GeoDANO: Geometric VLM with Domain Agnostic Vision Encoder

Anonymous ACL submission

Abstract

We introduce GeoDANO, a geometric visionlanguage model (VLM) with a domain-agnostic vision encoder, for solving plane geometry problems. Although VLMs have been employed for solving geometry problems, their ability to recognize geometric features remains insufficiently analyzed. To address this gap, we propose a benchmark that evaluates the recognition of visual geometric features, including primitives such as dots and lines, and relations such as orthogonality. Our preliminary study shows that vision encoders often used in general-purpose VLMs, e.g., OpenCLIP, fail to detect these features and struggle to generalize across domains. We develop GeoCLIP, a CLIP-based model trained on synthetic geo-016 metric diagram-caption pairs to overcome the 017 limitation. Benchmark results show that Geo-CLIP outperforms existing vision encoders in recognizing geometric features. We then propose our VLM, GeoDANO, which augments 021 GeoCLIP with a domain adaptation strategy 022 for unseen diagram styles. GeoDANO outperforms specialized methods for plane geometry problems and GPT-40 on MathVerse.

1 Introduction

037

041

Large language models (LLMs) have achieved remarkable success in automated math problem solving, particularly through code-generation capabilities integrated with proof assistants (Moura and Ullrich, 2021; Nipkow et al., 2002; Chen et al., 2023; Wu et al., 2022; Hendrycks et al., 2021). Although LLMs excel at generating solution steps and correct answers in algebra and calculus (Zhou et al., 2024), their unimodal nature limits performance in plane geometry, where solution depends on both diagram and text (Zhou et al., 2024).

Specialized vision-language models (VLMs) have accordingly been developed for plane geometry problem solving (PGPS) (Chen et al., 2021, 2022; Lu et al., 2021; Zhang et al., 2023; Zhang



Figure 1: Examples of diagram-caption pairs and their solution steps written in formal languages from GeoQA and PGPS9k datasets. In the problem description, the visual geometric premises and numerical variables are highlighted in green and red, respectively. A significant difference in the style of the diagram and formal language can be observable.

and Moshfeghi, 2024; Li et al., 2024b; Xia et al., 2024). Yet, it remains unclear whether these models genuinely leverage diagrams or rely almost exclusively on textual features. This ambiguity arises because existing PGPS datasets typically embed sufficient geometric details within problem statements, potentially making the vision encoder unnecessary (Zhang and Moshfeghi, 2024). Fig. 1 illustrates example questions from GeoQA and PGPS9K, where solutions can be derived without referencing the diagrams.

We propose a new benchmark created via a synthetic data engine, which systematically evaluates the ability of VLM vision encoders to recognize geometric premises. Our empirical findings reveal that previously suggested self-supervised learning (SSL) approaches, e.g., vector quantized variataional auto-encoder (VQ-VAE) (Liang et al., 2023) and masked auto-encoder (MAE) (Ning

et al., 2023; Xia et al., 2024), and widely adopted encoders, e.g., OpenCLIP (Radford et al., 2021) and DinoV2 (Oquab et al., 2024), struggle to detect geometric features such as perpendicularity and degrees.

061

062

063

066

069

077

078

080

084

086

100

101

103

To this end, we propose GeoCLIP, a model pre-trained on a large corpus of synthetic diagram-caption pairs. By varying diagram styles (e.g., color, font size, resolution, line width), GeoCLIP learns robust geometric representations and outperforms prior SSL-based methods on our benchmark. Building on GeoCLIP, we introduce a few-shot domain adaptation technique that efficiently transfers the recognition ability to realworld diagrams. We further combine this domainadapted GeoCLIP with an LLM, forming a domainagnostic VLM for solving PGPS tasks in Math-Verse (Zhang et al., 2024a).

In our experiments on MathVerse (Zhang et al., 2024a), which encompasses diverse plane geometry tasks and diagram styles, our VLM with a domain-adapted GeoCLIP consistently outperforms both task-specific PGPS models and generalist VLMs. Ablation studies confirm the effectiveness of our domain adaptation strategy, showing improvements in optical character recognition (OCR)-based tasks and robust diagram embeddings across different styles.

We summarize the contributions as follows: We propose a novel benchmark for systematically assessing how well vision encoders recognize geometric premises in plane geometry diagrams (§3); We introduce GeoCLIP, a vision encoder capable of accurately detecting visual geometric premises (§4.1), and a few-shot domain adaptation technique that efficiently transfers this capability across different diagram styles (§4.2); We show that our VLM, incorporating domainadapted GeoCLIP, surpasses existing specialized PGPS VLMs and generalist VLMs on the Math-Verse benchmark (§5.2) and effectively interprets diverse diagram styles (§5.3).

2 Related Work

104In this section, we summarize the studies related to105the benchmarks proposed to evaluate plane geome-106try problem solving (PGPS), the models trained for107PGPS, and the contrastive learning methods used108to enhance PGPS performance.

2.1 PGPS benchmarks

Several studies have introduced benchmarks for PGPS, including a set of diagrams and corresponding problem and solution descriptions (Chen et al., 2021; Lu et al., 2021; Zhang et al., 2023; Chen et al., 2022). The problem and solution descriptions are provided in natural languages or formal languages. Often, the solution steps are provided in the form of formal language. Given the dataset, the goal of PGPS is to train a model that produces a valid solution as an executable program. 110

111

112

113

114

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

147

148

149

150

151

152

153

154

155

156

157

158

However, as recent work by Zhang et al. (2024a) shows, the problem description contains too much information such that the model produces a valid solution program without having the diagram information as shown in Fig. 1. MathVerse (Zhang et al., 2024a) introduces modifications to existing PGPS benchmarks by directly encoding the geometric properties and relations into the diagrams. Therefore, it is impossible to produce a valid solution without recognizing the necessary information from diagrams. Despite the effort, it is still unclear to what extent the vision encoder recognizes the geometric conditions in a diagram as models are evaluated in an end-to-end fashion.

