Theorem-Validated Reverse Chain-of-Thought Problem Generation for Geometric Reasoning

Anonymous ACL submission

Abstract

Large Multimodal Models (LMMs) face limitations in geometric reasoning due to insufficient Chain of Thought (CoT) image-text training data. While existing approaches leverage template-based or LLM-assisted methods for 007 geometric CoT data creation, they often face challenges in achieving both diversity and precision. To bridge this gap, we introduce a twostage Theorem-Validated Reverse Chain-of-Thought Reasoning Synthesis (TR-CoT) frame-011 work. The first stage, TR-Engine, synthesizes theorem-grounded geometric diagrams with structured descriptions and properties. The 014 015 second stage, TR-Reasoner, employs reverse reasoning to iteratively refine question-answer pairs by cross-validating geometric properties 017 and description fragments. Our approach expands theorem-type coverage, corrects longstanding misunderstandings, and enhances geo-021 metric reasoning. Fine-grained CoT improves theorem understanding and increases logical consistency by 24.5%. Our best models surpass the baselines in MathVista and GeoQA by 10.1% and 4.7%, outperforming advanced closed-source models like GPT-4o.

1 Introduction

037

041

Large Language Models (LLMs) (OpenAI, 2024; Guo et al., 2025) have revolutionized textual mathematical reasoning through advanced inferential mechanisms. While architectural innovations now enable these models to process multimodal inputs via parameter-efficient vision-language alignment (e.g., GPT-40 (Islam and Moushi, 2024), Gemini (Team et al., 2023)), achieving human-competitive VQA performance (Fan et al., 2024), their geometric reasoning remains constrained (Wang et al., 2025). This limitation stems from training data dominated by natural scenes, which lack the geometric specificity required for rigorous spatial problem-solving.

Current methods for generating geometric reasoning data through Chain-of-Thought (CoT) frameworks face three fundamental limitations. First, rephrasing approaches (Gao et al., 2023b) use LLM to transform the CoT format of existing problems, which requires scarce high-quality annotations and domain-specific expertise to ensure theorem consistency (Fig. 1 (a)). Second, templatebased methods (Kazemi et al., 2023a; Zhang et al., 2024b) generate geometrically oversimplified images by combining predefined polygons in rigid configurations, lacking theorem-aware element interactions, limiting their applicability to advanced reasoning, as shown in Fig. 1 (b). Thirdly, while LMM-based reasoning (Peng et al., 2024) ensures reasoning diversity, insufficient mathematical priors often lead to incorrect reasoning, e.g., misusing theorems in the wrong situation, leading to logically invalid chains of reasoning(Fig. 1 (c)).

042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

078

079

081

We introduce Theorem-Validated Reverse Chainof-Thought (TR-CoT), a two-stage framework designed to generate geometric reasoning data and verify logical flows, as shown in Fig. 1 (d). We first develop the theorem-driven image and property generation engine (TR-Engine) to create images paired with geometric properties, ensuring dependencies among elements. Then, TR-Reasoner derives questions from answers by segmenting image descriptions, generating single-step reasoning, and combining them into multi-step reasoning chains. Each step is verified against geometric properties, discarding pairs that violate mathematical rules, ensuring the logical rigor of generated data.

With TR-CoT, we create TR-GeoMM and TR-GeoSup, comprehensive datasets of diverse geometric theorems, which fully leverage CoT information. TR-CoT can bring notable and consistent improvements across a range of LMM baselines such as LLaVA, Qwen, and InternVL. Using the recent LMM baselines, we achieve a new performance record in 2B, 7B, and 8B settings for solving



Figure 1: Comparison of TR-CoT with existing CoT data generation approaches. (a) Rephrase existing Q&A pairs using LLMs, relying on existing CoT data. (b) Generate images and CoT data using pre-defined templates containing a limited number of theorems. (c) Generate CoT using LMMs, where accuracy is limited by the performance of the LMMs. (d) Design the TR-Engine to generate images, corresponding descriptions, and geometric properties from theorems. And input the descriptions and properties into TR-Reasoner to generate reliable CoT Q&A pairs.

geometry problems. The main advantages of our method are summarized as follows:

- Compared to traditional template-based methods, our approach covers twice the number of theorem types, effectively correcting long-standing theorem misunderstandings in models and enhancing their geometric reasoning.
- Generating geometric data with fine-grained CoTs enhances the model's understanding of theorems, increasing the proportion of logically consistent and clear outputs by 24.5%.
- Our most advanced models achieve a 10.1% performance gain on MathVista and 4.7% on GeoQA over the baseline, outperforming advanced closed-source models such as GPT-40.

2 Related Work

083

086

101

102

104

105

106

108

110

111

112

113

114

Enhancing Reasoning with CoT in Inference. Chain-of-thought (CoT) prompting has improved reasoning in math tasks. KQG-CoT (Liang et al., 2023a) selects logical forms from unlabeled data via CoT-based KBQG. In general math, code-based self-verification (Zhou et al., 2023) and SSC-CoT (Zhao et al., 2024b) enhance reliability by combining reasoning with structured knowledge. Other prompting strategies, including PEP (Liao et al., 2024), Plan-and-Solve (Wang et al., 2023), and incontext demonstrations (Didolkar et al., 2024), further refine inference. In geometry, visual-symbolic CoT methods (Zhao et al., 2024a; Hu et al., 2024) align reasoning with multimodal representations.

Enhancing Reasoning in Geometry Training. Training geometric solvers requires scalable and diverse data. Early symbolic systems (e.g., GeoS (Seo et al., 2015), Inter-GPS (Lu et al., 2021)) relied on small benchmarks, while neural approaches like UniGeo (Chen et al., 2022) and PGPS9K (Zhang et al., 2023a) scaled up with costly manual annotations. Recent methods automate data generation using visual-language models (e.g., G-LLaVA (Gao et al., 2023a)) or code-based engines (Kazemi et al., 2023b; Zhang et al., 2024b). Geo-Eval (Zhang et al., 2024a) provides fine-grained evaluation across diverse reasoning settings. LLMgenerated CoT traces (Peng et al., 2024) offer new avenues for training data synthesis.

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

Recently, reverse engineering has helped diagnose and refine LLM reasoning. Techniques such as condition-answer swapping (Jiang et al., 2024; Weng et al., 2023), error localization (Xue et al., 2023), and prompt optimization (Yuan et al., 2024) validate reasoning consistency without model updates. However, they often lack integration into training. Our approach embeds reverse reasoning into CoT generation, producing fine-grained, theorem-aware supervision for model training.

3 Theorem-Validated Reverse Chain-of-Thought

There are two key challenges for generating geometry reasoning data: (1) Direct generation of question-answer pairs often leads to errors or unsolvable problems due to oversimplified scenarios. (2) Single-step reasoning processes lack validation of intermediate steps, compromising reliability.

We propose Theorem-Validated Reverse Chain-



Figure 2: The TR-Engine generates diverse images, corresponding descriptions, and geometric properties step by step based on geometric theorems. Subsequently, the TR-Reasoner is utilized to obtain accurate geometric Q&A pairs from descriptions and properties.

of-Thought (TR-CoT), a two-stage framework for creating geometry reasoning data with verified logical flow, as shown in Fig. 2. The pseudo-code of TR-CoT is shown in Appendix A.

