent with the caption. For the inconsistent questions, we then switch to a different prompt, encouraging the model to incorporate additional information and formulate new questions accordingly. This process continues until a self-consistent question is generated for the first time. The detailed two-stage prompt design is provided in Appendix A and B.

3.3 RLVR FRAMEWORK FOR DATA REFINEMENT

Our proposed RLVR framework iteratively optimizes both the model and the dataset through a novel alternating paradigm. The method consists of two phases: (1) a cold-start supervised fine-tuning phase to bootstrap initial captioning capabilities, and (2) an RLVR phase with RAFT (Dong et al., 2023) that cyclically refines the dataset and model via reinforcement learning. The overall framework is shown in Figure 4.

3.3.1 COLD-START PHASE

To initialize the model’s ability to generate geometrically aligned captions, we first perform supervised fine-tuning (SFT) on the base vision language model using the GeoReasoning-10K dataset. This phase minimizes the standard cross-entropy loss:

$$\mathcal{L}_{\text{SFT}} = -\mathbb{E}_{(I, c^*) \sim \mathcal{D}_0} [\log P_{\theta_0}(c|I)] \quad (1)$$

where I denotes an input geometric image, c^* and c indicate the ground-truth caption and the predicted caption respectively, and \mathcal{D}_0 represents the initial dataset. The model parameter θ_0 is optimized to establish basic image-to-text mapping capabilities.

¹This convention is similar to <https://mathbitsnotebook.com/Geometry/BasicTerms/BTnotation2.html>.

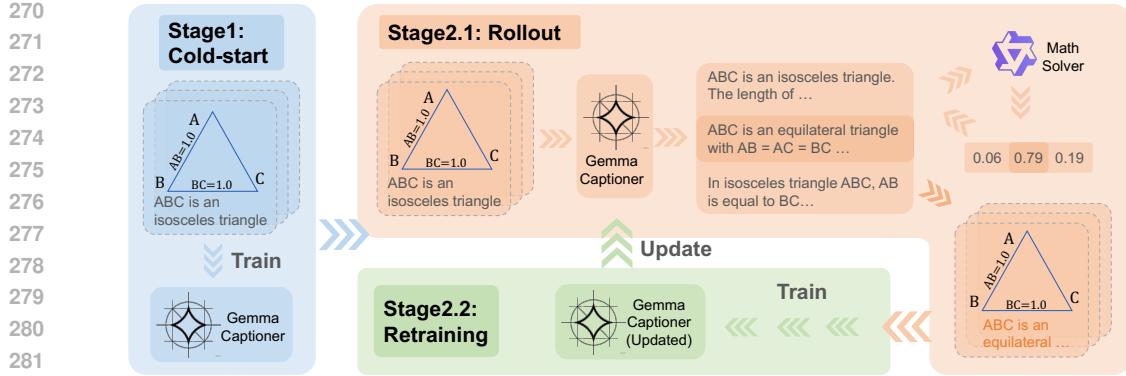


Figure 4: **The RLVR training framework.** In Stage 1, the model is trained to develop a preliminary ability to generate image captions. In Stage 2, an alternating optimization strategy is employed to jointly refine the generated captions and enhance the model’s overall performance. The data of Stage 1 comes from the rule-based image-caption generation pipeline illustrated in Figure 3.

3.3.2 RLVR PHASE

The RLVR phase with RAFT operates in alternating stages, as shown in Figure 4.

Rollout Experience Generation Suppose the current iteration is t . For each image I in the current dataset \mathcal{D}_t , we first generate N candidate captions $\{c_i\} (i = 1, 2, \dots, N)$ using the current vision language model with parameter θ_t . Then, we utilize a specifically designed reward function $R(c_i, Q_i, c^*)$ (detailed introduced in Section 3.3.3) to score each caption. Last, we retain the top-K caption $c_{\text{best}} = \arg \max_{c_i} R(c_i, Q_i, c^*)$ to update the current dataset and construct the refined dataset \mathcal{D}_{t+1} .

Model retraining We update the model by training on \mathcal{D}_{t+1} for one epoch, which is:

$$\theta_{t+1} = \arg \max_{\theta_t} \mathbb{E}_{(I, c_{\text{best}}) \sim \mathcal{D}_{t+1}} [\log P_{\theta_t}(c_{\text{best}}|I)] \quad (2)$$

This iterative process continues for $T = 5$ epochs, progressively enhancing both dataset quality and model performance.

3.3.3 REWARD FUNCTION

The composite reward $R(c, q, c^*)$ balances task correctness and caption-image alignment, as shown in Eq. 3. [And the overall reward modeling is shown in Figure 5.](#)

$$R(c, I) = \lambda_r \cdot R_{\text{reasoning}}(c, q) + (1 - \lambda_r) \cdot R_{\text{caption}}(c, c^*) \quad (3)$$

Reasoning reward To evaluate a candidate caption’s utility for solving downstream tasks, we leverage a frozen large language model of Qwen2.5-7B-Instruct (Yang et al., 2024) to generate an answer $a \sim P_{\text{LLM}}(a|q, a^*, c)$ where q, a^* is the geometric question and its groundtruth answer generated by a reasoning model (Gemini2.5 Flash) corresponding to the caption c^* in advance. As encouraged by mainstream RL process, we check both the format and correctness of the answer, which is:

$$R_{\text{Reasoning}} = s_c \cdot \mathbb{I}(a = a^*) + (1 - s_c) \cdot \mathbb{F}(a) \quad (4)$$

where $\mathbb{F}(\cdot)$ denotes the format checking function, and s_c indicate the weight of correctness, set as 0.9 in the experiments.

Caption reward To prevent reward sparsity during early training, we measure semantic relevance between c and the ground-truth caption c^* using ROUGE and BLEU-4, as shown in Eq. 5:

$$R_{\text{caption}} = w_r \cdot \text{ROUGE}(c, c^*) + (1 - w_r) \cdot \text{BLEU}(c, c^*) \quad (5)$$

where w_r represent the weight of ROUGE score, set as 0.7 in the experiments.