2.2 Program generation based PGPS

A core challenge in program generation-based PGPS is processing both diagrams and text to interpret geometric premises. One approach tackles the challenge by converting a diagram into alternative representations such as lists of geometric primitives and relations that can be represented as text (Seo et al., 2015; Sachan et al., 2017; Lu et al., 2021; Zhang and Moshfeghi, 2024; Zhang et al., 2022; Peng et al., 2023). Although reducing the problem to a single modality can be effective, building such converters typically requires labeled diagrams, which are expensive to collect and eventually limit generalization across diverse diagram styles.

Another line of research typically employs vision-language models (VLMs), where a VLM comprises a vision encoder and a language model (Zhang et al., 2023; Chen et al., 2021; Cao and Xiao, 2022; Ning et al., 2023; Chen et al., 2022; Liang et al., 2023; Xia et al., 2024; Li et al., 2024b). The vision encoder produces a visual embedding from the diagram, and the language model then generates solution steps in an autoregressive manner, conditioned on the textual description and the visual embedding. While the VLMs apply to various

diagram formats, the visual geometric premises per-159 ception of the VLMs remains underexplored due 160 to the abundance of textual information in existing benchmarks. Moreover, the VLMs are often 162 fine-tuned and tested on a single benchmark, leav-163 ing their domain generalization capabilities across 164 different diagram styles unexamined.

Contrastive learning in PGPS 2.3

161

165

166

188

190

191

195

197

199

201

202

206

Contrastive learning is applied in diverse domains 167 such as computer vision (Schroff et al., 2015) and 168 natural language processing (Gao et al., 2021). 169 In the context of PGPS, contrastive learning is 170 employed to address domain-specific challenges. 171 GeoX (Xia et al., 2024) applies contrastive learn-172 ing to the adapter layer of the VLM to enhance 173 formal language comprehension. Other approaches 174 train the vision encoder itself using the contrastive 175 language-image pre-training (CLIP) (Radford et al., 176 2021) objective: LANS (Li et al., 2024b) aligns 177 patch embeddings from a vision Transformer (ViT) 178 with text token embeddings if they describe the 179 same point, and MAVIS (Zhang et al., 2024b) employs diagram-caption pairs generated by a syn-181 thetic engine for CLIP. In this work, we examine how CLIP with varied caption styles influences the 183 visual geometric premises recognition of the vision encoder. In addition, a contrastive learning framework is introduced to strengthen robustness against 186 domain shifts in the styles of diagrams.

Visual Geometric Premises Recognition 3 **Benchmark for Vision Encoders**

In this section, we first develop a benchmark for evaluating a vision encoder's performance in recognizing geometric features from a diagram. We then report the performance of well-known vision encoders on this benchmark.

3.1 Benchmark preparation

We design our benchmark as simple classification tasks. By investigating PGPS datasets, we identify that recognizing geometric primitives, such as points and lines, and geometric properties representing relations between primitives, such as perpendicularity, is important for solving plane geometry problems. Recognized information forms geometric premises to solve the problem successfully. To this end, we carefully curate five classification tasks as follows:

• Concyclic: A circle and four points are given.

The task is to identify how many of those points lie on the circle.

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

- TwoLines: Two lines, AB and BC, are given alongside other geometric objects. The task is to determine whether AB and BC are perpendicular, collinear, or neither.
- **ObjectShape**: A given diagram includes one of the following geometric objects: a segment, triangle, square, or pentagon. The task is to classify which object is present.
- SquareShape: A diagram including a square ABCD and other geometric objects is given. The task is to classify whether the square is a trapezoid, parallelogram, or rectangle.
- AngleDetection: A diagram is given with at least three points: A, B, and C. The task is to classify the correct angle of ABC from $\{15^{\circ}, 20^{\circ}, \ldots, 75^{\circ}\}.$

An example of each task is provided in Fig. 2.

Our benchmark is built on top of AlphaGeometry (Trinh et al., 2024), which is designed to solve IMO-style plane geometry problems. The program provides useful functions such as formal language describing plane diagrams. The language predefines a set of geometric premises listed in Table 5, including all necessary properties to define our benchmark tasks. In addition, once a diagram description is given in formal language, the program renders a corresponding diagram with varying fonts, colors, widths, orientations, and resolutions, allowing us to have diagrams with diverse styles often observed in a real-world scenario.

We create question-and-answer pairs based on To sample a diverse set of AlphaGeometry. question-and-answers, we first establish a foundational geometric structure corresponding to the key problem of the task and then repeatedly add new points or lines with randomly selected geometric relationships to the existing diagram with the help of the formal language. The pseudo-code for the random question generation is presented in Algorithm 1. For each task, we generate 50,000, 10,000, and 10,000 question-and-answer pairs for training, validation, and testing, respectively.

3.2 Results

With the proposed benchmark, we evaluate four widely adopted vision encoders for the opensourced VLMs: OpenCLIP (Radford et al., 2021),



Figure 2: Illustration of the proposed visual feature perception benchmark. We introduce five different diagram classification tasks that require visual feature perception to answer geometry-related questions.

	Models	Object Shape	Con cyclic	Two Lines	Square Shape	Angle Detection
Baseline	OpenCLIP	100.00	99.13	86.57	85.20	64.81
	SigLIP	100.00	99.71	89.26	89.31	76.86
	DinoV2	100.00	98.01	85.30	91.24	22.43
	ConvNeXT	100.00	99.20	89.39	88.13	61.84
TSS	Jigsaw	86.11	63.85	49.98	61.88	11.44
	MAE	93.99	72.25	71.73	82.70	13.08
	VQ-VAE	63.05	60.97	48.10	57.35	9.22
GeoCLIP	$\begin{array}{l} \text{GeoCLIP}\left(F\times\right)\\ \text{GeoCLIP}\left(2K\right)\\ \text{GeoCLIP}\end{array}$	99.52 99.32 99.21	98.61 98.73 99.24	88.33 94.73 96.05	86.76 89.22 95.95	65.68 74.95 78.56

Table 1: Results on the proposed visual feature benchmark. We report the test accuracy of the models with the best validation performance.