147

148

149

150

151

152

153

155

157

158

159

160

161

162

163

Stage 1: Theorem-Driven Image & Property Generation. We collect 110 fundamental geometry theorems (Complete theorems and collection method are shown in Appendix I) and develop TR-Engine, a structured method to generate images paired with textual descriptions and geometric properties (e.g., angles, lengths). Unlike random image generation, TR-Engine guides image generation based on the sampled theorems and enforces dependencies between geometric elements across generation steps. Each current step must operate on the geometric primitives—such as lines, angles, and points—produced in the preceding step.

2) Stage 2: Q&A Generation with Stepwise 164 Validation. Using the descriptions and properties 165 from Stage 1, **TR-Reasoner** generates questions from answers through three steps: First, the image 167 description is divided into logical segments (e.g., "Triangle ABC is isosceles with AB = AC"). An 169 LLM processes these parts step-by-step, generating 170 individual inferences that are then combined into 171 multi-step reasoning chains. Secondly, for each reasoning step, the system creates corresponding 173 questions. For instance, the inference " $\angle B = \angle C$ " generates the question: "If triangle ABC is isosce-176 les with AB=AC, which angles are equal?" Finally, all Q&A pairs are cross-checked against the ge-177 ometric properties from Stage 1. Pairs violating 178 mathematical rules (e.g., claiming " $\angle A = 90^{\circ}$ " for 179 a non-right isosceles triangle) are discarded. 180

3.1 TR-Engine

TR-Engine is a theorem-guided framework for synthesizing geometrically valid images with rich relational structures, corresponding descriptions, and geometric properties. TR-Engine operates through four key components (Fig. 3): 181

182

183

184

185

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

1) Geometric Theorem Library. The 110 fundamental geometric theorems are classified into substrate-related theorems and line-element-related theorems. During the image generation process, 1 to 3 theorems from each category are sampled to guide the selection of geometric substrates and the addition of line elements.

2) Geometric Substrate Library. We curate 20 fundamental geometric shapes (substrates), such as triangles, circles, and quadrilaterals. Each substrate is paired with a set of relevant geometric theorems and description templates. During image generation, appropriate substrates are selected according to the sampled theorems. The description templates encode geometric conditions (e.g., "In triangle ABC, AB = 5 cm and BC = 6 cm") to anchor subsequent reasoning steps.

3) Theorem-Based Dynamic Element Injection. This component strategically injects elements to enable complex reasoning scenarios based on theorem requirements. For example: Adding parallel lines to invoke properties of alternate angles. Introducing auxiliary lines (e.g., medians, altitudes) to create congruent sub-shapes. Such operations expand reasoning opportunities while maintaining geometric validity. In addition, TR-Engine assigns line segment values and angle degrees using exact vertex coordinates, preventing numerical conflicts from geometric constraints.

4) Property Computation Module. As elements are added, the vertex coordinates are used to au-



Figure 3: Overview of the TR-Engine. Starting from a Geometric Substrate Library, dynamically injecting elements based on theorems, and integrating a property computation module to enable multi-step geometric reasoning and validation in image generation.

tomatically calculate: Metric properties: Lengths, angles, areas. Relational properties: Parallelism, congruence, symmetry. These properties serve as ground truth for verifying generated Q&A pairs. Additionally, we perform a visual fidelity check on geometric properties, filtering out distorted images with abnormal vertex spacing (the ratio of the maximum distance to the minimum distance exceeds a threshold) or extreme angles (less than 15 degrees or more than 160 degrees).

> By integrating theorem-driven construction with stepwise validation, TR-Engine ensures images inherently support multi-step geometric reasoning, which is a critical advance over prior generation methods in practice.

3.2 TR-Reasoner

218

219

221

228

231

237

241

242

243

244

245

Despite advances in LLMs, generating accurate and educationally viable geometric question-answer (Q&A) pairs remains challenging due to three persistent issues: (1) misapplication of geometric theorems in multi-step proofs, (2) diagram-text misalignment in problem formulation, and (3) inability to maintain answerability constraints during question generation. To address these limitations, we propose the TR-Reasoner to generate theoremgrounded Q&A pairs through coordinated interaction between geometric properties and structured reasoning chains (Fig. 4).

Description Patch Reasoning Fusion Building on the geometrically valid descriptions from TR-Engine, this module enforces logical coherence through causal dependencies between reasoning steps. Let $D = \{p_1, p_2, ..., p_x\}$ denote the x description patches extracted from an image, where each patch p_i corresponds to a geometrically meaningful component (e.g., "Circle O with chord AB and tangent CD"). The single-step reasoning r_i for patch p_i is generated through theorem-constrained transformation:

$$r_i = \mathcal{F}_{\text{LLM}}(p_i | r_{\langle i}, \mathcal{T}), \qquad (1)$$

252

253

254

255

256

257

258

260

261

262

264

265

266

267

269

270

271

272

273

274

275

276

277

278

279

280

281

where $r_{\langle i} = \{r_1, ..., r_{i-1}\}$ represents preceding reasoning states, and \mathcal{T} denotes the applicable theorem set (e.g., intersecting chords theorem for patch p_i describing chord intersections). This chained formulation ensures cumulative reasoning: later steps automatically inherit and extend prior conclusions (e.g., deriving arc lengths after establishing chord congruence).

Reverse Question Generation To prevent answerability drift, we implement *answerconstrained reverse generation* rather than open-ended question synthesis. Given a verified reasoning chain $R = \{r_1, r_2, ..., r_n\}$, each step r_i undergoes answerability assessment through a theorem-aware discriminator:

$$f_{aq}(r_i) = \begin{cases} f_q(r_i; \Phi_{\text{geo}}), & \text{if } \mathcal{V}(r_i, G_{\text{props}}) = \text{True} \\ \emptyset, & \text{otherwise} \end{cases}$$
(2)

where G_{props} denotes geometric properties from TR-Engine (e.g., coordinate-derived lengths), \mathcal{V} performs theorem-based validation (e.g., checking triangle congruence rules), and f_q generates questions using a geometry-specialized LLM with instruction prompt Φ_{geo} . This approach leverages the granular reasoning steps from the patch reasoning stage to generate theorem-aware Q&A pairs.



Figure 4: Overview of the TR-Reasoner. Image descriptions are segmented into patches to generate single-step reasoning results. Single-step reasoning results are fused progressively to get multi-step reasoning results. Then questions are generated based on the multi-step reasoning results. Finally, Q&A pairs that contradict geometric properties are filtered.

Error A&O Filtering The final verification stage applies bidirectional cross-validation to ensure Q&A quality. The forward validation aligns generated answers with deterministic geometric properties computed by TR-Engine's analytical algorithms, removing cases demonstrating: (1) final answer-property discrepancies, and (2) intermediate reasoning inconsistencies with verified properties. The reverse validation identifies ill-posed questions through semantic analysis, excluding those exhibiting answer ambiguity or logical indeterminacy. Both of the validation are conducted through single-round LLM inference, and only Q&A pairs that satisfy both verifications are reserved. Quantitative analysis revealed four main error patterns that were filtered out: Theorem Violation (36.3%): incorrect geometric principle application; Metric Discrepancy (24.9%): numerical inconsistency with problem constraints; Diagram-Text Mismatch (12.2%): references to non-existent diagram elements; and Answerability Ambiguity (26.6%): ill-defined problem statements.

285

290

296

297

299

302

311

Our proposed filtering mechanism can effectively reduce model hallucination and accumulate errors in previous reasoning steps. Among a sample of 200 generated Q&A pairs, the framework successfully suppresses reasoning error, reducing overall error rates from 16.0% (pre-validation) to 5.0% (post-validation). Showcases of invalid samples in Appendix D.

