MATHGLM-VISION: SOLVING MATHEMATICAL PROB-LEMS WITH MULTI-MODAL LARGE LANGUAGE MODEL

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language models (LLMs) have demonstrated significant capabilities in mathematical reasoning, particularly with text-based mathematical problems. However, current multi-modal large language models (MLLMs), especially those specialized in mathematics, tend to focus predominantly on solving geometric problems but ignore the diversity of visual information available in other areas of mathematics. Moreover, the geometric information for these specialized mathematical MLLMs is derived from several public datasets, which are typically limited in diversity and complexity. To address these limitations, we aim to construct a fine-tuning dataset named MathVL, and develop a series of specialized mathematical MLLMs termed MathGLM-Vision by conducting Supervised Fine-Tuning (SFT) on MathVL with various parameter-scale backbones. To extensively evaluate the effectiveness of MathGLM-Vision, we conduct experiments on several public benchmarks and our curated MathVL-test benchmark consisting of 2,000 problems. Experimental results demonstrate that MathGLM-Vision achieves significant improvements compared with some existing models, including backbone models and open-source mathematical MLLMs. These findings indicate the importance of diversity dataset in enhancing the mathematical reasoning abilities of MLLMs. Both MathGLM-Vision model (based on CogVLM2, GLM-4V-9B) and MathVL-test will be open-sourced.

028 029

031

004

010 011

012

013

014

015

016

017

018

019

020

021

024

025

026

027

1 INTRODUCTION

032 Recent advancements in computational linguistics have led to substantial progress in solving mathe-033 matical problems using Large Language Models (LLMs) with multi-step reasoning processes (Light-034 man et al., 2023). For example, models like GPT-4 (Achiam et al., 2023), Qwen (Bai et al., 2023a), GLM-4 (Team et al., 2024), LLaMA (Touvron et al., 2023a;b) have demonstrated impressive performance on mathematical datasets such as GSM8K (Cobbe et al., 2021) and MATH (Hendrycks 037 et al., 2021). Furthermore, the development of specialized mathematical models is expanding the 038 potential of LLMs in this domain. These models, specifically designed for mathematical problem solving, include notable contributions such as WizardMath (Luo et al., 2023), MAmmoTH (Yue et al., 2023), MathCoder (Wang et al., 2023a), MetaMath (Yu et al., 2023), DeepSeekMath (Shao et al., 040 2024), and others (Yang et al., 2023; Yuan et al., 2023; Gou et al., 2023; Yue et al., 2024b; Mitra 041 et al., 2024; Ying et al., 2024). These advancements highlight the growing proficiency of LLMs in 042 handling intricate mathematical reasoning and problem-solving tasks. 043

Despite significant advancements, the majority of models designed for mathematical problem solving
 still rely predominately on textual representations. This limits their effectiveness in scenarios that
 require visual information. Notably, approximately 63% of mathematics questions in Chinese K12 ed ucation include visual elements, highlighting the critical role of visual information in comprehending
 and solving mathematical problems.

Therefore, a crucial question arises: Is visual information essential for solving these mathematical
 problems that include visual elements? To verify this, we conduct a series of insightful experiments
 comparing the performance of these models such as GPT-40, Claude-3.5-Sonnet, Qwen-VL-Max,
 and Gemini-1.5-Pro on MathVL-test dataset, both with and without visual inputs. As shown in
 Figure 1, the results clearly demonstrate that the inclusion of visual elements significantly enhances
 the models' ability to accurately solve complex mathematical problems. Conversely, the exclusion of

066

067

068

081

082

084

085

087

088

090

091

092



Figure 1: Insight experiments demonstrates the significance of visual information in solving mathematical problems. (Left) A performance comparison of different models with and without visual inputs on MathVL-test dataset. (Right) The accuracies of MathGLM-Vision on MathVL-test with and without visual inputs.

visual information leads to a pronounced decrease in performance, emphasizing the essential role
 that visual context plays in solving mathematical problems that incorporate visual elements.

Currently, multi-modal large language models (MLLMs) are at the forefront of efforts to integrate 072 visual and textual information for solving mathematical problems. Close-source models such as 073 GPT-4V (OpenAI, 2023), Gemini (Team et al., 2023), Claude3 (Anthropic, 2024), Qwen-VL (Bai 074 et al., 2023b), along with several open-source MLLMs like CogVLM (Wang et al., 2023c; Hong 075 et al., 2024), MiniGPT Zhu et al. (2023), LLaVA-1.5 Liu et al. (2024a), SPHINX-MoE Gao et al. 076 (2024), and LLaVA-NeXT Liu et al., demonstrate substantial potential in addressing geometric 077 reasoning challenges. Additionally, specialized geometric MLLMs like G-LLaVA Gao et al. (2023a), 078 GeoGPT4V Cai et al. (2024) and Math-LLaVA Shi et al. (2024) are particularly focused on enhancing 079 capabilities in this domain. However, these models still face several challenges and limitations that need to be addressed. 080

- Current MLLMs, particularly those specialized in mathematics, predominantly focus on solving geometric problems and tend to overlook the diversity of visual information in mathematics. This visual information encompasses a broad spectrum of elements, including arithmetic, statistics, algebra and word problems, each integral to different mathematical domains beyond geometry.
 - Current fine-tuning dataset for specialized mathematical MLLMs, typically sourced from public datasets like GeoQA and Geometry3K, often lack diversity and complexity. This limitation restricts the models' ability to effectively solve a broader range of mathematical problems.

• Current specialized mathematical MLLMs are predominantly designed to process single-image inputs and lack the capability to handle multiple images simultaneously. This limitation hampers their ability to tackle complex problems that necessitate the integration of information from multiple visual sources.

In response to these challenges and limitations, we construct a fine-tuning dataset named MathVL, 094 which encompasses both open-source data and our specially curated Chinese data collected from 095 K12 education. The MathVL dataset is meticulously designed to incorporate a diverse range of 096 mathematical problems, consisting of textual and visual inputs. For textual information, the MathVL dataset covers a variety of mathematical subjects such as arithmetic, algebra, geometry, statistics, and 098 word problems. It includes various types of questions, including fill-in-the-blank, multiple-choice, 099 and free-form. For visual information, the MathVL dataset involves elements like functions, statistical data, graphs, charts, LaTeX expressions, and geometric figures, providing a comprehensive resource 100 for complex mathematical problem solving. 101

With our constructed MathVL dataset, we develop a series of specialized mathematical MLLMs, collectively referred to as MathGLM-Vision, with different parameter scales. Specifically, MathGLM-Vision-9B, MathGLM-Vision-19B and MathGLM-Vision-32B are fine-tuned on three backbone models: GLM-4V-9B, CogVLM2, and CogVLM-32B, respectively. Moreover, we establish a benchmark dataset named MathVL-test, which contains 2,000 problems designed to evaluate the ability of MathGLM-Vision and other MLLMs in solving mathematical problems involving multiple images. Through extensive evaluation experiments on three public benchmark datasets and one curated

108 80 80 MathGLM-Vision-32B Claude-3.5-Sonnet MathGLM-Vision-32B MathGLM-Vision-98 Claude3.5-sonnet 109 MathGLM-Vision-19B Gemini-1.5-pro MathGLM-Vision-19B GPT-40 Gemini-1.5-pro MathGLM-Vision-9B Claude-3-Opus 110 70 GPT-40 Owen-VL-Max 70 111 (%) (%) 112 Accuracy (Accuracy (99 113 114 115 **⊢** 50 50 Top. -doj 116 117 40 40 118 119 30 30

Figure 2: Performance comparison of the different multi-modal large language models. (Left) The accuracies of MathGLM-Vision and other MLLMs among three evaluation datasets. (Right) The accuracy of MathGLM-Vision and other MLLMs on MathVL-test across different categories.

MathVL-test, we validate the effectiveness of our MathGLM-Vision. The results in Figure 2 demonstrate that MathGLM-Vision exhibits superior performance in understanding and solving complex 126 mathematical problems with visual elements compared to existing MLLMs. For instance, on the geometry problem solving (GPS) minitest split of MathVista (Lu et al., 2023), MathGLM-Vision-9B 128 achieves a 39.68% relative improvement for GLM-4V-9B, MathGLM-Vision-19B achieves a 65.06% relative improvement for CogVLM2, and MathGLM-Vision-32B achieves a 51.05% relative improvement over CogVLM-32B. Last but not least, Both MathGLM-Vision model (based on CogVLM2, GLM-4V-9B) and MathVL-test will be open-sourced to facilitate the future development of this field.

We highlight our contributions as follows:

121

122

123 124 125

127

129

130

131

132

133 134

135

136 137

138

139

140

141

142 143 144

145

- **Data Perspective:** We construct MathVL, a diverse and comprehensive multi-modal mathematical supervised fine-tuning dataset that contains both textual and visual inputs.
- Model Perspective: We develop a suite of specialized mathematical multi-modal large language models, referred to as MathGLM-Vision, which demonstrates significant improvements on various mathematical benchmarks while maintaining general vision-language understanding capabilities.
- Benchmark Perspective: We establish a benchmark dataset called MathVL-test, which designed to evaluate the mathematical reasoning abilities of MLLMs using a multi-image format.
- 2 MATHVL: DATASET CURATION

146 To enhance the capabilities of MLLMs in solving mathematical problems, previous efforts (Chen 147 et al., 2021; 2022; Cao & Xiao, 2022; Gao et al., 2023a) focus on constructing high-quality datasets. 148 Nevertheless, the majority of these datasets fall into the category of Visual Question Answering 149 (VQA), which generally involves descriptive or identification tasks rather than conventional math-150 ematical problems. Furthermore, the answers in some public datasets like Geometry3K (Lu et al., 151 2021), GeoGPT4V (Cai et al., 2024), MathV360K (Shi et al., 2024) for standard mathematical 152 questions are often too simplistic, usually providing only the final answer without the intermediate 153 steps necessary for a thorough understanding. It is well-established that including step-by-step solutions can significantly enhance the reasoning capabilities of large language models (Wei et al., 154 2022; Lightman et al., 2023; Zhang et al., 2023; Wang et al., 2023b). Figure 3 demonstrates the 155 distribution of answer lengths in current open-source mathematical datasets. 156

157 To address these issues, we construct a fine-tuning dataset MathVL, including both several public 158 datasets and our curated Chinese dataset collected from K12 education levels. This dataset is meticulously crafted to encompass a diverse array of mathematical problems that incorporate visual 159 information. Each problem is presented with detailed step-by-step solutions, aiming to enhance 160 the problem-solving skills of MLLMs by providing them with both the context and the procedural 161 knowledge necessary for effective reasoning and comprehension.



