ERRORRADAR: BENCHMARKING COMPLEX MATHE MATICAL REASONING OF MULTIMODAL LARGE LAN GUAGE MODELS VIA ERROR DETECTION

Anonymous authors

Paper under double-blind review

Abstract

As the field of Multimodal Large Language Models (MLLMs) continues to evolve, their potential to handle mathematical reasoning tasks is promising, as they can handle multimodal questions via cross-modal understanding capabilities compared to text-only LLMs. Current mathematical benchmarks predominantly focus on evaluating MLLMs' problem-solving ability, yet there is a crucial gap in addressing more complex scenarios such as error detection, for enhancing reasoning capability in complicated settings. To fill this gap, we formally formulate the new task — multimodal error detection, and introduce ERRORRADAR, the first benchmark designed to assess MLLMs' capabilities in such a task. ERROR-RADAR evaluates two sub-tasks: error step identification and error categorization, providing a framework for evaluating MLLMs' complex mathematical reasoning ability. It consists of 2,500 high-quality multimodal K-12 mathematical problems, collected from real-world student interactions in an educational organization, with expert-based annotation and metadata such as problem type and error category. Through extensive experiments, we evaluated both open-source and closed-source representative MLLMs, benchmarking their performance against educational expert evaluators. Results indicate challenges still remain, as GPT-40 with best model performance is still around 10% behind human evaluation.

029 030 031

032

006

008 009 010

011 012 013

014

015

016

017

018

019

021

023

025

026

027

028

1 INTRODUCTION

033 On the path to Artificial General Intelligence, Large Language Models (LLMs) such as GPT-4 (Ope-034 nAI, 2023) have emerged as a central focus in both industry and academia (Minaee et al., 2024; Zhao 035 et al., 2023; Zhu et al., 2023). As the real world is inherently multimodal, the evolution of Multimodal Large Language Models (MLLMs) such as the latest GPT-40 (OpenAI, 2024a) and Gemini 037 1.5 (Reid et al., 2024), has become a rapidly growing area of interest, demonstrating remarkable 038 effectiveness in diverse applications (Xiao et al., 2024; He et al., 2024a; Yan et al., 2024; Hao et al., 2024). In particular, multimodal reasoning stands to significantly benefit education scenarios from the robust capabilities of MLLMs (Wang et al., 2024b; Li et al., 2024a), given its reliance on multi-040 modal inputs to comprehensively grasp users' intentions and needs. 041

Within the multimodal sphere, mathematical scenarios pose a significant challenge, demanding so-phisticated reasoning abilities from MLLMs (Lu et al., 2022; Ahn et al., 2024). These scenarios have attracted considerable research aimed at pushing the boundaries of MLLMs' reasoning capabilities (Hu et al., 2024; Jia et al., 2024; Lu et al., 2024c; Shi et al., 2024b; Zhuang et al., 2024). Besides, various representative benchmarks have been designed to measure MLLMs' performance in complex mathematical reasoning tasks, which involve multi-step reasoning and quantitative analysis within visual contexts (Lu et al., 2024b; Zhang et al., 2024; Qiao et al., 2024; Peng et al., 2024).

Scrutinizing the off-the-shelf mathematical reasoning benchmarks, there is a predominant focus
 on evaluating the problem-solving capabilities of MLLMs, prioritizing the accuracy with which
 MLLMs can solve mathematical problems (Wang et al., 2024a; Lu et al., 2024b; Zhang et al., 2024),
 as depicted in Figure 1 (a). However, in educational contexts, it is even more crucial to consider
 user-oriented needs, such as error detection. As indicated in Figure 1 (b), this involves not only
 pinpointing the first incorrect step in a student's step-by-step solution but also categorizing the types

067

068

069

071 072 073

074

075

054



Benchmarks	Venue	Modality	Student Ans.	Error Det.
TheoremQA (Chen et al., 2023a)	EMNLP	Т	-	-
MathBench (Liu et al., 2024b)	ACL	T	-	-
MR-GSM8K (Zeng et al., 2024)	arXiv	T	-	-
SciEval (Sun et al., 2024)	AAAI	T	-	-
EIC (Li et al., 2024b)	arXiv	T	-	\checkmark
CMMaTH (Li et al., 2024c)	arXiv	T, I	-	-
MathScape (Zhou et al., 2024)	arXiv	T, I	-	-
MATH-V (Wang et al., 2024a)	arXiv	T, I	-	-
QRData (Liu et al., 2024c)	ACL	T, I	-	-
IsoBench (Fu et al., 2024)	COLM	T, I	-	-
SciBench (Wang et al., 2024c)	ICML	T, I	-	-
MathVista (Lu et al., 2024b)	ICLR	T, I	-	-
MathVerse (Zhang et al., 2024)	ECCV	T, I	-	-
ERRORRADAR (Ours)	-	T, I	√	\checkmark

Table 1: Comparison between our proposed ERROR-RADAR benchmark vs. its relevant LLM-based mathematical reasoning benchmarks or datasets. Under the column of *Modality*, the letters T and I represent text and image, respectively. The column labeled as Student Ans. indicates whether the dataset contains real student data (i.e., students' incorrect answers); the column lavious work and our proposed ERRORRADAR bench- beled as Error Det. represents whether evaluation includes the complex reasoning task of error detection.

Figure 1: Comparison of research scope between premark on mathematical reasoning tasks.

of errors made, which is a multifaceted process that requires a deep understanding of both mathematical concepts and cognitive processes (Davies et al., 2021; Rabillas et al., 2023).

Towards this end, addressing the aforementioned research gap, we aim to formulate the new task 076 of evaluating MLLMs in the context of error detection scenarios, and therefore introduce the corre-077 sponding benchmark termed ERRORRADAR. We have designed two sub-tasks to comprehensively assess the performance: error step identification and error categorization. To construct a rich and 079 reliable dataset, we initially sourced a collection of multimodal K-12 level math problems from an educational organization and subsequently refined the dataset through rigorous manual annotation to 081 ensure quality. In particular, we also collect real students' answers for each multimodal question for a relatively robust experimental setting, compared to other relevant benchmarks (as shown in Table 083 1). Furthermore, we categorized the dataset to better align with diverse needs as follows: **Problem** 084 types: plane geometry, solid geometry, diagram, algebra, and mathematical common sense; and 085 Error categories: visual perception errors, calculation errors, reasoning errors, knowledge errors, and misinterpretation of the problem. In summary, the ERRORRADAR comprises 2,500 high-quality instances derived from real-life problem-solving data, providing a foundational dataset to enhance 087 the complex reasoning capabilities of MLLMs for the research community and industry. 880

For ERRORRADAR, we carry out an extensive experimental analysis to determine the proficiency 090 in complex mathematical reasoning of various MLLMs. The evaluation encompasses both the latest open-source MLLMs (e.g., InternVL2 (Chen et al., 2023b), LLaVA-NEXT (Liu et al., 2024a), 091 CogVLM2 (Wang et al., 2023a)), and closed-source MLLMs (e.g., GPT4-0 (OpenAI, 2024a), Gem-092 ini Pro 1.5 (Reid et al., 2024), Claude 3.5 (Anthropic, 2024b)). Our focus was on their error de-093 tection capabilities, specifically the identification of the erroneous step and the classification of the 094 error type. To establish a comparative human performance standard, we involved expert human ed-095 ucators who possess a graduate-level degree or higher qualifications. The results demonstrate that 096 ERRORRADAR, covering cutting-edge topics such as MLLMs' complex reasoning, poses a significant challenge, with human evaluation for two error detection tasks achieving less than 70%. 098

From in-depth evaluation of representative MLLMs, we obtain the following findings: **1** Closedsource MLLMs, particularly GPT-40, consistently outperform open-source MLLMs in both sub-100 tasks, and show more balanced accuracy across different error categories; @ Weaker MLLMs exhibit 101 an over-reliance on simpler categories, while stronger models handle complex tasks better; 3 Both 102 MLLMs and humans perform better on error step identification compared to error categorization, as 103 localizing specific errors is inherently simpler than categorizing errors. 104

- Our contributions can be summarized as follows: 105
- **0** We take the **first step to formulate the multimodal error detection task**, and introduce a mul-106 timodal benchmark termed ERRORRADAR for evaluation. This benchmark serves as a standard 107 operator for assessing the complex mathematical reasoning capabilities of the latest MLLMs.

We meticulously curate an extensive dataset comprising approximately 2,500 high-quality instances with rigorous annotation and rich metadata derived from real user interactions in an educational organization. To the best of our knowledge, this is the first attempt to use real-world student problem-solving data to evaluate MLLMs, providing a protocol for future research on MLLMs' complex mathematical reasoning.

- Our comprehensive experimental evaluation of more than 20 MLLMs, both proprietary and opensource, highlight the substantial room for improvement (*i.e.*, 7%-15% in accuracy) in the complex mathematical reasoning capabilities, underscoring the necessity for further research.
- 115 116

108

109

110

111

112

113

114

117 2

RELATED WORK

118 Benchmarks for Mathematical Reasoning. Recent advancements in mathematical reasoning 119 benchmarks have led to the development of both pure text and multimodal assessments (Lu et al., 120 2022; Wang et al., 2024a; Zheng et al., 2024; Huo et al., 2024). While datasets like GSM8K (Cobbe 121 et al., 2021), MATH (Hendrycks et al., 2021), SuperCLUE-Math (Xu et al., 2024), and MathBench 122 (Liu et al., 2024b) focus on text-based problems, the field has expanded to include multimodal 123 benchmarks that introduce visual elements, pushing the boundaries of AI's mathematical under-124 standing. For instance, MathVista (Lu et al., 2024b) evaluates AI's performance on visual math 125 OA tasks; MATH-V (Wang et al., 2024a) focuses on multimodal mathematical understanding with competition-derived questions; MathVerse (Zhang et al., 2024) assesses visual diagram compre-126 hension using CoT strategies; CMMU (He et al., 2024b) tests multi-disciplinary, multimodal math 127 understanding with a broad range of Chinese-language questions; MathScape (Zhou et al., 2024) 128 further advances the field by presenting longer, more complex, and open-ended multimodal prob-129 lems; and MMMU(Yue et al., 2024) covers college-level knowledge including interleaved mathe-130 matical questions. The aforementioned benchmarks assess the mathematical reasoning capabilities 131 of MLLMs by evaluating their problem-solving levels, but they overlook tasks based on the student's 132 perspective, such as error detection, and thus fail to comprehensively evaluate the more complex role 133 of current MLLMs. Therefore, we propose the ERRORRADAR benchmark, which is entirely based 134 on real student response data to evaluate the proficiency of MLLMs in error detection tasks.