SigLIP (Zhai et al., 2023), DinoV2 (Oquab et al., 2024), and ConvNeXT (Liu et al., 2022).

To evaluate the vision encoder, we adopt a linear probing approach. Specifically, we add a linear layer on top of each encoder as a prediction head and train the linear layer from scratch while freezing the parameters of the vision encoder. We use a training set to train the prediction head and report the test accuracy with the best validation performance. The details for the hyper-parameters are described in Appendix B.1.

261

262

265

271

272

273

275

As shown in Table 1, many existing vision encoders well recognize the shape of objects but fail to recognize the correct angle between two lines. The encoders also show some difficulties in recognizing the shape of a square and the relationship between two lines. Although the result may seem satisfactory at a glance, these errors will propagate to the downstream tasks when combined with LLMs. Hence, it is important to improve the recognition performance of the vision encoder further.

4 Improving the Vision Encoder Geometric Premises Recognition

In this section, we first propose GeoCLIP, a new vision encoder designed to recognize geometric premises from diverse styles of diagrams. To transfer the recognition to real-world PGPS benchmarks, we then propose a domain adaptation technique for GeoCLIP that leverages a small set of diagram–caption pairs from the target domains.

276

277

278

279

281

283

287

289

290

291

292

293

294

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

4.1 GeoCLIP

To make a vision encoder recognize geometric diagrams better, we propose a GeoCLIP, a vision encoder trained with CLIP objective with a newly developed 200,000 diagram-caption examples. From the random diagram generator developed in §3.1, we additionally sample 200,000 diagrams written in the formal language. Directly rendering these samples can result in a diagram that may not preserve the geometric properties. For example, the perpendicularity between two lines cannot be observed from the diagram without having the right angle sign, i.e., \models . Therefore, we ensure to render the images containing all necessary geometric premises from its visual illustration.

For the caption of a diagram, we filter out some geometric properties from the original description of a diagram used to render the image. Specifically, we only keep the following four properties, concyclic, perpendicularity, angle measures, and length measures, from the visual premises shown in Table 5. After that, we convert the remaining descriptions written in the formal language into natural language. We filter out some properties for two reasons. First, some properties are not recognizable from the rendered diagram without additional

407

408

information, e.g., congruency. These properties are listed as non-visual premises in Table 5. Second, collinearity and parallelity occur so frequently that they can marginalize others. Some examples of generated captions after filtering and translating are provided in the right-most column of Fig. 4. We call the filtered caption as *GeoCLIP-style caption*.

311

312

313

314

317

319

321

323

325

331

335

337

340

341

343

347

354

358

With this dataset, we fine-tune OpenCLIP (Radford et al., 2021) according to the CLIP objective which is formulated as:

$$\mathcal{L}_{\text{CLIP}}(\mathcal{D}, g, h) := \\ \mathbb{E}_{\mathcal{D}} \bigg[-\log \frac{\exp(g(D_i)^T h(X_i)/\tau)}{\sum_{X \in \{X_i\}_i} \exp(g(D_i)^T h(X)/\tau)} \bigg],$$
(1)

where $\mathcal{D} := \{(D_i, X_i)\}_{i=1}^N$ is the diagram-caption pairs, g is the vision encoder, h is the text encoder, and τ is a temperature parameter. We named the resulting vision encoder as GeoCLIP. Appendix B.1 provides the details, including hyper-parameters.

We compare the performance of GeoCLIP to other self-supervised approaches trained with the same dataset. We test three self-supervised approaches: Jigsaw (Chen et al., 2021; Cao and Xiao, 2022), MAE (Ning et al., 2023; Xia et al., 2024), and VQ-VAE (Liang et al., 2023) used in previous work to improve the recognition performance of plane diagrams. We use the same architecture used for GeoCLIP for Jigsaw and MAE with the hyper-parameters used in the previous works. For VQ-VAE, we follow the architecture of Liang et al. (2023). All model performances are measured through the linear probing used in §3.2.

As shown in Table 1, GeoCLIP recognizes geometric features better than existing baselines and self-supervised methods. The self-supervised approaches generally perform poorly for the benchmark, justifying the choice of the objective. We also compare the performance of GeoCLIP against other encoders such as OpenCLIP. Note that although we outperform the other encoders in difficult tasks such as SquareShape and AngleDetection, these results might be *unfair* since the training set of GeoCLIP is similar to the diagrams in the benchmark. The t-SNE plots of the embeddings from the vision encoders are illustrated at Fig. 5.

We further ablate the filtering process in Geo-CLIP. To this end, we compare GeoCLIP with its two variants: *GeoCLIP* ($F \times$), which uses the captions generated without filtering. We also test *Geo-CLIP* (2K), which is trained on only 2,000 pairs, to see the effectiveness of the large-scale dataset. The results in Table 1 imply both the filtering and the training set size matter in enhancing geometric properties recognition.

4.2 Domain adaptation of GeoCLIP

Although GeoCLIP enhances the geometric premises recognition on the benchmark set, the diagram styles in existing PGPS benchmarks differ, necessitating further adaptation. To overcome this challenge, we propose a domain adaptation method for GeoCLIP. To this end, we propose a few-shot domain adaptation method utilizing a few labeled diagrams.

A domain-agnostic vision encoder must match the same diagrams drawn in different styles. To do so, we need a target domain diagram translated into the source domain style or the source diagrams translated into the target domain style. With these translated images, we can guide the model to focus on key geometric information instead of irrelevant attributes, such as color and font family. However, in practice, it is difficult to obtain the same diagrams with different styles.