Figure 3: Analysis of answer lengths in several open-source mathematical datasets like MathV360K, GeoGPT4V, and Geometry3K.

Open-Source Data. We first collect open-source datasets from GeoQA+ (Cao & Xiao, 2022), 176 Geometry3K (Lu et al., 2021), ChartQA (Masry et al., 2022), and UniGEO-Calculation (Chen et al., 177 2022). These datasets commonly serve as seed data for constructing enhanced datasets. Through 178 observation and statistical analysis, we discover that 57% of the answers within these datasets are 179 comprised of fewer than 50 words, indicating that many questions are answered directly without 180 elaboration or explanation. To enrich these dataset with comprehensive step-by-step solutions, we 181 employ GPT-40 to generate the detailed solutions for each question, thereby enhancing the learning 182 and reasoning potential of these datasets. After generating the detailed answers, we perform a rigorous 183 judgement process to ensure the accuracy of the solutions provided by GPT-40. Additionally, we adopt a public instruction tuning dataset named Geo170K (Gao et al., 2023a), which is constructed 185 using GeoQA+ and Geometry3K as seed data and contains more than 110K geometric questionanswer pairs. We also incorporate another public dataset, GeomVerse (Kazemi et al., 2023), as part of our resources. In the end, the detailed statistics of the open-source datasets used in MathGLM-Vision 187 is provided in Table 1. 188

Datasets ChartQA	UniGeo-Calculation	Geometry3K	GeoQA+	Geo170K	GeomVerse	ALL
Samples 7,398	3,499	2,101	6,026	117,205	9,339	145,568

Table 1: The detailed statistics of the open-source datasets used in MathGLM-Vision.

194 **Chinese Data Collected from K12 Education.** We construct a dataset specifically focused on 195 K12 education, comprising 341,346 mathematical problems with textual and visual inputs. This 196 dataset is meticulously curated to encompass a board range of mathematical topics and difficulty levels tailored to the Chinese educational curriculum. It features various question types, such as 197 multiple-choice, fill-in-the-blank, and free-form questions, spanning disciplines including arithmetic, algebra, geometry, statistics, and word problems. Mathematically, this dataset can be represented as 199 $D_{\text{MathVL}}^{\text{zh}} = \{Q, A, I_s\}$, where Q represents the question, A represents the answer, and I_s represents 200 one or more images associated with each question. To build this dataset, we first process the images 201 by adding a white border around each image and enhancing their resolution to ensure that MLLMs 202 can effectively recognize and interpret these images. This modification is crucial for facilitating 203 the accurate extraction of visual information. Next, we extract 341,346 samples from a raw dataset 204 containing 685,670 samples by implementing a selective filtering process. This selection is based 205 on two specific criteria: (1) filtering out samples where the answer includes images or the question 206 is incomplete, and (2) eliminating samples with answer that are fewer than 50 words in length to 207 ensure the responses are sufficiently detailed for model training. After constructing this dataset, 208 we categorize and analyze it based on mathematical topics associated with each question. Detailed statistics about the distribution of these categories are presented in Table 2. Figure 4 demonstrates 209 some examples sampled from the constructed Chinese dataset, providing a visual representation of 210 the mathematical topics of questions included. More dataset cases are provided in Appendix A. 211

212

173

174 175

213	Types	Arithmetic	Geometry	Algebra	Statistics	Word Problems	ALL
214 215	Samples	7,207	291,879	20,111	18,284	3,865	341,346

Table 2: Detailed statistics regarding the distribution used in MathGLM-Vision.



Figure 4: Examples sampled from the constructed Chinese dataset.

MATHGLM-VISION: MODEL TRAINING 3

248 249

250 251

253

254

256

257

258

Model Architecture. We employ CogVLM2 (Hong et al., 2024; Wang et al., 2023c) and GLM-4V-9B (GLM et al., 2024) architectures as our backbone models, and conduct Supervised Fine-Tuning (SFT) on our constructed MathVL dataset. Specifically, we utilize three pre-trained multi-modal large language models (MLLMs) for the fine-tuning process: GLM-4V-9B, CogVLM2-19B, and CogVLM-32B. This results in the development of three distinct variants of MathGLM-Vision, designated as MathGLM-Vision-9B, MathGLM-Vision-19B, and MathGLM-Vision-32B, respectively. Further details about the abovementioned three pre-trained MLLMs are available in Appendix B.

259 **Model Training.** To maintain the general vision-language understanding skills of MathGLM-260 Vision, we incorporate 19 open-source visual question-answering datasets (VQA datasets) into the 261 MathVL dataset. More details about the task type and visual context of VQA datasets are provided in 262 Appendix C. These datasets are meticulously selected to challenge and enhance the model's ability to interpret and integrate visual and textual information, ensuring it retains a broad understanding across 264 various contexts. By merging these varied sources, we enhance MathGLM-Vision's specialized 265 capabilities for mathematical problem-solving and simultaneously preserve its robustness in general 266 vision-language tasks. In the end, we conduct supervised fine-tuning (SFT) across the combined VQA and MathVL datasets. The training process undergoes 35,000 iterations with a learning rate 267 of 1e-5 and a batch size of 128. To ensure the stability of the training, we activate the visual 268 encoder's parameters and adjust its learning rate to be one-tenth of that used for the remaining training parameters. The details of the SFT procedures are described in Appendix D.

²⁷⁰ 4 EXPERIMENTS

272 4.1 EXPERIMENTAL SETUP273

Evaluation Datasets. We assess our MathGLM-Vision using three well-established public benchmark datasets (MathVista (Lu et al., 2023), MathVerse (Zhang et al., 2024), and MATH-Vision (Wang et al., 2024) datasets) alongside our specially curated dataset MathVL-test benchmark. This benchmark comprises 2,000 sampled cases, distinct from those in the MathVL dataset, ensuring a rigorous and unbiased evaluation of MathGLM-Vision's capabilities. Additionally, we evaluate MathGLM-Vision's general vision-language understanding skills using the MMMU benchmark Yue et al. (2024a).
Detailed descriptions for these benchmark datasets are provided in Appendix E.

281 Compared Models. We compare MathGLM-Vision with other Multi-Modal Large Language Models 282 (MLLMs), including closed-source MLLMs such as Gemini (Team et al., 2023), GPT-4V (OpenAI, 283 2023), Claude3 (Anthropic, 2024), and Qwen-VL (Bai et al., 2023b), and open-source MLLMs 284 like mPLUG-Owl (Ye et al., 2023), LLaMA-Adapter-V2 (Gao et al., 2023b), InstrctBLIP (Dai 285 et al., 2024), LLaVA-1.5 (Liu et al.), ShareGPT4V (Chen et al., 2023), SPHINX (Gao et al., 2024), 286 InternLM-XC2 (Dong et al., 2024), and InternVL (Chen et al., 2024). Additionally, we compare 287 MathGLM-Vision with recent specialized mathematical MLLMs, including G-LLaVA (Gao et al., 288 2023a), LLaVA-1.5-G (Cai et al., 2024), ShareGPT4V-G (Cai et al., 2024), and Math-LLaVA (Shi 289 et al., 2024). 290

Evaluation Metrics. We adopt top-1 accuracy to evaluate our MathGLM-Vision across MathVista GPS, MathVista, MathVerse, MATH-V, and MathVL-test benchmarks. Our evaluation process
 follows the pipeline outlined in the aforementioned benchmark datasets, which involves using LLMs
 to extract predicted answers from the model's responses. Accuracy is then calculated by comparing
 these extracted answers against the ground truths.