135 Multimodal Large Language Models. Generative foundation models such as GPT-4 (OpenAI, 136 2023), Claude (Anthropic, 2024b), and Gemini (Pal & Sankarasubbu, 2024) have significantly ad-137 vanced various task solutions without fine-tuning (Cui et al., 2024; Yan & Lee, 2024; Zou et al., 138 2025; Zhong et al., 2024). Similarly, current open-source MLLMs, built on top of powerful LLMs, 139 have also demonstrated promising potential in multimodal tasks such as image captioning (Yang 140 et al., 2024) and visual question answering (Fan et al., 2024). For instance, LLaVA-NEXT (Liu 141 et al., 2024a) proposed projecting visual embeddings, extracted by a pretrained vision encoder, into 142 the word space through a single MLP layer, where LLMs like LLaMA, Vicuna, and Mistral are finetuned to understand these post-projection tokens. In a similar fashion, Phi3 (Abdin et al., 2024), 143 DeepSeek-VL (Lu et al., 2024a), MiniCPM-V (Yao et al., 2024), ChatGLM (GLM et al., 2024), 144 CogVLM (Wang et al., 2023a), Intern-VL (Chen et al., 2023b), Qwen-VL (Bai et al., 2023) and 145 Yi-VL (Young et al., 2024) also utilize a projector (or adapter, shared compression layer, etc.) to 146 align the visual embeddings extracted from a vision encoder with text embeddings, which are then 147 concatenated and fed into LLM. Therefore, we propose ERRORRADAR, a comprehensive bench-148 mark on a fine-grained evaluation of MLLMs' ability to detect errors based on students' answers 149 and reasoning steps, thereby advancing the development of complex multimodal system. 150

3 THE ERRORRADAR DATASET

3.1 TASK FORMULATION

Basic Setting. In this task, we assess the model's ability to detect errors in mathematical problemsolving processes across multiple samples. Let N denote the total number of samples in the evaluation set. For each sample $i \in \{1, 2, ..., N\}$, the input set \mathcal{I}_i is defined as:

$$\mathcal{I}_i = \{Q_{\text{text},i}, Q_{\text{image},i}, A_{\text{correct},i}, A_{\text{incorrect},i}, \{S_{k,i}\}_{k=1}^{n_i}\},\$$

159 where:

151

152

153 154

155

156

157 158

160

- $Q_{\text{text},i}$: the textual statement of the *i*-th problem.
 - $Q_{\text{image},i}$: the image representation of the *i*-th problem.



210		
217	Statistic	Number
218	Total multimodal questions	2,500
219	Problem Type	
220	- Plane Geometry	1559 (62.4%)
221	- Solid Geometry	191 (7.6%)
	- Diagram	233 (9.3%)
222	- Algebra	288 (11.5%)
223	- Math Commonsense	229 (9.2%)
224	Error Category	
225	- Visual Perception Error	395 (15.8%)
200	- Calculation Error	912 (36.5%)
226	 Reasoning Error 	951 (38.0%)
227	 Knowledge Error 	119 (4.8%)
228	- Misinterpretation of the Qns	123 (4.9%)
229	Average Reasoning Step	7.6
230	Maximum Reasoning Step	20
	Minimum Reasoning Step	3
231	Average Question Length	168
232	Maximum Question Length	719
233	Minimum Question Length	13

Table 2: Key statistics of ERRORRADAR.

234

235

237

238

239 240

241

242 243

244 245

246

247



Figure 3: Roadmap of ERRORRADAR dataset collection, annotation, and consistent update.



Figure 4: Dataset distribution of ERRORRADAR with respect to problem type and error category.

Performance Metric. The evaluation of both subtasks is conducted separately. The model's performance is evaluated using accuracy metrics for both subtasks:

• Error Step Identification Accuracy. Let $G_{\text{step},i}$ be the ground truth index of the first incorrect step for the *i*-th sample. The accuracy for this subtask is:

$$\operatorname{Acc}_{\operatorname{step}} = \frac{1}{N} \sum_{i=1}^{N} \mathbb{I}(x_i = G_{\operatorname{step},i}),$$

where $\mathbb{I}(\cdot)$ is indicator function, returning 1 if prediction matches ground truth, and 0 otherwise.

• Error Categorization Accuracy. Let G_{error,i} be the ground truth error category for the *i*-th sample. The accuracy for this subtask is:

$$\operatorname{Acc}_{\operatorname{cate}} = \frac{1}{N} \sum_{i=1}^{N} \mathbb{I}(C_{\operatorname{error},i} = G_{\operatorname{error},i}).$$

3.2 **DATA SOURCE & ANNOTATION**

253 Following the roadmap shown in Figure 3, this section includes how we collect and annotate ER-254 RORRADAR dataset to ensure the overall data quality. Different from the conventional benchmarks 255 that rely on public datasets or modified textbook collections (Lu et al., 2024b; Zhou et al., 2024), 256 ERRORRADAR dataset is uniquely sourced from the question bank of a global educational organi-257 zation. This repository encompasses a vast array of mathematical problems in K-12 levels, totaling 258 over a million entries. Initially, we curated approximately 180,000 math problems that feature a 259 single image-based question stem, aligning with our goal to target a multimodal assessment setup.

260 Subsequently, we refined our selection by evaluating the universality and articulation of the prob-261 lem content. For each problem, we identified multiple incorrect answers. To ensure the dataset's 262 relevance for error detection tasks, we selected the most frequently given incorrect answer as the 263 student's response. Additionally, we scrutinized cases where the most common incorrect answer was due to system input errors despite the answer being correct. In such instances, we amended the 264 dataset by incorporating the next most frequently incorrect answer. 265

266 Furthermore, since error detection tasks necessitate a step-by-step reasoning process, we enriched our dataset with new content through manual annotation. Specifically, we provided professional 267 annotators with the original multimodal QA data, student's incorrect answers, and the pedagogical 268 team's analysis of correct answer process. Based on this initial data, annotators delineated the 269 erroneous steps leading to the incorrect answers (More details in Appendix B.1 and B.2).

Our team of annotators, consisting of around ten educational experts with domain expertise, conducted two rounds of cross-checking to ensure the reliability of the annotations. In cases of inconsistency, the contentious question and related data were presented to the annotation lead for final adjudication. The annotators' results were subject to review and quality control by the educational organization from which the data originated, ensuring security, reliability, and consistent updates.

2752763.33.3DATASET DETAILS

As illustrated in Table 2, ERRORRADAR dataset comprises a substantial collection of 2,500 multi-277 modal math questions designed for error detection tasks. It predominantly includes plane geometry 278 problems, with solid geometry, diagram, algebra, and math commonsense questions making up the 279 remainder, highlighting its focus on diverse mathematical problems. it also categorizes errors into 280 visual perception, calculation, reasoning, knowledge, and question misinterpretation. Key statistics 281 indicate a diverse dataset with an average reasoning step of 7.6, a variety of question lengths, and 282 a wide range of reasoning steps (up to 38). Detailed distribution of ERRORRADAR, problem type 283 definition, and error category formulation can be seen in Figure 4, Appendix B.3 and B.4. 284

4 EXPERIMENTS AND ANALYSIS

287 4.1 EVALUATION PROTOCOLS

In ERRORRADAR benchmark, we propose an evaluation strategy using template matching rules.
 The evaluation process consists of three stages: *response generation, answer extraction,* and *performance calculation*. Initially, the MLLMs generate responses given the inputs, which incorporates the multimodal mathematical question, wrong answer, and its step-by-step reasoning, using the template from Appendix C.2. Subsequently, the short answer text can be extracted from the detailed response.
 Finally, the model performance is based on the detailed score calculation as shown in Section 3.1.
 The final score will be calculated by averaging the scores from three rounds of assessment.

297 4.2 EXPERIMENTAL SETUP

298 In our experimental setup, we meticulously categorized and evaluated a diverse array of MLLMs 299 into three distinct groups to assess their capabilities across error detection tasks. (i) The **Open**-300 Source MLLMs category encompassed models such as InternVL-2 (Chen et al., 2023b), Phi-3-301 vision (Abdin et al., 2024), Yi-VL (Young et al., 2024), DeepSeek-VL (Lu et al., 2024a), LLaVA-302 v1.6-Vicuna (Liu et al., 2024a), MiniCPM-LLaMA3-V2.5 (Yao et al., 2024), MiniCPM-V2.6 (Yao 303 et al., 2024), Qwen-VL (Bai et al., 2023), GLM-4v (GLM et al., 2024), and LLaVA-NEXT (Liu 304 et al., 2024a), each demonstrating their unique strengths and capabilities in handling different types of errors. (ii) The Closed-Source MLLMs featured proprietary models like Qwen-VL-Max (Bai 305 et al., 2023), Claude-3-Haiku (Anthropic, 2024a), Claude-3.5-Sonnet (Anthropic, 2024b), Gemini-306 Pro-1.5 (Reid et al., 2024), GPT-4o-mini (OpenAI, 2024b), and GPT-4o (OpenAI, 2024a), providing 307 a comparison point for the performance of models that are not publicly accessible. (iii) Lastly, the 308 **Human Performance** category served as a benchmark for natural intelligence, allowing us to gauge 309 how closely MLLMs can emulate human cognitive functions across tasks such as visual perception 310 (More details in Appendix C.1). We provide the prompts for MLLMs and sources of MLLMs in 311 Appendix C.2 and C.3, respectively.

312 313

285

286

296

4.3 EXPERIMENTAL RESULTS

4.3.1 MAIN RESULTS

316 Finding #1: Closed-source MLLMs generally outperform open-source MLLMs in both error 317 detection tasks, with GPT-40 demonstrating the strongest performance. Table 3 shows that 318 closed-source MLLMs generally outperform open-source MLLMs in both STEP and CATE tasks, 319 and they also exhibit relatively more balanced performance across the five error categories. This 320 superiority can likely be attributed to the proprietary datasets and advanced training resources avail-321 able to closed-source models, which allow for more robust fine-tuning (Shi et al., 2023; Yu et al., 2024; Wang et al., 2023b). Notably, GPT-40 stands out as the best model, achieving the highest 322 scores not only in STEP and CATE tasks, but also in the VIS, REAS, and MIS categories, demon-323 strating its overall versatility and strength. Given the current performance gap between open-source

Multimodal Large Language Models	Parameters	LLM	STEP	CATE	VIS	CAL	REAS	KNOW	MIS
	Open	-Source MLLMs							
InternVL2 (Chen et al., 2023b)	2B	InternLM-2	9.8	25.1	32.2	38.8	12.2	0.0	24.4
Phi-3-vision (Abdin et al., 2024)	4B	Phi-3	37.5	40.7	9.6	99.6	6.6	3.4	4.1
Yi-VL (Young et al., 2024)	6B	Yi	15.7	32.1	9.1	77.1	4.9	14.3	0.0
DeepSeek-VL (Lu et al., 2024a)	7B	DeepSeek	16.2	35.7	4.6	90.9	0.4	28.6	6.5
LLaVA-v1.6-Vicuna (Liu et al., 2024a)	7B	Vicuna-v1.5	30.3	17.7	40.3	14.9	8.3	0.0	55.3
InternVL-2 (Chen et al., 2023b)	8B	InternLM-2.5	44.2	44.1	12.4	99.6	13.6	10.9	2.4
MiniCPM-LLaMA3-V2.5 (Yao et al., 2024)	8B	LLaMA3	37.4	38.0	4.1	100.0	2.1	2.5	0.0
MiniCPM-V2.6 (Yao et al., 2024)	8B	Qwen2	17.0	39.8	11.4	87.8	12.1	10.1	17.9
Qwen-VL (Bai et al., 2023)	9B	Qwen	23.8	38.9	8.6	99.1	3.5	0.0	0.8
GLM-4v (GLM et al., 2024)	13B	GLM-4	44.6	44.1	2.5	92.9	25.8	0.0	0.0
LLaVA-v1.6-Vicuna (Liu et al., 2024a)	13B	Vicuna-v1.5	36.9	47.8	0.0	74.5	53.7	0.8	2.4
CogVLM2-LLaMA3 (Wang et al., 2023a)	19B	LLaMA3	15.0	20.1	43.3	33.8	0.7	13.4	0.0
InternVL2 (Chen et al., 2023b)	26B	InternLM-2	50.4	51.2	39.2	84.6	35.6	0.8	10.6
LLaVA-NEXT (Liu et al., 2024a)	72B	Qwen1.5	51.8	45.0	7.1	86.0	32.0	7.6	0.8
InternVL2 (Chen et al., 2023b)	76B	Hermes-2 Theta	54.4	49.5	33.4	92.4	25.1	10.9	8.1
	Closed	d-Source MLLMs							
Qwen-VL-Max (Bai et al., 2023)	-	-	48.7	52.9	15.2	78.9	50.5	14.3	36.6
Claude-3-Haiku (Anthropic, 2024a)	-	-	45.6	48.0	10.4	77.4	46.8	4.2	1.6
Claude-3.5-Sonnet (Anthropic, 2024b)	-	-	50.2	49.5	35.7	48.4	64.8	21.0	11.4
Gemini-Pro-1.5 (Reid et al., 2024)	-	-	55.0	52.7	43.5	55.7	63.1	18.5	13.0
GPT-4o-mini (OpenAI, 2024b)	-	-	52.0	44.5	9.1	46.8	62.7	31.9	13.0
GPT-4o (OpenAI, 2024a)	-	-	55.1	53.1	46.3	50.4	64.9	9.2	46.3
		Human							
Human performance	-	-	69.8	60.7	66.8	75.9	47.6	35.3	53.7