We develop a way to translate the target diagrams into source style. Thankfully, since wellknown PGPS datasets come with diagram captions written in formal languages (Lu et al., 2021), we can easily convert them to the AlphaGeometrystyle descriptions. Given the translated descriptions, we utilize the rendering engine of Alpha-Geometry to translate the target domain images into the source domain. With the translation, we can generate the same diagram in the source domain style. Fig. 6 provides examples of the diagram pairs with different styles. However, in some cases, the original description contains geometric premises that are unrecognizable from the diagram, such as $\angle ACB = 35.0$ in Fig. 1a. Therefore, we apply the same filtering process used in GeoCLIP to translate the AlphaGeometry-style descriptions into natural languages.

Formally, let $\mathcal{D}_S := \{(D_S^{(i)}, X_S^{(i)})\}_{i=1}^{N_S}$ be the diagram-caption pairs from source domain S, e.g., the synthetic diagrams, and let $\mathcal{D}_{T_j} := \{(D_{T_j}^{(i)}, X_{T_j}^{(i)})\}_{i=1}^{N_{T_j}}$ be the set of diagram-caption pairs of target domain T_j , e.g., the PGPS benchmarks. With the translation process described above, we can synthesize a style-transferred diagram-caption pair $(\hat{D}_{T_j}^{(i)}, \hat{X}_{T_j}^{(i)})$ for each diagram $D_{T_j}^{(i)}$ and caption $X_{T_j}^{(i)}$ in target domain T_j . We perform domain adaptation by fine-tuning the vision encoder through the style-transferred diagram-caption pairs. Let $\hat{\mathcal{D}}_{T_j}$ be a collection of the original diagram and style-transferred captions, i.e., $\hat{\mathcal{D}}_{T_j} = \{(D_{T_j}^{(i)}, \hat{X}_{T_j}^{(i)})\}_{i=1}^{N_{T_j}}$, and let $\hat{\mathcal{D}}_{T_jS}$ be a collection of the original and style transferred diagram pairs, i.e., $\hat{\mathcal{D}}_{T_jS} = \{(D_{T_j}^{(i)}, \hat{D}_{T_j}^{(i)})\}_{i=1}^{N_{T_j}}$. The cross-domain adaptation objective is written as

$$\mathcal{L}_{\text{CLIP-DA}}(\mathcal{D}_S, \{\mathcal{D}_{T_j}\}_j, g, h) := \mathcal{L}_{\text{CLIP}}(\mathcal{D}_S, g, h) + \Sigma_j \mathcal{L}_{\text{CLIP}}(\hat{\mathcal{D}}_{T_j}, g, h) + \mathcal{L}_{\text{CLIP}}(\hat{\mathcal{D}}_{T_jS}, g, g), \quad (2)$$

where g and h are the vision and text encoders of GeoCLIP, respectively. Note that we do not use the original captions from the target domain, since our goal is to adapt the vision encoder to the target domain, not the text encoder.

5 Experiments

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

In this section, we evaluate the PGPS performance of our VLM equipped with the domain adapted GeoCLIP on MathVerse (Zhang et al., 2024a). We compare its performance against established PGPS baselines. We also present ablation studies highlighting our VLM's strong visual feature recognition and resilience to domain shifts, both of which are facilitated by the adapted vision encoder.

5.1 Experimental settings and training details

Datasets. We use MathVerse (Zhang et al., 434 2024a) to measure the performance of VLMs. 435 MathVerse is a benchmark designed to evaluate 436 both the reasoning and visual-feature recognition 437 438 capabilities of VLMs, covering plane geometry, solid geometry, and function problems. 439 It is constructed by compiling problems from various 440 sources, including Geometry3K (Lu et al., 2021), 441 GeoQA (Chen et al., 2021), and GEOS (Seo et al., 442 2015). Each problem is presented in five variants: 443 text-dominant, which provides all essential tex-444 tual information for solving the problem; *text-lite*, 445 which omits descriptive details, e.g., object shapes, 446 from the text; vision-intensive, which removes cer-447 tain textual conditions that can be inferred from 448 remaining information; vision-dominant, which re-449 locates numerical measurements, such as angles 450 451 and lengths, from the text to the diagram; and vision-only, which offers only the diagram as input, 452 embedding all text within the diagram. In the fol-453 lowing experiments, we focus on plane geometry 454 problems and exclude the vision-only task. 455

Training details. We describe the construction of our **geo**metric VLM with **d**omain-**a**gnostic visio**n** encoder, named GeoDANO. Based on GeoCLIP developed in §4.1, we apply the domain adaptation to GeoQA and Geometry3K datasets. For the domain adaptation, we randomly sample 50 diagrams and translate the diagram and caption styles following the procedure described in §4.2. Finally, GeoCLIP is fine-tuned via Eq. (2). We name the GeoQA and Geometry3K adapted GeoCLIP as GeoCLIP-DA.

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

505

We combine LLama-3-8b-Instruct (Dubey et al., 2024) and GeoCLIP-DA to construct a VLM. The combined model is then fine-tuned again with the training set of GeoQA and PGPS9K to predict the solution program. For PGPS9K, we use the Geometry3K split. While previous works focusing on PGPS do not consider optical character recognition (OCR) from diagrams since the benchmark datasets, GeoQA and PGPS9K, provide necessary details in problem descriptions, numerical values can appear within diagrams in real-world settings. Therefore, we fine-tune GeoDANO with additional OCR capability by modifying the problem statements. Additional details about the modification process with hyper-parameter configurations can be found in Appendix D.

In addition, we unify the programming languages used in the solution programs of GeoQA and PGPS9K by converting GeoQA language into PGPS9K format. The unification makes the output of VLM consistent since both datasets use different types of formal languages.

Baselines. We use two different types of baseline models for the experiments: PGPS *specialist VLMs* and *generalist VLMs*. Specialist VLMs produce a solution program as an output of a given problem, and generalist VLMs produce a natural language solution as an output.