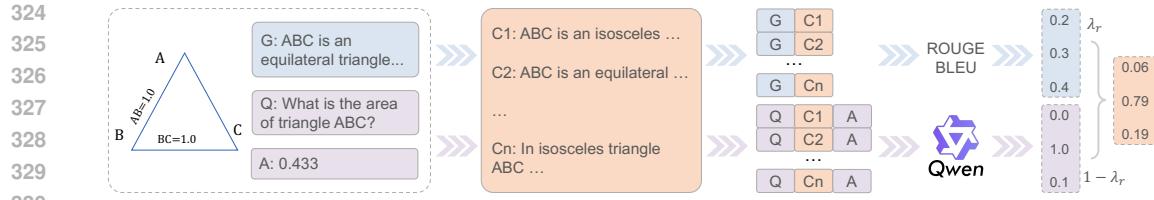


Figure 5: **The reward function.** Given the Generated caption (G), Question (Q), and Answer (A), the reward function measures the caption’s quality from two aspects: 1) its reasoning reward, and 2) caption reward, as formulated as $R(c, I) = \lambda_r \cdot R_{\text{reasoning}}(c, q) + (1 - \lambda_r) \cdot R_{\text{caption}}(c, c^*)$. **Reasoning reward** stands for the caption’s relevance to the image and question, especially the ability to capture the key reasoning information for solving the question. **Caption reward** is an auxiliary reward signal that measures the caption’s similarity to the ground truth caption.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

Our experiments utilize Gemma3-4B (Farabet & Warkentin, 2025), a commonly used lightweight multimodal architecture with strong mathematical reasoning capabilities, as our base model. All in-domain experiments are conducted on MathVista (Lu et al., 2023) and MathVerse (Zhang et al., 2024a), two established mathematical reasoning benchmarks focusing on visual and mathematical problem-solving.

4.2 IN-DOMAIN PERFORMANCE OF GEOREASONING-10K

In this section, we verify the effectiveness and scalability of our proposed dataset on in-domain benchmarks. Several commonly-used datasets in this field are chosen as baselines, including AutoGeo (Huang et al., 2025b), GeoPeP (Sun et al., 2025), GeoGPT4V (Cai et al., 2024), Geo170K (Gao et al., 2023), GeoQA (Chen et al., 2021), and MathVision Wang et al. (2024).

It can be observed from Table 1 that the model trained on GeoReasoning-10K obtains better mathematical reasoning performance compared to that trained on other caption datasets. This improvement mainly concentrates on in-domain mathematical subtasks, such as geometry, algebra, science, statistics, and most subtasks in MathVerse. The performance gain can be attributed to the symbolic synthesis process of our pipeline, which allows an infinite number of possible geometry problem types and offers diverse difficulty levels for the generated images.

Besides, the performance of models trained on the full datasets on MathVista and MathVerse is shown in Table 4 in Appendix D.

Table 1: **Better In-Domain Performance.** Accuracy of Gemma3-4B models trained on 10k random samples of each dataset over 4 trials. Results on several subtasks in MathVerse and MathVista are also reported here.

	MathVista (\uparrow)					MathVerse (\uparrow)			
	Overall	Geometry	Algebra	Science	Statistic	Overall	Vision-Dominant	Text-Dominant	Text-Lite
Base	46.2	60.7	59.1	53.3	43.2	25.2	24.0	32.0	25.9
AutoGeo	47.8 ± 0.8	62.3 ± 2.4	60.2 ± 1.9	52.5 ± 1.2	44.1 ± 0.9	24.6 ± 0.4	22.3 ± 1.4	35.2 ± 0.7	26.7 ± 1.3
GeoPeP	47.5 ± 0.4	61.0 ± 2.3	59.6 ± 1.8	54.1 ± 0.6	44.2 ± 0.8	24.2 ± 0.2	21.7 ± 0.9	33.7 ± 0.3	25.7 ± 1.3
GeoGPT4V	47.5 ± 0.2	60.5 ± 0.7	59.3 ± 1.3	54.1 ± 1.5	44.6 ± 1.0	25.2 ± 0.5	22.4 ± 0.8	36.4 ± 1.4	26.9 ± 1.0
Geo170K	47.6 ± 0.3	62.2 ± 1.5	60.6 ± 1.2	53.5 ± 1.5	43.7 ± 0.4	25.3 ± 0.1	22.5 ± 1.0	35.4 ± 1.7	26.9 ± 0.7
GeoReasoning	48.6 ± 0.3	62.8 ± 1.3	61.4 ± 1.4	54.3 ± 1.2	46.0 ± 0.5	25.8 ± 0.1	24.0 ± 0.8	36.8 ± 0.4	28.4 ± 0.5

GeoReasoning also demonstrates better scalability, as shown in Figure 6. The model trained on GeoReasoning improves progressively when the dataset sizes increase. Moreover, GeoReasoning

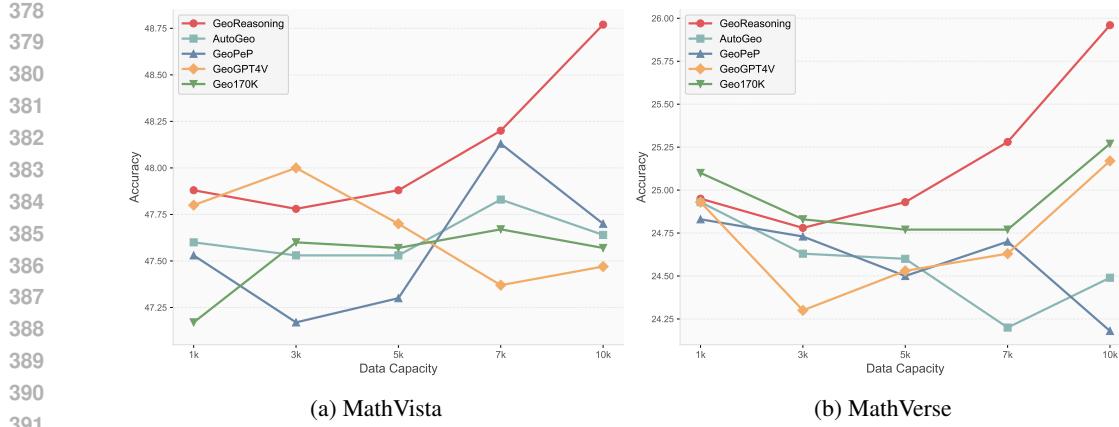


Figure 6: **Better Scalability.** The accuracy of models fine-tuned on different capacities and mathematical datasets on downstream evaluation benchmarks: (a) MathVista. (b) MathVerse.

outperforms existing datasets by a non-trivial margin at 10K size, verifying the effectiveness of the proposed cross-modal alignment method. In this setting, subsets of different sizes are randomly sampled from the original dataset, with the baseline models trained on these subsets.

4.3 OUT-OF-DOMAIN PERFORMANCE OF GEOREASONING-10K

Surprisingly, GeoReasoning-10K also demonstrates better out-of-domain generalization ability for non-geometric input images. Specifically, we evaluate the accuracy of the baseline models and the trained models on a commonly-used benchmark MMMU (Yue et al., 2023), as shown in Table 2.

Table 2: **Better Out-of-Domain Performance.** Accuracy of models evaluated on all subtasks of MMMU over 5 trials, where “A&D.”, “Busi.”, “Sci.”, “H&M.”, “Human.”, “Tech.” are short for “Art and Design”, “Business”, “Science”, “Health and Medicine”, “Humanities and Social Science”, “Tech and Engineering”, respectively.