312 **Context-Aware Prompt Engineering** We de-313 ploy an instruction-based context-aware prompting strategy to optimize reasoning. We construct a reasoning instruction template pool containing prototypical geometric problems with a corresponding reasoning process. For each input, 3-4 optimal templates that are most relevant to the theorem and content is selected and integrated into the prompt. Additionally, the pool also contains a series of geometric relationships that are easily misunderstood by LLM. We use the same sample strategy to integrate them into the prompt as well, referred to as Basic Knowledge. The sampled instruction templates and basic knowledge serve as examples to assist the LLM to perform correct reasoning. Such context-aware prompt engineering ensures a relatively ideal reasoning accuracy, improving the efficiency of data generation. More details about the prompt strategy in Appendix B.

314

315

316

317

318

319

320

321

322

323

324

325

327

328

329

330

331

332

333

334

335

336

338

340

341

344

3.3 TR-GeoMM

Through the TR-CoT pipeline, we construct the TR-GeoMM dataset to enhance LMM's geometric reasoning ability. From 15k figures, we obtain 45k high-quality Q&A pairs after error filtering, averaging 3.49 questions per figure. Detailed dataset statistics are shown in Fig. 6.

At the image level, TR-GeoMM covers 20 substrate shapes, mainly triangles, quadrilaterals, and circles. Unlike conventional polygon-based designs, TR-Engine builds figures from lines as primitive elements. It emphasizes key lines with geometric significance, e.g., midlines, angle bisectors, and radii, which frequently appear in theorems. As



Figure 5: Diversity analysis of TR-GeoMM.



Figure 6: Statistical information about TR-GeoMM.

illustrated in Fig. 5 (a), 1.7k unique patterns are formed through theorem-guided line combinations, where each addition must interact with existing elements(e.g., a new line's vertex must align with previously generated lines). At the text level, questions are categorized into four core types: side lengths, angles, areas, and geometric relationships. The hierarchical figure construction induces interdependent questions, where earlier solutions serve as prerequisites for subsequent ones. This subproblem design supports step-by-step learning of geometric concepts and reasoning. As shown in Fig. 5 (b), TR-GeoMM contains a theorem repository twice as large as existing synthetic datasets (MAVIS and GeomVerse). Furthermore, Fig. 5 (c) demonstrates superior data diversity through higher Q&A pair cosine distances. More information is provided in Appendix C and Appendix F.

3.4 TR-GeoSup

345

347

354

367

371

374

378

TR-CoT can not only generate reliable CoT geometric data but also be used to augment existing datasets. Real-world geometry CoTs often include key intermediate steps rich in problem-solving insights, yet these are typically implicit or oversimplified, relying on human prior knowledge. This lack of explicit reasoning may hinder model learning due to limited background knowledge and inference capability. Leveraging the TR-CoT pipeline, we decompose the original CoT process into explicit theorem-aware steps, then reverse generate new Q&A pairs with TR-Reasoner.

Specifically, our augmentation involves three steps: generating a comprehensive multi-step analysis of the geometric figure, segmenting it into essential problem-solving steps, and creating new Q&A pairs for each step. These fine-grained Q&A pairs explicitly guide the model with theorems and knowledge implicitly expressed in the original data, leading to improvement in comprehension and reasoning abilities. We applied TR-Reasoner to the GeoQA dataset, producing the TR-GeoSup dataset with 20k new entries. The final TR-GeoSup dataset does not contain the original GeoQA data. Examples of TR-GeoSup are shown in Appendix E.

380

381

383

385

386

387

390

391

392

393

394

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

During the augmentation, LLM receives the original question and its corresponding CoT answer to produce a more complete analysis, supplementing missing theorems and steps not explicitly stated in the original CoT. We sampled 200 examples from both the analysis and Q&A generation stages and observed no errors, confirming the reliability of our design. To streamline the data generation process, we did not introduce additional independent validation. After generation, 10% of the data was manually reviewed and corrected.

4 Experiments

4.1 Setup

We train multiple LMMs (Wang et al., 2024; Liu et al., 2024; Chen et al., 2024c) using existing geometric instruction datasets (Chen et al., 2021; Gao et al., 2023b) and our TR-CoT generated data (TR-GeoMM and TR-GeoSup). Both the projected linear layer and the language model are trainable. The models are trained for two epochs with a batch size of 128 on 16×64 G NPU, and learning rate set to 5e-6. For evaluation, we assess these models on the geometry problem solving on the testmini set of MathVista (Lu et al., 2023) and GeoQA (Chen et al., 2021) following Gao et al. (2023b). Top-1 accuracy serves as the metric, with predictions and ground truth evaluated via ERNIE Bot 4.0. Ablation experiments were done on Intern-VL-2.0-8B.

4.2 Ablation Study

Data generating procedures. To evaluate the contributions of TR-CoT components, we construct

ablated variants by removing specific modules, as 420 summarized in Tab. 1. Each variant is used to gener-421 ate training data, and the resulting models are eval-422 uated on MathVista and GeoQA. Generating Q&A 423 pairs from descriptions yields better performance 494 than from images, with gains of 5.3% on MathVista 425 and 6.3% on GeoQA. Incorporating reverse gener-426 ation further improves accuracy by 2.9% and 2.6% 427 on the two datasets, respectively. The full TR-CoT 428 pipeline achieves the best performance, confirming 429 the effectiveness of each component. 430

Table 1: Ablation study on the data generating procedures. 'Description' represents generation based on descriptions. 'Reverse' represents generating reasoning followed by reverse question generation. 'Filter' represents filtering errors based on geometric properties.

Configurations		MathVista	GeoOA	
Description	Reverse	Filter	iviatii vista	ULUQA
×	X	X	55.3	44.2
\checkmark	X	X	60.6	50.5
\checkmark	\checkmark	×	63.5	53.1
\checkmark	\checkmark	\checkmark	64.4	54.0

431

432

433

434

435 436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

Separate validity of synthetic and augmented data. We evaluated the impact of the TR-GeoSup and TR-GeoMM datasets on model performance, as shown in Tab. 2. Training with TR-GeoSup improved performance by 1.4% on MathVista and 7.9% on GeoQA compared to the baseline. Combining GeoQA with TR-GeoSup improves performance by 2.9% on MathVista and 3.9% on GeoQA compared to GeoQA alone, indicating their complementarity. It suggests TR-GeoSup effectively enhances in-domain performance with better extracted knowledge. A deeper understanding of knowledge may contribute to improved generalization on mixed out-of-domain datasets.

Table 2: Ablation study on the TR-CoT generated data.

Configurations		MathVisto	GaaOA	
GeoQA	TR-GeoSup	TR-GeoMM	Ivialii v ista	UEUQA
X	×	×	63.0	52.4
\checkmark	×	×	64.9	64.8
X	\checkmark	×	64.4	60.3
X	×	\checkmark	64.4	54.0
\checkmark	\checkmark	×	67.8	68.7
\checkmark	×	\checkmark	65.4	67.9
\checkmark	\checkmark	\checkmark	68.3	69.0

Second, training with TR-GeoMM improved performance by 1.4% on MathVista and 1.6% on GeoQA, confirming the strong generalization of TR-CoT synthetic data to real data. Moreover, joint training with GeoQA further improved performance, highlighting the effectiveness of synthetic data in supplementing real data. Finally, when jointly training on all three datasets (GeoQA, TR-GeoSup and TR-GeoMM). The model achieved the best performance, with improvements of 5.3% on MathVista and 6.6% on GeoQA over the baseline. These results support that TR-CoT-generated data compensate for the limitations of existing datasets and enhance the model's reasoning capability. 451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

We Compared with other synthesis datasets. train InternVL-2.0-8B using TR-GeoMM and two recent synthetic datasets for geometric problems, *i.e.* MAVIS (synthesis part) (Zhang et al., 2024b) and GeomVerse (Kazemi et al., 2023a), as summarized in Tab. 3. Compared to the baseline, models trained with GeomVerse or MAVIS show a slight performance gain on GeoQA and a decline on MathVista, both lower than TR-GeoMM. We attribute this to the limited diversity of image and Q&A pairs in these datasets, which benefits the simpler distribution of GeoQA but struggles with the diverse distributions in MathVista. In contrast, TR-GeoMM, with its diverse image and Q&A pairs, improves performance on both datasets.