For the specialist VLMs, we test PGP-SNet (Zhang et al., 2023), NGS (Chen et al., 2021), SCA-GPS (Ning et al., 2023), GeoFormer (Chen et al., 2022), UniMath-Flan-T5 (Liang et al., 2023), and GeoX (Xia et al., 2024). For GeoX, we use the two variants GeoX-Geo3K and GeoX-GeoQA, which are fine-tuned on Geometry3K and GeoQA, respectively.

For the generalist VLMs, we test two GPT-40 variants (Hurst et al., 2024): gpt-40-2024-11-20 and gpt-40-mini-2024-07-18, and the InternVL2.5 variants: 8B and 26B models (Chen et al., 2024).

Models	Text Dominant		Text Lite		Vision Intensive		Vision Dominant	
models	Completion \uparrow	Top-10 ↑	Completion ↑	Top-10 ↑	Completion ↑	Top-10 ↑	Completion \uparrow	Top-10 ↑
PGPSNet	4.37	14.55	2.08	12.06	2.08	11.02	-	-
NGS	6.45	34.57	6.64	28.52	5.86	26.37	-	-
SCA-GPS	6.84	18.16	5.66	16.80	3.52	15.23	-	-
GeoFormer	16.22	32.85	16.84	30.77	13.10	29.11	-	-
UniMath-Flan-T5	17.88	32.43	16.42	30.56	13.93	28.27	-	-
GeoX-Geo3K	5.41	9.98	4.16	6.86	3.53	5.61	-	-
GeoX-GeoQA	24.32	37.42	17.26	32.43	13.51	16.25	-	-
GeoDANO (OC)	19.13	40.12	16.63	34.72	13.31	31.81	1.25	8.12
GeoDANO (GC)	20.37	41.79	18.09	38.25	15.80	35.34	5.62	19.38
GeoDANO (GC-D)	22.66	43.45	21.00	38.46	18.30	35.76	6.67	20.42
GeoDANO	23.70	47.82	21.21	45.11	18.09	42.20	12.08	36.04

Table 2: PGPS accuracy on MathVerse benchmark. We compare the performance of GeoDANO against PGPS specialist models, which generate a solution program as an output. GeoDANO-OC, -GC, and -GCD are three variants of our model with different encoders. Further details about these variants can be found in §5.3.

Evaluation metric. For each plane geometry problem, both the specialist VLMs and GeoDANO generate 10 outputs via beam search. Following Zhang et al. (2023), we then use completion accuracy and top-10 accuracy as our primary evaluation metrics. The completion accuracy assesses whether the first successfully executed solution from the beam is correct; the solutions are reviewed in beam order, and success is recorded if the first executable solution produces the correct answer. Top-10 accuracy examines all ten beam outputs, counting a success if any of these solutions yield the correct result upon execution. Note that, as described before, the specialist VLMs do not have OCR capability. For the evaluation, we feed the correct values to the outputs of these models by using the parser developed in Zhang et al. (2023). For the models that are trained in Chinese, i.e., NGS and SCA-GPS, we use problem descriptions translated by GPT-40 (Hurst et al., 2024).

> To measure the performance of the generalist VLMs, we use multiple-choice questions instead of open-ended questions due to the difficulty in parsing the final answer from free-form text. We use the multiple-choice question provided in Math-Verse as an additional input to each problem. We ask VLMs to produce the answer in a pre-specified form. We report the top-1 accuracy of these models. To compare GeoDANO against the generalist models, we use the same protocol used in Zhang et al. (2023) to measure the accuracy.

5.2 Results

Performance against specialist VLMs. In Table 2, GeoDANO shows the best performance in almost all the problem variants and metrics except the completion accuracy in the text-dominant task.

	Text Dominant	Text Lite	Vision Intensive	Vision Dominant
GPT-40	40.35	39.18	38.01	36.95
GPT-4o-mini	41.12	39.53	35.59	30.50
InternVL2.5-8B	38.30	36.26	35.09	21.99
InternVL2.5-26B	42.40	40.06	38.01	38.71
GeoX-GeoQA	52.05	45.91	37.43	-
GeoDANO	48.54	49.71	41.81	39.30

Table 3: Comparison between GeoDANO and generalist VLMs on multiple choice questions. We report the accuracy of GeoDANO and GeoX following the same evaluation protocol suggested in Zhang et al. (2023).

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

Note that the specialist models cannot solve the vision-dominant problems since these problems do not contain variables representing numerical values, such as a length, in the problem description. When comparing the performance between text and vision-dominant tasks, the top-10 accuracy of Geo-DANO on vision-dominant task is higher than the top-10 accuracy of the specialist models on text-dominant task except for GeoX-GeoQA. Given that the two tasks use the same problem set, the result implies that GeoDANO performs better than the specialist models without having the geometric premises in the problem description. In other words, our vision encoder can extract geometric premises accurately from the visual information.

Performance against generalist VLMs. Table 3 reports the performance of generalist VLMs and GeoDANO on multiple choice questions. Geo-DANO outperforms proprietary closed models, i.e., GPT-40 variants, and open-sourced models, i.e., the InternVL2.5 variants. Especially, the performance gap between GeoDANO and InternVL2.5-26B reflects the parameter efficiency of our VLM. While GeoDANO shows impressive results among the variants, the performance of GeoX-GeoQA

532

534

536

537

538

539

540

541

506

507

508

509

510

Models	PGPS9K			GeoQA		
	$MR\downarrow$	$mAP\uparrow$	_	$MR\downarrow$	$mAP\uparrow$	
OpenCLIP	50.50	27.87		111.70	1.29	
GeoCLIP	88.99	17.61		128.73	1.05	
GeoCLIP-D	58.83	13.35		107.25	2.86	
GeoCLIP-DA	12.88	41.13		35.60	33.25	

Table 4: Domain adaptation analysis. We report the mean rank (MR) and mean average precision (mAP) of the test diagrams.

degrades dramatically as the visual information moves from the text to the diagram. Our work is the first to show that the specialist can compete with the generalist in MathVerse.

