

A Survey of Mathematical Reasoning in the Era of Multimodal Large Language Model: Benchmark, Method & Challenges

Anonymous ACL submission

Abstract

Mathematical reasoning, a core aspect of human cognition, is vital across many domains, from educational problem-solving to scientific advancements. As artificial general intelligence (AGI) progresses, integrating large language models (LLMs) with mathematical reasoning tasks is becoming increasingly significant. This survey provides the **first comprehensive analysis of mathematical reasoning in the era of multimodal large language models (MLLMs)**. We review over 200 studies published since 2021, and examine the state-of-the-art developments in Math-LLMs, with a focus on multimodal settings. We categorize the field into three dimensions: *benchmarks, methodologies, and challenges*. In particular, we explore multimodal mathematical reasoning pipeline, as well as the role of (M)LLMs and the associated methodologies. Finally, we identify five major challenges hindering the realization of AGI in this domain, offering insights into the future direction for enhancing multimodal reasoning capabilities. This survey serves as a critical resource for the research community in advancing the capabilities of LLMs to tackle complex multimodal reasoning tasks.

1 Introduction

Mathematical reasoning is a critical aspect of human cognitive ability, involving the process of deriving conclusions from a set of premises through logical and systematic thinking (Jonsson et al., 2022; Yu et al., 2024b). It plays an essential role in a wide range of applications, from problem-solving in education to advanced scientific discoveries. As artificial general intelligence (AGI) continues to advance (Zhong et al., 2024), the integration of large language models (LLMs) with mathematical reasoning tasks becomes increasingly significant. These models, with their impressive capabilities in language understanding, have the potential to simulate complex reasoning processes that were once

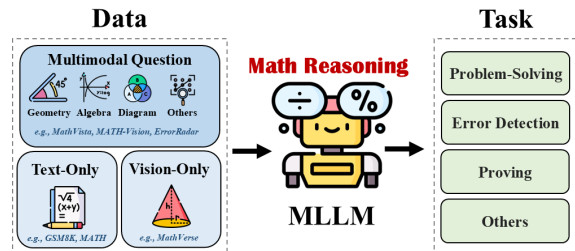


Figure 1: The illustration of our research scope (*i.e.*, investigating the MLLM’s math reasoning capability).

Survey	Venue & Year	Scope	Multimodal	LLM
(Lu et al., 2022b)	ACL’22	DL4Math		
(Li et al., 2023a)	arXiv’23	LLM4Edu		✓
(Liu et al., 2023b)	arXiv’23	LLM4Edu		✓
(Li et al., 2024f)	COLM’24	DL4TP		
(Ahn et al., 2024)	EACL’24	LLM4Math		✓
(Xu et al., 2024a)	IJMLC’24	LLM4Edu		✓
(Wang et al., 2024d)	arXiv’24	LLM4Edu		✓
Ours	-	MLLM4Math	✓	✓

Table 1: Comparisons between relevant surveys & ours.

thought to be inherently human. In recent years, both academia and industry have placed increasing emphasis on this direction (Wang et al., 2024d; Xu et al., 2024a; Lu et al., 2022b).

The inputs for mathematical reasoning tasks are diverse, extending beyond traditional text-only to multimodal settings, as illustrated in Figure 1. Mathematical problems often involve not only textual information but also visual elements, such as diagrams, graphs, or equations, which provide essential context for solving the problem (Wang et al., 2024e; Yin et al., 2024). In the past year, multimodal mathematical reasoning has emerged as a key focus for multimodal large language models (MLLMs) (Zhang et al., 2024c; Bai et al., 2024; Wu et al., 2023a). This shift is driven by the recognition that reasoning tasks in fields like mathematics require models capable of integrating and processing multiple modalities simultaneously to achieve human-like performance. However, multimodal mathematical reasoning poses significant

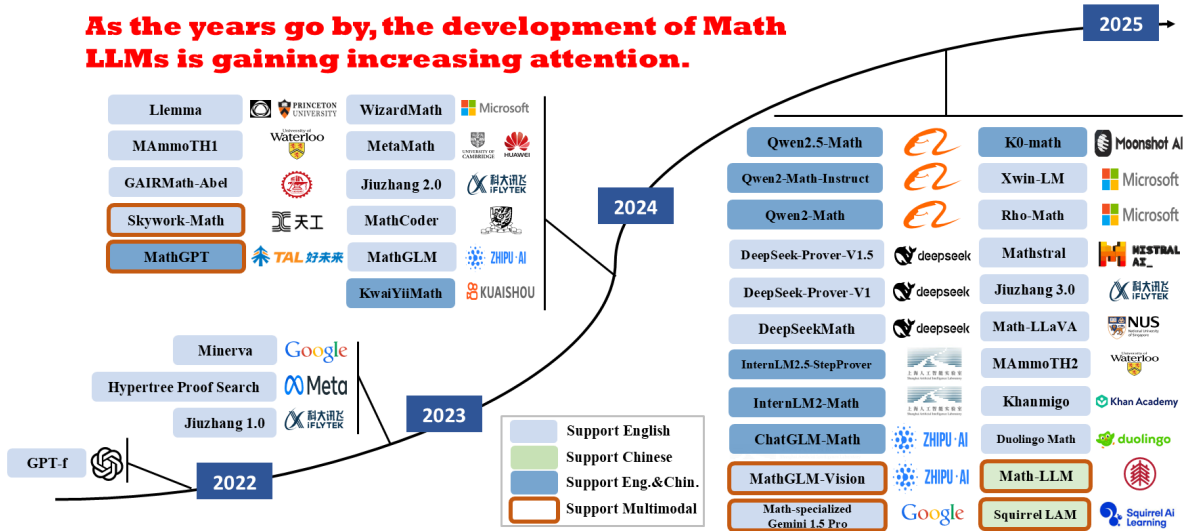


Figure 2: The release timeline of Math-LLMs in recent years.

challenges due to the complex interaction between different modalities, the need for deep semantic understanding, and the importance of context preservation across modalities (Liang et al., 2024a; Song et al., 2023; Fu et al., 2024b). These challenges are central to the realization of AGI, where models must integrate diverse forms of knowledge seamlessly to perform sophisticated reasoning tasks.

Math-LLM Progress. Figure 2 illustrates that, driven by the rapid development of LLMs since 2021, the number of math-specific LLMs (Math-LLMs) has grown steadily, alongside enhanced support for multilingual and multimodal capabilities (More details in Appendix A). The landscape was marked by the introduction of models like GPT-f (Polu and Sutskever, 2021) and Minerva (Lewkowycz et al., 2022), with Hypertree Proof Search (Lample et al., 2022) and Jiuzhang 1.0 (Zhao et al., 2022) highlighting advancements in theorem proving and mathematical question understanding capabilities, respectively. Year 2023 saw a surge in diversity and specialization, alongside multimodal support from models like Skywork-Math (Zeng et al., 2024). In year 2024, there was a clear focus on enhancing mathematical instruction (e.g., Qwen2.5-Math (Yang et al., 2024a)) and proof (e.g., DeepSeek-Proof (Xin et al., 2024a)) capabilities. The year also witnessed the emergence of Math-LLMs with a vision component, such as MathGLM-Vision (Yang et al., 2024b).

Scope. Previous surveys have not fully captured the progress and challenges of mathematical reasoning in the age of MLLMs. As indicated in Table

1, some works have concentrated on the application of deep learning techniques to mathematical reasoning (Lu et al., 2022b) or specific domains such as theorem proving (Li et al., 2024f), but they have overlooked the rapid advancements brought about by the rise of LLMs. Others have broadened the scope to include the role of LLMs in education (Wang et al., 2024d; Xu et al., 2024a; Li et al., 2023a) or mathematical fields (Ahn et al., 2024; Liu et al., 2023b), but have failed to explore the development and challenges of mathematical reasoning in multimodal settings in depth. Therefore, this survey aims to fill this gap by providing the **first-ever comprehensive analysis of the current state of mathematical reasoning in the era of MLLMs**, focusing on three key dimensions: *benchmark, methodology, and challenges*.