Table 3: Compared with other synthesis datasets.

Dataset	MathVista	GeoQA
/	63.0	52.4
GeomVerse(9k)	58.2	53.6
MAVIS(sample 48k)	57.2	53.2
TR-GeoMM(45k)	64.4	54.0
TR-GeoMM(sample 9k)	63.0	55.6

4.3 Comparison with Previous State-of-the-Art

With the proposed method, we train three special-476 ized models for geometry problem solving named 477 TR-CoT-InternVL-2.0-2B, TR-CoT-Qwen2.5-VL-478 7B, and TR-CoT-InternVL-2.5-8B on the joint 479 dataset of Geo170K and TR-CoT-generated data 480 (TR-GeoMM and TR-GeoSup). We compare our 481 models with both general and mathematical LMMs 482 on the geometry problems from testmini set of 483 MathVista and the test set of GeoQA. As shown 484 in Tab. 4, TR-CoT-InternVL-2.5-8B outperforms 485 GPT-40 by 17.3% on MathVista and TR-CoT-486 Qwen2.5-VL-7B outperforms GPT-40 by 17.8% 487 on GeoQA. Compared to mathematical LMMs, TR-488 CoT-InternVL-2.5-8B maintains a 11.1% lead on 489 MathVista, and TR-CoT-Qwen2.5-VL-7B achieves 490 a 12.5% advantage on GeoQA. For performance 491 analysis on more baselines, please refer to Ap-492 pendix J and Tab. 8. 493

_

Question: In the figure, KL is tangent to $M = \int_{17}^{17} OM $ at K. Find the value of x. Choices: (A) 6.00 (B) 9.45 (C) 18.9 (D) 37.8 L	1. To find the value of x, we can use the Pythagorean theorem in the right triangle KLM. 2. KL = 17, LM = 10, MK = x. Since KL is tangent to the circle at K, MK is the radius of the circle. Using the Pythagorean theorem in ΔKLM : KL/2 = $KM^2 + LM^2$. 3. Substitute the given values: $17^2 = x^2 + 10^2$, $x \approx 13.75$. Since none of the given choices match this value, it appears there might be an error in the problem.	Utilizing the Pythagorean theorem, we can express MK as $\sqrt{(ML^2 - KL^2)}$. For the equation $\sqrt{\{(x+10)^2 - 17^2\}} = x$, therefore, the value of x is 9.45. Answer: B
A 9x-1 B Question: Find the length of AC in the isosceles triangle ABC. Choices: (A) 1.5 (B) 7 (C) 11 (D) 12.5	To find the length of AC, we use the given information and the properties of isosceles triangles. AB = 9x - 1, AC = 4x + 1, BC = 5x - 0.5. Since ABC is an isosceles triangle, two sides are equal. We can set up the following equations based on the given information: 1. AB = AC 2. AB = BC. Step 1: Set AB = AC, 9x - 1 = 4x + 1, 9x = x + 2A = BC, 9x - 1 = 5x + 0.5. Since x must be consistent in both equations, $5/2 = 0.4, 0.375 \neq 0.4$	By equating the expressions for the sides of the triangle, we get $4x + 1$ = $5x - 0.5$. Solving for x, we find x = 1.5 , AC = $4(1.5) + 1 = 7$. Answer:C
	(a) pre-trained model	(b) post-trained model

Figure 7: Comparison of model problem solving before and after training.

Table 4: Top-1 Accuracy (%) on geometry problem solving on the testmini set of MathVista and the GeoQA test set. * represents the results from the existing papers.

Model	MathVista	GeoQA			
Closed-source LMMs					
GPT-40 (Islam and Moushi, 2024)	60.6	61.4			
GPT-4V	51.0*	43.4*			
Gemini Ultra (Team et al., 2023)	56.3*	-			
Open-source LMMs					
LLaVA2-13B (Liu et al., 2024)	29.3*	20.3*			
mPLUG-Owl2-7B (Ye et al., 2024)	25.5	21.4			
Qwen-VL-Chat-7B (Bai et al., 2023)	35.6	26.1			
Monkey-Chat-7B (Li et al., 2024)	24.5	28.5			
Deepseek-VL-7B (Lu et al., 2024)	34.6	33.7			
InternVL-2.0-2B (Chen et al., 2024c)	46.2	38.2			
InternLM-XC2-7B (Zhang et al., 2023b)	51.4	44.7			
InternVL-1.5-20B (Chen et al., 2024b)	60.1	49.7			
Qwen2-VL-7B (Wang et al., 2024)	55.1	55.7			
InternVL-2.0-8B (Chen et al., 2024c)	65.9	56.5			
InternVL-2.5-8B (Chen et al., 2024a)	67.8	59.0			
Qwen2.5-VL-7B (Wang et al., 2024)	71.6	74.5			
Open-source Mathematical I	MMs				
UniMath (Liang et al., 2023b)	-	50.0*			
Math-LLaVA-13B (Shi et al., 2024)	56.5*	47.8			
G-LLaVA-7B (Gao et al., 2023b)	53.4*	62.8*			
MAVIS-7B (Zhang et al., 2024b)	-	66.7*			
PUMA-Qwen2-7B (Zhuang et al., 2024)	48.1*	-			
MultiMath-7B (Peng et al., 2024)	66.8*	-			
TR-CoT-InternVL-2.0-2B	56.3	63.4			
TR-CoT-Qwen2.5-VL-7B	74.5	79.2			
TR-CoT-InternVL-2.5-8B	77.9	76.7			

5 Discussion

494

495

496

497

498

499

505

509

Fig. 7 highlights consistent improvements: posttrained models produce concise, logical CoTs with accurate conclusions, demonstrating robust geometric understanding. Pre-trained models show recurring errors (e.g., misdefining isosceles triangles as having two equal sides), reflecting foundational gaps in theorem comprehension. Our approach trains models on diverse theorems with structured reasoning, addressing these errors and enhancing general geometric problem-solving.

We use DeepSeek R1 and ERNIE Bot 4.0 to quantitatively evaluate model outputs before and after training, focusing on logical consistency, clarity, and lack of ambiguity (see Appendix K for detailed information). We use the average score



Figure 8: Comparison of model output quality and token length before and after training.

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

of the two models as the final score. As shown in Fig. 8 (a), the total mean score increased by 0.37 after training, the mean score for correct answers increased by 0.70, and outputs with scores of 8 or higher increased by 24.5%. We attribute these improvements to TR-CoT's explicit focus on the reasoning process, where step decomposition enhances the model's logical consistency and rigor.