295

4.2 MAIN RESULTS

297 298

299 **Results on public benchmark datasets.** To comprehensively assess the ability of MathGLM-Vision in solving mathematical problems, we evaluate its performance against other MLLMs across several 300 public benchmark datasets. Table 3 demonstrates the overall results from these evaluations. The ex-301 perimental results indicate that our constructed MathVL dataset can significantly improve MathGLM-302 Vision's mathematical reasoning capabilities. For example, MathGLM-Vision-9B achieves a 64.42% 303 accuracy on the MathVista-GPS dataset, marking a substantial 39.68% improvement over its backbone 304 model, GLM-4V-9B. Besides, across various parameter scales, MathGLM-Vision consistently surpass 305 all backbone models on different evaluation benchmarks, highlighting the significant enhancements 306 that MathVL brings to the MathGLM-Vision's problem-solving skills. Notably, MathGLM-Vision 307 outperforms all open-source specialized mathematical MLLMs across various benchmarks. The 308 superior performance suggests that the high-quality and diverse data, complete with detailed step-309 by-step solutions, are crucial for improving MLLM's mathematical reasoning capabilities. More importantly, MathGLM-Vision-32B outperforms even the advanced GPT-4V on the more challenging 310 MATH-V benchmark, demonstrating its superior capacity to tackle complex mathematical problems. 311 Detailed experimental results on public benchmark datasets across different task types can be found 312 in Appendix F. 313

314 Results on MathVL-test. We also evaluate MathGLM-Vision and several close-source MLLMs 315 using our specially constructed MathVL-test benchmark. As depicted in Table 4, the results clearly 316 demonstrate that MathGLM-Vision significantly outperforms both its backbone models and other 317 leading closed-source MLLMs across various model sizes. Specifically, our MathGLM-Vision-32B 318 outperforms the advanced GPT-40 with a significant margin, achieving an accuracy of 59.00% 319 compared to GPT-40's 51.05%. Compared to the backbone model, GLM-4V-9B, MathGLM-Vision-320 9B achieves an impressive accuracy of 57.05% with a significant improvement of 86.5%. This 321 superior performance suggests that MathGLM-Vision, when conducting SFT on the MathVL dataset, notably enhances its capability to tackle complex Chinese mathematical problems. Additionally, we 322 report the accuracy across various categories, as illustrated in Figure 2 (See Right). MathGLM-Vision 323 significantly outperforms other advanced MLLMs in the domains of geometry and statistics. In

Model	Input	LLM	MathVista (GPS)	MathVista	MathVerse	MATH-V
		Closed Sour	rce Models			
Gemini Pro	Q, I	-	40.40	45.20	36.80	17.66
Gemini-1.5-Pro	Q, I	-	53.85	63.90	51.08	19.24
GPT-4V	Q, I	-	50.50	49.90	50.80	22.76
GPT-4-turbo	Q, I	-	58.25	58.10	43.50	30.26
GPT-40	Q, I	-	64.71	63.80	56.65	30.39
Claude3-Opus	Q, I	-	52.91	50.50	31.77	27.13
Claude3.5-Sonnet	Q, I	-	64.42	67.70	48.98	37.99
Qwen-VL-Plus	Q, I	-	33.01	43.30	19.10	10.72
Qwen-VL-Max	Q, I	-	46.12	51.00	35.90	15.59
		Open Source	ce Models			
General Multi-modal LLN	As					
mPLUG-Owl	Q, I	LLaMA-7B	23.60	22.20	12.47	9.84
LLaMA-Adapter-V2	Q, I	LLaMA-7B	25.50	23.90	4.50	9.44
InstructBLIP	Q, I	Vicuna-7B	20.70	25.30	15.36	10.12
LLaVA-1.5	Q, I	Vicuna-13B	24.04	27.60	12.70	11.12
ShareGPT4V	Q, I	Vicuna-13B	38.35	29.30	16.20	11.88
SPHINX-MoE	Q, I	Mixtral 8*7B	31.20	42.30	19.60	14.18
SPHINX-Plus	Q, I	LLaMA2-13B	16.40	36.70	14.70	9.70
InternLM-XC2	Q, I	InternLM2-7B	63.00	57.60	24.40	14.54
InternVL-1.2-Plus	Q, I	Nous-Hermes-2-Yi-34B	61.10	59.90	21.70	16.97
Geo-Multi-modal LLMs						
G-LLaVA	Q, I	LLaMA2-7B	53.40	28.46	12.70	12.07
G-LLaVA	Q, I	LLaMA2-13B	56.70	35.84	14.59	13.27
LLaVA-1.5-G	Q, I	Vicuna-7B	32.69	45.22	13.96	14.13
LLaVA-1.5-G	Q, I	Vicuna-13B	36.54	48.34	15.61	14.88
ShareGPT4V-G	Q, I	Vicuna-7B	32.69	45.07	16.24	12.86
ShareGPT4V-G	Q, I	Vicuna-13B	43.27	49.14	16.37	14.45
Math-LLaVA	Q, I	Vicuna-13B	57.70	46.60	19.04	15.69
		MathGLM-Vision an	d Backbone Models			
GLM-4V-9B	Q, I	GLM-4-9B	46.12	46.70	35.66	15.31
MathGLM-Vision-9B	Q, I	GLM-4-9B	64.42	52.20	44.20	19.18
CogVLM2	Q, I	LLaMA-3-8B	39.61	40.85	25.76	13.20
MathGLM-Vision-19B	Q, I	LLaMA-3-8B	65.38	61.10	42.50	21.64
CogVLM-32B	Q, I	GLM2-32B	41.06	40.04	35.28	19.32
MathGLM-Vision-32B	Q, I	GLM2-32B	62.02	62.40	49.20	26.51

Table 3: Results on several public benchmark datasets. Comparison of model performance on the
testmini set of MathVista and geometry problem solving (GPS) of MathVista. For MathVerse dataset,
results are evaluated on Vision Dominant with CoT-E. For MATH-V dataset, all 3,040 samples
included in the data are evaluated.

363
364 contrast, Claude3.5-Sonnet excels in algebra and arithmetic, demonstrating superior performance.
365 Meanwhile, MathGLM-Vision-19B ranks second in performance in the domain of arithmetic, showing
366 its strong abilities in this area as well. GPT-40 exhibits the highest performance in word problems
367 domain, while MathGLM-Vision also exhibits robust performance, surpassing both Gemini-1.5-Pro
and Claude3.5-Sonnet in this category.

368 369 370

4.3 GENERALIZABILITY OF MATHGLM-VISION

In addition to its proficiency in mathematical reasoning, we further assess MathGLM-Vision's
 capabilities in general vision-language understanding by conducting experiments on the MMMU
 benchmark. This benchmark is specifically designed to evaluate the ability of models to comprehend
 and process information across a variety of academic and professional disciplines, providing a
 comprehensive test of general vision-language understanding. Table 5 shows the performance of
 MathGLM-Vision, a specific variant fine-tuned exclusively on MathVL without the inclusion of
 VQA datasets, and backbone models. Compared to CogVLM2, MathGLM-Vision-19B achieves
 comparable performance in terms of generalizability, underscoring its capacity for simultaneous

Model	Input	LLM Size	MathVL-test
Gemini-1.5-Pro	Q, I	-	52.03
GPT-4V	Q, I	-	35.89
GPT-4-turbo	Q, I	-	42.19
GPT-40	Q, I	-	51.05
Claude3.5-Sonet	Q, I	-	46.84
Claude3-Opus	Q, I	-	33.77
Qwen-VL-Plus	Q, I	-	28.50
Qwen-VL-Max	Q, I	-	35.61
GLM-4V-9B	Q, I	9B	30.59
MathGLM-Vision-9B	Q, I	9B	57.05
CogVLM2	Q, I	8B	27.47
MathGLM-Vision-19B	Q, I	8B	57.30
CogVLM-32B	Q, I	32B	30.86
MathGLM-Vision-32B	Q, I	32B	59.00

Table 4: **Results on MathVL-test.** A detailed comparison of the performance of MathGLM-Vision and various other leading close-source MLLMs on the MathVL-test benchmark.

multi-modal understanding and mathematical reasoning. However, MathGLM-Vision-32B shows a slight reduction in performance across multiple categories on the MMMU benchmark. Besides, MathGLM-Vision, when fine-tuned with VQA datasets, outperforms its variant lacking VQA datasets. This indicates that omitting VQA datasets from the fine-tuning process limits the general visionlanguage understanding abilities. Thus, the SFT process using our MathVL incorporated with VQA datasets not only enhances MathGLM-Vision's mathematical reasoning abilities but also preserves its generalizability.

Model	MMMU	Art & Design	Business	Sci.	Health & Med.	Human. & Social Sci.	Tech. & Eng.
CogVLM2	40.2	58.3	30.0	26.7	41.3	38.6	53.3
w/o VQA datasets	38.1	60.8	28.7	34.0	36.7	43.3	32.9
MathGLM-Vision-19B	40.2	63.3	37.3	27.3	36.0	46.7	37.6
CogVLM-32B	42.9	63.3	31.3	32.0	43.3	62.5	36.7
w/o VQA datasets	38.6	62.5	26.7	28.0	34.0	56.7	33.8
MathGLM-Vision-32B	40.0	60.0	28.7	34.0	38.7	52.5	34.8

Table 5: Generalizability of MathGLM-Vision on the MMMU benchmark.

4.4 FURTHER ANALYSIS

Effect of Chinese Dataset. To validate the effectiveness of the adopted Chinese dataset in MathVL, we conduct an extended experiment that involves fine-tuning GLM-4V-9B with open-source datasets, deliberately excluding Chinese data collected from K12 education. Table 6 shows a comparison of performance results. Compared to the backbone model GLM-4V-9B, a variant MathGLM-Vision-9B that undergoes SFT exclusively with open-source data exhibits significant improvement on the minitest of MathVista, particularly in geometry problem solving (GPS) and geometry reasoning (GEO). This indicates that fine-tuning on diverse open-source data can markedly enhance model performance in specific mathematical areas. MathGLM-Vision, incorporating both open-source data and Chinese data, outperforms the variant tuned only with open-source data on the minitest of MathVista, highlighting the significant value added by integrating the Chinese dataset in the training process. Notably, compared to the variant without Chinese data, MathGLM-Vision achieves a significantly higher accuracy on the MathVL-test benchmark. These findings confirm that the inclusion of the Chinese dataset not only enhances the model's capability in handling complex mathematical problems but also contributes significantly to its overall performance on a diverse set of tasks within MathVista.

431 Effect of VQA Datasets. To explore the effect of VQA datasets on the performance of MathGLM-Vision, an extended experiment can be designed where SFT is applied exclusively to mathematical

132 133 134	Model	MathVista GPS GEO ALL	MathVL-test	Model	MathVista GPS GEO ALL
435	GLM-4V-9B	46.12 44.35 46.70	30.59	GLM-4V-9B	46.12 44.35 46.70
436	+ SFT on Open-source Data	62.98 61.51 50.40	47.55	MathGLM-Vision-9B	64.42 62.34 52.20
437 438	MathGLM-Vision-9B	64.42 62.34 52.20	57.25	- SFT on VQA Datasets	61.54 58.58 41.34

Table 6: Effect of the constructed Chinese data.