5.3 Ablation studies

567

568

569

570

571

572

573

577

578

579

580

582

584

586

591

592

593

598

599

603

Variation of GeoCLIP. We perform a detailed empirical analysis to evaluate how effectively the GeoCLIP-style captions and the proposed domain adaptation technique improve GeoDANO's performance. Specifically, we compare GeoDANO against other VLMs trained on the GeoCLIP variants, including OpenCLIP (Radford et al., 2021) and the GeoCLIP without domain adaptation. We also test a variant of GeoCLIP trained with additional diagram-caption pairs from the target domains without having any filtering process. In this case, we utilize all the data in the training sets.

We show the experimental result in Table 2. GeoDANO-OC and GeoDANO-GC represent the VLM with OpenCLIP and GeoCLIP without domain adaptation, respectively. GeoDANO-GCD represents the GeoCLIP with additional unfiltered domain captions. GeoDANO outperforms other variants on most tasks, except the completion accuracy on the vision-intensive task.

OCR performance. We assess the accuracy of GeoDANO and its variants in OCR on the Math-Verse diagrams, focusing on the vision-dominant task. We evaluate the OCR performance of the first executable solution program in top-10 VLM predictions. GeoDANO-OC, GeoDANO-GC, GeoDANO-GCD, and GeoDANO achieve 1.84%, 20.26%, 13.95%, and 46.58% accuracy, respectively. The result explains the accuracy improvement of GeoDANO in the vision-dominant task against other variants.

Domain adaptation analysis. We examine how effectively GeoCLIP-DA generalizes to new domains with different diagram styles. For this experiment, we compare the embedding similarity



Figure 3: Visualization of OpenCLIP and GeoCLIP-DA embeddings. The orange, green, and blue dots represent PGPS9K, GeoQA, and synthetic diagrams, respectively. In the top row, the three diagrams on the left and right are those with the highest cosine similarities to the center under OpenCLIP and GeoCLIP-DA, respectively.

between two diagrams representing the same structure in different styles. To create the paired dataset, we use a similar process described in §4.2. Specifically, a total of 100 diagrams are sampled from the test sets of GeoQA and PGPS9K, and these samples are rendered in AlphaGeometry style through the diagram description.

For evaluation, we sample 100 diagrams from each of the target domain's training sets and compare the similarity against the original diagram via cosine similarity. We also compute the similarity between the style transferred diagram and the original diagram. We report two metrics for test diagrams: the mean rank (MR) and the mean average precision (mAP) of the style-transferred diagram.

As reported in Table 4, GeoCLIP-DA produces similar embeddings for structurally equivalent diagrams, regardless of their stylistic differences. Fig. 3 visualizes the diagram embeddings of Open-CLIP and GeoCLIP-DA. As one can observe, the OpenCLIP embeddings are largely separated by the domain of the diagrams, whereas those of GeoCLIP-DA appear to capture and align with underlying visual features more effectively.

6 Conclusion

In this work, we propose a domain-agnostic PGPS method, GeoDANO, by implementing a synthetic data engine and proposing a contrastive learning framework with domain adaptation. We demonstrate the effectiveness of GeoDANO in visual feature perception at both VLM and vision encoder levels by evaluating on the MathVerse and through a newly proposed geometric feature recognition benchmark for vision encoders. Eventually, the reasoning ability in plane geometry problems is enhanced with the improved perceptual capabilities.

642

607

608

609

610

611

643 Limitations

In this work, we present a domain-agnostic VLM for PGPS by refining the vision encoder. Although our VLM exhibits strong performance in recognizing visual features, its coverage remains limited to geometric premises. Building on the success of the synthetic data engine and contrastive learning, extending this combination to different kinds of visual features, e.g., sub-structures in molecular graphs (Kamoi et al., 2024), statistics from charts (Masry et al., 2022), and solid geometry, promises further improvements in recognition of VLM. Due to the limitations in the experimental environment, we are unable to test LLMs with more than 30B parameters.

References

659

666

667

671

672

673

674

675

678

679

685

690

691

696

- Jie Cao and Jing Xiao. 2022. An augmented benchmark dataset for geometric question answering through dual parallel text encoding. In *Proceedings of the* 29th International Conference on Computational Linguistics, pages 1511–1520, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- 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. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3313–3323, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- 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, Online. Association for Computational Linguistics.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2023. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *Transactions on Machine Learning Research*.
- Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. 2024. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple contrastive learning of sentence embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6894–6910, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. 697

698

699

700

701

704

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the MATH dataset. In *Thirtyfifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2).*
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.
- Ryo Kamoi, Yusen Zhang, Sarkar Snigdha Sarathi Das, Ranran Haoran Zhang, and Rui Zhang. 2024. Visonlyqa: Large vision language models still struggle with visual perception of geometric information. *arXiv preprint arXiv:2412.00947*.
- Diederik P Kingma. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. 2024a. Llavaonevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*.
- Zhong-Zhi Li, Ming-Liang Zhang, Fei Yin, and Cheng-Lin Liu. 2024b. LANS: A layout-aware neural solver for plane geometry problem. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 2596–2608, Bangkok, Thailand. Association for Computational Linguistics.
- Zhenwen Liang, Tianyu Yang, Jipeng Zhang, and Xiangliang Zhang. 2023. UniMath: A foundational and multimodal mathematical reasoner. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7126–7133, Singapore. Association for Computational Linguistics.
- Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. 2022.
 A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11976–11986.
- Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In *International Conference on Learning Representations*.

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. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 6774–6786, Online. Association for Computational Linguistics.