We further compare the token usage for correct answers before and after training. As shown in Fig. 8 (b), the model after training requires fewer tokens on average, with the percentage of correct answers within 200 tokens increasing by 35%. We assume this improvement results from the data diversity, which enables the model to find more efficient solutions across different theorems, while a deeper understanding of the theorems allows for more concise reasoning.

6 Conclusion

We propose TR-CoT, a novel theorem-based reverse generation pipeline that enhances theorem coverage and supports fine-grained theorem understanding in geometric datasets. Models trained on TR-CoT data demonstrate a significant improvement in geometric problem solving with more concise and rigorous reasoning. We will extend this approach to other mathematical domains to further analyze the impact of theorem mastery on problemsolving, offering insights for future research.

591 592

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

Limitations 539

For our method, one major constraint is that there is still room for further improvement in the generation 541 efficiency. The overall efficiency can be divided 542 into time efficiency and data efficiency. First, in 543 our process, LLM is called multiple times for reasoning generation. The limited reasoning speed of 545 LLM becomes the bottleneck of time efficiency. In 546 addition, although we have adopted various meth-547 ods to improve the reasoning accuracy of LLM, due to the limitations of model performance, there 549 is still a certain proportion of errors in the direct output of the model. We observe that about 10%551 of the direct output is deleted in the Error A&Q 553 Filtering stage.

References

554

555

557

559

560

561

564

566

567 568

569

570

571

576

577

579

581

582

583

586

- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. arXiv preprint arXiv:2308.12966.
- Jiaqi Chen, Tong Li, Jinghui Qin, Pan Lu, Liang Lin, Chongyu Chen, and Xiaodan Liang. 2022. Unigeo: Unifying geometry logical reasoning via reformulating mathematical expression. Preprint, arXiv:2212.02746.
- Jiaqi Chen, Jianheng Tang, Jinghui Qin, Xiaodan Liang, Lingbo Liu, Eric Xing, and Liang Lin. 2021. Geoqa: A geometric question answering benchmark towards multimodal numerical reasoning. In Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021, pages 513-523.
- Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, and 1 others. 2024a. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. arXiv preprint arXiv:2412.05271.
- Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, and 1 others. 2024b. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. arXiv preprint arXiv:2404.16821.
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, and 1 others. 2024c. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 24185–24198.

- Jacob Devlin. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.
- Aniket Didolkar, Anirudh Goyal, Nan Rosemary Ke, Siyuan Guo, Michal Valko, Timothy Lillicrap, Danilo Rezende, Yoshua Bengio, Michael Mozer, and Sanjeev Arora. 2024. Metacognitive capabilities of llms: An exploration in mathematical problem solving. Preprint, arXiv:2405.12205.
- Yue Fan, Jing Gu, Kaiwen Zhou, Qianqi Yan, Shan Jiang, Ching-Chen Kuo, Yang Zhao, Xinze Guan, and Xin Wang. 2024. Muffin or chihuahua? challenging multimodal large language models with multipanel vqa. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 6845–6863.
- Jiahui Gao, Renjie Pi, Jipeng Zhang, Jiacheng Ye, Wanjun Zhong, Yufei Wang, Lanqing Hong, Jianhua Han, Hang Xu, Zhenguo Li, and Lingpeng Kong. 2023a. G-llava: Solving geometric problem with multi-modal large language model. Preprint, arXiv:2312.11370.
- Jiahui Gao, Renjie Pi, Jipeng Zhang, Jiacheng Ye, Wanjun Zhong, Yufei Wang, Lanqing Hong, Jianhua Han, Hang Xu, Zhenguo Li, and 1 others. 2023b. G-llava: Solving geometric problem with multi-modal large language model. arXiv preprint arXiv:2312.11370.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. arXiv preprint arXiv:2501.12948.
- 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. Preprint, arXiv:2406.09403.
- Raisa Islam and Owana Marzia Moushi. 2024. Gpt-4o: The cutting-edge advancement in multimodal llm. Authorea Preprints.
- Weisen Jiang, Han Shi, Longhui Yu, Zhengying Liu, Yu Zhang, Zhenguo Li, and James T. Kwok. 2024. Forward-backward reasoning in large language models for mathematical verification. Preprint, arXiv:2308.07758.
- Mehran Kazemi, Hamidreza Alvari, Ankit Anand, Jialin Wu, Xi Chen, and Radu Soricut. 2023a. Geomverse: A systematic evaluation of large models for geometric reasoning. arXiv preprint arXiv:2312.12241.
- Mehran Kazemi, Hamidreza Alvari, Ankit Anand, Jialin Wu, Xi Chen, and Radu Soricut. 2023b. Geomverse: A systematic evaluation of large models for geometric reasoning. Preprint, arXiv:2312.12241.

Zhang Li, Biao Yang, Qiang Liu, Zhiyin Ma, Shuo Zhang, Jingxu Yang, Yabo Sun, Yuliang Liu, and Xiang Bai. 2024. Monkey: Image resolution and text label are important things for large multi-modal models. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pages 26763–26773.

653

660

664

669

673

675

676

677

684

694

- Yuanyuan Liang, Jianing Wang, Hanlun Zhu, Lei Wang, Weining Qian, and Yunshi Lan. 2023a. Prompting large language models with chain-of-thought for fewshot knowledge base question generation. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 4329– 4343, Singapore. Association for Computational Linguistics.
- Zhenwen Liang, Tianyu Yang, Jipeng Zhang, and Xiangliang Zhang. 2023b. Unimath: A foundational and multimodal mathematical reasoner. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7126–7133.
 - Haoran Liao, Jidong Tian, Shaohua Hu, Hao He, and Yaohui Jin. 2024. Look before you leap: Problem elaboration prompting improves mathematical reasoning in large language models. *Preprint*, arXiv:2402.15764.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024. Visual instruction tuning. *Advances in neural information processing systems*, 36.
- Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, and 1 others. 2024. Deepseek-vl: Towards real-world vision-language understanding. *CoRR*.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2023.
 Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *The 3rd Workshop on Mathematical Reasoning and AI at NeurIPS'23*.
- 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.
- OpenAI. 2024. Openai o1 system card. preprint.
- Shuai Peng, Di Fu, Liangcai Gao, Xiuqin Zhong, Hongguang Fu, and Zhi Tang. 2024. Multimath: Bridging visual and mathematical reasoning for large language models. *arXiv preprint arXiv:2409.00147*.
- Minjoon Seo, Hannaneh Hajishirzi, Ali Farhadi, Oren Etzioni, and Clint Malcolm. 2015. Solving geometry problems: Combining text and diagram interpretation. In Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, pages 1466–1476, Lisbon, Portugal. Association for Computational Linguistics.