Table 3: Comparison of open-source and closed-source MLLM performance (accuracy in percentage) across
 error detection tasks. We denote STEP and CATE for the performance of error step identification task (*i.e.*, Acc_{step}) and error categorization task (*i.e.*, Acc_{cate}), respectively. We also denote VIS, CAL, REAS, KNOW,
 and MIS for visual perception error, calculation error, reasoning error, knowledge error, and misinterpretation
 of the question. The highest and second highest scores (except for exceptional values) among MLLMs in each
 column are highlighted in red and blue, respectively. Exceptional values in CAL column are highlighted in
 grey, as more than 70% categories predicted by the MLLM are CAL (More analysis on Sec 4.3.1 Finding #2).

351 352

353

and closed-source MLLMs, open-source MLLMs can further enhance themselves by distilling the error detection capabilities of closed-source ones (Hsieh et al., 2023).

354 Finding #2: Weak open-source MLLMs 355 tend to predict CAL category, leading 356 to unusually high performance. Table 357 3 indicates that MLLMs with relatively 358 low performance in the CATE task tend to exhibit unusually high performance in 359 the CAL category. Specifically, open-360 source models like MiniCPM-LLaMA3-361 v2.5 even achieve a 100% accuracy in 362 CAL, while Phi-3-vision and InternVL-2-363 8B reach 99.6%. Upon analyzing the cat-364 egory prediction proportions of CAL from Figure 5 (See details of all MLLMs in Ap-366 pendix C.4), it becomes clear that open-367 source MLLMs with the top five CAL ac-368 curacy predict over 80% of instances as



Figure 5: The proportion of CAL predictions of closedsource and open-source MLLMs with top-5 CAL accuracy.

CAL category, suggesting an over-reliance on this category. In contrast, closed-source MLLMs 369 with top-five CAL accuracy do not exhibit this extreme trend of prediction bias. This phenomenon 370 likely arises from weaker MLLMs attempting to overfit on the CAL category, a relatively simpler 371 classification, to compensate for their inability to handle more complex scenarios (Tirumala et al., 372 2022; Xu et al., 2021). Models exhibiting this phenomenon can assign different weights to samples 373 of different categories during training to reduce the model's preference for a particular category. 374 This can be achieved by adjusting the weight parameters in the loss function (e.g., Focal Loss & 375 AdaFocal) (Li et al., 2022; Ghosh et al., 2022). 376

Finding #3: MLLMs with strong overall performance tend to handle STEP easier than CATE. From Table 3, the best open-source MLLMs, such as InternVL2-76B, and the best closed-source



MLLMs, like GPT-40, exhibit a tendency where their STEP performance surpasses that of CATE. 392 This trend holds even for human performance, where accuracy on STEP is higher (69.8%) compared 393 to CATE (60.7%). The reason for this disparity is likely that identifying the error step is inherently 394 easier, as it involves localizing a specific point of failure. On the other hand, categorizing the error 395 requires more complex reasoning and contextual understanding to classify the nature of the error, 396 which adds difficulty. This mirrors the settings in object detection, where localization (*i.e.*, pre-397 dicting where an object is) is relatively simpler than classification (*i.e.*, predicting what an object 398 is) (Zou et al., 2023; Jiao et al., 2021). To improve the performance of error categorization tasks, 399 MLLMs need to better understand the relationship between the problem itself and the steps where 400 errors occur. Thus, modeling this part of the relationship can be a focus in the design of training 401 data (Ling et al., 2023; Shi et al., 2024a).

402 Finding #4: CAL is the easiest category for MLLMs, while KNOW is the most difficult. CAL 403 is the category with the highest performance among most MLLMs (excluding those with excep-404 tional values), which could be attributed to the structured and deterministic nature of calculations, 405 where errors often result in clear, quantifiable deviations from expected outcomes, making them 406 more straightforward to detect (Lewkowycz et al., 2022; Kojima et al., 2022). Conversely, KNOW 407 stands out as the most challenging category, suggesting that MLLMs struggle significantly with tasks requiring deep factual understanding and contextual reasoning. The complexity of knowledge 408 errors likely stems from the need for comprehensive domain expertise, which current MLLMs may 409 not fully encapsulate yet. Even human performance reflects this trend, with knowledge error scoring 410 notably lower than other categories, albeit with higher accuracy than MLLMs, highlighting the in-411 herent difficulty of this task for both humans and AI (Kandpal et al., 2023; Feng et al., 2023). Thus, 412 adding domain-specific knowledge to the dataset of MLLM is a direct solution (Ling et al., 2023). 413

414 Finding #5: MLLMs still have a gap to close to reach human-level intelligence in error detection. Human performance significantly outperforms the best MLLMs in both the STEP and CATE 415 tasks, with accuracy scores of 69.8% and 60.7% respectively, compared to the highest MLLM scores 416 of 55.1% and 53.1%. Notably, the detection of VIS by humans is markedly superior to the best 417 MLLMs, with a difference of nearly 20%. This substantial lead may be attributed to the sophisti-418 cated pattern recognition inherent to human visual processing (Doerig et al., 2022), which MLLMs, 419 despite their advancements, have yet to fully emulate. Besides, it is interesting to note that human 420 performance in REAS detection is lower than all closed-source MLLMs but higher than almost 421 all open-source MLLMs. This may suggest that closed-source MLLMs benefit from proprietary 422 datasets and algorithms that better capture the nuances of logical reasoning (Wang et al., 2024d). To 423 achieve human-level performance in error detection tasks, we can further introduce a reinforcement 424 learning from human feedback (RLHF) approach, enabling the model to align with human thinking 425 mechanisms in understanding error causes (Liu et al., 2023).

426 427 428

4.3.2 VISUAL PERCEPTION ANALYSIS

Finding #1: Closed-source MLLMs are most likely to misjudge VIS as REAS in error catego rization task. Taking the best-performing GPT-40 model as an example, as shown in the Figure 6,
 48% of VIS are misclassified as REAS, followed by 30% being misjudged as MIS. In multimodal mathematical scenarios, where the MLLM needs to handle information involving both visual and

linguistic elements simultaneously, particularly in problems related to plane and solid geometry, the
 complexity of the diagrams makes it difficult for the model to accurately extract certain features,
 leading to the frequent misclassification of VIS as REAS. For instance, if an erroneous response to a
 mathematical query originates from VIS (*e.g.*, misinterpreting a diagram), MLLM may mistakenly
 attribute this to a flaw in logical reasoning that occurs subsequent to initial visual misinterpretation.

437 Finding #2: Open-source MLLMs are more likely to misclassify VIS as CAL. Taking the open-438 source model CogVLM2-LLaMA3, which performs best in identifying VIS, as an example, CAL 439 accounts for 64% of misclassified category, as illustrated in the Figure 7. When handling complex 440 visual information, especially in geometry problems, the MLLM often struggles to accurately ex-441 tract key features. Due to the open-source MLLM's weaker multimodal integration capabilities, it 442 simplifies visual issues into numerical calculation problems. The lack of sufficient training and data for visual-related errors is also a key reason behind this phenomenon (Wichmann & Geirhos, 2023). 443 More analysis on misclassification for each category can be seen in Appendix C.5. 444

445 446

4.3.3 RELATION BETWEEN ERROR CATEGORY AND ERROR STEP

Finding #1: There is a close relationship between different error category and their distribution in the reasoning steps. As shown in Figure 8, VIS tends to occur in the earlier to mid-stages, accounting for a median proportion of 0.5 of total steps. In contrast, MIS, REAS, CAL, and KNOW are more likely to arise in the later stages, with their median proportions ranging from 0.7 to 0.9. More analysis of this relationship across MLLMs can be seen in Appendix C.6.

Finding #2: VIS occurs in the earlier stages of problem-solving reasoning. This finding could be closely linked to the sequence in which students approach the task (Binz & Schulz, 2023; Kennedy & Romig, 2024). Since image content often serves as a key reference early on, any misinterpretation of this visual information directly impacts the subsequent problem-solving steps. Students typically first examine the image, and then integrate the information before proceeding to reasoning or calculation. As a result, visual perception errors arise earlier compared to other types of errors.

Finding #3: Other error categories are primarily in later stages of problem-solving reasoning. 459 This may be linked to the increasing cognitive load students encounter during problem-solving. 460 Cognitive Load Theory posits that information complexity ranges from low to high interactivity 461 (Paas et al., 2010; Binz & Schulz, 2023). While low-interactivity information can be understood 462 independently, high-interactivity information requires simultaneous processing of related elements, 463 thus increasing cognitive load (Kennedy & Romig, 2024; Abbad-Andaloussi et al., 2023). In later 464 stages, students must integrate complex information from multiple sources. For instance, calculating 465 the distance between two points needs increasing interactivity heightens cognitive load, leading to errors like forgetting to take the square root or miscalculating differences. Consequently, as cognitive 466 load rises, the frequency of errors in later steps also increases. 467

468 4.3.4 SCALING ANALYSIS

469 Finding #1: The performance of MLLMs on STEP 470 task increases with the scale of parameters. We ob-471 serve a phenomenon similar to the scaling law (Kaplan 472 et al., 2020) in our experiments. As shown in Figure 473 9, when the size of the InternVL2 model increases from 474 Tiny to Huge, the accuracy of STEP task rises from 9.8% 475 to 54.4%, showing an improvement of 44.6%. Similarly, 476 as the size of LLaVA-NEXT increases from Small to Large, its accuracy of STEP also improves from 30.3% to 477 51.8%, demonstrating that larger MLLMs exhibit greater 478 reasoning ability in localizing erroneous steps. 479

Finding #2: CATE task is relatively more difficult to
improve through scaling. While the accuracy of CATE
shows a trend of improvement for both the InternVL2 and
LLaVA-NEXT models as their size increases from Tiny
(Small) to Middle, a slight decrease is also observed when





the model size reaches Large. We presume that this is because CATE is a more challenging task compared to STEP, and merely increasing the model size without fine-tuning is insufficient for





512

513

sustained improvement and may even introduce bias (Aghajanyan et al., 2023; Muennighoff et al., 2024). This phenomenon can also be seen in Table 3, consistent with Section 4.3.1 Finding #2.