753

754

763

767

771

773

774

775

779

781

783

784

785

790

793

794

795

796

799

801

802

804

805

- Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*.
- Leonardo de Moura and Sebastian Ullrich. 2021. The lean 4 theorem prover and programming language.
 In Automated Deduction – CADE 28: 28th International Conference on Automated Deduction, Virtual Event, July 12–15, 2021, Proceedings, page 625–635, Berlin, Heidelberg. Springer-Verlag.
- Maizhen Ning, Qiu-Feng Wang, Kaizhu Huang, and Xiaowei Huang. 2023. A symbolic characters aware model for solving geometry problems. In Proceedings of the 31st ACM International Conference on Multimedia, MM '23, page 7767–7775, New York, NY, USA. Association for Computing Machinery.
- Tobias Nipkow, Markus Wenzel, and Lawrence C. Paulson. 2002. *Isabelle/HOL: a proof assistant for higher-order logic*. Springer-Verlag, Berlin, Heidelberg.
- Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, Mido Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. 2024. DINOv2: Learning robust visual features without supervision. *Transactions on Machine Learning Research*. Featured Certification.
- Shuai Peng, Di Fu, Yijun Liang, Liangcai Gao, and Zhi Tang. 2023. GeoDRL: A self-learning framework for geometry problem solving using reinforcement learning in deductive reasoning. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 13468–13480, Toronto, Canada. Association for Computational Linguistics.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR.

Mrinmaya Sachan, Kumar Dubey, and Eric Xing. 2017. From textbooks to knowledge: A case study in harvesting axiomatic knowledge from textbooks to solve geometry problems. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 773–784, Copenhagen, Denmark. Association for Computational Linguistics. 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

855

856

857

858

859

860

861

862

863

864

- Florian Schroff, Dmitry Kalenichenko, and James Philbin. 2015. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823.
- 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.
- Trieu H. Trinh, Yuhuai Wu, Quoc V. Le, He He, and Thang Luong. 2024. Solving olympiad geometry without human demonstrations. *Nature*, 625(7995):476–482.
- Yuhuai Wu, Albert Qiaochu Jiang, Wenda Li, Markus Norman Rabe, Charles E Staats, Mateja Jamnik, and Christian Szegedy. 2022. Autoformalization with large language models. In Advances in Neural Information Processing Systems.
- Renqiu Xia, Mingsheng Li, Hancheng Ye, Wenjie Wu, Hongbin Zhou, Jiakang Yuan, Tianshuo Peng, Xinyu Cai, Xiangchao Yan, Bin Wang, et al. 2024. Geox: Geometric problem solving through unified formalized vision-language pre-training. *arXiv preprint arXiv:2412.11863.*
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11975–11986.
- Jiaxin Zhang and Yashar Moshfeghi. 2024. GOLD: Geometry problem solver with natural language description. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 263–278, Mexico City, Mexico. Association for Computational Linguistics.
- Ming-Liang Zhang, Fei Yin, Yi-Han Hao, and Cheng-Lin Liu. 2022. Plane geometry diagram parsing.
 In Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22, pages 1636–1643. International Joint Conferences on Artificial Intelligence Organization. Main Track.
- Ming-Liang Zhang, Fei Yin, and Cheng-Lin Liu. 2023. A multi-modal neural geometric solver with textual clauses parsed from diagram. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence*, IJCAI '23.

866 Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun 867 Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, 868 Kai-Wei Chang, Yu Qiao, Peng Gao, and Hongsheng 869 Li. 2024a. Mathverse: Does your multi-modal llm 870 truly see the diagrams in visual math problems? In Computer Vision - ECCV 2024: 18th European Con-871 ference, Milan, Italy, September 29-October 4, 2024, 872 Proceedings, Part VIII, page 169-186, Berlin, Hei-873 874 delberg. Springer-Verlag.

875

876

877

878

879

880

881 882

883

884

885

- Renrui Zhang, Xinyu Wei, Dongzhi Jiang, Ziyu Guo, Shicheng Li, Yichi Zhang, Chengzhuo Tong, Jiaming Liu, Aojun Zhou, Bin Wei, et al. 2024b. Mavis: Mathematical visual instruction tuning with an automatic data engine. arXiv preprint arXiv:2407.08739.
- Aojun Zhou, Ke Wang, Zimu Lu, Weikang Shi, Sichun Luo, Zipeng Qin, Shaoqing Lu, Anya Jia, Linqi Song, Mingjie Zhan, and Hongsheng Li. 2024. Solving challenging math word problems using GPT-4 code interpreter with code-based self-verification. In *The Twelfth International Conference on Learning Representations*.

Visual premises	Non-visual premises			
 Perpendicularity Collinearity Concyclicity Parallelity Angle measure Length measure 	 Middle point Congruency in degree Congruency in length Congruency in ratio Triangle similarity Triangle congruency Circumcenter 			
	• Foot			

Table 5: Geometric premises used in AlphaGeometry. *Visual premises* denotes the geometric premises which can be directly perceived from the diagram. *Non-visual premises* requires reasoning to be recognized.

Algorithm 1 Sampling process of the synthetic data engine

Input Geometric relations R, geometric objects O, number of clauses n_c

Output AlphaGeometry program c

1: Initialize points and clauses with the sampled object: $P, C \sim O$

2: for $i \leftarrow 1$ to n_c do

- 3: Generate points: P_{new}
- 4: Sample relation and points: $r, P_{old} \sim R, P$
- 5: Construct clause: $C_{\text{new}} = r(P_{\text{new}}, P_{\text{old}})$
- 6: Update points and clauses: $P, C \leftarrow P \cup P_{\text{new}}, C \cup C_{\text{new}}$
- 7: Generate program with points and clauses: $c \leftarrow \text{Clauses2Program}(P, C)$
- 8: **return** *c*

Appendix

887

892

895

896

900

901

902

A Synthetic Data Engine

In this section, we provide the details of our synthetic data engine. Based on AlphaGeometry (Trinh et al., 2024), we generate synthetic diagram and caption pairs by randomly sampling a AlphaGeometry program with Algorithm 1.

Examples for randomly sampled AlphaGeometry problems and their corresponding diagrams and lists of geometric premises are described in Fig. 4. The types of geometric premises that appear in our synthetic data engine are listed in Table 5.