- Wenhao Shi, Zhiqiang Hu, Yi Bin, Junhua Liu, Yang Yang, See-Kiong Ng, Lidong Bing, and Roy Ka-Wei Lee. 2024. Math-llava: Bootstrapping mathematical reasoning for multimodal large language models. *arXiv preprint arXiv:2406.17294*.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, and 1 others. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.
- Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. 2023. Planand-solve prompting: Improving zero-shot chain-ofthought reasoning by large language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2609–2634, Toronto, Canada. Association for Computational Linguistics.
- Peijie Wang, Zhong-Zhi Li, Fei Yin, Xin Yang, Dekang Ran, and Cheng-Lin Liu. 2025. Mv-math: Evaluating multimodal math reasoning in multi-visual contexts. *arXiv preprint arXiv:2502.20808*.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, and 1 others. 2024. Qwen2vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.
- Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Shengping Liu, Bin Sun, Kang Liu, and Jun Zhao. 2023. Large language models are better reasoners with self-verification. *Preprint*, arXiv:2212.09561.
- Tianci Xue, Ziqi Wang, Zhenhailong Wang, Chi Han, Pengfei Yu, and Heng Ji. 2023. Rcot: Detecting and rectifying factual inconsistency in reasoning by reversing chain-of-thought. *Preprint*, arXiv:2305.11499.
- Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Anwen Hu, Haowei Liu, Qi Qian, Ji Zhang, and Fei Huang. 2024. mplug-owl2: Revolutionizing multimodal large language model with modality collaboration. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pages 13040–13051.
- Jiahao Yuan, Dehui Du, Hao Zhang, Zixiang Di, and Usman Naseem. 2024. Reversal of thought: Enhancing large language models with preferenceguided reverse reasoning warm-up. *Preprint*, arXiv:2410.12323.
- Jiaxin Zhang, Zhong-Zhi Li, Ming-Liang Zhang, Fei Yin, Cheng-Lin Liu, and Yashar Moshfeghi. 2024a. GeoEval: Benchmark for evaluating LLMs and multimodal models on geometry problem-solving. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 1258–1276, Bangkok, Thailand. Association for Computational Linguistics.

Ming-Liang Zhang, Fei Yin, and Cheng-Lin Liu. 2023a. A multi-modal neural geometric solver with textual clauses parsed from diagram. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence*, pages 3374–3382.

759

760

762

763

764

765

771 772

773

774

775

776

778

779

780 781

789 790

791

793

- Pan Zhang, Xiaoyi Dong Bin Wang, Yuhang Cao, Chao Xu, Linke Ouyang, Zhiyuan Zhao, Shuangrui Ding, Songyang Zhang, Haodong Duan, Hang Yan, and 1 others. 2023b. Internlm-xcomposer: A vision-language large model for advanced text-image comprehension and composition. *arXiv preprint arXiv:2309.15112*.
- Renrui Zhang, Xinyu Wei, Dongzhi Jiang, Yichi Zhang, Ziyu Guo, Chengzhuo Tong, Jiaming Liu, Aojun Zhou, Bin Wei, Shanghang Zhang, and 1 others. 2024b. Mavis: Mathematical visual instruction tuning. arXiv preprint arXiv:2407.08739.
- Xueliang Zhao, Xinting Huang, Tingchen Fu, Qintong Li, Shansan Gong, Lemao Liu, Wei Bi, and Lingpeng Kong. 2024a. Bba: Bi-modal behavioral alignment for reasoning with large vision-language models. *Preprint*, arXiv:2402.13577.
- Zilong Zhao, Yao Rong, Dongyang Guo, Emek Gözlüklü, Emir Gülboy, and Enkelejda Kasneci. 2024b. Stepwise self-consistent mathematical reasoning with large language models. *Preprint*, arXiv:2402.17786.
- Aojun Zhou, Ke Wang, Zimu Lu, Weikang Shi, Sichun Luo, Zipeng Qin, Shaoqing Lu, Anya Jia, Linqi Song, Mingjie Zhan, and Hongsheng Li. 2023. Solving challenging math word problems using gpt-4 code interpreter with code-based self-verification. *Preprint*, arXiv:2308.07921.
- Wenwen Zhuang, Xin Huang, Xiantao Zhang, and Jin Zeng. 2024. Math-puma: Progressive upward multimodal alignment to enhance mathematical reasoning. *arXiv preprint arXiv:2408.08640*.

A Pseudo Code

We have written pseudo-code for the overall flow of TR-CoT, the details of which are given in Algor. 1.

Algorithm 1: Pseudo-code of TR-CoT

	Input: Geometry substrates sampling rounds <i>n</i> , plot					
	function f , image-description pair sets S , line					
	sampling rounds k , geometric property					
	calculation module \mathcal{V} , large language model					
	\mathcal{M}					
	Output: Generated Image \mathcal{I} , Description \mathcal{D} ,					
	Geometric Properties \mathcal{T} , Question \mathcal{Q} ;					
	Answer \mathcal{A}					
1	Initialization: $\mathcal{I} \leftarrow \emptyset$, $\mathcal{D} \leftarrow \emptyset$, $\mathcal{T} \leftarrow \emptyset$, vertex					
	coordinate $\mathcal{C} \leftarrow \emptyset, r_s \leftarrow \emptyset$					
2	for $i \leftarrow 1$ to n do					
3	Sample geometry substrate G_i and description					
	\mathcal{D}_i from image-description pair sets \mathcal{S}					
4	Refresh \mathcal{I} using plot function: $\mathcal{I} \leftarrow f(\mathcal{I}, \mathcal{G}_i)$					
5	Refresh corresponding description:					
	$\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_i$					
6	Refresh vertex coordinate: $C \leftarrow C \cup C_i$					
7	end					
8	for $j \leftarrow 1$ to k do					
9	Select line drawing position \mathcal{P}_j					
10	Draw line and label length: $\mathcal{I} \leftarrow f(\mathcal{I}, \mathcal{P}_j)$					
11	Refresh corresponding description:					
	$\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{P}_j$					
12	Refresh vertex coordinate: $C \leftarrow C \cup C_i$					
13	if $j = k$ then					
14	Calculate all angle information \mathcal{R}					
15	Draw angles and label degrees:					
	$\mathcal{I} \leftarrow f(\mathcal{I}, \mathcal{R})$					
16	Refresh corresponding description:					
	$ \mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{R}$					
17	end					
18	end					
19	Refresh Geometric Properties: $\mathcal{T} \leftarrow \mathcal{V}(C)$					
20	20 Produce single-step reasoning result r_c using prompt					
	$P_s: r_c \leftarrow \mathcal{M}(\mathcal{D}, P_s)$					
21	Generate answer A_e and its corresponding question					
	Q_e using prompt P_q : $A_e, Q_e \leftarrow \mathcal{M}(r_c, P_q)$					
22	22 Filtering for correct answer A and its corresponding					
	question Q using prompt P_e :					
$A, Q \leftarrow \mathcal{M}(A_e, \bar{Q}_e, \bar{T}, P_e)$						
23	Return: $\mathcal{I}, \mathcal{D}, \mathcal{Q}, \mathcal{A}$					

B Details of prompt in TR-Reasoner

We used ERNIE Bot 4.0 to implement TR-Reasoner. We describe the prompts used in TR-Reasoner, including the prompts for the Description Patch Reasoning Fusion (Fig. 9), the Reverse Question Generation (Fig. 10), and the Error A&Q Filtering (Fig. 11). In these figures, the texts in blue represent the Task Description, while the texts in orange represent the input information. Each prompt includes three contextual examples, and we show only one of them, with the remaining examples replaced by ellipses. In addition to the examples, some prompts also include an instruction section



Figure 9: The prompt of the Description Patch Reasoning Fusion.

that specifies more detailed requirements, some incorporate additional basic knowledge, and others outline more specific goals that must be achieved. 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

In preliminary experiments, we observed that LLMs often failed to accurately interpret certain geometric relationships. To systematically identify such issues, we selected 50 representative instances per geometric substrate from the TR-GeoMM dataset and applied the TR-Reasoner framework for Patch Reasoning. We analyzed the most frequently misinterpreted relationships and formalized their correct representations into a base knowledge library. During formal generation, prompts are dynamically constructed by retrieving relevant geometric relationships from this library based on the target substrate.