514 4.3.5 VISUAL PERCEPTION CASE STUDY

515 Visual perception errors are critical in multimodal error detection tasks, as they impact the accurate 516 comprehension of mathematical problems presented with both text and diagrams. As illustrated in 517 Figure 10 and Appendix C.7, the five primary categories of visual errors observed in GPT-40 (the 518 MLLM with best overall and VIS performance) include distance perception, diagram percep-519 tion, spatial perception, flip/fold perception, and shape perception. These categories differ in 520 their cognitive demands: distance perception focuses on point identification; diagram perception on 521 quantitative estimation; spatial perception on geometric visualization; flip/fold perception on mental 522 rotation; and shape perception on object classification (Lu et al., 2024b; Zhang et al., 2024). Detecting such errors is challenging because they often require both intricate visual processing and precise 523 interpretation of mathematical relations, which can be difficult to encode in current MLLMs. To 524 overcome these challenges, future MLLMs should incorporate more advanced visual reasoning capabilities, possibly through enhanced alignment between vision and language modalities, enabling 526 better detection and correction of complex perception errors (Song et al., 2023). This could signifi-527 cantly improve the robustness of MLLMs in mathematical and other perception-heavy tasks. 528

529

530 5 CONCLUSION

531 In conclusion, this work introduces ERRORRADAR, the first multimodal benchmark designed specif-532 ically for evaluating MLLMs's reasoning in mathematical error detection scenarios. By focusing on 533 both error step identification and error categorization, ERRORRADAR bridges a critical research 534 gap in assessing MLLMs' capabilities in complex mathematical reasoning. The dataset's construction, based on real-world student interactions, ensures a robust evaluation framework that reflects genuine user needs. Our extensive experimental analysis, comparing leading open-source and pro-536 prietary MLLMs, reveals significant challenges in error detection, highlighting the need for con-537 tinued advancements in this domain. As MLLMs continue to evolve, ERRORRADAR serves as an 538 essential benchmark for driving improvements in the effectiveness of multimodal reasoning systems in real-world applications, on the path to Artificial General Intelligence.

540 REFERENCES

549

556

560

561 562

563

564 565

566

567 568

569

570

- Amine Abbad-Andaloussi, Andrea Burattin, Tijs Slaats, Ekkart Kindler, and Barbara Weber. Complexity in declarative process models: Metrics and multi-modal assessment of cognitive load. *Expert Systems with Applications*, 233:120924, 2023.
- Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.
- Armen Aghajanyan, Lili Yu, Alexis Conneau, Wei-Ning Hsu, Karen Hambardzumyan, Susan Zhang, Stephen Roller, Naman Goyal, Omer Levy, and Luke Zettlemoyer. Scaling laws for generative mixed-modal language models. In *International Conference on Machine Learning*, pp. 265–279. PMLR, 2023.
- Janice Ahn, Rishu Verma, Renze Lou, Di Liu, Rui Zhang, and Wenpeng Yin. Large language models
 for mathematical reasoning: Progresses and challenges. *arXiv preprint arXiv:2402.00157*, 2024.

557 Anthropic. Claude 3, 2024a. URL https://www.anthropic.com/news/ claude-3-haiku.

- Anthropic. Claude 3.5, 2024b. URL https://www.anthropic.com/news/ claude-3-5-sonnet.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023.
- Marcel Binz and Eric Schulz. Turning large language models into cognitive models. *arXiv preprint arXiv:2306.03917*, 2023.
- Wenhu Chen, Ming Yin, Max Ku, Pan Lu, Yixin Wan, Xueguang Ma, Jianyu Xu, Xinyi Wang, and Tony Xia. Theoremqa: A theorem-driven question answering dataset. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 7889–7901, 2023a.
- 571
 572
 573
 574
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 574
 575
 575
 575
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
 576
- 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.
- 579
 580
 581
 581
 582
 583
 583
 584
 585
 585
 585
 586
 586
 587
 588
 588
 588
 588
 588
 588
 588
 588
 589
 589
 580
 580
 581
 581
 582
 583
 583
 584
 584
 584
 585
 585
 586
 586
 587
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
- Alex Davies, Petar Veličković, Lars Buesing, Sam Blackwell, Daniel Zheng, Nenad Tomašev,
 Richard Tanburn, Peter Battaglia, Charles Blundell, András Juhász, et al. Advancing mathematics
 by guiding human intuition with ai. *Nature*, 600(7887):70–74, 2021.
- Adrien Doerig, Tim C Kietzmann, Emily Allen, Yihan Wu, Thomas Naselaris, Kendrick Kay, and Ian Charest. Visual representations in the human brain are aligned with large language models. *arXiv preprint arXiv:2209.11737*, 2022.
- Yue Fan, Jing Gu, Kaiwen Zhou, Qianqi Yan, Shan Jiang, Ching-Chen Kuo, Yang Zhao, Xinze Guan, and Xin Wang. Muffin or chihuahua? challenging multimodal large language models with multipanel vqa. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 6845–6863, 2024.

594	Zhangyin Feng, Weitao Ma, Weijiang Yu, Lei Huang, Haotian Wang, Qianglong Chen, Weihua
595	Peng, Xiaocheng Feng, Bing Qin, et al. Trends in integration of knowledge and large language
596	models: A survey and taxonomy of methods, benchmarks, and applications. arXiv preprint
597	arXiv:2311.05876, 2023.
598	

- Deqing Fu, Ghazal Khalighinejad, Ollie Liu, Bhuwan Dhingra, Dani Yogatama, Robin Jia, and Willie Neiswanger. Isobench: Benchmarking multimodal foundation models on isomorphic rep-600 resentations. arXiv preprint arXiv:2404.01266, 2024. 601
- 602 Arindam Ghosh, Thomas Schaaf, and Matthew Gormley. Adafocal: Calibration-aware adaptive focal loss. Advances in Neural Information Processing Systems, 35:1583–1595, 2022.
- Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu 605 Feng, Hanlin Zhao, Hanyu Lai, et al. Chatglm: A family of large language models from glm-130b 606 to glm-4 all tools. arXiv preprint arXiv:2406.12793, 2024. 607
- 608 Xixuan Hao, Wei Chen, Yibo Yan, Siru Zhong, Kun Wang, Qingsong Wen, and Yuxuan Liang. Urbanvlp: A multi-granularity vision-language pre-trained foundation model for urban indicator 609 prediction. arXiv preprint arXiv:2403.16831, 2024. 610
- 611 Jinlong He, Pengfei Li, Gang Liu, Zixu Zhao, and Shenjun Zhong. Pefomed: Parameter efficient 612 fine-tuning on multimodal large language models for medical visual question answering. arXiv 613 preprint arXiv:2401.02797, 2024a. 614
- Zheqi He, Xinya Wu, Pengfei Zhou, Richeng Xuan, Guang Liu, Xi Yang, Qiannan Zhu, and Hua 615 Huang. Cmmu: A benchmark for chinese multi-modal multi-type question understanding and 616 reasoning. arXiv preprint arXiv:2401.14011, 2024b. 617
- 618 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, 619 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. arXiv 620 preprint arXiv:2103.03874, 2021. 621
- Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alexander Rat-622 ner, Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. Distilling step-by-step! outperform-623 ing larger language models with less training data and smaller model sizes. arXiv preprint 624 arXiv:2305.02301, 2023. 625
- 626 Yushi Hu, Weijia Shi, Xingyu Fu, Dan Roth, Mari Ostendorf, Luke Zettlemoyer, Noah A Smith, 627 and Ranjay Krishna. Visual sketchpad: Sketching as a visual chain of thought for multimodal language models. arXiv preprint arXiv:2406.09403, 2024. 628
- 629 Jiahao Huo, Yibo Yan, Boren Hu, Yutao Yue, and Xuming Hu. Mmneuron: Discovering 630 neuron-level domain-specific interpretation in multimodal large language model. arXiv preprint 631 arXiv:2406.11193, 2024. 632
- 633 Mengzhao Jia, Zhihan Zhang, Wenhao Yu, Fangkai Jiao, and Meng Jiang. Describe-then-reason: Improving multimodal mathematical reasoning through visual comprehension training. arXiv 634 preprint arXiv:2404.14604, 2024. 635
- 636 Licheng Jiao, Ruohan Zhang, Fang Liu, Shuyuan Yang, Biao Hou, Lingling Li, and Xu Tang. New 637 generation deep learning for video object detection: A survey. IEEE Transactions on Neural 638 *Networks and Learning Systems*, 33(8):3195–3215, 2021. 639
- Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. Large language 640 models struggle to learn long-tail knowledge. In International Conference on Machine Learning, 641 pp. 15696-15707. PMLR, 2023. 642
- 643 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, 644 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language 645 models. arXiv preprint arXiv:2001.08361, 2020. 646
- Michael J Kennedy and John Elwood Romig. Cognitive load theory: An applied reintroduction for 647 special and general educators. TEACHING Exceptional Children, 56(6):440-451, 2024.