	Overall	A&D.	Busi.	Sci.	H&M.	Human.	Tech.
Base	43.3 \pm 0.7	57.8 \pm 4.0	44.1 \pm 0.6	34.3 \pm 0.9	46.8 \pm 2.2	59.2 \pm 2.1	29.0 \pm 1.3
AutoGeo	43.5 \pm 0.5	59.3 \pm 1.4	43.3 \pm 1.1	34.9 \pm 1.3	47.4 \pm 1.1	58.9 \pm 1.5	30.7 \pm 2.6
GeoPeP	43.7 \pm 0.9	59.2 \pm 1.1	40.4 \pm 1.4	34.0 \pm 1.7	45.1 \pm 0.9	59.6 \pm 1.0	32.6 \pm 0.7
GeoGPT4V	44.0 \pm 0.7	60.2\pm1.1	43.1 \pm 1.5	34.5 \pm 0.7	46.0 \pm 0.7	58.3 \pm 1.6	30.8 \pm 2.0
Geo170K	42.9 \pm 1.0	58.5 \pm 0.8	43.6 \pm 1.4	30.9 \pm 2.0	46.8 \pm 2.2	59.9 \pm 2.2	30.9 \pm 1.6
GeoReasoning	44.9\pm0.7	60.2\pm2.0	44.5\pm2.5	36.0\pm2.0	46.7\pm1.1	60.0\pm0.5	32.9\pm1.3

It can be observed from Table 2 that the trained model outperforms the baseline on most of the domains, especially subtasks involving line art or drawing, such as Art & Design, and Tech & Engineering. These types of improvement can be mainly attributed to the proposed RLVR training process, whose reasoning reward forces the model to focus on the key element in the geometric image for problem-solving. This key-element capturing ability can help model generalize better to scenarios other than geometric problems.

We also conduct qualitative analyses on representative examples from MathVista and MMMU to understand which types of reasoning abilities our method generalizes especially well. As shown in Figure 1, 11, and 12, GeoReasoning is in particular favorable for arithmetic reasoning, line-art reasoning, and spatial reasoning. Additional examples are also available in Appendix C.

4.4 ABLATION STUDY

To understand the individual contribution of Cold-Start and RLVR phases, additional ablation studies are conducted. Specifically, we implement the Cold Start and RLVR pipelines on Gemma3-4B, generating refined models and datasets at each optimization stage, as shown in Table 3.

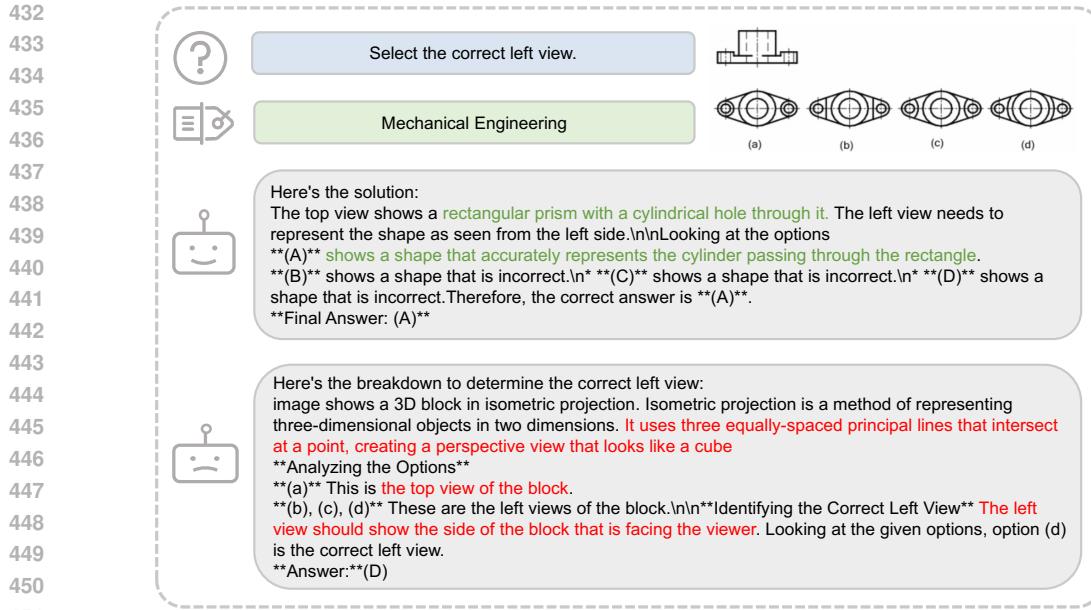


Figure 7: **Cast Study.** An engineering case, where the base model’s answer is relatively general and the analysis of shape is not rigorous enough, while the model after training on GeoReasoning is more detailed and accurate in observing shape and has spatial reasoning ability.

Table 3: **Ablation Study.** Accuracy of Gemma3-4B models with various stages.

Stage 1: Cold Start	Stage 2: RAFT	MathVerse	MathVista
✗	✗	25.2	46.2
✓	✗	25.9 (+0.7)	47.6 (+1.4)
✗	✓	26.1 (+0.9)	49.4 (+3.2)
✓	✓	27.4 (+2.2)	50.0 (+3.8)

It can be observed that RLVR benefits both the base model and the one after cold start, demonstrating that accuracy-guided captioning enhances the model’s general geometric problem-solving ability on in-domain questions.

When it comes to out-of-domain generalization, such as most non-geometric samples in MathVista, RAFT plays a more important role, contributing most of the improvements. This implies that the RLVR stage enables the captioner to focus more on key aspects of the geometric image for problem-solving, resulting in higher-quality captions for generalization. For example, captions like “the line’s color is blue” or “the left number is 26” are mostly irrelevant, whereas a caption like “the length of line AB equals line BC” is more useful for generalization purposes.

5 CONCLUSION

In this paper, we propose **Geo-Image-Textualization**, a novel reinforcement learning-based framework designed to symbolically synthesize high-quality, geometry-centered multimodal data. Leveraging this framework, we construct **GeoReasoning-10K**, a new dataset aimed at bridging the gap between visual and linguistic modalities in the geometry domain. Extensive experiments on the MathVista and MathVerse benchmarks demonstrate that our dataset significantly enhances the cross-modal reasoning capabilities of trained MLLMs, with improvements generalizing to non-geometric domains. Detailed analysis shows that geometric image caption datasets are especially beneficial for generalized skills of arithmetic reasoning, line-art reasoning, and spatial reasoning.