Table 7: Effect of the VQA datasets.

datasets, deliberately excluding VQA datasets. Table 7 demonstrates the performance comparison achieved by different models on MathVista. Compared to the backbone model GLM-4V-9B, a 442 variant of MathGLM-Vision-9B achieves significant improvements on geometry problem solving 443 (GPS) and geometry Reasoning (GEO). However, it exhibits a decline in the overall accuracy on 444 the minitest of MathVista (ALL). The decline can be attributed to the composition of MathVista, 445 which comprises five tasks, with question-answering types (such as graphical question-answering, textbook question-answering, and visual question-answering) comprising up to 60.6% of the tasks. 446 Omitting VQA training in MathGLM-Vision impacts the model's ability to effectively process and 447 respond to these multi-modal questions. Notably, within specific subsets of MathVista, such as GPS 448 and GEO, a variant of MathGLM-Vision-9B slightly below the standard MathGLM-Vision-9B. This 449 observation suggests that VQA datasets are crucial for preserving overall multi-modal understanding, 450 their impact may vary depending on different task types. Besides, VQA datasets can indirectly bolster mathematical reasoning skills, which in turn enhances image recognition capabilities.

451 452 453

454

439 440

441

4.5 ERROR ANALYSIS

455 We meticulously analyze the causes of errors in MathGLM-Vision-32B on the MathVL-test bench-456 mark and illustrate the distribution of these errors in Figure 5. We summarize these errors in 457 MathGLM-Vision-32B into five types: reasoning error, knowledge error, vision recognition error, 458 calculation error, and question misunderstood error. The most prevalent type of errors, accounting for 69.1% of the total, is identified as Reasoning Error. This indicates a significant challenge in 459 the MathGLM-Vision-32B's logical deductions and inferential reasoning. Improving these capa-460 bilities can dramatically enhance the MathGLM-Vision-32B's overall performance. Knowledge 461 Error, which made up 12.7% of the errors, relates to the model's misapplication or lack of specific 462 factual information. Vision Recognition Error accounts for 11.4% of the total errors and involves 463 inaccuracies in interpreting visual data. This type of error can be reduced through the implementation 464 of more advanced vision encoders. Furthermore, the fact that Calculation Error constitutes only 465 4.3% of the errors suggests that MathGLM-Vision-32B demonstrates considerable robustness in 466 numerical and computational tasks. Lastly, Question Misunderstood Error, which constitutes 2.5% of 467 the total, occurs when the model fails to correctly interpret question. Enhancing natural language 468 processing capabilities and refining context understanding can significantly reduce these types of 469 errors. Addressing these identified error types through targeted enhancements can significantly boost the overall effectiveness of MathGLM-Vision-32B. Figure 6 demonstrate some cases of the 470 Calculation Error category. More detailed examples of these errors can be found in Appendix G. 471

472

5 **RELATED WORKS**

473 474 475

Multi-Modal Language Model. The development of Multi-Modal Language Models (MLLMs) 476 have emerged as a significant area of research, which are designed to integrate information from 477 multiple modalities-typically text and images-to perform tasks that require a holistic understanding 478 of both visual and linguistic inputs. Pioneering efforts such as ViLBERT Lu et al. (2019) and 479 LXMERT Tan & Bansal (2019) have advanced this field by conducting the joint pre-training on 480 image-text datasets. They process text and image inputs separately before fusing them for final 481 task layers, significantly improving performance on tasks like image captioning and visual question 482 answering. The continues evolution of MLLMs has lead to innovations in data fusion techniques. 483 Notable models such as CLIP Radford et al. (2021), ALIGN Jia et al. (2021), and BLIP Li et al. (2022) have adopted contrastive learning paradigms to align visual and language information from billions 484 of image-text pairs. Concurrently, the success of LLMs Brown et al. (2020); Du et al. (2021); Zeng 485 et al. (2022); Le Scao et al. (2023); Bai et al. (2022); Touvron et al. (2023a); Ouyang et al. (2022);



Figure 5: Error distribution of MathGLM-Vision-32B.

Figure 6: Cases of Calculation Error category.

Hoffmann et al. (2022); Smith et al. (2022); Chowdhery et al. (2023) facilitates the integration of LLMs into multi-modal tasks by utilizing pre-training alignment and visual instruction tuning, leading 504 to the emergence of multi-modal language models (MLLMs) Liu et al. (2024b); Liu et al.; Wang et al. 505 (2023c); Li et al. (2023); Dai et al. (2024); Bai et al. (2023a). Despite MLLMs have demonstrated 506 remarkable capabilities on tasks such as image caption and visual question answering, they stall face 507 significant challenges in solving mathematical problems that involve visual information Yue et al. 508 (2024a); Lu et al. (2023); Zhang et al. (2024); Wang et al. (2024). 509

510 Mathematical Reasoning. Recently, math-specific LLMs Azerbayev et al. (2023); Wang et al. 511 (2023a); Yue et al. (2024b); Ying et al. (2024); Yu et al. (2023); Yue et al. (2023); Yuan et al. 512 (2023); Luo et al. (2023) have demonstrated remarkable abilities in handling mathematical reasoning 513 tasks that only involve textual information. These models have been specifically trained on web-514 scale instruction mathematical dataset or fine-tuned on specialized mathematical problem sets. For 515 instance, WizardMath Luo et al. (2023) and MetaMath Yu et al. (2023) have implemented data augmentation methods to enhance the models' ability to understand and solve mathematical problems 516 by enriching the MATH Hendrycks et al. (2021) and GSM8K Cobbe et al. (2021) datasets. Recent 517 research has also focused on creating specialized MLLMs for mathematical tasks. UniGeo Chen 518 et al. (2022) and UniMath Liang et al. (2023) have demonstrated enhanced datasets and conventional 519 deep learning approaches for geometric problem solving. MLLMs like G-LLaVA Gao et al. (2023a), 520 GeoGPT4V Cai et al. (2024), and Math-LLaVA Shi et al. (2024) are tailored for mathematical 521 problem solving, incorporating both geometric understanding and algebraic reasoning. Additionally, 522 several benchmark datasets Yue et al. (2024a); Lu et al. (2023); Zhang et al. (2024); Wang et al. 523 (2024) are proposed to evaluate the multi-modal mathematical reasoning abilities of MLLMs. 524

- 6 CONCLUSION
- 526 527

525

500

501

In this paper, we attempt to address the issues in current mathematical MLLMs. We construct a fine-528 tuning dataset named MathVL, upon which we conduct a Supervised Fine-Tuning (SFT) process. This 529 initiative results in the development of a series of enhanced MLLMs, designated as MathGLM-Vision. 530 Specially, MathGLM-Vision contains three variations: MathGLM-Vision-9B, MathGLM-Vision-19B, 531 and MathGLM-Vision-32B, each fine-tuned on different backbone models: GLM-4-V, CogVLM2, 532 and CogVLM-32B, respectively. These developed MathGLM-Vision significantly improve the 533 capabilities of mathematical reasoning, achieving substantial performance improvements. Relative to their respective backbone models, MathGLM-Vision-9B, MathGLM-Vision-19B, and MathGLM-534 Vision-32B show improvements of 39%, 65%, and 53.7% on the Geometry Problem Solving (GPS) minitest split of MathVista, demonstrating the effectiveness of MathVL in enhancing the mathematical 536 problem-solving abilities of MLLMs. Additionally, we evaluate the effectiveness of MathGLM-Vision 537 on our curated MathVL-test benchmark. Experimental results reveal that MathGLM-Vision not only 538 surpass their backbone models in specialized mathematical tests but also preserve the generalizability capabilities in general vision-language understanding domains.

540 REFERENCES

552

553

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Anthropic. The claude 3 model family: Opus, sonnet, haiku. 2024. URL https://www-cdn. anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_ Card_Claude_3.pdf.
- Zhangir Azerbayev, Hailey Schoelkopf, Keiran Paster, Marco Dos Santos, Stephen McAleer, Albert Q
 Jiang, Jia Deng, Stella Biderman, and Sean Welleck. Llemma: An open language model for
 mathematics. arXiv preprint arXiv:2310.10631, 2023.
 - Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023a.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang
 Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities.
 arXiv preprint arXiv:2308.12966, 2023b.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain,
 Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with
 reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Shihao Cai, Keqin Bao, Hangyu Guo, Jizhi Zhang, Jun Song, and Bo Zheng. Geogpt4v: Towards geometric multi-modal large language models with geometric image generation. *arXiv preprint arXiv:2406.11503*, 2024.
- Jie Cao and Jing Xiao. An augmented benchmark dataset for geometric question answering through
 dual parallel text encoding. In *Proceedings of the 29th International Conference on Computational Linguistics*, pp. 1511–1520, 2022.
- Jiaqi Chen, Jianheng Tang, Jinghui Qin, Xiaodan Liang, Lingbo Liu, Eric P Xing, and Liang Lin.
 Geoqa: A geometric question answering benchmark towards multimodal numerical reasoning. *arXiv preprint arXiv:2105.14517*, 2021.
- Jiaqi Chen, Tong Li, Jinghui Qin, Pan Lu, Liang Lin, Chongyu Chen, and Xiaodan Liang. Unigeo:
 Unifying geometry logical reasoning via reformulating mathematical expression. *arXiv preprint arXiv:2212.02746*, 2022.
- Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*, 2023.
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong
 Zhang, Xizhou Zhu, Lewei Lu, et al. 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*, pp. 24185–24198, 2024.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113, 2023.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve
 math word problems. *arXiv preprint arXiv:2110.14168*, 2021.

- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. Instructblip: Towards general-purpose visionlanguage models with instruction tuning. *Advances in Neural Information Processing Systems*, 36, 2024.
- Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Xilin Wei, Songyang Zhang, Haodong Duan, Maosong Cao, et al. Internlm-xcomposer2: Mastering free-form text-image composition and comprehension in vision-language large model. *arXiv preprint arXiv:2401.16420*, 2024.
- ⁶⁰³ Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang.
 ⁶⁰⁴ Glm: General language model pretraining with autoregressive blank infilling. *arXiv preprint* arXiv:2103.10360, 2021.
- Jiahui Gao, Renjie Pi, Jipeng Zhang, Jiacheng Ye, Wanjun Zhong, Yufei Wang, Lanqing Hong,
 Jianhua Han, Hang Xu, Zhenguo Li, et al. G-llava: Solving geometric problem with multi-modal
 large language model. *arXiv preprint arXiv:2312.11370*, 2023a.
- Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan Lu, Conghui He, Xiangyu Yue, et al. Llama-adapter v2: Parameter-efficient visual instruction model. *arXiv preprint arXiv:2304.15010*, 2023b.
- Peng Gao, Renrui Zhang, Chris Liu, Longtian Qiu, Siyuan Huang, Weifeng Lin, Shitian Zhao, Shijie
 Geng, Ziyi Lin, Peng Jin, et al. Sphinx-x: Scaling data and parameters for a family of multi-modal large language models. *arXiv preprint arXiv:2402.05935*, 2024.
- Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu
 Feng, Hanlin Zhao, Hanyu Lai, et al. Chatglm: A family of large language models from glm-130b
 to glm-4 all tools. *arXiv preprint arXiv:2406.12793*, 2024.
- ⁶²⁰
 ⁶²¹ Zhibin Gou, Zhihong Shao, Yeyun Gong, Yujiu Yang, Minlie Huang, Nan Duan, Weizhu Chen, et al. Tora: A tool-integrated reasoning agent for mathematical problem solving. *arXiv preprint arXiv:2309.17452*, 2023.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv* preprint arXiv:2103.03874, 2021.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- Wenyi Hong, Weihan Wang, Ming Ding, Wenmeng Yu, Qingsong Lv, Yan Wang, Yean Cheng,
 Shiyu Huang, Junhui Ji, Zhao Xue, et al. Cogvlm2: Visual language models for image and video
 understanding. *arXiv preprint arXiv:2408.16500*, 2024.
- Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung,
 Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with
 noisy text supervision. In *International conference on machine learning*, pp. 4904–4916. PMLR, 2021.
- Mehran Kazemi, Hamidreza Alvari, Ankit Anand, Jialin Wu, Xi Chen, and Radu Soricut. Geomverse:
 A systematic evaluation of large models for geometric reasoning. *arXiv preprint arXiv:2312.12241*, 2023.
- Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. Bloom: A 176bparameter open-access multilingual language model. 2023.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre training for unified vision-language understanding and generation. In *International conference on machine learning*, pp. 12888–12900. PMLR, 2022.

648 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image 649 pre-training with frozen image encoders and large language models. In International conference 650 on machine learning, pp. 19730–19742. PMLR, 2023. 651 Zhenwen Liang, Tianyu Yang, Jipeng Zhang, and Xiangliang Zhang. Unimath: A foundational and 652 multimodal mathematical reasoner. In Proceedings of the 2023 Conference on Empirical Methods 653 in Natural Language Processing, pp. 7126–7133, 2023. 654 655 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan 656 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. arXiv preprint 657 arXiv:2305.20050, 2023. 658 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 659 Llava-next: Improved reasoning, ocr, and world knowledge (january 2024). URL https://llava-vl. 660 github. io/blog/2024-01-30-llava-next, 1(8). 661 662 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 663 pp. 26296–26306, 2024a. 664 665 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. Advances in 666 neural information processing systems, 36, 2024b. 667 668 Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. Advances in neural information processing systems, 669 32, 2019. 670 671 Pan Lu, Ran Gong, Shibiao Jiang, Liang Qiu, Siyuan Huang, Xiaodan Liang, and Song-Chun Zhu. 672 Inter-gps: Interpretable geometry problem solving with formal language and symbolic reasoning. 673 *arXiv preprint arXiv:2105.04165*, 2021. 674 Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, 675 Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning 676 of foundation models in visual contexts. arXiv preprint arXiv:2310.02255, 2023. 677 678 Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, 679 Qingwei Lin, Shifeng Chen, and Dongmei Zhang. Wizardmath: Empowering mathematical 680 reasoning for large language models via reinforced evol-instruct. arXiv preprint arXiv:2308.09583, 681 2023. 682 Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A bench-683 mark for question answering about charts with visual and logical reasoning. arXiv preprint 684 arXiv:2203.10244, 2022. 685 686 Arindam Mitra, Hamed Khanpour, Corby Rosset, and Ahmed Awadallah. Orca-math: Unlocking the 687 potential of slms in grade school math. arXiv preprint arXiv:2402.14830, 2024. 688 Gpt-4v(ision) system card. In technical report, 2023. URL https://api. OpenAL. 689 semanticscholar.org/CorpusID:263218031. 690 691 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong 692 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow 693 instructions with human feedback. Advances in neural information processing systems, 35:27730– 27744, 2022. 694 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, 696 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual 697 models from natural language supervision. In International conference on machine learning, pp. 8748-8763. PMLR, 2021. 699 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, YK Li, Yu Wu, 700 and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language 701 models. arXiv preprint arXiv:2402.03300, 2024.

702 703 704	Wenhao Shi, Zhiqiang Hu, Yi Bin, Junhua Liu, Yang Yang, See-Kiong Ng, Lidong Bing, and Roy Ka-Wei Lee. Math-llava: Bootstrapping mathematical reasoning for multimodal large language models, 2024.
705 706 707 708 709	Shaden Smith, Mostofa Patwary, Brandon Norick, Patrick LeGresley, Samyam Rajbhandari, Jared Casper, Zhun Liu, Shrimai Prabhumoye, George Zerveas, Vijay Korthikanti, et al. Using deepspeed and megatron to train megatron-turing nlg 530b, a large-scale generative language model. <i>arXiv</i> preprint arXiv:2201.11990, 2022.
710 711	Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from trans- formers. <i>arXiv preprint arXiv:1908.07490</i> , 2019.
712 713 714 715	Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> , 2023.
716 717 718	GLM Team, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, et al. Chatglm: A family of large language models from glm-130b to glm-4 all tools. <i>arXiv e-prints</i> , pp. arXiv–2406, 2024.
719 720 721	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> , 2023a.
723 724 725	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> , 2023b.
726 727 728	Ke Wang, Houxing Ren, Aojun Zhou, Zimu Lu, Sichun Luo, Weikang Shi, Renrui Zhang, Linqi Song, Mingjie Zhan, and Hongsheng Li. Mathcoder: Seamless code integration in llms for enhanced mathematical reasoning. <i>arXiv preprint arXiv:2310.03731</i> , 2023a.
729 730 731 732	Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Mingjie Zhan, and Hongsheng Li. Measuring multimodal mathematical reasoning with math-vision dataset. <i>arXiv preprint arXiv:2402.14804</i> , 2024.
733 734 735	Peiyi Wang, Lei Li, Zhihong Shao, RX Xu, Damai Dai, Yifei Li, Deli Chen, Y Wu, and Zhifang Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations. CoRR, abs/2312.08935, 2023b.
736 737 738 739	Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. Cogvlm: Visual expert for pretrained language models. arXiv preprint arXiv:2311.03079, 2023c.
740 741 742	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837, 2022.
743 744 745 746	Zhen Yang, Ming Ding, Qingsong Lv, Zhihuan Jiang, Zehai He, Yuyi Guo, Jinfeng Bai, and Jie Tang. Gpt can solve mathematical problems without a calculator. <i>arXiv preprint arXiv:2309.03241</i> , 2023.
747 748 749	Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. mplug-owl: Modularization empowers large language models with multimodality. <i>arXiv preprint arXiv:2304.14178</i> , 2023.
750 751 752 753	Huaiyuan Ying, Shuo Zhang, Linyang Li, Zhejian Zhou, Yunfan Shao, Zhaoye Fei, Yichuan Ma, Jiawei Hong, Kuikun Liu, Ziyi Wang, et al. Internlm-math: Open math large language models toward verifiable reasoning. arXiv preprint arXiv:2402.06332, 2024.
754 755	Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for large language models. <i>arXiv preprint arXiv:2309.12284</i> , 2023.

756 757 758	Zheng Yuan, Hongyi Yuan, Chengpeng Li, Guanting Dong, Chuanqi Tan, and Chang Zhou. Scaling relationship on learning mathematical reasoning with large language models. <i>arXiv preprint arXiv:2308.01825</i> , 2023.
760 761 762	Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. Mammoth: Building math generalist models through hybrid instruction tuning. <i>arXiv preprint</i> <i>arXiv:2309.05653</i> , 2023.
763 764 765 766	Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal under- standing and reasoning benchmark for expert agi. In <i>Proceedings of the IEEE/CVF Conference on</i> <i>Computer Vision and Pattern Recognition</i> , pp. 9556–9567, 2024a.
767 768 769	Xiang Yue, Tuney Zheng, Ge Zhang, and Wenhu Chen. Mammoth2: Scaling instructions from the web. <i>arXiv preprint arXiv:2405.03548</i> , 2024b.
770 771 772	Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. Glm-130b: An open bilingual pre-trained model. <i>arXiv preprint arXiv:2210.02414</i> , 2022.
773 774 775	Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Peng Gao, et al. Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems? <i>arXiv preprint arXiv:2403.14624</i> , 2024.
777 778	Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. Multimodal chain-of-thought reasoning in language models. <i>arXiv preprint arXiv:2302.00923</i> , 2023.
779 780 781 782 783 784 785 786 787 788 789 790 791 792 793 794 795 796 797 798 799 800 801 802 803 804 805 806 807 808	Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. <i>arXiv preprint arXiv:2304.10592</i> , 2023.
807 808 809	

A DATASET CASES

811 812 813

814

815

816

817

In this section, we provide a detailed overview of specific cases from our constructed MathVL dataset. These cases demonstrate the variety of mathematical disciplines covered by MathVL, including arithmetic, geometry, algebra, statistics, and word problems. Figure 7, Figure 8, Figure 9, Figure 10, and Figure 11 depict the types of problems that MathGLM-Vision is designed to tackle in each respective category. Each of these categories is critical for assessing the comprehensive mathematical capabilities of MathGLM-Vision. By tackling a wide range of problems, MathGLM-Vision demonstrates its versatility and robustness in addressing diverse mathematical tasks.