684

686

687

688

689

- 648 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large 649 language models are zero-shot reasoners. Advances in neural information processing systems, 650 35:22199–22213, 2022.
- Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ra-652 masesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative 653 reasoning problems with language models. Advances in Neural Information Processing Systems, 654 35:3843-3857, 2022. 655
- 656 Hang Li, Tianlong Xu, Chaoli Zhang, Eason Chen, Jing Liang, Xing Fan, Haoyang Li, Jiliang Tang, 657 and Qingsong Wen. Bringing generative ai to adaptive learning in education. arXiv preprint 658 arXiv:2402.14601, 2024a. 659
- Xiang Li, Chengqi Lv, Wenhai Wang, Gang Li, Lingfeng Yang, and Jian Yang. Generalized focal 660 loss: Towards efficient representation learning for dense object detection. IEEE transactions on 661 pattern analysis and machine intelligence, 45(3):3139–3153, 2022. 662
- 663 Xiaoyuan Li, Wenjie Wang, Moxin Li, Junrong Guo, Yang Zhang, and Fuli Feng. Evaluating mathe-664 matical reasoning of large language models: A focus on error identification and correction. arXiv 665 preprint arXiv:2406.00755, 2024b. 666
- 667 Zhong-Zhi Li, Ming-Liang Zhang, Fei Yin, Zhi-Long Ji, Jin-Feng Bai, Zhen-Ru Pan, Fan-Hu Zeng, Jian Xu, Jia-Xin Zhang, and Cheng-Lin Liu. Cmmath: A chinese multi-modal math skill evalua-668 tion benchmark for foundation models. arXiv preprint arXiv:2407.12023, 2024c. 669
- 670 Chen Ling, Xujiang Zhao, Jiaying Lu, Chengyuan Deng, Can Zheng, Junxiang Wang, Tanmoy 671 Chowdhury, Yun Li, Hejie Cui, Xuchao Zhang, et al. Domain specialization as the key to make 672 large language models disruptive: A comprehensive survey. arXiv preprint arXiv:2305.18703, 673 2023. 674
- 675 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 676 Llava-next: Improved reasoning, ocr, and world knowledge, January 2024a. URL https:// llava-vl.github.io/blog/2024-01-30-llava-next/. 677
- 678 Hongwei Liu, Zilong Zheng, Yuxuan Qiao, Haodong Duan, Zhiwei Fei, Fengzhe Zhou, Wen-679 wei Zhang, Songyang Zhang, Dahua Lin, and Kai Chen. Mathbench: Evaluating the theory 680 and application proficiency of llms with a hierarchical mathematics benchmark. arXiv preprint 681 arXiv:2405.12209, 2024b. 682
- 683 Wenhao Liu, Xiaohua Wang, Muling Wu, Tianlong Li, Changze Lv, Zixuan Ling, Jianhao Zhu, Cenyuan Zhang, Xiaoqing Zheng, and Xuanjing Huang. Aligning large language models with human preferences through representation engineering. arXiv preprint arXiv:2312.15997, 2023. 685
 - Xiao Liu, Zirui Wu, Xueqing Wu, Pan Lu, Kai-Wei Chang, and Yansong Feng. Are Ilms capable of data-based statistical and causal reasoning? benchmarking advanced quantitative reasoning with data. arXiv preprint arXiv:2402.17644, 2024c.
- 690 Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, 691 Zhuoshu Li, Yaofeng Sun, et al. Deepseek-vl: towards real-world vision-language understanding. 692 arXiv preprint arXiv:2403.05525, 2024a.
- 693 Pan Lu, Liang Qiu, Wenhao Yu, Sean Welleck, and Kai-Wei Chang. A survey of deep learning for 694 mathematical reasoning. arXiv preprint arXiv:2212.10535, 2022. 695
- 696 Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-697 Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of 698 foundation models in visual contexts. arXiv preprint arXiv:2310.02255, 2024b.
- Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, 700 and Jianfeng Gao. Chameleon: Plug-and-play compositional reasoning with large language mod-701 els. Advances in Neural Information Processing Systems, 36, 2024c.

732

- Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. Large language models: A survey. *arXiv preprint arXiv:2402.06196*, 2024.
- Niklas Muennighoff, Alexander Rush, Boaz Barak, Teven Le Scao, Nouamane Tazi, Aleksandra Piktus, Sampo Pyysalo, Thomas Wolf, and Colin A Raffel. Scaling data-constrained language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- 709 OpenAI. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 711 OpenAI. GPT-4V(ision) system card, 2024a. URL https://openai.com/index/ 712 gpt-4o-system-card/.
- 713 OpenAI. Gpt-40 mini: advancing cost-efficient intelligence, 2024b. URL https://openai. com/index/gpt-40-mini-advancing-cost-efficient-intelligence/.
- Fred Paas, Alexander Renkl, and John Sweller. Cognitive load theory and instructional design:
 Recent developments. *Educational Psychologist*, 2010.
- Ankit Pal and Malaikannan Sankarasubbu. Gemini goes to med school: exploring the capabilities
 of multimodal large language models on medical challenge problems & hallucinations. *arXiv preprint arXiv:2402.07023*, 2024.
- Shuai Peng, Di Fu, Liangcai Gao, Xiuqin Zhong, Hongguang Fu, and Zhi Tang. Multimath: Bridging visual and mathematical reasoning for large language models. *arXiv preprint arXiv:2409.00147*, 2024.
- Runqi Qiao, Qiuna Tan, Guanting Dong, Minhui Wu, Chong Sun, Xiaoshuai Song, Zhuoma GongQue, Shanglin Lei, Zhe Wei, Miaoxuan Zhang, et al. We-math: Does your large multi-modal model achieve human-like mathematical reasoning? *arXiv preprint arXiv:2407.01284*, 2024.
- Annabelle Rabillas, Osias Kit Kilag, Neil Cañete, Maria Trazona, Mery Lou Calope, and Jacqueline
 Kilag. Elementary math learning through piaget's cognitive development stages. *Excellencia: International Multi-disciplinary Journal of Education (2994-9521)*, 1(4):128–142, 2023.
- Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jeanbaptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- Haizhou Shi, Zihao Xu, Hengyi Wang, Weiyi Qin, Wenyuan Wang, Yibin Wang, Zifeng Wang,
 Sayna Ebrahimi, and Hao Wang. Continual learning of large language models: A comprehensive
 survey. *arXiv preprint arXiv:2404.16789*, 2024a.
- Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi
 Chen, and Luke Zettlemoyer. Detecting pretraining data from large language models. *arXiv* preprint arXiv:2310.16789, 2023.
- Wenhao Shi, Zhiqiang Hu, Yi Bin, Junhua Liu, Yang Yang, See-Kiong Ng, Lidong Bing, and Roy
 Ka-Wei Lee. Math-Ilava: Bootstrapping mathematical reasoning for multimodal large language
 models. *arXiv preprint arXiv:2406.17294*, 2024b.
- Shezheng Song, Xiaopeng Li, and Shasha Li. How to bridge the gap between modalities: A comprehensive survey on multimodal large language model. *arXiv preprint arXiv:2311.07594*, 2023.
- Liangtai Sun, Yang Han, Zihan Zhao, Da Ma, Zhennan Shen, Baocai Chen, Lu Chen, and Kai Yu. Scieval: A multi-level large language model evaluation benchmark for scientific research. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 19053–19061, 2024.
- Kushal Tirumala, Aram Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. Memorization
 without overfitting: Analyzing the training dynamics of large language models. Advances in Neural Information Processing Systems, 35:38274–38290, 2022.

778

784

- Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Mingjie Zhan, and Hongsheng Li. Measuring multimodal mathematical reasoning with math-vision dataset. *arXiv preprint arXiv:2402.14804*, 2024a.
- Shen Wang, Tianlong Xu, Hang Li, Chaoli Zhang, Joleen Liang, Jiliang Tang, Philip S Yu, and
 Qingsong Wen. Large language models for education: A survey and outlook. *arXiv preprint arXiv:2403.18105*, 2024b.
- 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, 2023a.
- Xiao Wang, Guangyao Chen, Guangwu Qian, Pengcheng Gao, Xiao-Yong Wei, Yaowei Wang,
 Yonghong Tian, and Wen Gao. Large-scale multi-modal pre-trained models: A comprehensive
 survey. *Machine Intelligence Research*, 20(4):447–482, 2023b.
- Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R
 Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. Scibench: Evaluating college-level scientific problem-solving abilities of large language models. *arXiv preprint arXiv:2307.10635*, 2024c.
- Yiqi Wang, Wentao Chen, Xiaotian Han, Xudong Lin, Haiteng Zhao, Yongfei Liu, Bohan Zhai, Jianbo Yuan, Quanzeng You, and Hongxia Yang. Exploring the reasoning abilities of multimodal large language models (mllms): A comprehensive survey on emerging trends in multimodal reasoning. *arXiv preprint arXiv:2401.06805*, 2024d.
- Felix A Wichmann and Robert Geirhos. Are deep neural networks adequate behavioral models of human visual perception? *Annual Review of Vision Science*, 9(1):501–524, 2023.
- Hanguang Xiao, Feizhong Zhou, Xingyue Liu, Tianqi Liu, Zhipeng Li, Xin Liu, and Xiaoxuan Huang. A comprehensive survey of large language models and multimodal large language models in medicine. *arXiv preprint arXiv:2405.08603*, 2024.
- Liang Xu, Hang Xue, Lei Zhu, and Kangkang Zhao. Superclue-math6: Graded multi-step math
 reasoning benchmark for llms in chinese. *arXiv preprint arXiv:2401.11819*, 2024.
- Runxin Xu, Fuli Luo, Zhiyuan Zhang, Chuanqi Tan, Baobao Chang, Songfang Huang, and Fei Huang. Raise a child in large language model: Towards effective and generalizable fine-tuning. *arXiv preprint arXiv:2109.05687*, 2021.
- Yibo Yan and Joey Lee. Georeasoner: Reasoning on geospatially grounded context for natural language understanding. *arXiv preprint arXiv:2408.11366*, 2024.
- Yibo Yan, Haomin Wen, Siru Zhong, Wei Chen, Haodong Chen, Qingsong Wen, Roger Zimmermann, and Yuxuan Liang. Urbanclip: Learning text-enhanced urban region profiling with contrastive language-image pretraining from the web. In *Proceedings of the ACM on Web Conference*2024, pp. 4006–4017, 2024.
- Xu Yang, Yongliang Wu, Mingzhuo Yang, Haokun Chen, and Xin Geng. Exploring diverse incontext configurations for image captioning. *Advances in Neural Information Processing Systems*, 36, 2024.
- Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li,
 Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv preprint arXiv:2408.01800*, 2024.
- Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, et al. Yi: Open foundation models by 01. ai. *arXiv preprint arXiv:2403.04652*, 2024.
- Yue Yu, Yuchen Zhuang, Jieyu Zhang, Yu Meng, Alexander J Ratner, Ranjay Krishna, Jiaming Shen,
 and Chao Zhang. Large language model as attributed training data generator: A tale of diversity
 and bias. Advances in Neural Information Processing Systems, 36, 2024.

810 811 812 813	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 multi-modal understanding and reasoning benchmark for expert agi. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 9556–9567, 2024.
814 815 816 817	Zhongshen Zeng, Yinhong Liu, Yingjia Wan, Jingyao Li, Pengguang Chen, Jianbo Dai, Yuxuan Yao, Rongwu Xu, Zehan Qi, Wanru Zhao, et al. Mr-ben: A comprehensive meta-reasoning benchmark for large language models. <i>arXiv preprint arXiv:2406.13975</i> , 2024.
818 819 820	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.
821 822 823	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. <i>arXiv</i> preprint arXiv:2303.18223, 2023.
825 826 827	Kening Zheng, Junkai Chen, Yibo Yan, Xin Zou, and Xuming Hu. Reefknot: A comprehensive benchmark for relation hallucination evaluation, analysis and mitigation in multimodal large language models. <i>arXiv preprint arXiv:2408.09429</i> , 2024.
828 829 830	Siru Zhong, Xixuan Hao, Yibo Yan, Ying Zhang, Yangqiu Song, and Yuxuan Liang. Urban- cross: Enhancing satellite image-text retrieval with cross-domain adaptation. <i>arXiv preprint</i> <i>arXiv:2404.14241</i> , 2024.
831 832 833 834	Minxuan Zhou, Hao Liang, Tianpeng Li, Zhiyu Wu, Mingan Lin, Linzhuang Sun, Yaqi Zhou, Yan Zhang, Xiaoqin Huang, Yicong Chen, et al. Mathscape: Evaluating mllms in multimodal math scenarios through a hierarchical benchmark. <i>arXiv preprint arXiv:2408.07543</i> , 2024.
835 836 837	Yutao Zhu, Huaying Yuan, Shuting Wang, Jiongnan Liu, Wenhan Liu, Chenlong Deng, Haonan Chen, Zhicheng Dou, and Ji-Rong Wen. Large language models for information retrieval: A survey. <i>arXiv preprint arXiv:2308.07107</i> , 2023.
838 839 840	Wenwen Zhuang, Xin Huang, Xiantao Zhang, and Jin Zeng. Math-puma: Progressive upward multimodal alignment to enhance mathematical reasoning. <i>arXiv preprint arXiv:2408.08640</i> , 2024.
842 843 844	Xingchen Zou, Yibo Yan, Xixuan Hao, Yuehong Hu, Haomin Wen, Erdong Liu, Junbo Zhang, Yong Li, Tianrui Li, Yu Zheng, et al. Deep learning for cross-domain data fusion in urban computing: Taxonomy, advances, and outlook. <i>Information Fusion</i> , 113:102606, 2025.
845 846 847 848	Zhengxia Zou, Keyan Chen, Zhenwei Shi, Yuhong Guo, and Jieping Ye. Object detection in 20 years: A survey. <i>Proceedings of the IEEE</i> , 111(3):257–276, 2023.
849 850 851	
852 853 854	
855 856	
857 858	
859 860 861	
862 863	