486 ETHICS STATEMENT
487488 After carefully reviewing the ethical regulations of the conference, to the best of our knowledge, this
489 work does not present any foreseeable ethical concerns. No negative societal or ethical impacts are
490 anticipated for the contribution of this work. The proposed GeoReasoning-10K dataset only contains
491 geometric images and corresponding captions, does not involve anything about human subjects, po-
492 tentially harmful insights, potential conflicts of interest and sponsorship, discrimination/bias/fairness
493 concerns, privacy and security issues, legal compliance, and research integrity issues.
494495 REPRODUCIBILITY STATEMENT
496497 We have made efforts to ensure that our work is reproducible, with details provided in Section 4 and
498 Appendix A, B. The code for data generation and RLVR will be released upon the acceptance of the
499 paper.
500501 REFERENCES
502503 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang,
504 Shijie Wang, Jun Tang, et al. Qwen2.5-vl technical report. [arXiv preprint arXiv:2502.13923](https://arxiv.org/abs/2502.13923), 2025.505 Shihao Cai, Keqin Bao, Hangyu Guo, Jizhi Zhang, Jun Song, and Bo Zheng. GeoGPT4V: Towards
506 Geometric Multi-modal Large Language Models with Geometric Image Generation. [arXiv e-prints](https://arxiv.org/abs/2406.11503),
507 art. [arXiv:2406.11503](https://arxiv.org/abs/2406.11503), June 2024. doi: 10.48550/arXiv.2406.11503.
508509 Jiaqi Chen, Jianheng Tang, Jinghui Qin, Xiaodan Liang, Lingbo Liu, Eric Xing, and Liang Lin.
510 GeoQA: A geometric question answering benchmark towards multimodal numerical reasoning.
511 In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), [Findings of the](https://aclanthology.org/2021.findings-acl.46)
512 [Association for Computational Linguistics: ACL-IJCNLP 2021](https://aclanthology.org/2021.46.pdf), pp. 513–523, Online, August
513 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.findings-acl.46. URL
514 <https://aclanthology.org/2021.46.pdf>.515 Zhe Chen, Jiannan Wu, Wenhui Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong
516 Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning
517 for generic visual-linguistic tasks. In [Proceedings of the IEEE/CVF Conference on Computer](https://openaccess.thecvf.com/content/CVPR2024/papers/Internvl_CVPR2024.pdf)
518 [Vision and Pattern Recognition](https://openaccess.thecvf.com/content/CVPR2024/papers/Internvl_CVPR2024.pdf), pp. 24185–24198, 2024.519 Yihe Deng, Hritik Bansal, Fan Yin, Nanyun Peng, Wei Wang, and Kai-Wei Chang. Openvlthinker:
520 An early exploration to complex vision-language reasoning via iterative self-improvement. [arXiv](https://arxiv.org/abs/2503.17352)
521 [preprint arXiv:2503.17352](https://arxiv.org/abs/2503.17352), 2025.
522523 Justin K Dimmel and Patricio G Herbst. The semiotic structure of geometry diagrams: How textbook
524 diagrams convey meaning. [Journal for Research in Mathematics Education](https://doi.org/10.1080/10720813.2015.1000000), 46(2):147–195, 2015.
525526 Hanze Dong, Wei Xiong, Deepanshu Goyal, Rui Pan, Shizhe Diao, Jipeng Zhang, Kashun Shum,
527 and T. Zhang. Raft: Reward ranked finetuning for generative foundation model alignment.
528 [ArXiv](https://arxiv.org/abs/2304.06767), abs/2304.06767, 2023. URL <https://api.semanticscholar.org/CorpusID:258170300>.
529530 Clement Farabet and Tris Warkentin. Introducing gemma 3: The most capable model you can run on
531 a single gpu or tpu. <https://blog.google/technology/developers/gemma-3/>,
532 2025.533 Jiahui Gao, Renjie Pi, Jipeng Zhang, Jiacheng Ye, Wanjun Zhong, Yufei Wang, Lanqing Hong,
534 Jianhua Han, Hang Xu, Zhenguo Li, and Lingpeng Kong. G-LLaVA: Solving Geometric Problem
535 with Multi-Modal Large Language Model. [arXiv e-prints](https://arxiv.org/abs/2312.11370), art. [arXiv:2312.11370](https://arxiv.org/abs/2312.11370), December 2023.
536 doi: 10.48550/arXiv.2312.11370.
537538 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
539 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms
via reinforcement learning. [arXiv preprint arXiv:2501.12948](https://arxiv.org/abs/2501.12948), 2025.

540 Wenzuan Huang, Bohan Jia, Zijie Zhai, Shaosheng Cao, Zheyu Ye, Fei Zhao, Zhe Xu, Yao Hu, and
 541 Shaohui Lin. Vision-r1: Incentivizing reasoning capability in multimodal large language models.
 542 [arXiv preprint arXiv:2503.06749](https://arxiv.org/abs/2503.06749), 2025a.

543

544 Zihan Huang, Tao Wu, Wang Lin, Shengyu Zhang, Jingyuan Chen, and Fei Wu. Autogeo: Automating
 545 geometric image dataset creation for enhanced geometry understanding. *IEEE Transactions on
 546 Multimedia*, 2025b.

547 Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec
 548 Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. [arXiv preprint
 549 arXiv:2412.16720](https://arxiv.org/abs/2412.16720), 2024.

550

551 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee.
 552 Llava-next: Improved reasoning, ocr, and world knowledge, January 2024. URL <https://llava-vl.github.io/blog/2024-01-30-llava-next/>.

553

554 Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng,
 555 Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning
 556 of foundation models in visual contexts. [arXiv preprint arXiv:2310.02255](https://arxiv.org/abs/2310.02255), 2023.

557

558 Yiting Lu, Jiakang Yuan, Zhen Li, Shitian Zhao, Qi Qin, Xinyue Li, Le Zhuo, Licheng Wen,
 559 Dongyang Liu, Yuewen Cao, et al. Omnicaptioner: One captioner to rule them all. [arXiv preprint
 560 arXiv:2504.07089](https://arxiv.org/abs/2504.07089), 2025.

561

562 Yingzhe Peng, Gongrui Zhang, Miaozen Zhang, Zhiyuan You, Jie Liu, Qipeng Zhu, Kai Yang,
 563 Xingzhong Xu, Xin Geng, and Xu Yang. Lmm-r1: Empowering 3b lmms with strong reasoning
 564 abilities through two-stage rule-based rl. [arXiv preprint arXiv:2503.07536](https://arxiv.org/abs/2503.07536), 2025.

565

566 Renjie Pi, Jianshu Zhang, Jipeng Zhang, Rui Pan, Zhekai Chen, and Tong Zhang. Image textualization:
 567 An automatic framework for creating accurate and detailed image descriptions. [arXiv preprint
 568 arXiv:2406.07502](https://arxiv.org/abs/2406.07502), 2024.

569

570 Hao Shao, Shengju Qian, Han Xiao, Guanglu Song, Zhuofan Zong, Letian Wang, Yu Liu, and
 571 Hongsheng Li. Visual cot: Advancing multi-modal language models with a comprehensive dataset
 572 and benchmark for chain-of-thought reasoning. *Advances in Neural Information Processing
 573 Systems*, 37:8612–8642, 2024.