C More information of TR-GeoMM

Through the TR-CoT, we construct a high-quality geometric dataset, TR-GeoMM. In Fig. 12, we provide a detailed overview of specific cases from TR-GeoMM. These cases demonstrate the variety of mathematical geometry question types covered by TR-GeoMM, including solving for lengths, angles, areas, and geometry elemental relations. Each of these categories is critical for improving the geometric reasoning ability of LMMs.

For Cosine distance based data diversity, we first randomly sample 5000 instances from each

799



Figure 10: The prompt of the Reverse Question Generation.

Error A&Q Filtering Prompt
All the correct answers are given in Ground Truth, if the result obtained in Answer contradicts the Ground Truth, return Score: 0, otherwise return Score: 1
Junnet.
Input: Ground Truth: Length of segment $OA = 2.00$. Length of segment $OB = 2.00$. Length of segment $AC = 1.04$. Length of segment $BC = 4.00$. Length of segment $BD = 1.04$. Angle OAC measures 75.0 degrees. Angle AOB measures 150.0 degrees. Angle ACB measures 75.0 degrees. Angle OBD measures 75.0 degrees. Angle CBD measures 75.0 degrees. Question: There is a circle O with a radius of 2. The central angle AOB measures 150 degrees. Line BD is parallel to line CA. What is the relationship between angle DBC and angle ACB?Answer: Since line BD is parallel to line CA, according to the properties of parallel lines, the corresponding angles are equal. Therefore, angle DBC is equal to angle
ACB.
Output: Analysis: From the ground truth we can get angle ACB measures 75.0 degrees and angle CBD measures 75.0 degrees. So angle ACB = angle DBC.Score: 1
Input: [GROUND TRUTH AND ANSWER]
What must be achieved:
 When Ground Truth is the same as the content of the Answer, Score: 1. When Ground Truth contradicts the content of the Answer, Score: 0. When Ground Truth is not related to the content of the Answer, Score: 4. You don't need to focus on the solution process, as long as the result agrees with the Ground Truth, return Score: 1.

Figure 11: The prompt of the Error A&Q Filtering.

dataset(MAVIS, GeomVerse, and TR-GeoMM), then we encode the instances into embedding features using pretrained BERT model (Devlin, 2018). Finally, we calculate the average cosine distance of each dataset using the BERT output features. Higher distance score indicates better diversity, and our TR-GeoMM has the highest distance score

841

845



Figure 12: Examples of TR-GeoMM dataset.

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

among the three datasets.

We further report the computation and generation cost of the generation pipeline here. Geometric images, descriptions, and properties are all generated simultaneously by the TR-Engine, our Python-built graphics rendering engine. A total of 15,000 geometric images were generated during the geometric image generation phase, and it could be done in approximately 10 minutes by a 12th-gen Intel Core i7-12700. TR-Reasoner, our LLM-based module for Patch Reasoning, Q&A generation, and error filtering, runs via cloud APIs with 32 parallel processes. These stages take approximately 6 hours in total (2 hours per stage). As a one-time cost, the generated data can be reused across model training tasks.

D Qualitative examples of filtered errors

Fig. 13 presents four representative types of errors identified by the Error A&Q Filtering module. Theorem Violation refers to cases where conclusions or assumptions contradict established mathematical theorems. Metric Discrepancies involve inconsistencies between the given numerical values or angles and the geometric properties. Diagramtext Mismatches occur when elements described

Table 5: Statistic comparison between Geo170K, Geom-Verse, and our data. En. Exis means Enhance Existing data. Fully syn. means Fully synthesized. '/'indicates the same number as GeoQA

Dataset name	Data type	Img.	Q&A	Theorem
Geo170K	En. Exis.	6.4K	110K	/
GeomVerse	Fully syn.	9.3K	9.3K	60
TR-GeoMM	Fully syn.	15K	45K	110
TR-GeoSup	En. Exis.	6.4K	20K	/

in the problem statement are either absent from the diagram or inconsistent with it. Ambiguous Answerability denotes problems in which the information provided is insufficient to derive a unique solution, or essential data is not explicitly stated in the question.

871

874

875

876

877

878

879

884

886

890

894

901

902

903

904

905

906

907

909

E Examples of TR-GeoSup dataset

Fig. 14 illustrates an example from the TR-GeoSup dataset, showcasing the transformation of a multistep reasoning problem from the original GeoQA dataset. In the original Q&A pair, the reasoning process is condensed and lacks explicit intermediate steps, relying on implicit knowledge. TR-GeoSup decomposes the original reasoning process into three hierarchical sub-questions, each accompanied by a detailed and theorem-aware reasoning chain. This augmentation not only clarifies the implicit knowledge embedded in the original data but also provides a step-by-step guide for model training.

F Statistical Comparison With Related Datasets

Here we make a brief comparison between our proposed dataset with some related academic datasets.

- GeomVerse is a representative template-based method that generates geometrically oversimplified images by combining predefined polygons in fixed configurations. These images only contain polygon compositions and lack theorem-aware elements (e.g., midlines and angle bisectors). It has 9.3k synthetic images accompanied by Q&A pairs, but their richness was limited by the absence of theorem-aware elements, covering only 60 theorems.
- Geo170K represents an augmented version of the existing GeoQA dataset. It primarily focuses on rephrasing Q&A pairs, such as altering wording, swapping conditions and answers, or scaling numerical values while keeping the underlying

theorems identical. This approach does not enhance the diversity of theorems covered in the dataset.

910

911

912

913

914

915

916

917

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

• TR-CoT enables theorem-driven multimodal reasoning by designing substrates and embedding theorem-aware elements based on theorem conditions, allowing generated images to support complex Q&A construction. Unlike prior approaches, TR-CoT is not limited by existing data coverage and can expand a model's geometric knowledge. The framework supports both new data synthesis and augmentation of existing datasets, and covers 110 theorems through its structured generation process.

As shown in Tab. 5, our data possess a notable diversity in theorem coverage, image distribution, and Q&A quantity. Ablation studies in Section 4.2 further discuss the training effectiveness of our proposed data.

G Detail of polygon distribution

We conducted robustness experiments for different polygon distributions, where the details of the polygon distributions are shown in Tab. 6. From top to bottom, the percentage of triangles and quads gradually decreases, and the percentage of pentagons and hexagons gradually increases. There is also a clear difference in the percentage of circles.

Similar quantitative results within 0.6% in Tab. 7 show that the impact of polygon distributions is almost negligible, demonstrating the strong robustness of our method to different polygon distributions. Therefore, the performance gain is mainly attributed to the diverse geometry representation and reasoning knowledge provided by our method.

Table 6: Details of polygon distribution for distribu-tional robust ablation studies.

Method	Polygon Distribution				
Wiethou	triangle	quad	circle	pentagon	hexagon
Group I	29%	46%	17%	5%	3%
Group II	32%	40%	14%	8%	6%
Group III	25%	33%	21%	12%	8%

Table 7: Ablation study on the robustness to polygonal distributions.

Polygon Distribution	MathVista	GeoQA
Group I	64.4	54.0
Group II	64.4	53.7
Group III	63.9	53.4



Figure 13: Examples of filtered errors.



Figure 14: Examples of TR-GeoSup dataset.

H The Case of Direct Generation and TR-Reasoner Generation

945

946The core idea of the TR-Reasoner is to improve947the accuracy of Q&A pairs by simplifying the rea-948soning based on descriptions and then generating949corresponding questions from the answers in a re-950versed manner. A straightforward approach is di-951rectly prompting ERNIE Bot 4.0 to generate Q&A952pairs from the input image description. However,953as shown on the left of Fig. 16, this approach often954fails to determine the correct answer. In contrast,955the Q&A pairs produced by TR-Reasoner are cor-956rect for all three instances with our design.