Structure. In this paper, we survey over 200 publications from the AI community since 2021 related to (M)LLM-based mathematical reasoning, and summarize the progress of Math-LLMs. We first approach the field from the benchmark perspective, analyzing the LLM-based mathematical reasoning task through three key aspects: dataset, task, and evaluation (Section 2). Subsequently, we explore the roles that (M)LLMs play in mathematical reasoning, categorizing them as enhancers, reasoners, and planners (Section 3). Finally, we identify five core challenges that the mathematical reasoning faces in the era of MLLMs (Section 4). This survey aims to provide the community with comprehensive insights for advancing the multimodal complex reasoning capabilities of LLMs.

2 Benchmark Perspective

2.1 Overview

Benchmarking for mathematical reasoning plays a crucial role in advancing LLM research, as it provides standardized, reproducible pipeline for assessing the performance on reasoning tasks. While previous benchmarks such as GSM8K (Cobbe et al., 2021) and MathQA (Amini et al., 2019) were instrumental in the pre-LLM era, our scope is centered on those relevant to (M)LLMs. In this section, we present a comprehensive analysis of recent benchmarks for mathematical reasoning in the context of (M)LLMs (Shown in Table 2). The section is organized into three subsections: Datasets (Sec.2.2), Tasks (Sec.2.3) & Evaluation (Sec.2.4).

2.2 Dataset

Basic Format. In a math reasoning task (taking problem-solving as a basic setting), the goal is to solve a mathematical problem given a specific format of input and output. The input consists of a statement that describes the problem to be solved. As illustrated in Figure 3, this can be presented in either a textual format (text-only) or a multimodal format (text accompanied by visual elements, such as figures or diagrams). The output is the correct or predicted solution to the mathematical problem, often represented as a numerical or symbolic result.

Language & Size. The majority of benchmarks are available in English, with a few exceptions like Chinese (Li et al., 2024h) or Romanian (Cosma et al., 2024) datasets. This predominance of English datasets underscores the challenges of multilingual representation in the mathematical reasoning domain, suggesting an opportunity for future work to diversify datasets across languages, especially those in underrepresented regions. Moreover, the size of these datasets varies widely, from smaller sets (e.g., QRData (Liu et al., 2024d) with 411 questions) to massive corpora (e.g., OpenMathInstruct-1 (Toshniwal et al., 2024) with 1.8 million problem-solution pairs). Larger datasets are more likely to support robust model training and evaluation, but their size can also present challenges in terms of computational requirements and quality control.

Source. The sources of datasets predominantly consist of public (i.e., derived from public repositories or datasets) and private sources. The private datasets typically offer specialized problem types and tasks, and may present unique challenges, such

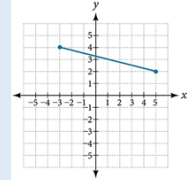
(a) Text-only Math Reasoning Setting

[Qns] Find the distance between the two endpoints using the distance formula. The two endpoints of the line are (-3, 4) and (5, 2), respectively.

[Ans] 8.246

(b) Multimodal Math Reasoning Setting

[Qns] Find the distance between the two endpoints.



[Ans] 8.246

Figure 3: Typical data format of math reasoning task for text-only & multimodal settings. Examples are derived from MathVerse (Zhang et al., 2024f), which assess whether and how much MLLMs can truly understand the visual diagrams for mathematical reasoning.

as restricted access or ethical considerations. On the other hand, public datasets foster wider community collaboration, though they may suffer from limitations in diversity and task coverage. Some works have also leveraged LLMs to generate the datasets tailored to specific needs. For instance, GeomVerse constructs synthetic datasets to evaluate the multi-hop reasoning abilities required in geometric math problems (Kazemi et al., 2023).

Educational Level. The benchmarks span various educational levels, ranging from elementary school to university-level problems. Besides, there has also been a surge in datasets focused on competition-level problems (Tsoukalas et al.), offering insights into the current limitations of LLMs in comparison to the upper bound of human cognitive abilities. Future directions could involve more focused datasets targeting specific educational levels to enable models to specialize in handling particular age groups or skill sets.

2.3 Task

Model Choice. The choice of models in these benchmarks spans open-source and closed-source models, with a growing interest in Math-LLMs. This trend indicates an increasing recognition of the need for models tailored to mathematical reasoning, which often require specialized training and handling of structured knowledge. Additionally, with the recent release of GPT-4o (OpenAI, 2024) and Gemini-Pro-1.5 (Reid et al., 2024), which have

Benchmarks	Venue	Language	Size	Source	Level(s)	Evaluation	Model(s)	Task(s)
GSM-Plus (Li et al., 2024d)	ACL'24	English	10,552	P	E M H U	Generative	Closed/Open/Math	S
MuggleMath (Li et al., 2024c)	ACL'24	English	37,365	P	E H	Discriminative	Open	S
OlympiadBench (He et al., 2024) ★	ACL'24	English/Chinese	8,476	S	H C	Generative	Closed/Open/Math	S
MathBench (Liu et al., 2024b)	ACL Findings'24	English/Chinese	3,709	P S	E M H U	Generative	Closed/Open/Math	S
GeoEval (Zhang et al., 2024d) ★	ACL Findings'24	English	5,050	P G	E M H U	Discriminative	Closed/Open/Math	S
QRData (Liu et al., 2024d)	ACL Findings'24	English	411	S	U	Discriminative	Closed/Open/Math	S
EIC-Math (Li et al., 2024e)	ACL Findings'24	English	1,800	P	E M H	Discriminative	Closed/Open	H D
Srivastava et al. (2024)	ACL Findings'24	English	270	P	H	Discriminative	Closed/Open	S
CHAMP (Mao et al., 2024)	ACL Findings'24	English	270	S	H	Generative	Closed/Open	S
IMO-AG-30 (Trinh et al., 2024)	Nature'24	English	30	S	C	Discriminative	Closed	P
PutnamBench (Tsoukalas et al.)	NeurIPS'24	English	1,697	S	C	Generative	Closed	S P
MATH-Vision (Wang et al., 2024a) ★	NeurIPS'24	English	3,040	S	E M H U	Discriminative	Closed/Open	S
CARP (Zhang et al., 2024a)	NeurIPS'24	Chinese	4,886	S	C	Discriminative	Closed	S
SMART-840 (Cherian et al., 2024) ★	NeurIPS'24	English	840	S	E M H U	Discriminative	Closed/Open	S
OpenMathInstruct-1 (Toshiwal et al., 2024)	NeurIPS'24	English	1,800,000	P	E M H C	Generative	Closed/Open/Math	S
Didolkar et al. (2024)	NeurIPS'24	English	8,600	P	E	Discriminative	Closed	S O
Scibench (Wang et al., 2023b) ★	ICML'24	English	869	S	U	Discriminative	Closed/Open	S
GeomVerse (Kazemi et al., 2023) ★	ICML workshop'24	English	1,000	G	U	Discriminative	Closed	S
MathVista (Lu et al., 2023) ★	ICLR'24	English	6,141	S P	E M H U	Discriminative	Closed/Open	S
MMMU _{math} (Yue et al., 2024a) ★	CVPR'24	English	540	U	U	Discriminative	Closed/Open	S
MathVerse (Zhang et al., 2024) ★	ECCV'24	English	2,612	S P	H	Generative	Closed/Open	S
Mathador-LM (Kurtic et al., 2024)	EMNLP'24	English	-	S	E	Both	Closed/Open	S D
MM-MATH (Sun et al., 2024a) ★	EMNLP Findings'24	English	5,929	S	M H	Discriminative	Closed/Open	S D
Scieval (Sun et al., 2024b)	AAAI'24	English	15,901	S P	H	Both	Closed/Open	S
ArqMATH (Satpute et al., 2024)	SIGIR'24	English	450	U	U	Generative	Closed/Open/Math	S
IsoBench (Fu et al., 2024a) ★	COLM'24	English	1,887	S	E M H U	Discriminative	Closed/Open	S
MMMU-Pro _{math} (Yue et al., 2024b) ★	arXiv'24	English	60	S	U	Discriminative	Closed/Open	S
MathOdyssey (Fang et al., 2024)	arXiv'24	English	387	S	H U C	Both	Closed/Open/Math	S
MathScape (Zhou et al., 2024b) ★	arXiv'24	Chinese	1,325	S	E M H U	Generative	Closed/Open	S
U-Math (Chernyshev et al., 2024) ★	arXiv'24	English	1,100	S	U	Discriminative	Closed/Open/Math	S D
MathHay (Wang et al., 2024b)	arXiv'24	English	673	S P	H	Both	Closed/Open	S
MathCheck (Zhou et al., 2024d) ★	arXiv'24	English/Chinese	4,536	P	E M H U	Discriminative	Closed/Open/Math	S
ErrorRador (Yan et al., 2024a) ★	arXiv'24	English	2,500	S	E M H U	Discriminative	Closed/Open	D
FaultyMath (Rahman et al., 2024) ★	arXiv'24	English	363	G	E M H	Discriminative	Closed/Open/Math	D
MathChat (Liang et al., 2024c)	arXiv'24	English	1,319	P	E	Both	Closed/Open/Math	S D O
E-GSM (Xu et al., 2024e)	arXiv'24	Chinese	4,500	P	E	Both	Closed/Open/Math	S O
Tangram (Tang et al., 2024) ★	arXiv'24	English	4,320	S	E M H C	Discriminative	Closed/Open	S
GSM-Symbolic (Mirzadeh et al., 2024)	arXiv'24	English	5,000	P	E	Discriminative	Closed/Open	S
CMM-Math (Liu et al., 2024c) ★	arXiv'24	Chinese	28,069	S	E M H U	Both	Closed/Open/Math	S
CMMaTH (Li et al., 2024b) ★	arXiv'24	English/Chinese	23,856	S	E M H U	Both	Closed/Open/Math	S
EAGLE (Li et al., 2024g) ★	arXiv'24	English	170,000	P	E M H U	Discriminative	Closed/Open/Math	S
VisAidMath (Ma et al., 2024) ★	arXiv'24	English	1,200	S	M H C	Discriminative	Closed/Open	S
AutoGeo (Huang et al., 2024d) ★	arXiv'24	English	100,000	S	E M H U	Both	Closed/Open	O
NTKEval (Guo et al., 2024a)	arXiv'24	English	1,860	P G	H	Discriminative	Open	S
Mamo (Huang et al., 2024b)	arXiv'24	English	1,209	S G	U	Generative	Closed/Open/Math	O
RoMath (Cosma et al., 2024)	arXiv'24	Romanian	70,000	S	M H C	Discriminative	Closed/Open/Math	S
MaTT (Davoodi et al., 2024)	arXiv'24	English	1,958	S	U	Discriminative	Closed/Open	S
Li et al. (2024a)	arXiv'24	English	15,000	P	E M H U	Generative	Closed/Open/Math	S
DynaMath (Zou et al., 2024) ★	arXiv'24	English	5,010	S P G	E M H U	Both	Closed/Open	S
Polymath (Gupta et al., 2024) ★	arXiv'24	English	5,000	S	M H U	Discriminative	Closed/Open	S
SuperCLUE-Math6 (Xu et al., 2024b)	arXiv'24	English/Chinese	2,144	S	E	Generative	Closed/Open	S
TheoremQA (Chen et al., 2023)	EMNLP'23	English	800	S	U	Discriminative	Closed/Open	S
LLA (Mishra et al., 2022)	EMNLP'22	English	133,815	P	H	Discriminative	Closed	S
MATH (Hendrycks et al., 2021)	NeurIPS'21	English	12,500	S	C	Discriminative	Closed	S