818 819 820 821

822

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846 847

848

B BACKBONE MODELS

We utilize the following multi-modal large language models as our backbone models for conducting
Specialized Fine-Tuning (SFT) on the constructed MathVL. The detailed description of each backbone
model can be represented as follows:

- **GLM-4V-9B** is a bilingual (Chinese and English) multi-modal large language model, developed collaboratively by Zhipu.AI and Tsinghua University. It is built upon the foundational architecture of GLM-4-9B, enhancing its capabilities to handle complex multi-modal interactions. GLM-4V-9B takes a high resolution of 1120 * 1120 images as visual inputs. In comprehensive evaluations that test various capabilities including combined language skills, perceptual reasoning, text recognition, and chart understanding, GLM-4V-9B consistently outperforms competitors such as GPT-4-turbo-2024-04-09, Gemini 1.0 Pro, Qwen-VL-Max, and Claude 3 Opus, demonstrating its superior performance across multiple modalities.
- **CogVLM2** is a series of open-source multi-modal large language models derived from Meta-Llama-3-8B-Instruct, developed by Zhipu.AI and Tsinghua University. This series contains two models: cogvlm-llama3-chat-19B and cogvlm2-llama3-chinese-chat-19B. The former is a monolingual language model focused on English and the latter is a bilingual model supporting both English and Chinese. CogVLM2 is designed to handle extended content lengths up to 8K and accepts high-resolution images up to 1344 * 1344. Here, we choose cogvlm2-llama3-chinese-chat-19B as our backbone to pre-train our MathGLM-Vision-19B.
- **CogVLM-32B** is a close-source multi-modal large language model, developed by Zhipu.AI and Tsinghua University. It is based on the GLM-32B architecture and is optimized for handling complex multi-modal tasks. CogVLM-32B is engineered to process visual inputs at a high resolution of 1120 * 1120, enabling detailed image analysis and enhanced interaction with visual data.

Table 8 demonstrates an overview of the MathGLM-Vision series, detailing the different model parameters and configurations.

Model	LIM Size	ToTal Size	Language Model Image Encoder					
Widder			Layers	Hidden Size	Heads	Layers	Hidden Size	Heads
MathGLM-Vision-9B	9B	20B	40	4096	32	63	1792	16
MathGLM-Vision-19B	8B	19B	32	4096	32	63	1792	16
MathGLM-Vision-32B	32B	43B	58	6656	52	63	1792	16

Table 8: An overview of MathGLM-Vision series along with model parameters and configurations.

C DESCRIPTIONS OF VQA DATASETS

860 861

858 859

Here, we provide a detailed description of collected visual question answering datasets (VQA)
 datasets. Table 9 demonstrates details 19 different VQA datasets, including task types and visual context.





Figure 8: Cases of geometric problems in our MathVL dataset.



In the same Cartesian coordinate system, the graphs of the linear function $y = k_1x + b$ and the proportional function $y = k_2x$ are shown in the figure. The range of x values that satisfy $k_1x + b > k_2x$ is

Answer



[Solution]When $x \le -3$, the line l_1 : $y_1 = k_1x + b$ is above the line l_2 : $y_2 = k_2x$, meaning $k_1x + b > k_2x$. Therefore, the range of x values that satisfy $k_1x + b > k_2x$ is x < -3. [Answer]The answer is x < -3.

Question

As shown in the figure, the graph represents the quadratic function $y = ax^2 + bx + c$. Among the following statements: (1) ac > 0; (2) a - b + c < 0; $(3) 4ac < b^2$; (4) 2a + b > 0; (5) When x > 0, y decreases as x increases. The number of correct statements is () A. 1 B. 2 C. 3 D. 4

Answer

[Solution] ① From the graph, we know that a > 0 and c < 0, so ac < 0. Thus, ① is incorrect. ② From the graph, we know that when x = -1, y = a - b + c > 0, so ② is incorrect. ③ Since the parabola has two intersection points with the x-axis, $\triangle = b^2 - 4ac > 0$. Thus, ③ is correct. ④ From the axis of symmetry, we know that -b/2a < 1, so 2a + b > 0. Thus, ④ is correct. ⑤ When x > -b/2a, y increases as x increases, so ⑤ is incorrect.

[Answer]The answer is B.

Question

The graph of the quadratic function $y = ax^2 + bx + c$ is shown in the figure. Among the following conclusions: (1) ab > 0; (2) a + b - 1 = 0; (3) a > 1; (4) One root of the quadratic equation $ax^2 + bx + c = 0$ is 1, and the other root is -1/a. The correct conclusions are ______



Answer



Figure 9: Cases of algebraic problems in our MathVL dataset.



Figure 10: Cases of statistical problems in our MathVL dataset.



Figure 11: Cases of word problems in our MathVL dataset.

Dataset	Task	Visual Context
DocVQA	Figure Question Answering (FQA)	Document Image
DVQA	Figure Question Answering (FQA)	Bar Chart
FigureQA	Figure Question Answering (FQA)	Charts and Plots
PlotQA	Figure Question Answering (FQA)	Bar, Line, Scatter
MapQA	Figure Question Answering (FQA)	Map Chart
IconQA	Math Word Problem (MWP)	Abstract Scene
TabMWP	Math Word Problem (MWP)	Table
CLEVR-Math	Math Word Problem (MWP)	Synthetic Scene
TQA	Textbook Question Answering (TQA)	Scientific Figure
AI2D	Textbook Question Answering (TQA)	Scientific Figure
ScienceQA	Textbook Question Answering (TQA)	Scientific Figure
A-OKVQA	Visual Question Answering (VQA)	Natural Image
VQA2.0	Visual Question Answering (VQA)	Natural Image
PMC-VQA	Visual Question Answering (VQA)	Medical Image
VizWiz	Visual Question Answering (VQA)	Natural Image
Super-CLEVR	Visual Question Answering (VQA)	Synthetic Scene
VQA-AS	Visual Question Answering (VQA)	Abstract Scene
VQA-RAD	Visual Question Answering (VQA)	Medical Image
TextVQA	Visual Question Answering (VQA)	Natural Image

Table 9: Summary of VQA datasets.

IMPLEMENTATION DETAILS D

We provide a detailed overview of the Specialized Fine-Tuning (SFT) process applied to our MathGLM-Vision. The specific hyperparameters used during this process are outlined in Table 10.

1163 1164	parameters	MathGLM-Vision-9B	MathGLM-Vision-19B	MathGLM-Vision-32B
1165	Total steps	35,000	35,000	35,000
1166	Global Batch Size	128	128	128
1167	Learning Rate	$1e^{-5}$	$1e^{-5}$	$1e^{-5}$
1168	Learning Rate Schedule	cosine decay	cosine decay	cosine decay
1169	Warmup Ratio	0.01	0.01	0.01
1170	Weight Decay	$5e^{-2}$	$5e^{-2}$	$5e^{-2}$
1171	Optimizer	AdamW	AdamW	AdamW
1172	Input Resolution	1120 * 1120	1344 * 1344	1120 * 1120
1173	Image Length	1600	2304	1600
1174				

Table 10: The detailed setup of the SFT procedures.

Ε THE DETAILED DESCRIPTION OF BENCHMARK DATASETS

In this section, we provide an in-depth description of the benchmark datasets used to evaluate the performance of MathGLM-Vision. These benchmark datasets have been carefully curated to test the MLLMs' capabilities. The detailed description of benchmark datasets is provides as follows.

• MathVista

MathVista is a comprehensive benchmark dataset designed to rigorously evaluate the math-metical reasoning capabilities of language models (LMs), especially in varied visual contexts. This dataset offers a comprehensive evaluation benchmark designed to integrate mathemat-ical reasoning with visual understanding, focusing on five primary tasks: figure question

1188	answering (FOA) , geometry problem solving (GPS) , math word problem (MWP) , teythook
1189	question answering (TOA); and visual question answering (VOA).
1190	
1191	• MathVista-GPS
1192	MathVista-GPS, a subset of the MathVista Dataset, specifically focuses on the domain of
1193	geometry problem solving. The questions in this subset range from basic shape recognition
1194	to more advanced problems involving theorems, calculations and reasoning.
1195	• MathVerse
1196	MathVerse is designed to provide a fair and comprehensive assessment of MLLMs' capa-
1197	bilities in visual mathematics. The benchmark comprises 2,612 high-quality, multi-subject
1198	math problems, each featuring diagrams and converted into six different versions by human
1199	annotators. These versions offer varying levels of multi-modal information, allowing for a
1200	thorough evaluation of MLLMs' understanding of visual diagrams.
1201	• MATH Vision

• MATH-Vision

The MATH-Vision (MATH-V) dataset comprises 3,040 high-quality mathematical problems, each featuring a visual context and sourced from 19 real math competitions. This extensive and diverse collection allows for a comprehensive evaluation of LMMs' ability to interpret and reason with visual information in mathematical contexts.

MMMU

The Massive Multi-discipline Multi-modal Understanding and Reasoning (MMMU) benchmark encompasses 11.5K questions across six disciplines, including Art, Business, Health & Medicine, Science, Humanities & Social Science, and Tech & Engineering. The tasks in MMMU challenge models to perform sophisticated multi-modal analysis and apply domain-specific knowledge, demanding a higher level capability in comprehension and integration.