574

575 Yanpeng Sun, Shan Zhang, Wei Tang, Aotian Chen, Piotr Koniusz, Kai Zou, Yuan Xue, and Anton
 576 van den Hengel. Mathglance: Multimodal large language models do not know where to look in
 577 mathematical diagrams. [arXiv preprint arXiv:2503.20745](https://arxiv.org/abs/2503.20745), 2025.

578

579 Trieu H Trinh, Yuhuai Wu, Quoc V Le, He He, and Thang Luong. Solving olympiad geometry
 580 without human demonstrations. *Nature*, 625(7995):476–482, 2024.

581

582 Junxiao Wang, Ting Zhang, Heng Yu, Jingdong Wang, and Hua Huang. Magicgeo: Training-free
 583 text-guided geometric diagram generation. [arXiv preprint arXiv:2502.13855](https://arxiv.org/abs/2502.13855), 2025.

584

585 Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Houxing Ren, Aojun Zhou, Mingjie Zhan, and
 586 Hongsheng Li. Measuring multimodal mathematical reasoning with math-vision dataset. *Advances
 587 in Neural Information Processing Systems*, 37:95095–95169, 2024.

588

589 Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Mingjie Zhan, and Hongsheng Li. Measuring Multi-
 590 modal Mathematical Reasoning with MATH-Vision Dataset. [arXiv e-prints](https://arxiv.org/abs/2402.14804), art. arXiv:2402.14804,
 591 February 2024. doi: 10.48550/arXiv.2402.14804.

592

593 Le Xue, Manli Shu, Anas Awadalla, Jun Wang, An Yan, Senthil Purushwalkam, Honglu Zhou, Viraj
 594 Prabhu, Yutong Dai, Michael S Ryoo, et al. xgen-mm (blip-3): A family of open large multimodal
 595 models. [arXiv preprint arXiv:2408.08872](https://arxiv.org/abs/2408.08872), 2024.

596

597 An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li,
 598 Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. [arXiv preprint
 599 arXiv:2412.15115](https://arxiv.org/abs/2412.15115), 2024.

594 Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Anwen Hu, Haowei Liu, Qi Qian, Ji Zhang, and Fei
 595 Huang. *mplug-owl2: Revolutionizing multi-modal large language model with modality collabora-*
 596 *tion*. In *Proceedings of the ieee/cvf conference on computer vision and pattern recognition*, pp.
 597 13040–13051, 2024.

598 Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruochi Liu, Ge Zhang, Samuel Stevens, Dongfu
 599 Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin,
 600 Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen.
 601 *MMMU: A Massive Multi-discipline Multimodal Understanding and Reasoning Benchmark for*
 602 *Expert AGI*. *arXiv e-prints*, art. arXiv:2311.16502, November 2023. doi: 10.48550/arXiv.2311.
 603 16502.

604 Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan
 605 Lu, Kai-Wei Chang, Yu Qiao, et al. *Mathverse: Does your multi-modal lilm truly see the diagrams*
 606 *in visual math problems?* In *European Conference on Computer Vision*, pp. 169–186. Springer,
 607 2024a.

608 Renrui Zhang, Xinyu Wei, Dongzhi Jiang, Ziyu Guo, Shicheng Li, Yichi Zhang, Chengzhuo Tong,
 609 Jiaming Liu, Aojun Zhou, Bin Wei, et al. *Mavis: Mathematical visual instruction tuning with an*
 610 *automatic data engine*. *arXiv preprint arXiv:2407.08739*, 2024b.

611 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. *Minigpt-4: En-*
 612 *hancing vision-language understanding with advanced large language models*. *arXiv preprint*
 613 *arXiv:2304.10592*, 2023.

614
 615
 616
 617
 618
 619
 620
 621
 622
 623
 624
 625
 626
 627
 628
 629
 630
 631
 632
 633
 634
 635
 636
 637
 638
 639
 640
 641
 642
 643
 644
 645
 646
 647

648 A QUESTION-ANSWER PAIR GENERATION PROMPT
649650 The rule-and-LLM-based pipeline contains two stages. We first design a prompt that satisfies all the
651 above conditions, using a relatively low temperature (0.2 in our experiments) to encourage the large
652 model (Gemini 2.5 Flash) to generate initial questions based on the caption, while also labeling those
653 that are inconsistent with the caption. For the inconsistent questions, we then switch to a different
654 prompt, encouraging the model to incorporate additional information and formulate new questions
655 accordingly, while increasing the temperature to 0.8. This process continues until a self-consistent
656 question is generated for the first time.657 The prompt of the first question generation stage is set as:
658659 **Prompt1**
660

```

661 You are a helpful dataset processor. Please:
662 1. Generate a mathematical question according to the
663 given description of a geometric image with the following
664 requirements:
665 1.1 The problem should base on the given description, i.e., you
666 must **NOT** generate problems that are unrelated to the given
667 description.
668 1.2 The problem difficulty should not be too low, such as
669 determining some information in the description.
670 1.3 It should also not be too hard, like introducing too much
671 extra information, but anyway you can introduce some extra
672 information to form a good geometric problem.
673 1.4 You should **NOT** include or repeat any information in the
674 description, and just contain the real question. For example,
675 if the description says: 'Line segment AB is present. The
676 length of BA is 1.24.', then when you generate the question,
677 you should not include the length of AB.
678 1.5 If the question is inconsistent with the given description,
679 the final answer should be 'None'.
680 2. Answer the question you just provided, and express the
681 final answer to two decimal places. The final answer should
682 be in \boxed{{answer}}.
683
684 Description:
685 {description}
686 Generated Question:
687 {question}
688 Generated Response:
689 {response}
690 Final Answer:
691 \boxed{{answer}}
```

692
693
694
695
696
697
698
699
700
701

702 The prompt of the question re-generation stage is set as:
 703

704
 705 **Prompt1**

706 You are a helpful dataset processor. Please: 1. Generate a
 707 mathematical question according to the given description of a
 708 geometric image with the following requirements:
 709 1.1 The problem should base on the given description, i.e., you
 710 must **NOT** generate problems that are unrelated to the given
 711 description.
 712 1.2 You can introduce some extra information to form a good
 713 geometric problem.
 714 1.3 If you find that it is hard to generate some difficult
 715 questions, just give a simple question. For example, when the
 716 lengths of all four sides of a quadrilateral are given, you can
 717 no longer assume it is a parallelogram or rectangle. In such
 718 cases, the problem may only allow for questions like asking for
 719 the perimeter, or determining the length of a segment when a
 720 certain point divides a side into an n-equal part, etc.
 721 1.4 You should **NOT** include or repeat any information in the
 722 description, and just contain the real question. For example,
 723 if the description says: 'Line segment AB is present. The
 724 length of BA is 1.24.', then when you generate the question,
 725 you should not include the length of AB.
 726 1.5 If the question is inconsistent with the given description,
 727 the final answer should be 'None'.
 728 2. Answer the question you just provided, and express the
 729 final answer to two decimal places. The final answer should
 730 be in `\boxed{{answer}}`.
 731 Description:
 732 `{description}`
 733 Generated Question:
 734 `{question}`
 735 Generated Response:
 736 `{response}`
 737 Final Answer:
 738 `\boxed{{answer}}`