Table 2: **Overview of LLM-based benchmarks for mathematical reasoning.** ★ refers to those designed to evaluate the multimodal mathematical setting. Different colors indicate different types for the following columns: **Source:** **S**= Self-Sourced, **P**= Collected from Public Dataset, **G**= Generated by LLM **Level:** **E**= Elementary, **M**= Middle School, **H**= High School, **U**= University, **C**= Competition, **H**= Hybrid **Task:** **S**= Problem-Solving, **D**= Error Detection, **P**= Proving, **O**= Others

demonstrated significant advancements in multi-modal reasoning capabilities, the latest benchmarks have begun to include them in the evaluations. For example, ErrorRadar, in its initial formulation of multimodal error detection setting, incorporates these state-of-the-art MLLMs to highlight the real-world performance gap between AI systems and human-level reasoning (Yan et al., 2024a).

Reasoning Task. Problem-solving tasks typically dominate, reflecting the emphasis on students’ ability to apply knowledge and reasoning skills in real-world contexts. This also serves as the core objective of current Math-LLMs. In addition, a growing proportion of error detection tasks suggests an increasing focus on helping students rec-

ognize and correct mistakes (Li et al., 2024e; Yan et al., 2024a; Kurtic et al., 2024). Meanwhile, proving tasks, often associated with higher-order thinking, highlight a shift towards cultivating logical reasoning and systematic problem-solving abilities (Tsoukalas et al.). Moreover, a smaller portion of work has addressed tasks that align with real-world educational needs but lack systematic formulation. For instance, Li et al. (2024e) further introduces error correction (which goes beyond simple error detection); Didolkar et al. (2024) explores automated skill discovery for problem-solving; and MathChat (Liang et al., 2024c) focuses on reasoning in multi-turn settings (such as follow-up QA and problem generation). Given the higher demands on rea-

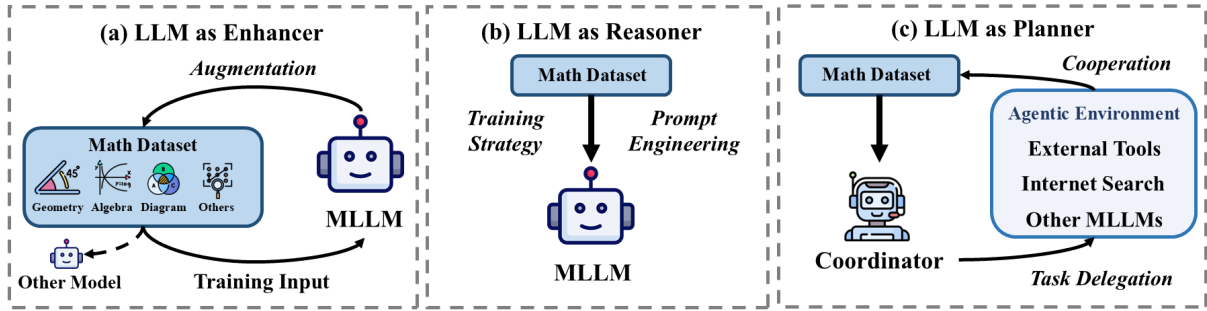


Figure 4: The illustration of the comparisons among three paradigms of (M)LLM-based mathematical reasoning.

soning capabilities in multimodal settings, many studies have also evaluated the aforementioned reasoning tasks in image-text problem settings. These efforts aim to provide the LLM community with more diverse, real-world task scenarios, catering to the needs of multimodal learning environments.

2.4 Evaluation

Discriminative Evaluation is a common approach, focusing on the ability of M(LLM)s to correctly classify or choose the correct answer (Hendrycks et al., 2021; Mishra et al., 2022; Li et al., 2024c). Based on specific motivations, some works also build their metrics upon accuracy for further expansion. For example, GSM-PLUS, a new adversarial benchmark for evaluating the robustness of LLMs in mathematical reasoning, develops performance drop rate (PDR) to measure the relative decline in performance on question variations compared to the original questions (Li et al., 2024d). Error-Radar uses error step accuracy and error category accuracy together to evaluate the multimodal error detection of MLLMs (Yan et al., 2024a).

Generative Evaluation, on the other hand, measures a M(LLM)’s ability to produce detailed explanations or solve problems from scratch. This evaluation type is gaining traction, particularly for complex mathematical tasks where step-by-step solutions are required. For instance, MathVerse, which modifies problems with varying degrees of information content in multi-modality, employs GPT-4 to score each key step in the reasoning process generated by MLLMs (Zhang et al., 2024f). CHAMP proposes a solution evaluation pipeline where GPT-4 is utilized as a grader for the answer summary, given the ground truth answer (Mao et al., 2024).

Due to page limit, more details of both types of evaluation metrics can be seen in Appendix B.