1212 1213 1214

1202

1203

1205

1207

1208

1209

1210

1211

F DETAILED EXPERIMENTAL RESULTS ON PUBLIC BENCHMARK DATASETS

1215 1216

1217 **Results on the testmini subset of MathVista.** To comprehensively evaluate the performance of 1218 MathGLM-Vision across various task types featured in the MathVista dataset, we systematically 1219 evaluate it on the testmini subset. This subset has been carefully selected to represent a diverse range 1220 of mathematical problem types, ensuring a robust assessment of our model's capabilities. Table 11 shows the evaluation results on the testmini subset of MathVista across various task types. Notably, 1221 MathGLM-Vision-19B and MathGLM-Vision-32B surpass human performance in overall accuracy, 1222 highlighting the advanced capabilities of these models in handling complex mathematical problems. 1223 In particular, MathGLM-Vision excels significantly in geometry problem solving (GPS) and geometry 1224 reasoning (GEO), demonstrating its superior proficiency in mathematical reasoning. 1225

1226															
1227	Model	Input	ALL	FQA	GPS	MWP	TQA	VQA	ALG	ARI	GEO	LOG	NUM	SCI	STA
1228	Human Performance	Q, I	60.30	59.70	48.40	73.00	63.20	55.90	50.90	59.20	51.40	40.70	53.80	64.90	63.90
1229	2-shot CoT GPT-4 2-shot PoT GPT-4	$\begin{array}{c} Q, I_c, I_t \\ Q, I_c, I_t \end{array}$	30.50 31.74	27.21 27.58	35.91 37.35	21.30 23.87	43.13 43.00	28.17 30.27	35.72 37.15	25.17 27.93	35.80 37.48	24.74 22.68	15.41 15.83	47.28 44.47	31.29 31.87
1230	GPT-4V	Q, I	49.90	43.10	50.50	57.50	65.20	38.00	53.00	49.00	51.00	21.60	20.10	63.10	55.80
1231	LLaVA-LLaMA-2-13B	Q, I	25.40	22.86	24.57	18.15	35.82	29.69	26.93	22.47	24.45	19.07	19.05	34.71	21.61
1232	MathGLM-Vision-9B MathGLM-Vision-19B	Q, I Q, I	52.20 61.10	46.10 59.85	64.42 65.38	58.60 68.28	55.70 53.80	37.43 55.31	59.79 59.79	43.91 59.21	62.34 63.18	10.81 18.92	37.50 59.03	54.10 53.28	54.82 68.44
1233	MathGLM-Vision-32B	Q, I	62.40	62.83	62.02	69.35	62.03	54.19	60.50	60.62	61.92	16.22	52.08	60.66	72.09

1235 Table 11: Accuracy scores on the *testmini* subset of MathVista. Input: Q: question, I: image, I_c : image caption, I_t : OCR texts detected from the image. ALL: overall accuracy. Task types: 1236 FQA: figure question answering, GPS: geometry problem solving, MWP: math word problem, TQA: 1237 textbook question answering, VQA: visual question answering. Mathematical reasoning types: ALG: 1238 algebraic reasoning, ARI: arithmetic reasoning, GEO: geometry reasoning, LOG: logical reasoning, 1239 NUM: numeric common sense, SCI: scientific reasoning, STA: statistical reasoning. The highest 1240 accuracy among all baseline MLLMs is marked in red, while the highest accuracy among various 1241 variants of MathGLM-Vision is marked bold.

1242 **Results on the testmini set of MathVerse.** To thoroughly evaluate the performance of MathGLM-1243 Vision across 12 detailed subjects within the MathVerse dataset, we conduct comprehensive exper-1244 iments and report the results in Table 12. This analysis delves into the model's ability to address 1245 a broad spectrum of mathematical challenges, ranging from geometry to functions. As shown in 1246 Table 12, MathGLM-Vision surpasses all open-source MLLMs and most close-source MLLMs. However, it still falls short by 14% compared to the performance of GPT-4V. In some subjects such as 1247 Angle, Analytic, and Property, MathGLM-Vision achieves better performance compared the advanced 1248 GPT-4V. For example, MathGLM-Vision-32B shows remarkable performance in plane geometry, 1249 particularly in handling angle-related problems, where it achieves a 60.1% accuracy, showcasing its 1250 strong geometric reasoning capabilities. 1251

	Model	A 11	Plane Geometry						Solid Geometry				Functions				
		All	All	Len	Area	Angle	Anal	Apply	All	Len	Area	Vol	All	Coord	Prop	Exp	Apply
							Closed-	source M	LLMs								
	Qwen-VL-Plus Gemini-Pro Qwen-VL-Max GPT-4V	21.3 35.3 37.2 54.4	17.3 33.0 38.4 56.9	19.1 32.2 41.7 60.8	16.4 42.6 46.4 63.4	16.1 28.4 32.6 52.6	23.6 30.2 40.6 48.5	13.2 32.3 38.7 60.9	24.8 33.4 33.7 50.2	18.1 35.0 25.4 54.8	18.7 29.3 28.3 39.9	33.4 36.1 42.6 56.8	31.3 28.3 38.4 52.8	52.5 25.7 43.7 72.3	25.1 26.6 35.5 47.1	10.8 10.8 13.6 30.9	50.3 51.3 61.0 70.1
Open-source MLLMs																	
	LLaMA-Adapter V2 ImageBind-LLM mPLUG-Owl2	5.8 10.0	5.9 9.7	4.0 12.1 8.2	5.9 9.9	6.6 9.2	13.4 10.2	3.3 4.8	4.6	5.3 4.9	3.1 3.5 6.7	5.7 5.3	6.2 14.9	6.7 12.3 22.8	6.1 13.8 18.6	4.5 4.6	7.9 25.9 22.2
	MiniGPT-v2 LLaVA-1 5	10.9	11.6	10.0	9.8 15.1	14.3 9.7	9.1 9.4	11.8	1.7	2.2	1.6	0.5	11.2	4.2	15.7	4.0 9.5	21.1 23.7
	SPHINX-Plus	14.0	14.4	14.2	10.5	14.1	16.5	16.8	7.0	7.2	6.1	7.6	17.9	11.1	19.1	6.3	27.7
	LLaVA-NeXT	17.2	15.9	14.8	13.0	16.3	17.7	17.8	19.6	33.3	11.7	12.6	23.1	24.5	23.4	8.0	33.1
	SPHINX-MoE InternLM-XC2	22.8 25.9	24.5 26.2	26.3 27.1	28.4 29.7	21.1 20.6	26.6 18.5	21.1 24.4 22.2	15.8 20.1	9.4 34.5	10.9 10.7 14.1	26.3 25.2	19.5 23.7	23.5 24.4	19.3 24.9	6.4 9.2 10.6	23.8 30.3 36.3
							Math	GLM-Visi	ion								
	MathGLM-Vision-9B MathGLM-Vision-19B	44.2 42.5	45.3	43.7 34.8	48.9 55.3	41.5 38.9	53.5 46.5	52.2 53.6	42.0 51.3	54.2 66.7	50.0 52.3	29.4 43.1	42.1	25.0 18.8	42.3 38.0	43.8 28.1	47.5 55.0
	MathGLM-Vision-32B	49.2	49.0	42.4	59.6	51.3	48.8	50.7	45.4	62.5	40.9	41.2	52.8	43.8	59.2	40.6	55.0

1269 Table 12: Mathematical Evaluation on Different Subjects and Subfields in MathVerse's testmini 1270 Set. Len: Length; Anal: Analytic; Apply: Applied; Vol: Volume; Coord: Coordinate; Prop: Property; 1271 Exp: Expressio. The highest accuracy among all baseline MLLMs is marked in red, while the highest 1272 accuracy among various variants of MathGLM-Vision is marked bold. 1273

1274

Results on Math-Vision datasets. To effectively assess MathGLM-Vision's ability across diverse 1275 subjects and difficulty levels within the Math-Vision dataset, we conduct a series of detailed evaluation 1276 experiments and report results in Table 13. Specifically, GPT-4V leads the close-source models with 1277 an overall accuracy of 22.76%, yet it remains significantly below the human performance benchmark 1278 of 75.66%. MathGLM-Vision shows competitive performance across a variety of mathematical 1279 disciplines compared to most of close-source MLLMs, with MathGLM-Vision-32B achieving the 1280 overall accuracy of 26.5%, closely approaching that of GPT-4V. Notably, MathGLM-Vision-32B 1281 excels in solid geometry with a accuracy of 29.1%, significantly outperforming the accuracy of 1282 23.8% on GPT-4V. This superior performance in solid geometry highlights MathGLM-Vision-32B's 1283 advanced spatial reasoning and geometric processing capabilities, which are essential for tackling 1284 complex three-dimensional problems.

- 1285 1286
- G ERROR CASES

1287

Figure 12, Figure 13, Figure 14, and Figure 15 show examples of errors made by MathGLM-Vision-32B on the MathVL-test dataset. Each figure highlights a specific type of error, providing valuable 1290 insights into the model's limitations and areas for improvement. 1291

- - Η CASE STUDY
- 1293 1294
- Figure 16, Figure 17, Figure 18, and Figure 19 present several case studies from MathGLM-Vision-1295 32B. These figures showcase the model's performance in various scenarios, highlighting its strengths



Figure 13: An example of vision recognition error. MathGLM-Vision-32B incorrectly interpreted the geometric properties of the diagram, leading to a vision recognition error.