739
 740
 741
 742
 743
 744
 745
 746
 747
 748
 749
 750
 751
 752
 753
 754
 755

756 **B EXPERIMENTAL DETAILS**
757758 **B.1 EXPERIMENTAL SETUP**
759760 The training and optimization pipeline contains two stages:
761

- 762 **Cold-Start phase:** we train each base model on the initial GeoReasoning-10K dataset for 1
763 epoch using standard supervised fine-tuning (SFT). The peak learning rate is 10^{-5} , with a
764 cosine learning rate scheduler and a 0.03 warm-up ratio for linear warm-up.
- 765 **RLVR phase:** we run RAFT (Dong et al., 2023) for 5 epochs, alternating between two
766 sub-phases:
 - 767 2.1) Caption Refinement: The model generates 8 candidate captions for each image,
768 and the top-1 caption per image is selected based on a composite reward with reasoning
769 reward weight of $\lambda_r = 0.7$ and caption reward weight of $1 - \lambda_r = 0.3$.
 - 770 2.2) Model Retraining: Fine-tune the model on the selected dataset for 1 epoch using
771 the same hyperparameters as the cold start phase.

772 To ensure consistency, we adopt the official evaluation codebases of both MathVerse and MathVista,
773 using the GPT-4o-mini API to evaluate the performance of our MLLM. Specifically, following each
774 benchmark’s official setup, we use GPT-4o-mini to extract and assess the correctness of answers for
775 MathVerse, and to extract answers for MathVista.
776

777 We evaluate MLLMs on MathVerse, MathVista and MMMU using A100 by VLLM. We employ
778 Gemma3-4B as our base model and fine-tune it on Georeasoning-10K using 4 L20 GPUs. The
779 training process is distributed using torchrun with the DeepSpeed ZeRO-3 optimization strategy.
780

781 **B.2 DATA SOURCE**

782 GeoReasoning-10K dataset is generated through rule-based methods and further refined using the
783 RAFT framework. The question-answer pairs are generated by Gemini 2.5-Flash with a specific
784 prompt.
785

786 **B.3 LICENSE**

- 787 • **GeoReasoning-10K** is released under the MIT License.
- 788 • **MathVerse** and **MathVista** are evaluated using their official codebases, which are publicly
789 available under the Apache 2.0 License and MIT License, respectively.
- 790 • Our use of the **GPT-4o-mini API** for evaluation complies with OpenAI’s API usage policies.
- 791 • All third-party datasets and models used in this work are under their respective licenses, and
792 we ensure compliance with their terms of use.

795 **C CASE STUDIES**

797 Additional mathematical cases that further demonstrate the generalization capability of our RAFT
798 method are presented here. Figure 8, Figure 9, and Figure 10 correspond to the geometric, arithmetic,
799 and numeric domains, respectively.
800

801 Figure 11 and Figure 12 exhibit non-mathematical cases of the economics and physics domains,
802 respectively.

803 All these examples indicate that training on geometric caption tasks stimulates the reasoning capacity
804 of models.
805

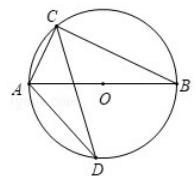
806 **D COMPARISON TO OTHER DATASETS**
807

808 We test the model performance trained on total sizes of AutoGeo, GeoPeP, GeoGPT4V, Geo170K,
809 GeoQA, MathVision, with the accuracy of MathVista and MathVerse shown in Table 4:
810

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

As shown in the figure, AB is the diameter of $\odot O$, CD is the chord of $\odot O$, $\angle ADC = 26.0$, then the degree of $\angle CAB$ is ()

Geometry



Let's analyze the given information. We are given that AB is a diameter of circle O, and CD is a chord of the circle. We are also given that $\angle ADC = 26^\circ$.
...
Since $OA = OC$, $\angle OAC = \angle OCA = \frac{180^\circ - 52^\circ}{2} = \frac{128^\circ}{2} = 64^\circ$.
Therefore, $\angle CAB = 64^\circ$.
Final Answer: The final answer is (C)

Let's analyze the given information. We are given that AB is a diameter of circle O, and CD is a chord of the circle. We are also given that $\angle ADC = 26^\circ$.
...
 $\angle AOC = 180^\circ - 52^\circ = 128^\circ$.
In $\triangle AOC$, $OA = OC$, so $\angle OAC = \angle OCA = \frac{180^\circ - 128^\circ}{2} = \frac{52^\circ}{2} = 26^\circ$.

Figure 8: A geometric case.

833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854

Is the number of brown suvs less than the number of brown rubber school buss?

Arithmetic



Let's analyze the image. We can identify the following vehicles:
* **Brown SUV:** There is one brown SUV.
...
The number of brown SUVs is 1.
The number of brown buses is 1. Since 1 is less than 1 is false, the number of brown SUVs is not less than the number of brown buses.
Therefore, the answer is (B).
...
Let's analyze the image. We need to count the number of brown SUVs and brown school buses.
* **Brown SUVs:** There is one brown SUV
...
Since the number of brown SUVs (1) is less than the number of brown school buses (1), the statement "the number of brown suvs is less than the number of brown rubber school buss" is true.
...
A

Figure 9: An arithmetic case.

E ADDITIONAL ABLATION STUDIES

This section serves as a complement of Section 4.4, exhibiting ablation studies on various domains, and hyperparameters of the reward function.