3 Methodology Perspective

3.1 Overview & Findings

MLLMs have been leveraged in various ways to tackle the broad spectrum of mathematical reasoning tasks. Based on our comprehensive review of recent methodologies (summarized in Table 3), we classify the works into three distinct paradigms: LLM as Enhancer (Sec.3.2), LLM as Reasoner (Sec.3.3), and LLM as Planner (Sec.3.4).

Findings. First, single-modality settings dominate the current landscape of method-oriented research, with the majority focusing solely on algebraic tasks. However, since 2024, multimodal approaches have been increasingly incorporated, expanding the scope of mathematical reasoning to include geometry, diagrams, and even broader mathematical concepts. This shift signals a growing interest in enhancing model robustness through multimodal learning, which can address the diverse nature of mathematical problems. Second, regarding the evaluated tasks, problem-solving and proving are gaining prominence, while some research also focuses on error detection or others (e.g., RefAug includes error correction and follow-up QA as evaluation tasks (Zhang et al., 2024i)). Finally, in terms of the role of LLMs, Reasoner is the most common role, followed by Enhancer, while Planner remains less explored but holds promise due to recent advancements in multi-agent intelligence.

3.2 LLM as Enhancer

Definition. In the *Enhancer* paradigm, M(LLM)s are primarily used to augment data, thereby enabling improvements in mathematical reasoning, as illustrated in Figure 4 (a). This can be achieved by synthesizing new training data, refining existing datasets, or introducing new variations that target specific problem-solving abilities (Li et al., 2022). Data augmentation can include paraphrasing math-

Methods	Venue	Evaluated Math Dataset(s)	Task(s)	Scope(s)	LLM as Enhancer	LLM as Reasoner	LLM as Planner
AlphaGeometry (Trinh et al., 2024) ★	Nature'24	IMO-AG-30	S P	G	✓	✓	
Masked Thought (Chen et al., 2024)	ACL'24	GSM8K/MATH/GSM8K-RFT/MetaMathQA/MathInstruct	S	A		✓	
MathGenie (Lu et al., 2024b)	ACL'24	GSM8K/MATH/SVAMP/SimulEq/Mathematics	S	A	✓	✓	
MATH-SHEPHERD (Wang et al., 2024c)	ACL'24	GSM8K/MATH	S P	A	✓	✓	
SEGO (Zhao et al., 2024)	ACL'24	GSM8K/MATH	S P	A	✓	✓	
Deng et al. (2023)	ACL Workshop'24	GSM8K/SVAMP/MultiArith/MathQA/CSQA	S	A		✓	
MathCoder (Wang et al., 2023a)	ICLR'24	GSM8K/MATH	S P	A		✓	
ToRA (Gou et al., 2023)	ICLR'24	GSM8K/MATH	S P	A		✓	
Visual Sketchpad (Hu et al., 2024) ★	NeurIPS'24	Geometry3K/ IsoBench	S P	G			✓
Minimo (Poesia et al., 2024)	NeurIPS'24	-	S P	A		✓	
Sinha et al. (2024) ★	NeurIPS Workshop'24	IMO-AG-30	S P	G		✓	
SBIRAG (Dixit and Oates, 2024)	NeurIPS Workshop'24	GSM8K	S	A		✓	
VerityMath (Han et al., 2023)	ICML Workshop'24	GSM8K	S	A		✓	
RefAug (Zhang et al., 2024)	EMNLP'24	GSM8K/MATH/Mathematics/MAWPS/SVAMP/MMLU-Math/SAT-Math/MathChat-FQA/MathChat-EC/MINI-Math	S P D	A	✓	✓	
Math-LLaVA (Shi et al., 2024) ★	EMNLP Findings' 24	MathVista/Math-V	S P	M	✓	✓	
COPRA (Thakur et al., 2024)	COLM'24	mini2F-test	S	A			✓
FRP (Wu et al., 2024b)	AAAI'24	MAWPS/ASDivA/Math23k/SVAMP/Un-biasedMWP	S	A		✓	
PERC (Jin et al., 2024)	L@S'24	PERC	S	A		✓	
Math-PUMA (Zhuang et al., 2024) ★	arXiv'24	MathVerse/MathVista/WE-MATH	S P	M		✓	
MultiMath (Peng et al., 2024) ★	arXiv'24	MathVista/MathVerse/MultiMath-300K	S P	M		✓	
MathAttack (Zhou et al., 2024e)	arXiv'24	GSM8K/MultiArith	S	A		✓	
MiniT (Liang et al., 2023b)	arXiv'24	GSM8K/MathQA/CM17k/Ape210k	S	A		✓	
DotaMath (Li et al., 2024b)	arXiv'24	GSM8K/MATH/Mathematics/SVAMP/TabMWP/ASDiv	S	A		✓	
DFE-GPS (Zhang et al., 2024h)	arXiv'24	FORMALGEO7k	S	G		✓	
JiuZhang 3.0 (Zhou et al., 2024a)	arXiv'24	GSM8K/MATH/SVAMP/ASDiv/MAWPS/CARP	S P	A		✓	
PGPSNet-v2 (Zhang et al., 2024e) ★	arXiv'24	Geometry3K/PGPS9K	S	G D		✓	
LLaMA-Berry (Zhang et al., 2024b)	arXiv'24	GSM8K/MATH/GaoKao2023En/OlympiadBench/College Math/MMLU STEM	S P	A		✓	
Skywork-Math (Zeng et al., 2024) ★	arXiv'24	GSM8K/MATH	S P	A	✓	✓	
SLaM (Yu et al., 2024a)	arXiv'24	GSM8K/CMATH	S P	A		✓	
InternLM-Math (Ying et al., 2024)	arXiv'24	GSM8K/MATH	S P	A		✓	
MathGLM-Vision (Yang et al., 2024b) ★	arXiv'24	MathVista/MathVerse/MathVision	S P	M	✓	✓	
Qwen2.5-Math (Yang et al., 2024a) ★	arXiv'24	GSM8K/MATH/MMLU-STEM/CMATH/GaoKao-Math-Cloze/GaoKao-Math-QA	S P	A	✓	✓	
S3c-Math (Yan et al., 2024c)	arXiv'24	GSM8K/MATH/SVAMP/Mathematics	S P	A	✓	✓	
BMA-MTIPL (Xiong et al., 2024)	arXiv'24	GSM8K/MATH	S P	A		✓	
SIRP (Wu et al., 2024a)	arXiv'24	CSQA/GSM8K/MATH/MBPP	S P	A		✓	
AIPS (Wei et al., 2024)	arXiv'24	MO-INT-20	S	G		✓	
DART-Math (Tong et al., 2024)	arXiv'24	MATH/GSM8K/College/DM/Olympiad/Theorem	S P	A	✓	✓	
DeepSeekMath (Shao et al., 2024)	arXiv'24	GSM8K/MATH/OCW/SAT/MMLU STEM/CMATH-Gaokao MathCloze/Gaokao MathQA	S	A		✓	
MMIQC (Liu et al., 2024a)	arXiv'24	MATH/MMIQC	S	A	✓	✓	
LANS (Li et al., 2023c) ★	arXiv'24	Geometry3K/PGPS9K	S	G D		✓	
VCAR (Jia et al., 2024) ★	arXiv'24	MathVista/MathVerse	S	M		✓	
KPDDS (Huang et al., 2024c)	arXiv'24	GSM8K/MATH/SVAMP/TabMWP/ASDiv/MAWPS	S P	A	✓	✓	
HGR (Huang et al., 2024a) ★	arXiv'24	Geometry3K	S	G		✓	
InfiMM-Math (Han et al., 2024) ★	arXiv'24	GSM8K/MMLU/MathVerse/We-Math	S P	A		✓	
CoSC (Han et al., 2024)	arXiv'24	GSM8K/MATH	S P	A	✓	✓	
Math-TSMC (Feng et al., 2024)	arXiv'24	GSM8K/MATH500	S P	A		✓	
SICCV (Liang et al., 2024b)	arXiv'24	GSM8K/MATH500	S P	A		✓	
BEATS (Sun et al., 2024c)	arXiv'24	GSM8K/MATH/SVAMP/SimulEq/NumGLUE	S P	A		✓	
MindStar (Kang et al., 2024)	arXiv'24	GSM8K/MATH	S P	A		✓	
UMM (Zhang et al., 2024g)	arXiv'24	MMLU/GSM8K-COT/GSM8K-Coding/MATH-COT/MATH-Coding/HumanEval/InfiBench	S P	A		✓	
STIC (Deng et al., 2024) ★	arXiv'24	ScienceQA/TextVQA/ChartQA/LLaVA-Bench/MMBench/MM-Vet/MathVista	S	M		✓	
SPMWP (Zhang et al., 2023)	ACL'23	GSM8K	S	A		✓	
CoRe (Zhu et al., 2022)	ACL'23	GSM8K/ASDiv-A/SingleOp/SimulEq/MultiArith	S	A		✓	
TabMWP (Lu et al., 2022a)	ICLR'23	TabMWP	S	A D		✓	
Chameleon (Lu et al., 2024a) ★	NeurIPS'23	ScienceQA/TabMWP	S	A D		✓	
ATHENA (Kim et al., 2023)	EMNLP'23	MAWPS/ASDivA/Math23k/SVAMP/Un-biasedMWP	S P	A		✓	
UniMath (Liang et al., 2023a) ★	EMNLP'23	SVAMP/GeoQA/TabMWP/MathQA/UniGeo-Proving	S P	A D		✓	
Jiuzhang 2.0 (Zhao et al., 2023)	KDD'23	MCQ/BFQ/CAG/BAG/KPC/QRC/JCAG/JBAG	S	A		✓	
TCDP (Qin et al., 2023)	TNNLS'23	Math23k/CM17K	S	A		✓	
UniGeo (Chen et al., 2022)	ACL'22	GeoQA	S	G		✓	
LogicSolver (Yang et al., 2022)	EMNLP Findings' 22	InterMWP/Math23K	S	A	✓	✓	
Jiuzhang (Zhao et al., 2022)	KDD'22	KPC/QRC/QAM/SQR/QAR/MCQ/BFQ/CAG/BAG	S	A		✓	
MWP-BERT (Liang et al., 2021)	NAACL'22	Math23k/MathQA/Ape-210k	S	A		✓	
Inter-GPS (Lu et al., 2021) ★	ACL'21	Geometry3K/GEOS	S	G		✓	