I Details of the theorems

The support of mathematical theorems is crucial for the accuracy of TR-Engine. In Tab. 9, we present the geometric theorems and properties that we used. These define the rules for combining elements, establishing a logically coherent chain throughout the figure construction process. They serve as the foundation for extending reasoning scenarios and also assist in the computation and verification of question-answer pairs. 957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

We collect and organize geometric theorems through three main approaches: (1) Systematic **Textbook Mining:** We analyzed standard textbooks and online educational resources to compile core geometric axioms and theorems from primary and secondary school mathematics curricula in Mainland China. (2) Alignment with Public Academic Datasets: We extracted theorems referenced in public academic datasets (e.g., PGPS9K, MAVIS, GeomVerse) to ensure consistency with commonly used training corpora. (3) Expert Consultation: We consulted primary and secondary school educators to identify important theorems and conclusions grounded in real-world teaching practices.

J Effectiveness of TR-CoT

As shown in Fig. 8, models jointly trained on Geo170K and TR-CoT-generated data (TR-GeoMM and TR-GeoSup) consistently outperform those trained solely on Geo170K ('Geo-'). InternVL2.5-8B receives a 1.5% improvement on MathVista and GeoQA, and Qwen2.5-VL-7B

Quality Judgment Prompt				
You are provided with a language model's response to a geometric question. Your mission is to judge the quality of the response based on the following				
standards, and give a score between 0 to 10.				
Judging Standards:				
1.Logic consistency. Assess whether the response is self-consistent,				
logically coherent, and free from contradictions or illogical reasoning.				
2. Clarity. Evaluate whether the response is clear and easy to understand,				
avoiding ambiguity or vague expressions,				
3. Output format: Score: your score(from 0 to 10)				

Figure 15: Comparison of model problem solving before and after training.

improves by 1.0% and 2.0% on MathVista and GeoQA, respectively. These results indicate that TR-CoT-generated data can supplement existing datasets and is widely effective in various LMMs.

Table 8: TR-CoT generated data effectiveness validation on different models. 'Geo-' indicates the model is finetuned only with geometric instruction data of Geo170K. Consistent and significant improvement without adding any additional parameters.

Model	MathVista	GeoQA
Geo-InternVL-2.0-2B	51.9	62.5
TR-CoT-InternVL-2.0-2B	56.3 (4.4^)	63.4 (0.9↑)
Geo-LLaVA-1.5-7B	27.9	47.6
TR-CoT-LLaVA-7B	29.3 (1.4↑)	51.7 (4.1↑)
Geo-Qwen2-VL-7B	59.9	69.1
TR-CoT-Qwen2-VL-7B	67.6 (7.7↑)	70.4 (1.3↑)
Geo-InternVL-2.0-8B	70.2	74.9
TR-CoT-InternVL-2.0-8B	72.1 (1.9↑)	76.7 (1.8↑)
Geo-InternVL-2.5-8B	76.4	75.2
TR-CoT-InternVL-2.5-8B	77.9 (1.5↑)	76.7 (1.5↑)
Geo-Qwen2.5-VL-7B	73.5	77.2
TR-CoT-Qwen2.5-VL-7B	74.5 (1.0↑)	79.2 (2.0↑)

K Details of CoT quality evaluation

We used ERNIE Bot 4.0 and DeepSeek R1 to evaluate model outputs. For each response, the evaluation model gives a score between 0 and 10 to judge the logical consistency, clarity, and lack of ambiguity. We use the average score of the two models as the final score. To ensure more accurate evaluation, we include specific judging standards. The prompts used are shown in Fig. 15. The blue part represents the Task Description.

989 990 991



Figure 16: The Case of Direct Generation and TR-Reasoner Generation.

Catagory	Droportion	Critorio
Category Denallal Lines	Company and a sough Alternate inte	Equal company on alogy Symptomer
Parallel Lilles	corresponding angles equal; Alternate inte-	tary appaautive angles, Supplement
	gles supplementary	gles; Parallel to the same line
General Triangles	Interior angles sum to 180°	AA similarity; SSS/SAS/ASA/AAS/HL congruence
Isosceles Trian-	Equal base angles: Three-line coincidence	Two equal angles : Two equal sides
gles	(angle bisector, median, altitude) :Base an-	
0	gles are 45° in right-isosceles case	
Equilateral Trian- gles	All angles are 60° ; Three - line coincidence	Three equal sides ; Three equal angles ; Isosceles triangle with a 60° angle
Right Triangles	Acute angles are complementary ; Side op-	Contains a right angle ; HL congruence for
0 0	posite 30° angle is half of the hypotenuse ;	right - triangles
	Median on the hypotenuse is half of the hy-	
	potenuse ; Pythagorean theorem: $a^2 + b^2 =$	
	c^2	
Angle Bisector	Points on the perpendicular bisector are	A ray that divides an angle into two equal
	equidistant from the endpoints	parts
Triangle Midline	Parallel to the third side and half of its	Connects the mid-points of two sides
	length	
Parallelogram	Opposite sides are equal ; Diagonals bisect	Both pairs of opposite sides are parallel;
	each other ; Area = $base \times height$	Diagonals bisect each other; Opposite sides
		are equal
Rectangle	All angles are 90° ; Diagonals are equal	A parallelogram with a right angle; A
		quadrilateral with three right angles
Rhombus	All sides are equal; Diagonals are perpen-	A parallelogram with adjacent sides equal;
	dicular to each other	A quadrilateral with four equal sides
Square	All sides and angles are equal; Diagonals are equal and perpendicular	Prove it is both a rectangle and a rhombus
Isosceles Trape-	Legs are equal; Base angles on the same	Two equal legs; Equal base angles on the
zoid	base are equal	same base
Trigonometric	$\sin 30^{\circ} = \frac{1}{2}$; $\sin 45^{\circ} = \frac{\sqrt{2}}{2}$; $\sin 60^{\circ} =$	/
Functions	$\frac{\sqrt{3}}{2}$; $\sin 90^\circ = 1$; $\cos 30^\circ = \frac{\sqrt{3}}{2}$;	
	$\cos 45^{\circ} = \frac{\sqrt{2}}{2}$; $\cos 60^{\circ} = \frac{1}{2}$; $\cos 90^{\circ} =$	
	$0 ; \tan 30^\circ = \frac{\sqrt{3}}{3} ; \tan 45^\circ = 1 ;$	
	$\tan 60^\circ = \sqrt{3}$	
Circle	The perpendicular bisector of a chord is per-	/
	pendicular to the chord; The perpendicular	
	bisector of a chord passes through the cen-	
	ter	
Central Angle	Equal central angles subtend equal chords	/
	and arcs	
Inscribed Angle	An inscribed angle is half of the central	/
	angle subtended by the same arc; An angle	
	subtended by a diameter is a right angle	
Cyclic Quadrilat-	Opposite angles are supplementary	/
eral		
Tangent	A tangent is perpendicular to the radius at	A line perpendicular to the radius at the
	the point of contact; Tangents from an ex-	endpoint on the circle is a tangent
	ternal point to a circle are equal in length	
Regular Polygon	For an equilateral triangle inscribed in a	/
	circle of radius R, side length $a = R\sqrt{3}$;	
	For a square inscribed in a circle of radius	
	R, side length $a = R\sqrt{2}$	

Table 9: Summary of Geometric Theorems and Properties