MORE MULTIMODAL QUESTION EXAMPLES А

Question (Visual Part):	Question (Textual Part):
	The diagram has small squares.
 Student's Step-by-Step Solution Step 1: Read the question and the Step 2: Identify that there are 5 g Step 3: Identify that there are 9 in Step 5: Combine the 5 hundreds, Step 6: Output the result as 529. 	 ion (Reasoning Process): e text in the diagram "The diagram below has small blocks." groups of small blocks in the diagram, each containing 100, representing 5 hundreds. 2 strips of small blocks in the diagram, each containing 10, representing 2 tens. ndividual small blocks in the diagram, representing 9 ones. , 2 tens, and 9 ones to write: 529.
True Answer: 549 Student's Answer: 529	ErrorRadar Error Category Error Step GPT-40 Visual Perception Error Step 3
Problem Type: Counting	GPT-40-mini Knowledge Error Step 2
Error Category: Visual Perception	Error Gemini-Pro-1.5 Visual Perception Error Step 2
T (1) (1)	

Figure 11: Multimodal mathematical example one (type: counting) from ERRORRADAR dataset.







Figure 13: Multimodal mathematical example three (type: plane geometry) from ERRORRADAR dataset.



Question (Visual Part):	Question (Textual P	'art):		
	As shown in the figu	re, a large rectang	le is divided into 6 id	entical smaller
24cm	rectangles. If the leng area of the large recta	gth of a smaller reangle is sq	ctangle is 24 centime juare centimeters.	ters, then the
dent's Sten-by-Sten Solution	n (Regsoning Process).		
ep 1: Read the problem statemer	nt and the text in the diagr	am, "Divide a large i	rectangle into 6 identical	smaller
ectangles." Step 2: Determine that the length o	of each smaller rectangle	is 24 centimeters.		
Step 3: From the diagram, it is kno	own that the length of the	smaller rectangle is	times its width.	
Step 4: Calculate the width of the s	maller rectangle, 24 ÷ 4 maller rectangle, 24 × 6 =	= 6 (centimeters). = 144 (square centime	eters).	
• Step 6: Calculate the area of the	large rectangle, which i	s the sum of the are	as of 6 smaller rectang	les, 144 × 6 = 764
Step 7: The output result is 764.	ita			
rue Answer: 864		ErrorRadar	Error Category	Error Step
udent's Answer: 764		GPT-40	Calculation Error	Step 6
lem Type: Plane Geometry		GPT-40-mini	Reasoning Error 🥮	Step 6
Category: Calculation Error	G	emini-Pro-1.5	Calculation Error	Step 3 🚳
rror Step: Step 6		LLaVA-v1.6	Reasoning Error 🏁	Step 3
Figure 15: Multimodal mathe	matical example five	(type: plane geor	netry) from ERRORI	RADAR dataset.

1134 B ADDITIONAL DATASET DETAILS

1136 B.1 ANNOTATION DETAILS

To ensure the quality and relevance of the ERRORRADAR dataset for error detection tasks, we employed a rigorous manual annotation process, involving professional educational experts as annotators. This section outlines the details of the annotation procedure, focusing on how the data was enriched with step-by-step reasoning processes, identification of erroneous steps, and error categorization.

Annotator Selection and Training. Given the complexity of the task, we recruited a group of ten annotators with specialized knowledge in educational theory and mathematics, particularly in K-12 pedagogy. These annotators were trained extensively to familiarize themselves with the structure and expectations of the task. The training covered the specifics of multimodal problem-solving in mathematics, typical student error patterns, and the need for precise identification of reasoning steps that led to incorrect answers. The annotators were also briefed on using the provided tools and the quality assurance process.

Annotation Process. Each mathematical problem in the dataset was annotated with a step-by-step
 reasoning process, capturing both correct and incorrect approaches to problem-solving. Annotators
 were provided with:

1153

1158

1160

1161

1162

1163

- 1. The original question stem (comprising both text and image components).
- 11541. The original question stem (comprising both t11552. The student's most frequent incorrect answer.
- 1156 3. The correct answer to the question.
- 1157 4. The pedagogical analysis of the correct reasoning process, prepared by educational experts.
- 1159 Based on these inputs, annotators were tasked with:
 - 1. **Step-by-Step Reasoning Annotation**: For each problem, annotators mapped out the logical steps that students should ideally follow to arrive at the correct answer. This involved identifying key stages in the problem-solving process, such as formula application, arithmetic operations, or logical deductions.
- 2. Error Step Identification: For problems where students provided incorrect answers, annotators identified the exact steps where the reasoning went wrong. These error steps were explicitly marked and linked to the incorrect responses, ensuring that they could be traced back to specific problem-solving mistakes.
- 3. Error Categorization: Once the erroneous step was identified, annotators assigned an appropriate error category based on a predefined schema. These categories included common types of errors such as misinterpretation of the question (More details can be seen in Section 3.1). The categorization was designed to align with known student error patterns in mathematical learning.

Quality Control and Cross-Validation. To ensure annotation accuracy and consistency, each problem underwent two rounds of cross-checking:

- First Round of Cross-Validation: After the initial annotation, another annotator independently reviewed the annotations. Any discrepancies between the first and second annotators were flagged for further analysis.
- 1179
 1180
 1180
 1181
 1182
 1183
 2. Second Round of Cross-Validation: In the second round, if the errors or discrepancies persisted, the problem was escalated to a senior educational expert who acted as the annotation lead. The annotation lead adjudicated these contentious cases, ensuring that the final decision was both pedagogically sound and aligned with the dataset's goals.

1184 Dataset Refinement. Throughout the annotation process, we worked closely with the educational
1185 organization from which the dataset originated. This collaboration ensured that the annotations were
1186 not only reliable but also adhered to the standards of the organization's question bank. Addition1187 ally, ongoing feedback and updates from the organization helped refine the dataset, making it more
accurate and relevant for multimodal error detection tasks.

1188 Annotation Duration and Effort. The annotation process for the ERRORRADAR dataset spanned 1189 over a period of at least two months. During this time, the annotators, comprised of both profes-1190 sional educational experts and domain specialists, worked meticulously through several stages of 1191 preparation, annotation, and validation. Each annotator dedicated significant time to understanding 1192 the dataset, reviewing the provided pedagogical analyses, and applying their domain knowledge to identify and categorize errors. The first phase, involving step-by-step reasoning annotation, took ap-1193 proximately six weeks, while the subsequent cross-validation and quality control efforts accounted 1194 for the remaining two weeks. Given the complexity of the tasks and the necessity for high precision, 1195 the team's sustained efforts ensured that the final dataset was of the highest quality. 1196

By incorporating these annotations, ERRORRADAR provides a robust foundation for studying student errors in mathematical reasoning and enables the development of advanced models for error detection and correction.

1200

1207

1209

1210

1211

1212

1213

1216

1217

1218

1219

1220

1222

1225

1226 1227

1228

1229

1230

1231

1232

1233

1237

1239

1240

1201 B.2 DETAILS OF HANDLING INCONSISTENT ANNOTATIONS

To ensure the quality and reliability of our dataset for the multimodal mathematical error detection task, we established a systematic approach to resolve annotation inconsistencies. This process balances annotator independence with rigorous quality control, ensuring that the dataset is both accurate and representative.

- 1208 B.2.1 ANNOTATION AGREEMENT PRINCIPLES
 - 1. **Guided Consensus**: Annotations must align with clear, predefined guidelines covering the five error categories. Annotators are trained extensively to reduce subjective biases.
 - 2. Cross-Checking and Agreement Threshold: Each instance is annotated by at least three annotators. Disagreements are flagged for further review.
- **3. Systematic Review Process**: For inconsistent cases, a multi-step resolution process is applied:
 - (a) **Initial Review**: Annotators discuss disagreements, referencing annotation guidelines and the specific problem context.
 - (b) **Expert Arbitration**: For unresolved cases, a domain expert (e.g., an educational professional) reviews and finalizes the annotation.
 - (c) **Consensus-Driven Decisions**: When possible, annotations are harmonized based on majority opinion or shared agreement after discussions.
- 1224 B.2.2 CASE RESOLUTION FRAMEWORK

Case Example 1: Visual Perception vs. Reasoning Error

- **Example**: A problem presents a bar chart requiring students to determine the highest value. A student misidentifies the tallest bar and selects the wrong answer.
 - Annotator A labels this as a Visual Perception Error, arguing the mistake stems from misreading the chart.
 - Annotator B classifies it as a Reasoning Error, interpreting the mistake as a failure to compare values logically.

• Resolution: Annotators revisit the problem:

- If evidence shows the student misunderstood the chart format (e.g., interpreting height as quantity but misjudging due to poor visualization), it is classified as a Visual Perception Error.
- If the student correctly interprets the chart but misapplies logical comparisons (e.g., failing to compare values explicitly), it is categorized as a Reasoning Error.
- For persistent disagreement, an expert examines the student's work, including any notes or intermediate steps, to determine the correct annotation.

1242	Case E	xample 2: Knowledge vs. Misinterpretation of the Question
1243		
1244	•	• Example: A problem asks for the perimeter of a rectangle, but the student calculates the
1245		area instead.
1246		- Annotator A identifies this as a Knowledge Error, attributing the mistake to a lack of
1247		understanding of perimeter concepts.
1248		- Annotator B labels it as a Misinterpretation of the Question, asserting that the student
1249		misunderstood what was being asked.
1250		• Resolution:
1251		- Did the student's work demonstrate understanding of the concept but apply it incor-
1252		rectly (Misinterpretation of the Ouestion)?
1253		- Did the mistake reveal a fundamental gap in knowledge about perimeter (Knowledge
1204		Error)?
1200		If disagreement persists, the annotators consult the expert, who may analyze additional
1250		context (e.g., previous responses or annotations).
1257		
1250	B.2.3	HANDLING IRRECONCILABLE DISAGREEMENTS
1260	TO 11	
1261	It discr	epancies persist despite review and arbitration, the affected data points are excluded from the
1262	dataset.	I have a semilar maintain high reliability
1263	Tetamet	i samples mantani ingri renaority.
1264	B 2 4	MONITORING AND FEEDBACK
1265	D.2.4	MONTORING AND I LEDDACK
1266	Periodi	c feedback sessions are conducted to recalibrate annotators and refine guidelines based on
1267	observe	ed patterns of disagreement. This iterative approach minimizes future inconsistencies and
1268	enhance	es annotator alignment over time.
1269		
1270		
1271		
1272		
1273		
1274		
1275		
1276		
1277		
1278		
1279		
1280		
1201		
1202		
1203		
1204		
1286		
1287		
1288		
1289		
1290		
1291		
1292		
1293		
1294		
1295		

1296 B.3 DEFINITION OF PROBLEM TYPE CATEGORY

1298 The ERRORRADAR dataset distinguishes five primary types of multimodal mathematical problems, 1299 each characterized by unique features:

 Plane Geometry Problems: These involve two-dimensional shapes and figures, requiring knowledge of properties such as angles, lines, and polygons. Solving these problems often depends on understanding basic geometric principles and theorems about plane figures.