Figure 15: An example of premature conclusion error. MathGLM-Vision-32B prematurely 1399 concluded that AB is perpendicular to CD without proper reasoning, leading to a premature conclusion 1400 error. 1401

						Humar	Perfor	nance									
Model	Overall	Alg	AnaG	Ari	CombG	Comb	Cnt	DescG	GrphT	Log	Angle	Area	Len	SolG	Stat	Торо	TransG
Human (testmini)	75.66	57.9	79.0	100.0	100.0	47.4	94.7	89.5	63.2	63.2	36.8	52.6	73.7	89.5	89.5	100.0	73.7
						Open-s	ource M	LLMs									
LLaVA-v1.5-7B	8.52	7.0	7.1	10.7	7.1	4.8	10.5	7.7	10.0	9.2	15.6	10.2	9.8	5.3	8.6	4.4	4.8
SPHINX (V2) ShareGPT4V-7B	9.70 10.53	6.7	7.1	12.9	7.5 10.1	7.7 4.8	6.0 7.5	9.6 11.5	16.7 14.4	10.1	11.0 16.2	11.8	12.5	8.2 9.8	8.6 15.5	8.7 17.4	6.0 11.3
LLaVA-v1.5-13B	11.12	7.0	14.3	14.3	9.1	6.6	6.0	13.5	5.6	13.5	10.4	12.6	14.7	11.5	13.8	13.0	10.7
ShareGPT4V-13B	11.88	7.5	15.5	16.4	10.7	8.9	9.0	11.5	8.9	7.6	11.6	13.0	17.4	10.3	8.6	8.7	12.5
InternLM-XComposer2-VL	14.13	9.3	15.5	12.1	15.3	11.3	10.5	14.4	22.2	19.3	19.7	15.6	15.0	11.9	15.5	26.1	15.5
						Closed-s	source N	ILLMs									
Qwen-VL-Plus	10.72	11.3	17.9	14.3	12.7	4.8	10.5	15.4	8.9	14.3	11.6	6.4	10.0	14.3	6.9	8.7	11.31
Qwen-VL-Max Gemini Pro	15.59	10.7	19.1	20.0	16.9	12.5	17.9	16.4	12.2	21.0	13.3	14.2	19.8	11.5	20.7	13.0	17.3
GPT4V	22.76	27.3	32.1	35.7	21.1	16.7	13.4	22.1	14.4	16.8	22.0	22.2	20.9	23.8	24.1	21.7	25.6
						Math	GLM-V	ision									
MathGLM-Vision-9B	19.2	18.6	20.2	19.3	15.3	18.5	20.9	26.0	18.9	15.1	23.1	20.4	18.3	23.8	19.0	17.4	14.3
MathGLM-Vision-19B MathGLM-Vision-32B	21.6 26.5	22.0 22.9	29.8 20.2	23.6 24.3	22.4 23.1	18.5 28.0	25.4 20.9	25.0 34.6	17.8 27.8	16.0 23.5	20.2 31.2	22.0 26.8	20.3 30.1	21.3 29.1	20.7 22.4	30.4 17.4	23.2 26.2
															-		

Table 13: Comparison of model performances across various mathematical subjects. Subjects:
Alg: algebra, AnaG: analytic geometry, Ari: arithmetic, CombG: combinatorial geometry, Comb:
combinatorics, Cnt: counting, DescG: descriptive geometry, GrphT: graph theory, Log: logic, Angle:
metric geometry - angle, Area: metric geometry - area, Len: metric geometry - length, SolG: solid
geometry, Stat: statistics, Topo: topology, TransG: transformation geometry. The highest accuracy
among all baseline MLLMs is marked in red, while the highest accuracy among various variants of
MathGLM-Vision is marked bold.

1424 1425

1429

in providing concise and clear answers with logical mathematical reasoning. Compared to other close source MLLMs, MathGLM-Vision-32B stands out for its ability to deliver precise and understandable
 solutions.



Figure 16: An example of solid geometry problem. MathGLM-Vision-32B correctly calculated the shortest path with shorter steps, while GPT-40 made a reasoning error, leading to a different conclusion.



Figure 17: An example of a planar geometry problem. MathGLM-Vision-32B correctly utilized axial symmetry to determine the angle sum, while Claude-3.5-Sonnet arrived at the correct answer through an erroneous calculation process. Despite the correct final answer, the calculation process was incorrect and overly complex.

1539

1540



Figure 18: **An example of a planar geometry problem.** MathGLM-Vision-32B correctly used the properties of similar triangles to find the area ratio, while GPT-40 misinterpreted the geometric relationships and misapplied the formula for the area ratio, leading to an erroneous conclusion.



Figure 19: **An example of a geometric transformation problem.** MathGLM-Vision-32B correctly used the properties of rotation to determine the angle, while QWen-VL-Max misunderstood the problem requirements and incorrectly calculated the rotation angle as 90°.

1566 I MODEL EVALUATION

1568

1577

Evaluation on public benchmarks The existing public benchmarks for evaluating a wide array

of open-source and close-source models are neither timely nor comprehensive enough. To compare our MathGLM-Vision with the state-of-the-art open-source and close-source LLMs, we have supplemented the evaluations for some models missing from the public benchmark leaderboard.

We generate LLMs' responses through API access (for closed-source models) and local inference
(for open-source models). The evaluation was then conducted following the official evaluation code
from each benchmark's GitHub repository. The source of the models used in the evaluation can be
found in Table 14.

Model	Input LLM Size		Source				
	С	losed Source	Models				
Multi-modal LLMs							
Gemini Pro	Q, I	-	gemini-pro				
Gemini 1.5 Pro	Q, I	-	gemini-1.5-pro				
GPT-4V	Q, I	-	gpt-4-vision-preview				
GPT-4-turbo	Q, I	-	gpt-4-turbo				
GPT-40	Q, I	-	gpt-4o				
Claude-3-Opus	Q, I	-	claude-3-opus-20240229				
Claude-3.5-Sonnet	Q, I	-	claude-3-5-sonnet-2024620				
Qwen-VL-Plus	Q, I	-	qwen-vl-plus				
Qwen-VL-Max	<i>Q</i> , <i>I</i>	-	qwen-vl-max				
	(Open Source N	Nodels				
General Multi-modal	LLMs						
mPLUG-Owl	Q, I	7B	mPLUG-Owl				
LLaMA-Adapter-V2	Q, I	7B	LLaMA-Adapter V2				
InstructBLIP	Q, I	7B	InstructBLIP				
LLaVA-1.5	Q, I	13B	LLaVA-v1.5-13B				
ShareGPT-4V	Q, I	13B	ShareGPT4V-13B				
SPHINX-MoE	Q, I	8*7B	SPHINX-MoE				
SPHINX-Plus	Q, I	13B	SPHINX-Plus				
InternLM-XC2	Q, I	7B	InternLM-XComposer2-VL-7B				
InternVL-1.2-Plus	Q, I	34B	InternVL-Chat-V1-2-Plus				
Geo-Multi-modal LLM	ls						
G-LLaVA	Q, I	7B	G-LLaVA-7B				
G-LLaVA	Q, I	13B	G-LLaVA-13B				
LLaVA-1.5-G	Q, I	7B	LLaVA-1.5-7B-GeoGPT4V				
LLaVA-1.5-G	Q, I	13B	LLaVA-1.5-13B-GeoGPT4V				
ShareGPT4V-G	Q, I	7B	ShareGPT4V-7B-GeoGPT4V				
ShareGPT4V-G	Q, I	13B	ShareGPT4V-1.5-13B-GeoGPT4V				
Math-LLaVA	Q, I	13B	Math-LLaVA-13B				

1619

Table 14: The source of the models used in the evaluation.

Evaluation on MathVL-test We evaluate MathGLM-Vision and several close-source MLLMs using our specially constructed MathVL-test. The evaluation process of our MathVL-test is conducted through 3 key-step: generation, extraction, and scoring.

For the generation step, the model responses are generated by providing the model with queries which incorporate the Chain of Thought (CoT) template, questions, and diagram information. The reponses of close-source MLLMs is generated through API access. For the extraction step, we use GPT-3.5-turbo to extract the model's answer based on the reponses of first step. Finally, in the scoring step, the score for each question is determined by GLM-4 based on the comparison between the extracted answer and the standard answer.

The prompts used to guide the LLM in response generation, answer extraction and scoring can be found in Table 15.

Task	Prompt
Response Generation	You are a very skilled math teacher. Please provide a detailed, step-by- step solution to the question, following a step-by-step format. Be sure to conclude with a summary that states "The answer to this question is" followed by the final result.
Answer Extraction	Please read the following example. Then extract the answer from the model response and type it at the end of the prompt. Hint: Please answer the question requiring an integer answer and provide the final value, e.g., 1, 2, 3, at the end. Question: Which number is missing? Model response: The number missing in the sequence is 14. Extracted answer: 14 Hint: Please answer the question requiring a floating-point number with one decimal place and provide the final value, e.g., 1.2, 1.3, 1.4, at the end. Question: What is the fraction of females facing the camera? Model response: The fraction of females facing the camera. Extracted answer: 0.6
Scoring	Please determine if the extracted_answer correctly answers the question. The correct answer needs to be extracted from the answer without re- calculating it, and the answer in the answer should be considered the final answer. Also, do not judge whether the answer is correct. The question may contain multiple sub-questions, and correctly answering the question includes correctly answering every sub-question and every result within each sub-question. A relative error divided by the absolute value of the original answer of less than 0.01 is allowed. If the prediction does not contain an answer, it is considered wrong. If the answer is not numerical, determine the equivalence of the expression, not just the value. If there is one mistake, the answer is wrong. Only if all results given in the prediction are correct is it considered correct. There is no need to consider whether the solution process of the prediction is com- plete. Please first extract the answers given by the prediction, determine the relative error, check if each sub-question is answered correctly, and finally give the judgment in a single line (output only "yes" or "no" in a single line).
Table 15: Pr	rompts for response generation, answer extraction and scoring.

1670 1671

1672