Table 3: **Overview of LLM-based methods for mathematical reasoning.** ★ refers to those specifically designed to tackle the multimodal mathematical setting. Different colors indicate different types for the following columns: **Task:** **S**= Problem-Solving, **D**= Error Detection, **P**= Proving, **O**= Others **Scope:** **G**= Geometry, **A**= Algebra, **D**= Diagram, **M**= General Math

emathical problems, adding noise to mathematical expressions, or generating problem variants for underrepresented cases.

Examples. A typical example of a single-modality enhancement approach is Masked Thought, which introduces perturbations to the input and randomly masks tokens within the chain of thought during training (Chen et al., 2024). MathGenie, which aims to generate diverse and reliable math problems and solution from a small-scale dataset, leverages a solution augmentation model to iteratively create new solutions from existing

ones (Lu et al., 2024b). For multimodal methods, AlphaGeometry proves most olympiad-level mathematical theorems, via trained from scratch on large-scale synthetic data guiding the symbolic deuction (Trinh et al., 2024); LogicSolver introduces interpretable formula-based tree-structure for each solution equation (Yang et al., 2022); InfiMM-Math achieves the exceptional performance as it is trained on a large-scale multimodal interleaved math dataset developed and validated by LLMs such as LLaMA3-70B-Instruct (Han et al., 2024); DFE-GPS constructs its synthetic training

set, which integrates visual features and geometric formal language (Zhang et al., 2024h).

Summary & Outlook. This paradigm offers substantial performance improvements by enriching the training set. However, challenges remain in ensuring the diversity and relevance of the generated data. Moreover, while text-based augmentation methods have proven effective, the potential for multimodal augmentation is still underexplored. Future research should focus on advancing multimodal data augmentation techniques, especially for tasks that require interaction between visual and textual modalities (Xiao et al., 2023).

3.3 LLM as Reasoner

Definition. In the *Reasoner* paradigm, M(LLM)s harness their inherent reasoning capabilities to solve mathematical problems, as shown in Figure 4 (b). This can either involve fine-tuning existing LLMs on task-specific datasets or utilizing zero-shot or few-shot learning strategies. These models utilize advanced semantic understanding and reasoning techniques, such as symbolic manipulation, logical deduction, and multi-step reasoning.

Examples. Deng et al. (2023) develops a unified framework for answer calibration that integrates step-level and path-level strategies on multi-step reasoning of LLMs. MATH-SHEPHERD serves as a process-oriented math verifier, which assigns a reward score to each step of the LLM’s outputs on math questions (Wang et al., 2024c). As for multimodal approaches, Math-PUMA introduces progressive upward multimodal alignment strategy for reasoning-enhanced training (Zhuang et al., 2024); Math-LLaVA, a LLaVA-1.5-based model, directly bootstraps mathematical reasoning via fine-tuned on 360K high-quality math QA pairs, which can ensure the depth and breadth of multimodal mathematical problems (Shi et al., 2024); STIC develops a two-stage self-training pipeline (consisting of Image Comprehension Self-Training phase & Description-Infused Fine-Tuning phase) for enhancing visual comprehension (Deng et al., 2024); VCAR emphasizes on the visual-centric supervision, thus proposing a similar two-step training pipeline which handles the visual description generation task first, followed by mathematical rationale generation task (Jia et al., 2024).

Summary & Outlook. This paradigm has shown significant promise, particularly in solving problems requiring multiple steps of reasoning. However, despite improvements, issues with robust-

ness remain, particularly with zero-shot reasoning tasks. Future work should focus on combining reasoning with structured knowledge retrieval systems and enhancing models’ ability to reason effectively across diverse domains, especially in multimodal contexts (Fan et al., 2024; Pan et al., 2023).

3.4 LLM as Planner

Definition. In the *Planner* paradigm, M(LLM)s are treated as coordinators that guide the solution of complex mathematical problems by delegating tasks to other models or tools, as illustrated in Figure 4 (c). This includes scenarios where multiple agents or models collaborate to achieve a single objective, thereby enhancing the performance of mathematical problem-solving through cooperative interactions. These models often work in environments with multiple steps or require iterative refinement of solutions.

Examples. A notable tool-integrated agent is ToRA, which plans the sequential use of natural language rationale and program-based tools synergistically to solve mathematical problems in an optimal manner (Gou et al., 2023). Additionally, COPRA simulates a single agent-like reasoning mechanism where GPT-4 proposes tactic applications within a stateful backtracking search, leveraging feedback from the proof environment (Thakur et al., 2024). This can also extend to multimodal scenarios, as seen in Chameleon, which serves as an AI system that augments MLLMs with plug-and-play modules for compositional reasoning, leveraging an LLM-based planner to assemble tools for complex tasks (Lu et al., 2024a). Furthermore, Visual Sketchpad presents the concept of sketching as a ubiquitous tool used by humans for communication, ideation, and problem-solving. Hence, MLLMs can enable external tools (*e.g.*, matplotlib) to generate intermediate sketches to aid in reasoning, which includes an iterative interaction process with an environment (Hu et al., 2024). Although there has been much work on Compositional Visual Reasoning in the past (Gupta and Kembhavi, 2023; Surís et al., 2023; Yao et al., 2022), Visual Sketchpad is the first work that integrates the planning capabilities of MLLMs with the real gap of mathematical reasoning settings (*i.e.*, sketch-based reasoning involving visuo-spatial concepts).

Summary & Outlook. While the Planner paradigm introduces significant improvements, particularly for complex tasks that require multi-agent collaboration, it remains a relatively under-

explored area (Xi et al., 2023; Guo et al., 2024b). There is potential for further improvement in task decomposition, agent cooperation strategies, and integration of diverse computational tools. Future work will likely focus on refining these planning strategies, especially for multimodal systems that can jointly leverage visual and textual knowledge to solve more intricate problems (Xie et al., 2024; Durante et al., 2024; Li et al., 2023b).

4 Challenges

In the realm of MLLMs for mathematical reasoning, the following key challenges persist that hinder their full potential. Addressing these challenges is essential for advancing MLLMs toward more robust and flexible systems that can better support mathematical reasoning in real-world settings.