* Solid Geometry Problems: In contrast to plane geometry, solid geometry involves three-dimensional objects, such as cubes, cylinders, and spheres. These problems require spatial visualization and understanding of volume, surface area, and the relationships between different three-dimensional shapes.

- * Diagram-Based Problems: These require analysis of provided visual information, such as graphs, charts, or diagrams, to solve mathematical queries. Interpreting visual data correctly is crucial, as these problems test the ability to extract and analyze quantitative information from visual aids.
- Algebra Problems: Algebra problems focus on abstract symbols and variables to represent numbers and relationships. These include tasks like solving equations, manipulating algebraic expressions, and understanding functions. The problem-solving process typically involves logical reasoning and manipulation of mathematical symbols.
- Math Commonsense Questions: These encompass a variety of problem types, including time judgment, direction judgment, counting, and pattern recognition. Unlike the other categories, math commonsense challenges rely on everyday mathematical reasoning and problem-solving strategies that do not necessarily require formal mathematical knowledge, testing intuitive understanding rather than procedural skills.

These problem types highlight the ERRORRADAR dataset's diverse nature, with each category presenting distinct challenges and requiring specific reasoning abilities.

1350 B.4 DEVELOPMENT AND VALIDATION PROCESS OF ERROR CATEGORY

1352 1. CROSS-TEAM COLLABORATION TO ALIGN TASK NEEDS

The process began with close collaboration between the research team and the education team to ensure that the error categories aligned with the unique requirements of the multimodal math error detection task. The research team provided insights into the task's technical objectives, focusing on precision and comprehensive error coverage. Simultaneously, the education team contributed their understanding of real-world educational scenarios, emphasizing the practical relevance and applicability of the error taxonomy to students' and teachers' needs.

1360 Key Outcomes:

1361

1362

1363

1364

1373

1374

1375

1376

1384

1385 1386

1387 1388

1389

1390

- Initial consensus that the categories must address both multimodal challenges and real-life classroom scenarios.
- Recognition of the need to balance academic rigor with user-friendly categorization.

1365 2. BENCHMARK SURVEY AND FOCUS ANALYSIS

The research team conducted an extensive survey of representative benchmarks, focusing on error
analysis frameworks in existing datasets. Examples included studies on problem-solving steps in
educational AI and cognitive error modeling in multimodal tasks. The aim was to identify gaps in
current frameworks and understand how existing taxonomies handle errors specific to visual, textual,
and logical reasoning elements.

1372 Key Outcomes:

- Identification of inadequacies in current benchmarks, particularly in addressing multimodal interactions like visual misinterpretations and reasoning errors tied to diagram-based tasks.
- Validation of the necessity for distinct categories to capture errors unique to multimodal math problems.

3. COLLECTION OF FEEDBACK FROM STUDENTS AND TEACHERS

The education team collected qualitative and quantitative feedback from students and teachers to
 ensure that the proposed error categories were grounded in real-world educational needs. Focus
 groups, surveys, and interviews were used to gather perspectives on common error patterns encountered during classroom activities and assessments.

Key Insights:

- Teachers highlighted frequent calculation errors (CAL) and reasoning errors (REAS) as significant roadblocks to effective problem-solving.
- Students often reported confusion stemming from visual misinterpretations (VIS) and misunderstanding the question intent (MIS).
- Feedback emphasized the importance of separating reasoning-based errors from knowledge-based errors (KNOW) for better diagnostic support.
- 1391 1392 1393

1399

1400

1401

1402

1403

4. SECOND ROUND OF DISCUSSION AND ALIGNMENT

Following the feedback collection, the research and education teams reconvened to refine and align the error taxonomy. This phase involved iterative discussions to ensure that each category was distinct, comprehensive, and intuitive for annotators and end-users.

1398 Adjustments Made:

- Clarified the scope of **Reasoning Errors** (**REAS**) to focus on improper logical application rather than factual knowledge gaps.
 Strangthaned the definition of Visual **Percention Errors** (VIS) to address multimodal.
 - Strengthened the definition of **Visual Perception Errors (VIS)** to address multimodalspecific challenges, such as interpreting diagrams or image-based data.
 - Enhanced examples for each category to support annotation clarity.

1404 1405	5. INITIAL FINALIZATION AND FEEDBACK FROM EDUCATIONAL ORGANIZATION
1406	The refined error categories were presented to a partner educational organization for feedback. This
1407	organization, which specializes in global education assessments, conducted an independent review
1408	and provided expert input.
1409	Key Outcomes:
1410	
1411	 Positive validation of the categories' relevance and comprehensiveness.
1412	• Minor recommendations, such as specifying units and signs in the Calculation Errors
1413	(CAL) category, were integrated.
1414	
1415	6. FINAL VALIDATION AND ALIGNMENT WITH ANNOTATION TEAM
1416	
1417	After incorporating feedback, the final set of error categories was finalized. The annotation team,
1418	cation of the taxonomy during the annotation process. Mock annotations were conducted to test the
1419	clarity and usability of the categories.
1420	
1421	Final Adjustments:
1422	• Annotators highlighted the need for clearer boundaries between Reasoning Errors
1423	(REAS) and Knowledge Errors (KNOW), leading to additional examples and decision
1425	rules in the annotation guidelines.
1426	• Alignment meetings ensured that all discrepancies and ambiguities were resolved before
1427	the dataset's official annotation began.
1428	
1429	The aforementioned development process ensured that the five categories are comprehensive, robust,
1430	and applicable to both multimodal tasks and real-world educational scenarios.
1431	
1432	
1433	
1434	
1435	
1436	
1437	
1438	
1439	
1440	
1441	
1442	
1443	
1445	
1446	
1447	
1448	
1449	
1450	
1451	
1452	
1453	
1454	
1455	
1456	
1457	

¹⁴⁵⁸ C Additional Experimental Details

1460 C.1 HUMAN PERFORMANCE EVALUATION

In the Human Performance section, the evaluation involved three educational expert evaluators, each independently performing the two subtasks — error step identification and error categorization — on a set of multimodal math problems. To ensure the validity of their assessments, a rigorous cross-checking procedure was implemented. After the initial independent evaluations, the results from all three experts were compared for both the identification of error steps and the categorization of those errors. When discrepancies arose, particularly in cases where the experts disagreed on which step contained an error or how an error should be classified, a structured conflict resolution process was followed.

The cross-check process began with identifying areas of disagreement between the evaluators. These conflicts were discussed in a series of consensus meetings, where the evaluators would review the conflicting steps or categorizations in detail. Each expert provided their rationale, referencing the mathematical principles involved as well as the multimodal representations of the problems. Through open dialogue, the evaluators aimed to reach a consensus on the correct interpretation of the error.

In cases where consensus could not be easily achieved, a majority-vote system was employed. However, for particularly complex or ambiguous cases, an additional adjudicator — who did not participate in the initial evaluations but had equivalent expertise — was consulted to provide a final judgment. This adjudicator reviewed the contentious cases along with the evaluators' justifications, ensuring an unbiased final decision. The outcome of this process was the creation of a refined ground truth dataset that balanced expert knowledge with the goal of consistent and reliable error identification and categorization.

¹⁵¹² C 1513	2.2 PROMPT FOR MLLM EVALUATION
1514	Teck Definition. Vou are an advantion expert preficient in K 12 methametics. Your teck
1516	is to identify the first step where the mistake occurred in the incorrect answer reasoning
1517	steps based on the following mathematical question (including the textual and visual parts)
1518	reference answer, and incorrect answer.
1519	
1520	Output format:
1521	Error Step: Step X
1522	Below is the reference content you need to identify the error sten:
523	Question Image: {image}
24	Question text: {content}
25	Correct Answer: {answer}
6	Incorrect Answer: {user_answer}
7	Incorrect Answer Reasoning Steps:{user_answer_steps}
3	Instruction : Please provide the corresponding error step identification in the format "Error
9	Step: Step X", without any additional content.
Ľ	
))	Figure 16: Prompt for error step identification task.

1619

1567 1568 1569 1570 1571 1572 1573 1574 1575 1576 1577 Task Definition: You are an education expert proficient in K-12 mathematics. Your task 1578 is to identify the category of error for the incorrect answer based on the following question 1579 (including the textual and visual parts), reference answer, and incorrect answer. The error 1580 should belong to one of the following categories: Visual Perception Error, Reasoning Error, 1581 Knowledge Error, Calculation Error, or Misinterpretation of the Question. **Output format:** 1584 Error Category: Clearly indicate which error category it belongs to. 1585 The definitions of the error categories are as follows: 1586 \star Visual Perception Error: Failure to accurately obtain information from the images or charts 1587 in the question due to visual issues, leading to errors. \star Reasoning Error: Improper reasoning during the problem-solving process, failure to cor-1589 rectly apply logical relationships or draw conclusions, leading to errors 1590 *Knowledge Error: Errors occur when applying relevant knowledge points due to incom-1591 plete or incorrect understanding of knowledge. 1592 \star Calculation Error: Errors occur in the calculation process, such as addition, subtraction, 1593 multiplication, division mistakes, or unit conversion errors, or errors in numerical symbols 1594 between multiple steps. *Misinterpretation of the Question: Failure to correctly understand the requirements of the 1595 question or misinterpreting the meaning of the question stem, leading to an irrelevant an-1596 swer, such as answering with numbers when letters are required, and vice versa. 1597 1598 Below is the reference content you need to identify the error step: Question Image: {image} Question text: {content} Correct Answer: {answer} Incorrect Answer: {user_answer} Incorrect Answer Reasoning Steps: {user_answer_steps} 1604 **Instruction:** Please provide the corresponding error category in the format "Error Category: X", without any additional content. 1608 Figure 17: Prompt for error categorization task. 1609 1610 1611 1612 1613 1614 1615 1616 1617 1618

1620 C.3 MODEL SOURCES

1622Table 4 details specific sources for the various MLLMs we evaluated. The hyperparameters for the1623experiments are set to their default values unless specified otherwise.