E.1 ABLATION STUDY ON VARIOUS DOMAINS

We record the skills across diverse domains like geometry and arithmetic on various training stages in Table 5 and Table 6:

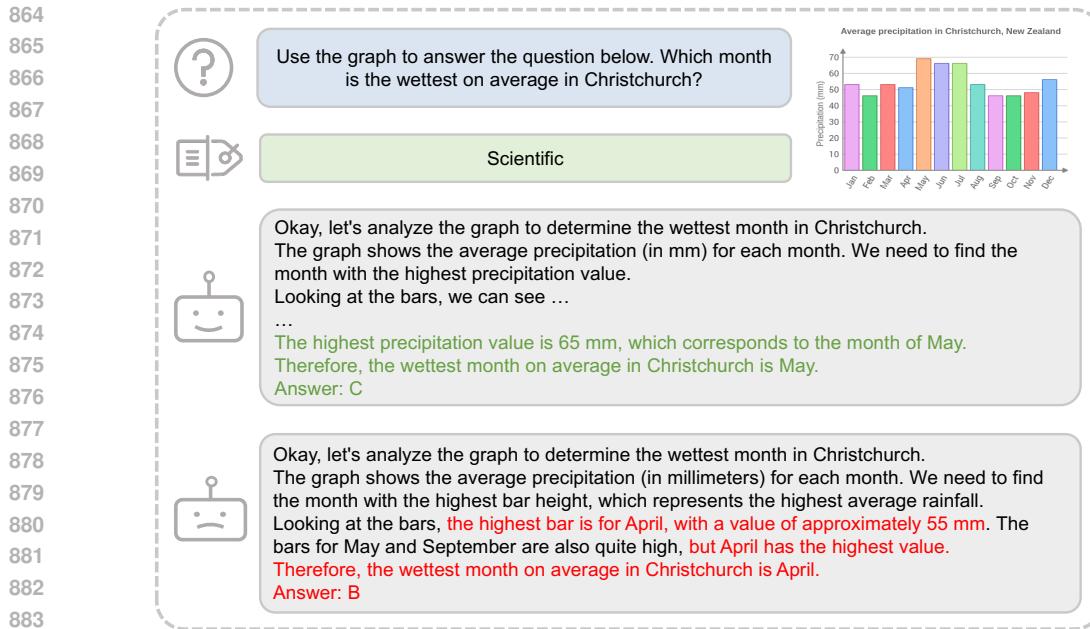


Figure 10: A numeric case.

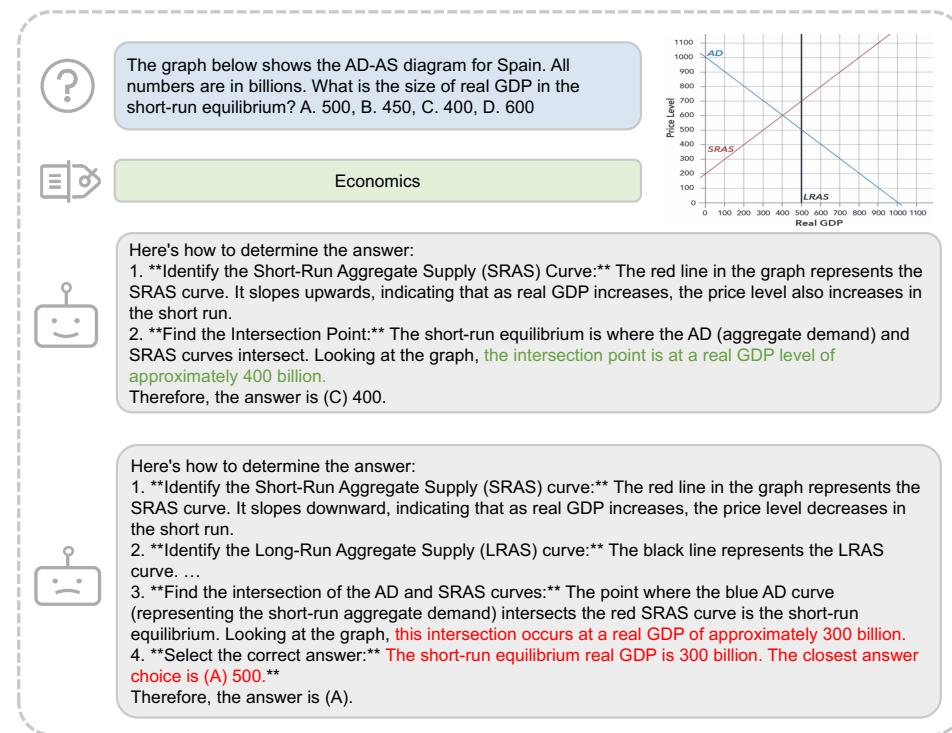
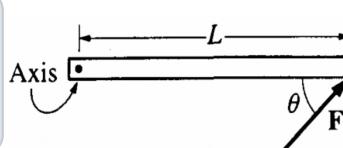


Figure 11: An economics case.

914 It can be observed from Table 5 and Table 6 that the model after RAFT stages outperforms the base
915 model across all domains. Specifically, the model achieves significant performance improvements
916 across the arithmetic, algebraic, and numeric domains, with respective gains of 5.1%, 4.2%, and
917 5.5%. These results demonstrate the effectiveness of our approach in enhancing model performance
as well as its generalization capability across different domains.

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

A rod on a horizontal tabletop is pivoted at one end and is free to rotate without friction about a vertical axis, as shown. A force F is applied at the other end, at an angle θ to the rod. If F were to be applied perpendicular to the rod, at what distance from the axis should it be applied in order to produce the same torque? A. $L \sin \theta$, B. $L \cos \theta$, C. L , D. $L \tan \theta$



View from Above

Physics

Here's how to solve this problem step-by-step:
 1. Understand Torque Torque (τ) is the rotational equivalent of force. It's calculated as: $\tau = r \times F \times \sin(\theta)$...
 2. Calculate the Torque with the Given Force In the given scenario, the force F is applied at a distance L from the axis, at an angle θ to the rod. Therefore: $\tau_{given} = L \times F \times \sin(\theta)$
 3. Calculate the Torque with the Perpendicular Force If the force F were applied perpendicular to the rod, the distance from the axis would be ' x '. The angle between the force and the lever arm would be 90 degrees ($\sin(90^\circ) = 1$). So: $\tau_{perp} = x \times F \times \sin(90^\circ) = x \times F \times 1 = x \times F$
 4. Set the Torques Equal To produce the same torque, $\tau_{given} = \tau_{perp}$: $L \times F \times \sin(\theta) = x \times F$
 **5. Solve for x ** Notice that ' F ' appears on both sides of the equation, so we can cancel it out: $L \times \sin(\theta) = x$
 = x
 Therefore, $x = L \sin(\theta)$
 Answer: (A) $L \sin \theta$

Here's how to solve this problem step-by-step:
 1. Understand Torque Torque (τ) is the rotational force that causes an object to rotate. It's calculated as...
 2. Calculate the Torque with the Given Force In the given scenario, the force F is applied at a distance L from the axis, at an angle θ to the rod. Therefore: $\tau_{given} = L \times F \times \sin(\theta)$
 3. Calculate the Torque with the Perpendicular Force If the force F were applied perpendicular to the rod, the distance from the axis would be L , and the angle θ would be 90 degrees ($\sin(90^\circ) = 1$). The torque would be: $\tau_{perp} = L \times F \times \sin(90^\circ) = L \times F \times 1 = L \times F$
 4. Set the Torques Equal To produce the same torque, $\tau_{given} = \tau_{perp}$: $L \times F \times \sin(\theta) = L \times F$
 5. Solve for the New Distance Divide both sides by $L \times F$: $\sin(\theta) = 1$ Therefore, $\theta = 90^\circ$. However, the question asks for the distance from the axis where the force should be applied to produce the same torque when the force is applied perpendicular to the rod. Let ' x ' be the distance from the axis where the force should be applied. $\tau_{given} = x \times F \times \sin(\theta)$
 We want $\tau_{given} = L \times F \times \sin(\theta)$, so: $x \times F \times \sin(\theta) = L \times F \times \sin(\theta)$ Divide both sides by $F \times \sin(\theta)$: $x = L$
 Answer: (C) L

Figure 12: A physics case.