❶ **Insufficient Visual Reasoning.** Many math problems require extracting and reasoning over visual content, such as charts, tables, or geometric diagrams. Current models struggle with intricate visual details, such as interpreting three-dimensional geometry or analyzing irregularly structured tables (Zhang et al., 2024f). Hence, it may be beneficial to introduce enhanced visual feature extraction modules and integrate scene graph representations for better reasoning over complex visual elements (Ibrahim et al., 2024; Guo et al., 2024c).

❷ **Reasoning Beyond Text and Vision.** While the current research focus on the combination of text and vision, mathematical reasoning in real-world applications often extends beyond these two modalities. For instance, audio explanations, interactive problem-solving environments, or dynamic simulations might play a role in some tasks. Current models are not well-equipped to handle such diverse inputs (Abrahamson et al., 2020; Jusslin et al., 2022). To address this, datasets should be expanded to include more diverse modalities, such as audio, video, and interactive tools. MLLMs should also be designed with flexible architectures capable of processing and reasoning over multiple types of inputs, allowing for a richer representation of mathematical problems (Dasgupta et al., 2023).

❸ **Limited Domain Generalization.** Mathematical reasoning spans many domains, such as algebra, geometry, diagram and commonsense, each with its own specific requirements for problem-solving (Liu et al., 2023b; Lu et al., 2022b). Math-LLMs that perform well in one domain often fail to generalize across others, which can limit their

utility. By pretraining and fine-tuning Math-LLMs on a wide array of problem types, models may handle cross-domain tasks more effectively, improving their ability to generalize across different mathematical topics and problem-solving strategies.

❹ **Error Feedback Limitations.** Mathematical reasoning involves various types of errors, such as calculation mistakes, logical inconsistencies, and misinterpretations of the problem. Currently, MLLMs lack mechanisms to detect, categorize, and correct these errors effectively, which can result in compounding mistakes throughout the reasoning process (Yan et al., 2024a; Li et al., 2024e). A potential solution is to integrate error detection and classification modules that can identify errors at each step of the reasoning process. Besides, multi-agent collaboration mechanism could be introduced, via involving multiple agents collaborating by exchanging feedback and collectively refining the reasoning process (Xu et al., 2024d).

❺ **Integration with Real-world Educational Needs.** Existing benchmarks and models often overlook real-world educational contexts, such as how students use draft work, like handwritten notes or diagrams, to solve problems (Xu et al., 2024c; Wang et al., 2024d). These real-world elements are crucial for understanding how humans approach mathematical reasoning (Mouchere et al., 2011; Gervais et al., 2024). By incorporating draft notes, handwritten calculations, and dynamic problem-solving workflows into the training data, MLLMs can be tailored to provide more accurate and contextually relevant feedback for students.

5 Conclusion

In this survey, we have provided a comprehensive overview of the progress and challenges in mathematical reasoning within the context of MLLMs. We highlighted the significant advances in the development of Math-LLMs and the growing importance of multimodal integration for solving complex reasoning tasks. We identified five key challenges that are crucial for the continued development of AGI systems capable of performing sophisticated mathematical reasoning tasks. As research continues to advance, it is essential to focus on these challenges to unlock the full potential of LLMs in multimodal settings. We hope this survey provides insights to guide future LLM research, ultimately leading to more effective and human-like mathematical reasoning capabilities in AI systems.

540 **Limitations**

541 Despite our best efforts to ensure comprehensive
542 coverage of the published works, it is possible that
543 some relevant studies were overlooked. Addition-
544 ally, human errors could have occurred during the
545 categorization or referencing of papers in the sur-
546 vey. To minimize such errors, we made a con-
547 certed effort to gather studies from multiple sources
548 and performed a multiple-round checking process.
549 While minor inconsistencies or omissions may still
550 exist, we believe this survey represents the most
551 comprehensive review of MLLM-based mathemat-
552 ical reasoning to date, effectively capturing key re-
553 search trends and highlighting ongoing challenges.

References

- 555 Dor Abrahamson, Mitchell J Nathan, Caro Williams-
556 Pierce, Candace Walkington, Erin R Ottmar, Hortensia
557 Soto, and Martha W Alibali. 2020. The future
558 of embodied design for mathematics teaching and
559 learning. In *Frontiers in Education*, volume 5, page
560 147. Frontiers Media SA.
- 561 Janice Ahn, Rishu Verma, Renze Lou, Di Liu, Rui
562 Zhang, and Wenpeng Yin. 2024. Large language
563 models for mathematical reasoning: Progresses and
564 challenges. [arXiv preprint arXiv:2402.00157](#).
- 565 Aida Amini, Saadia Gabriel, Peter Lin, Rik Koncel-
566 Kedzioriski, Yejin Choi, and Hannaneh Hajishirzi.
567 2019. Mathqa: Towards interpretable math word
568 problem solving with operation-based formalisms.
569 [arXiv preprint arXiv:1905.13319](#).
- 570 Zhangir Azerbayev, Hailey Schoelkopf, Keiran Paster,
571 Marco Dos Santos, Stephen McAleer, Albert Q Jiang,
572 Jia Deng, Stella Biderman, and Sean Welleck. 2023.
573 Llemma: An open language model for mathematics.
574 [arXiv preprint arXiv:2310.10631](#).
- 575 Tianyi Bai, Hao Liang, Binwang Wan, Yanran Xu, Xi Li,
576 Shiyu Li, Ling Yang, Bozhou Li, Yifan Wang, Bin
577 Cui, et al. 2024. A survey of multimodal large lan-
578 guage model from a data-centric perspective. [arXiv
579 preprint arXiv:2405.16640](#).
- 580 Changyu Chen, Xiting Wang, Ting-En Lin, Ang Lv,
581 Yuchuan Wu, Xin Gao, Ji-Rong Wen, Rui Yan, and
582 Yongbin Li. 2024. Masked thought: Simply masking
583 partial reasoning steps can improve mathematical rea-
584 soning learning of language models. [arXiv preprint
585 arXiv:2403.02178](#).
- 586 Jiaqi Chen, Tong Li, Jinghui Qin, Pan Lu, Liang Lin,
587 Chongyu Chen, and Xiaodan Liang. 2022. Uni-
588 geo: Unifying geometry logical reasoning via re-
589 formulating mathematical expression. [arXiv preprint
590 arXiv:2212.02746](#).
- 591 Wenhui Chen, Ming Yin, Max Ku, Pan Lu, Yixin
592 Wan, Xueguang Ma, Jianyu Xu, Xinyi Wang, and
593 Tony Xia. 2023. Theoremqa: A theorem-driven
594 question answering dataset. In *Proceedings of the
595 2023 Conference on Empirical Methods in Natural
596 Language Processing*, pages 7889–7901.
- 597 Anoop Cherian, Kuan-Chuan Peng, Suhas Lohit, Joanna
598 Matthiesen, Kevin Smith, and Joshua B Tenenbaum.
599 2024. Evaluating large vision-and-language mod-
600 els on children’s mathematical olympiads. [arXiv
601 preprint arXiv:2406.15736](#).
- 602 Ethan Chern, Haoyang Zou, Xuefeng Li, Jiewen Hu, Ke-
603 hua Feng, Junlong Li, and Pengfei Liu. 2023. Gen-
604 erative ai for math: Abel. [https://github.com/
605 GAIR-NLP/abel](https://github.com/GAIR-NLP/abel).
- 606 Konstantin Chernyshev, Vitaliy Polshkov, Ekaterina
607 Artemova, Alex Myasnikov, Vlad Stepanov, Alexei
608 Miasnikov, and Sergei Tilga. 2024. U-math: A
university-level benchmark for evaluating mathemat- 609
ical skills in llms. [arXiv preprint arXiv:2412.03205](#). 610
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, 611
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias 612
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro 613
Nakano, et al. 2021. Training verifiers to solve math 614
word problems. [arXiv preprint arXiv:2110.14168](#). 615
- Adrian Cosma, Ana-Maria Bucur, and Emilian Radoi. 616
2024. Romath: A mathematical reasoning bench- 617
mark in romanian. [arXiv preprint arXiv:2409.11074](#). 618
- Jiaxi Cui, Zongjian Li, Yang Yan, Bohua Chen, and 619
Li Yuan. 2023. Chatlaw: Open-source legal large 620
language model with integrated external knowledge 621
bases. [arXiv preprint arXiv:2306.16092](#). 622
- Ishita Dasgupta, Christine Kaeser-Chen, Kenneth 623
Marino, Arun Ahuja, Sheila Babayan, Felix Hill, 624
and Rob Fergus. 2023. Collaborating with language 625
models for embodied reasoning. [arXiv preprint
626 arXiv:2302.00763](#). 627
- Arash Gholami Davoodi, Seyed Pouyan Mousavi 628
Davoudi, and Pouya Pezeshkpour. 2024. Llms are 629
not intelligent thinkers: Introducing mathematical 630
topic tree benchmark for comprehensive evaluation 631
of llms. [arXiv preprint arXiv:2406.05194](#). 632
- Shumin Deng, Ningyu Zhang, Nay Oo, and Bryan 633
Hooi. 2023. Towards a unified view of answer cal- 634
ibration for multi-step reasoning. [arXiv preprint
635 arXiv:2311.09101](#). 636
- Yihe Deng, Pan Lu, Fan Yin, Ziniu Hu, Sheng Shen, 637
James Zou, Kai-Wei Chang, and Wei Wang. 2024. 638
Enhancing large vision language models with self- 639
training on image comprehension. [arXiv preprint
640 arXiv:2405.19716](#). 641
- Aniket Didolkar, Anirudh Goyal, Nan Rosemary Ke, 642
Siyuan Guo, Michal Valko, Timothy Lillicrap, Danilo 643
Rezende, Yoshua Bengio, Michael Mozer, and San- 644
jeev Arora. 2024. Metacognitive capabilities of llms: 645
An exploration in mathematical problem solving. 646
[arXiv preprint arXiv:2405.12205](#). 647
- Prakhar Dixit and Tim Oates. 2024. Sbi-rag: Enhanc- 648
ing math word problem solving for students through 649
schema-based instruction and retrieval-augmented 650
generation. [arXiv preprint arXiv:2410.13293](#). 651
- Duolingo. 2024. [Duolingo official platform](#). 652
- Zane Durante, Qiuyuan Huang, Naoki Wake, Ran Gong, 653
Jae Sung Park, Bidipta Sarkar, Rohan Taori, Yusuke 654
Noda, Demetri Terzopoulos, Yejin Choi, et al. 2024. 655
Agent ai: Surveying the horizons of multimodal in- 656
teraction. [arXiv preprint arXiv:2401.03568](#). 657
- Wenqi Fan, Yujian Ding, Liangbo Ning, Shijie Wang, 658
Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing 659
Li. 2024. A survey on rag meeting llms: To- 660
wards retrieval-augmented large language models. In 661
Proceedings of the 30th ACM SIGKDD Conference 662