MLLMs	Source	URL
InternVL2-2B	local checkpoint	https://huggingface.co/ OpenGVLab/InternVL2-2B
InternVL2-8B	local checkpoint	https://huggingface.co/ OpenGVLab/InternVL2-8B
InternVL2-26B	local checkpoint	https://huggingface.co/ OpenGVLab/InternVL2-26B
InternVL2-76B	local checkpoint	https://huggingface. co/OpenGVLab/ InternVL2-Llama3-76B
Phi-3-vision-4B	local checkpoint	<pre>https://huggingface. co/microsoft/ Phi-3-vision-128k-instr</pre>
Yi-VL-6B	local checkpoint	https://huggingface.co/(Yi-VL-6B
DeepSeek-VL-7B	local checkpoint	https://huggingface. co/deepseek-ai/ deepseek-vl-7b-chat
LLaVA-v1.6- Vicuna-7B	local checkpoint	https://huggingface. co/llava-hf/llava-v1. 6-vicuna-7b-hf
LLaVA-v1.6- Vicuna-13B	local checkpoint	https://huggingface. co/llava-hf/llava-v1. 6-vicuna-13b-hf
LLaVA-NEXT-72B	local checkpoint	https://huggingface.co/ llava-hf/llava-next-72b
MiniCPM-V2.5-8B	local checkpoint	https://huggingface.co/ openbmb/MiniCPM-Llama3-V
MiniCPM-V2.6-8B	local checkpoint	https://huggingface.co/ openbmb/MiniCPM-V-2_6
Qwen-VL-9B	local checkpoint	https://huggingface.co/Ç Qwen-VL-Chat
GLM-4v-13B	local checkpoint	https://huggingface.co/1 glm-4v-9b
CogVLM2-19B	local checkpoint	https://huggingface.co/I cogvlm2-llama3-chat-19B
Qwen-VL-Max	qwen-vl-max-0809	https://modelscope.cn/ studios/qwen/Qwen-VL-Max
Claude-3-Haiku	claude-3-haiku	https://www.anthropic.co
Claude-3.5-Sonnet	claude-3-5-sonnet	https://www.anthropic.co
Gemini-Pro-1.5	gemini-1.5-pro-latest	https://deepmind.google/ technologies/gemini/pro/
GPT-4o-mini	gpt-4o-mini-2024-07-18	https://platform.openai. docs/models/gpt-4o-mini
GPT-40	gpt-40-2024-08-06	https://platform.openai.

1673

 Table 4: Sources of our evaluated MLLMs.

C.4 CAL AND NON-CAL DISTRIBUTION OF MLLMS

In this section, we indicate the distribution of CAL and non-CAL category predictions of 21 MLLMs we evaluate, as shown in Figure 18. It can be seen that there is a bias towards CAL category among most open-source MLLMs, while closed-source ones except for Claude-3-Haiku and Qwen-VL-Max do not have such a bias for error categorization task.



Figure 18: Distribution of CAL and non-CAL category predictions of all MLLMs we evaluate.

1728 C.5 ANALYSIS OF CONFUSION MATRIX FOR CATE TASK

Figures 19 and 20 present the confusion matrices for InternVL2-76B and GPT-40, two MLLMs
evaluated on five error categories. The matrices show the count of predictions for each category,
with diagonal entries representing correct predictions and off-diagonal entries indicating misclassifications. These visualizations provide insights into each model's strengths and weaknesses.

InternVL2-76B shows strong performance in detecting CAL, with 843 correct predictions, indicating its robust numerical reasoning capability. However, the model struggles to distinguish between
REAS and CAL, misclassifying 626 REAS instances as CAL. This confusion suggests an overreliance on numerical features and an inability to separate logical reasoning tasks from computational ones. Additionally, there is significant misclassification of VIS into CAL, with 244 cases,
highlighting a potential weakness in integrating visual and textual modalities. These trends may
stem from InternVL2-76B's limited domain-specific reasoning ability.

GPT-40, on the other hand, demonstrates relatively good performance in VIS, with 183 correct pre-dictions, significantly outperforming InternVL2-76B. Its capability in REAS is also notable, with 617 correct predictions, suggesting a more balanced reasoning ability. However, GPT-40 struggles more with CAL, achieving only 460 correct predictions, and shows significant confusion between CAL and REAS, with 299 CAL instances misclassified as REAS. Furthermore, the model has diffi-culty with MIS, misclassifying 45 MIS cases as REAS, pointing to challenges in identifying nuanced interpretational issues. These trends suggest that GPT-4o's emphasis on multimodal alignment and contextual understanding contributes to its strengths in VIS and REAS but comes at the expense of CAL performance.

Comparing the two models reveals distinct strengths and weaknesses. GPT-40 significantly outper-forms InternVL2-76B in VIS, likely due to superior multimodal visual-text alignment capabilities. Both models exhibit confusion between REAS and CAL, but GPT-40 shows a more balanced clas-sification ability in REAS. MIS remains a challenging category for both models, though GPT-40 struggles slightly more in distinguishing it from REAS. These differences may arise from varia-tions in model architecture and training objectives. This analysis underscores the complementary strengths of these models: InternVL2-76B excels in numerical reasoning, while GPT-40 performs better in visual perception and logical reasoning. Future research could explore ways to integrate their strengths for a more robust multimodal error detection system.



Figure 19: The confusion matrix of five error categories predicted by InternVL2-76B, the open-source MLLM with the best overall performance on error detection.



Figure 20: The confusion matrix of five error categories predicted by GPT-40, the closed-source MLLM with the best overall performance on error detection.

1836 C.6 COGNITIVE LOAD ANALYSIS ACROSS MLLMS

1838 In analyzing the error step distribution for the multimodal error detection task using InternVL2-76B (see Figure 21) and GPT-40 (see Figure 22), we observe a consistency in the pattern of error category distribution across both MLLM's predictions and those in ERRORRADAR (see Figure 8). In 1840 particular, VIS tends to occur in the earlier stages of problem-solving for both MLLMs, which aligns 1841 with the sequence in which students typically approach tasks. Since visual content often serves 1842 as a key reference at the outset, any misinterpretation of this information can significantly impact 1843 subsequent steps. Students generally examine the image first and then integrate the information 1844 before proceeding to reasoning or calculation, leading to visual perception errors arising earlier 1845 compared to other types of errors. 1846

Other error categories, such as REAS, CAL, MIS, and KNOW, are more likely to emerge in the later 1847 stages of problem-solving. This pattern is linked to the increasing cognitive load students encounter 1848 as they progress. According to Cognitive Load Theory, information complexity ranges from low 1849 to high interactivity. Low-interactivity information can be understood independently, whereas high-1850 interactivity information requires the simultaneous processing of related elements, thereby increas-1851 ing cognitive load. In the later stages, students must integrate complex information from multiple 1852 sources, which can lead to errors like forgetting to take the square root or miscalculating differ-1853 ences when calculating distances, for example. Consequently, the frequency of errors in later steps 1854 increases with the rising cognitive load.

1855
1856 Despite the overall pattern being consistent, there may be subtle differences between InternVL2-76B and GPT-40 in terms of error step distribution, especially for MIS category. These differences could be attributed to the models' distinct architectures and training data, which might influence their approaches to error detection. As an open-source MLLM, InternVL2-76B might not have been optimized for specific types of questions or educational contexts, which could lead to a higher variability in MIS.



Figure 21: The error step distribution (in percentage) of error categories predicted by InternVL2-76B, the open-source MLLM with the best overall performance on error detection.



Figure 22: The error step distribution (in percentage) of error categories predicted by GPT-40, the closed-source MLLM with the best overall performance on error detection.

888

1889

1884 1885

1875

1876

1890 C.7 VISUAL BAD CASES PREDICTED BY GPT-40 1891

Figures 23, 24 and 25 show bad cases where GPT-40 predicts visual perception errors incorrectly, in terms of diagram, spatial and flip & fold perception.

94			
95	Diagram Devention		
96	Diagram Perception		
97	Question: To understand the growth situation of a certain type of flower in a botanical		
98	garden, a sample of 200 flowers was randomly selected from a plot with approximately 3000		
99	such flowers, and their heights were measured. After organizing the statistical results, the data		
0	the maximum value) Estimate the number of such flowers in the plot with a height less than		
1	55 cm and not less than 45 cm to be about flowers.		
2			
3	If ue Answer: 900 Student's Answer: 1440 Height(cm) 40~45. 45~50. 50~55. 55~60. 60~65.		
4	Frequency 33. 42. 22. 24. 43. 36.		
	Reasoning Steps:		
	◆ step 1: Read the problem statement and chart text "Estimate the number of such nowers in the garden with a neight less than 55 cm and not less than 45 cm is about plants."		
	 vstep :: Determine the total sample size as 200 plants. ♦ Step 3: Determine the total number of flowers as 3000 plants. ♦ Step 4: Dead for the chort that the numbers of flowers with a bailed of 45.50 cm is 42 plants. 		
	 Step 5: Read from the chart that the number of flowers with a height of 50-55 cm is 32 plants. ♦ Step 5: Read from the chart that the number of flowers with a height of 55 cm is 32 plants. 		
	• Step 5: Calculate the proportion of flowers with a height between 45-55 cm in the sample as $72 / 200 = 0.36$. • Step 7: Calculate the proportion of flowers with a height between 45-55 cm in the sample as $72 / 200 = 0.36$. • Step 8: Multiply the total number of flowers. 3000 plants by the proportion in the sample as $72 / 200 = 0.36$.		
	♦ Step 9: Output the result as 1440. End Calculation Frame		
	Err. Cale. by Gr 1-40: Calculation Err.		
	Figure 23: Diagram bad case where GPT-40 predicts visual perception errors incorrectl		
	Spatial Perception		
	Question: Using small cubes with an edge length of 1 cm to form the geometric shape shown in the figure place this shape on a table and paint the exposed surfaces with color (evoluting the		
	bottom surface), then the total area of the painted surfaces is cm ² .		
	True Answer: 27		
	Student's Answer: 41		
	Keasoning Steps: • Step 1: Read the problem statement and the text in the diagram, "The geometric shape shown is made up of small cubes with an edge length of 1 cm. Place this		
	geometric shape on the table and paint the exposed surfaces (excluding the bottom surface), then the total area of the painted surfaces is $_1_$." \blacklozenge Step 2: Observe the structure of the geometric shape and confirm it is composed of small cubes.		
	 ♦ Step 3: From the front view, count the number of small squares, totaling 6. ♦ Step 4: From the left view, count the number of small squares, totaling 6. 		
	 ♦ Step 5: From the top view, count the number of small squares, totaling 8. ♦ Step 6: From the right view, count the number of small squares, totaling 6. 		
	 ♦ Step 7: From the back view, count the number of small squares, totaling 6. ♦ Step 8: Add the number of small squares from each face, totaling 6 + 6 + 8 + 6 + 6 = 32. 		
	 step 9: count the number of small squares on the bottom face, totaling 9. Step 10: Add the number of small squares on the bottom face to the previous total, 32 + 9 = 41. 		
	 ▼ step 11: vacuate the total area of the painted surfaces, which is the area of 41 small squares. ♦ Step 12: Output the result as 41. 		
	Err. Cate. by GPT-40: Misinterpretation of the Question		
	Figure 24: Spatial bad case where GPT-40 predicts visual perception errors incorrectly		
	Flip/Fold Perception		
	Question: A number is shown in the mirror as illustrated, then the actual number is		
	Student's Answer: 2505		
	Reasoning Steps:		
	 ♦ Step 1: Read the problem statement and the text in the image: "The number in the mirror is shown in the figure." ♦ Step 2: Observe that the number in the mirror is "5025." 		
	 Step 3: Based on the properties of mirror symmetry, the number seen in the mirror is actually flipped horizontally and vertically. Step 4: Flip the number "2502" horizontally to get the actual number "2505." 		
	• Step 5: Output the result as 2505. Frr Cate by GPT-40: Reasoning Frr		
	Figure 25: Flip & fold bad case where GPT-40 predicts visual perception errors incorrec		
	Figure 25: Flip & fold bad case where GPT-40 predicts visual perception errors incorrect		
	Figure 25: Flip & fold bad case where GPT-40 predicts visual perception errors incorrect		