B Details of Benchmark

B.1 Training details

To evaluate the visual feature perception of the vision encoder, we utilize a linear probing approach, which involves freezing the vision encoder parameters and training a simple linear classifier on top of its features. 903

904

905

906

907

908

909

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

944

945

946

947

948

949

We train the linear classifier on the training set of each task for 50 epochs with batch size 128 and learning rate 1e-4. We use Adam optimizer for optimization.

B.2 Visualization of the vision encoders

We visualize the embeddings of the vision encoders used in §3.2 at Fig. 5.

C GeoCLIP-DA

C.1 Domain adaptation data

We adopt GeoCLIP to the two PGPS benchmarks: GeoQA (Chen et al., 2021) and PGPS9K (Zhang et al., 2023). For PGPS9K, we use the Geometry3K split. Fig. 6 shows the pairs used to adapt the domain of GeoCLIP.

C.2 Training details

We start from OpenCLIP (Radford et al., 2021), a pre-trained model where the architecture is ViT-L/14 with image resolution 336×336 . To train OpenCLIP, we use total of 200,000 diagramcaption pairs generated with our synthetic data engine. For the domain adaptation to GeoQA and Geometry3K datasets, we randomly sample 50 diagrams and translate the diagram and caption styles following the procedure described in §4.2. Finally, GeoCLIP is fine-tuned via Eq. (2). We name the GeoQA and Geometry3K adopted GeoCLIP as GeoCLIP-DA.

We set the batch size for the source domain diagram-caption pairs to 256. For the domain adaptation parts, i.e., applying CLIP on the diagram-caption pairs and the diagram pairs of target domains, we vary the batch size to 32. We set weight decay to 0.2. We optimize for 50 epochs using Adam optimizer (Kingma, 2014) and a cosine annealing scheduler with 2,000 warmup steps and the maximum learning rate is set to be 1e-4. We train the model with eight RTX3090 GPUs for approximately 24 hours.

D GeoDANO

D.1 Modification of training data

Our fine-tuning strategy differs slightly from previous works (Chen et al., 2022; Zhang et al., 2023; Xia et al., 2024). Here, we clarify the difference between our approach and previous approaches.



Figure 4: Example of randomly sampled AlphaGeometry problems. For each row, the first element describes the randomly sampled AlphaGeometry problem and the others are the geometric premises, diagram, and GeoCLIP-style caption that can be obtained from the AlphaGeometry problem. Note that the GeoCLIP-style caption can be obtained by filtering certain geometric properties, e.g., angle measure, perpendicularity, and concyclicity, from the geometric premises.



Figure 5: The embeddings of the vision encoders on the diagrams of TwoLines task. We visualize the embeddings of the vision encoders on the diagrams of TwoLines task. The blue, orange, and green dots are the diagrams where the two lines AB and BC are collinear, perpendicular, and otherwise, respectively.



Figure 6: Examples of diagram pairs curated for domain adaptation. For each row, the first diagram is from the target domain, and the remaining diagrams are from the source domain. To generate source domain diagrams, we translate the target diagram by our diagram generator with the textual description of the target image.

In previous works, the VLM is trained to produce the solution program given diagram and problem description as shown in Fig. 1. An interesting observation from GeoQA and PGPS9K datasets is that the numerical measurements, such as angles, lengths, and volumes, are not written in the problem description but given as additional conditions, and the numerals are substituted as a variable in the problem description as shown in Fig. 1a. Therefore, the VLM only needs to produce the solution program without having optical character recognition (OCR) from the diagram. The variables are automatically substituted by the actual numbers when the program is executed. Therefore, the vision encoders do not need to learn OCR from the image.

950

951

952

953

954

955

956

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

However, this approach cannot be generalized to a wider class of problems where the numerals are embedded in the diagram instead of written in the problem description. Some variants of MathVerse, such as the vision-dominant problems, fall into this category as well. To incorporate OCR into the solution of the problem, we modify some problem statements in the training set, such that the numerical measurements are only shown in the diagram and not in the statements. We further modify the solution problem so that the solution contains OCR results as a part of the final output. Finally, we unify the language of the solution programs used in GeoQA and PGPS9K by converting GeoQA programs into PGPS9K format. The unification makes the output of VLM consistent since both datasets use different types of formal languages.

Fig. 7 shows examples of the modified input pairs and solutions, where the first problem state-



Figure 7: Examples of the training data for GeoDANO. While previous PGPS models require the only to predict the solution steps and assume the numerical values are explicitly given, GeoDANO is trained to predict both the solution steps and the numerical values in the diagram and text.

ment does not have numerical measurements and the OCR results are in the part of the output solution program.

D.2 Training details

We begin by summarizing the architecture of our VLM, a combination of a vision encoder and a language model. For the vision encoder, we use GeoCLIP-DA, with a two-layer MLP of GeLU activation as the projection layers following LLaVA-OneVision (Li et al., 2024a). For the language model, we employ LLama-3-8B-Instruct (Dubey et al., 2024). For a given diagram and question pair in PGPS, we feed the vision encoder with the given diagram, and then the output of the encoder is used as an input token of LLM through the projection layer. The question text is then fed into the LLM, followed by the diagram embedding.

With the modified training data, we apply supervised fine-tuning on the VLM, i.e., the gradient only flows through the prediction of numerical values and solution steps, not the diagram and text.

We train the VLM with AdamW optimizer (Loshchilov and Hutter, 2019) and cosine annealing scheduler with warmp up ratio 0.03 and maximum learning rate 1e-5. We use LoRA (Hu et al., 2022) with rank 128. We set the batch size to 16 and train with 5 epochs. We train the VLM with four A100-80GB GPUs for approximately 24 hours.

985

987

988

989

991

993

994

995

997

999

1000

1001

1002

1003

1005

1006

1007

1008

1009

1010

1011

1012