663	on Knowledge Discovery and Data Mining , pages	Himanshu Gupta, Shreyas Verma, Ujjwala Anan-	716
664	6491–6501.	theswaran, Kevin Scaria, Mihir Parmar, Swaroop	717
665	Meng Fang, Xiangpeng Wan, Fei Lu, Fei Xing, and	Mishra, and Chitta Baral. 2024. Polymath: A chal-	718
666	Kai Zou. 2024. Mathodyssey: Benchmarking math-	lenging multi-modal mathematical reasoning bench-	719
667	ematical problem-solving skills in large language	mark. arXiv preprint arXiv:2410.14702 .	720
668	models using odyssey math data. arXiv preprint	Tanmay Gupta and Aniruddha Kembhavi. 2023. Vi-	721
669	arXiv:2406.18321 .	sual programming: Compositional visual reason-	722
670	Shengyu Feng, Xiang Kong, Shuang Ma, Aonan Zhang,	ing without training. In Proceedings of the	723
671	Dong Yin, Chong Wang, Ruoming Pang, and Yim-	IEEE/CVF Conference on Computer Vision and	724
672	ing Yang. 2024. Step-by-step reasoning for math	Pattern Recognition , pages 14953–14962.	725
673	problems via twisted sequential monte carlo. arXiv	Vernon Toh Yan Han, Ratish Puduppully, and Nancy F	726
674	preprint arXiv:2410.01920 .	Chen. 2023. Veritymath: Advancing mathematical	727
675	Deqing Fu, Ruohao Guo, Ghazal Khalighine-	reasoning by self-verification through unit consis-	728
676	jad, Ollie Liu, Bhuwan Dhingra, Dani Yo-	tenency. arXiv preprint arXiv:2311.07172 .	729
677	gatama, Robin Jia, and Willie Neiswanger. 2024a.	Xiaotian Han, Yiren Jian, Xuefeng Hu, Haogeng Liu,	730
678	Isobench: Benchmarking multimodal foundation	Yiqi Wang, Qihang Fan, Yuang Ai, Huaibo Huang,	731
679	models on isomorphic representations. arXiv	Ran He, Zhenheng Yang, et al. 2024. Infimm-	732
680	preprint arXiv:2404.01266 .	webmath-40b: Advancing multimodal pre-training	733
681	Jiayi Fu, Lei Lin, Xiaoyang Gao, Pengli Liu, Zhengzong	for enhanced mathematical reasoning. arXiv preprint	734
682	Chen, Zhirui Yang, Shengnan Zhang, Xue Zheng,	arXiv:2409.12568 .	735
683	Yan Li, Yuliang Liu, et al. 2023. Kwaiiimath: Tech-	Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu,	736
684	nical report. arXiv preprint arXiv:2310.07488 .	Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han,	737
685	Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu	Yujie Huang, Yuxiang Zhang, et al. 2024. Olympiad-	738
686	Wang, Xudong Lin, Dan Roth, Noah A Smith, Wei-	bench: A challenging benchmark for promoting agi	739
687	Chiu Ma, and Ranjay Krishna. 2024b. Blink: Multi-	with olympiad-level bilingual multimodal scientific	740
688	modal large language models can see but not perceive.	problems. arXiv preprint arXiv:2402.14008 .	741
689	In European Conference on Computer Vision , pages	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul	742
690	148–166. Springer.	Arora, Steven Basart, Eric Tang, Dawn Song, and Ja-	743
691	Philippe Gervais, Asya Fadeeva, and Andrii Maksai.	cob Steinhart. 2021. Measuring mathematical prob-	744
692	2024. Mathwriting: A dataset for handwritten	lem solving with the math dataset. arXiv preprint	745
693	mathematical expression recognition. arXiv preprint	arXiv:2103.03874 .	746
694	arXiv:2404.10690 .	Yushi Hu, Weijia Shi, Xingyu Fu, Dan Roth, Mari Os-	747
695	Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen,	tendorf, Luke Zettlemoyer, Noah A Smith, and Ran-	748
696	Yujiu Yang, Minlie Huang, Nan Duan, and Weizhu	jay Krishna. 2024. Visual sketchpad: Sketching as	749
697	Chen. 2023. Tora: A tool-integrated reasoning agent	a visual chain of thought for multimodal language	750
698	for mathematical problem solving. arXiv preprint	models. arXiv preprint arXiv:2406.09403 .	751
699	arXiv:2309.17452 .	Litian Huang, Xinguo Yu, Feng Xiong, Bin He, Sheng-	752
700	Siyuan Guo, Aniket Didolkar, Nan Rosemary Ke,	bing Tang, and Jiawen Fu. 2024a. Hologram rea-	753
701	Anirudh Goyal, Ferenc Huszár, and Bernhard	soning for solving algebra problems with geometry	754
702	Schölkopf. 2024a. Learning beyond pattern match-	diagrams. arXiv preprint arXiv:2408.10592 .	755
703	ing? assaying mathematical understanding in llms.	Xuhan Huang, Qingning Shen, Yan Hu, Anningzhe Gao,	756
704	arXiv preprint arXiv:2405.15485 .	and Benyou Wang. 2024b. Mamo: a mathematical	757
705	Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang,	modeling benchmark with solvers. arXiv preprint	758
706	Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xian-	arXiv:2405.13144 .	759
707	gliang Zhang. 2024b. Large language model based	Yiming Huang, Xiao Liu, Yeyun Gong, Zhibin Gou,	760
708	multi-agents: A survey of progress and challenges.	Yelong Shen, Nan Duan, and Weizhu Chen. 2024c.	761
709	arXiv preprint arXiv:2402.01680 .	Key-point-driven data synthesis with its enhance-	762
710	Tiezheng Guo, Qingwen Yang, Chen Wang, Yanyi	ment on mathematical reasoning. arXiv preprint	763
711	Liu, Pan Li, Jiawei Tang, Dapeng Li, and Yingyou	arXiv:2403.02333 .	764
712	Wen. 2024c. Knowledgenavigator: Leveraging	Zihan Huang, Tao Wu, Wang Lin, Shengyu Zhang,	765
713	large language models for enhanced reasoning over	Jingyuan Chen, and Fei Wu. 2024d. Autogeo: Au-	766
714	knowledge graph. Complex & Intelligent Systems ,	tomating geometric image dataset creation for en-	767
715	10(5):7063–7076.	hanced geometry understanding. arXiv preprint	768
		arXiv:2409.09039 .	769