Table 4: Performance of Gemma3-4B models trained on the total capacity of our dataset and counterpart datasets.

	Capacity	MathVerse	MathVista
Geo170K	117k	22.0	46.8
GeoPeP	100k	22.7	47.1
GeoGPT4V	23k	24.6	46.0
AutoGeo	100k	24.7	46.1
MathVision	3k	24.7	46.9
GeoQA	5k	24.9	46.0
GeoReasoning	10k	25.8	48.6

E.2 ABLATION STUDY ON HYPERPARAMETERS

We evaluated the RAFTed models with various hyperparameters on MathVista and MathVerse, as shown in Table 7:

As shown in Table 7, the reasoning reward plays a more important role in MathVista than MathVerse, indicating that the gain of generalization comes more from the helpfulness in solving the question other than comparison with captions.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
Table 5: Accuracy of Gemma3-4B models at various stages tested on MathVista

	baseline	cold-start	raft-1	raft-2	raft-3	raft-4	raft-5
all	46.2	47.6	48.7	48.1	49.2	49.0	50.0
geometry	60.7	62.3	63.2	64.0	63.6	60.3	64.0
arithmetic	42.5	45.0	44.8	45.3	45.9	47.6	46.5
algebraic	59.1	60.5	62.3	62.3	62.3	59.1	63.3
numeric	26.4	31.9	29.9	31.3	31.3	31.9	31.9

Table 6: Accuracy of Gemma3-4B models at various stages tested on MathVerse

	baseline	cold-start	raft-1	raft-2	raft-3	raft-4	raft-5
all	25.2	25.9	25.7	25.8	25.5	26.5	27.4
text dominant	32.0	35.5	35.1	35.2	35.1	36.5	36.5
text lite	25.9	27.4	28.2	28.5	27.4	26.6	26.3
vision intensive	24.0	24.8	24.4	24.4	23.1	26.1	26.5

Table 7: Accuracy of RAFTed models with various hyperparameters evaluated on MathVista and MathVerse, where λ_r stands for the weight of reasoning reward.

	MathVista	MathVerse
$\lambda_r = 1$	49.8	27.5
$\lambda_r=0.7$	50.0	27.4
$\lambda_r=0$	48.9	27.5

In addition, it is observed in the result that the performance is not very sensitive to the selection of this hyperparameter, indicating the robustness of our RAFT method.

F THE SCALING OF GOREASONING

To further understand the scaling property of GeoReasoning, we scale the dataset to 20K samples. We train the base model on this expanded dataset using exactly the same hyperparameters as in the prior experiments, and then validate its performance on MathVista and MathVerse.

As shown in Figure 13, the generalization performance keeps improving from 10K to 20K. This property is different from AutoGeo and GeoPeP (Table 1 and Table 4), which show performance degradation when scaling up the training samples, demonstrating the higher quality of GeoReasoning compared to them.

Besides, we evaluate the out-of-domain performance on MMMU, as shown in Table 8.

Table 8: **Scaling from 10K to 20K.** The performance of models trained on GeoReasoning with 10K or 20K samples on MMMU.

	MMMU
GeoReasoning-10K	44.9 ± 0.7
GeoReasoning-20K	45.2 ± 0.5

As shown in Table 8, the out-of-domain performance improves from 10K to 20K, indicating the strong power of GeoReasoning to improve the generalization capacity of models.

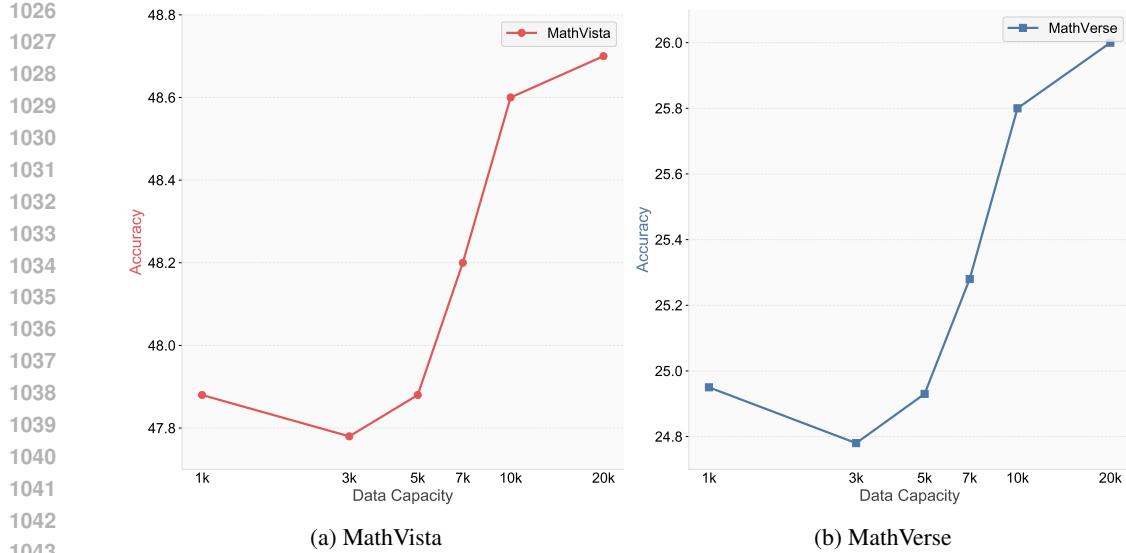


Figure 13: **Scalability with more samples.** The accuracy of models fine-tuned on different capacities and mathematical datasets on downstream evaluation benchmarks: (a) MathVista. (b) MathVerse.

G THE USE OF LARGE LANGUAGE MODELS

ChatGPT and GPT-5 were adopted to polish the writing of the paper, where all revised sentences were double-checked by the authors. ChatGPT was also utilized to write external parallelization scripts to speed up the image generation process, and the scripts were carefully reviewed by the authors.

H BROADER IMPACTS

The provided dataset pipeline and the generated dataset contribute to enhancing the generalizable reasoning abilities of multimodal large language models (MLLMs). In narrow domains, they are particularly effective for improving the geometric problem-solving capabilities of MLLMs, while in broader domains, they support the development of mathematical reasoning skills applicable to everyday scenarios. As the dataset is limited to geometric mathematical problems, it is considered safe for release and is unlikely to pose direct negative social impacts.