770	Nourhan Ibrahim, Samar Aboulela, Ahmed Ibrahim, and Rasha Kashef. 2024. A survey on augmenting knowledge graphs (kgs) with large language models (llms): models, evaluation metrics, benchmarks, and challenges. <i>Discover Artificial Intelligence</i> , 4(1):76.	825
771		826
772		827
773		
774		
775	Mengzhao Jia, Zhihan Zhang, Wenhao Yu, Fangkai Jiao, and Meng Jiang. 2024. Describe-then-reason: Improving multimodal mathematical reasoning through visual comprehension training. <i>arXiv preprint arXiv:2404.14604</i> .	828
776		829
777		830
778		831
779		832
780	Hyoungwook Jin, Yoonsu Kim, Yeon Su Park, Bekzat Tilekbay, Jinho Son, and Juho Kim. 2024. Using large language models to diagnose math problem-solving skills at scale. In <i>Proceedings of the Eleventh ACM Conference on Learning@ Scale</i> , pages 471–475.	833
781		834
782		835
783		836
784		837
785		
786	Bert Jonsson, Julia Mossegård, Johan Lithner, and Linnea Karlsson Wirebring. 2022. Creative mathematical reasoning: Does need for cognition matter? <i>Frontiers in Psychology</i> , 12:797807.	838
787		839
788		840
789		841
790	Sofia Jusslin, Kaisa Korpinen, Niina Lilja, Rose Martin, Johanna Lehtinen-Schnabel, and Eeva Anttila. 2022. Embodied learning and teaching approaches in language education: A mixed studies review. <i>Educational Research Review</i> , 37:100480.	842
791		843
792		844
793		845
794		
795	Jikun Kang, Xin Zhe Li, Xi Chen, Amirreza Kazemi, Qianyi Sun, Boxing Chen, Dong Li, Xu He, Quan He, Feng Wen, et al. 2024. Mindstar: Enhancing math reasoning in pre-trained llms at inference time. <i>arXiv preprint arXiv:2405.16265</i> .	846
796		847
797		848
798		849
799		850
800	Mehran Kazemi, Hamidreza Alvari, Ankit Anand, Jialin Wu, Xi Chen, and Radu Soricut. 2023. Geomverse: A systematic evaluation of large models for geometric reasoning. <i>arXiv preprint arXiv:2312.12241</i> .	851
801		852
802		853
803		854
804	KhanAcademy. 2024. <i>Khanmigo official platform</i> .	855
805	JB Kim, Hazel Kim, Joonghyuk Hahn, and Yo-Sub Han. 2023. Athena: Mathematical reasoning with thought expansion. <i>arXiv preprint arXiv:2311.01036</i> .	856
806		857
807		858
808	Eldar Kurtic, Amir Moeini, and Dan Alistarh. 2024. Mathador-1m: A dynamic benchmark for mathematical reasoning on large language models. <i>arXiv preprint arXiv:2406.12572</i> .	859
809		860
810		861
811		
812	Guillaume Lample, Timothee Lacroix, Marie-Anne Lachaux, Aurelien Rodriguez, Amaury Hayat, Thibaut Lavril, Gabriel Ebner, and Xavier Martinet. 2022. Hypertree proof search for neural theorem proving. <i>Advances in neural information processing systems</i> , 35:26337–26349.	862
813		863
814		864
815		865
816		866
817		
818	Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. 2022. Solving quantitative reasoning problems with language models. <i>Advances in Neural Information Processing Systems</i> , 35:3843–3857.	867
819		868
820		869
821		870
822		
823		
824		
	Bohan Li, Yutai Hou, and Wanxiang Che. 2022. Data augmentation approaches in natural language processing: A survey. <i>Ai Open</i> , 3:71–90.	871
		872
		873
		874
		875
	Chen Li, Weiqi Wang, Jingcheng Hu, Yixuan Wei, Nanning Zheng, Han Hu, Zheng Zhang, and Houwen Peng. 2024a. Common 7b language models already possess strong math capabilities. <i>arXiv preprint arXiv:2403.04706</i> .	876
		877
		878
		879
		880
		881
	Chengpeng Li, Guanting Dong, Mingfeng Xue, Ru Peng, Xiang Wang, and Dayiheng Liu. 2024b. Dotamath: Decomposition of thought with code assistance and self-correction for mathematical reasoning. <i>arXiv preprint arXiv:2407.04078</i> .	
	Chengpeng Li, Zheng Yuan, Hongyi Yuan, Guanting Dong, Keming Lu, Jiancan Wu, Chuanqi Tan, Xiang Wang, and Chang Zhou. 2024c. Mugglemath: Assessing the impact of query and response augmentation on math reasoning. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 10230–10258.	
	Qingyao Li, Lingyue Fu, Weiming Zhang, Xianyu Chen, Jingwei Yu, Wei Xia, Weinan Zhang, Ruiming Tang, and Yong Yu. 2023a. Adapting large language models for education: Foundational capabilities, potentials, and challenges. <i>arXiv preprint arXiv:2401.08664</i> .	
	Qintong Li, Leyang Cui, Xueliang Zhao, Lingpeng Kong, and Wei Bi. 2024d. Gsm-plus: A comprehensive benchmark for evaluating the robustness of llms as mathematical problem solvers. <i>arXiv preprint arXiv:2402.19255</i> .	
	Wenhua Li, Tao Zhang, Rui Wang, Shengjun Huang, and Jing Liang. 2023b. Multimodal multi-objective optimization: Comparative study of the state-of-the-art. <i>Swarm and Evolutionary Computation</i> , 77:101253.	
	Xiaoyuan Li, Wenjie Wang, Moxin Li, Junrong Guo, Yang Zhang, and Fuli Feng. 2024e. Evaluating mathematical reasoning of large language models: A focus on error identification and correction. <i>arXiv preprint arXiv:2406.00755</i> .	
	Zhaoyu Li, Jialiang Sun, Logan Murphy, Qidong Su, Zenan Li, Xian Zhang, Kaiyu Yang, and Xujie Si. 2024f. A survey on deep learning for theorem proving. <i>arXiv preprint arXiv:2404.09939</i> .	
	Zhihao Li, Yao Du, Yang Liu, Yan Zhang, Yufang Liu, Mengdi Zhang, and Xunliang Cai. 2024g. Eagle: Elevating geometric reasoning through llm-empowered visual instruction tuning. <i>arXiv preprint arXiv:2408.11397</i> .	
	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. 2024h. Cmath: A chinese multi-modal math skill evaluation benchmark for foundation models. <i>arXiv preprint arXiv:2407.12023</i> .	

882	Zhong-Zhi Li, Ming-Liang Zhang, Fei Yin, and Cheng-Lin Liu. 2023c. Lans: A layout-aware neural solver for plane geometry problem. arXiv preprint arXiv:2311.16476 .	936
883		937
884		938
885		939
		940
886	Paul Pu Liang, Yun Cheng, Xiang Fan, Chun Kai Ling, Suzanne Nie, Richard Chen, Zihao Deng, Nicholas Allen, Randy Auerbach, Faisal Mahmood, et al. 2024a. Quantifying & modeling multimodal interactions: An information decomposition framework. <i>Advances in Neural Information Processing Systems</i> , 36.	941
887		942
888		943
889		944
890		945
891		946
892		
893	Zhenwen Liang, Ye Liu, Tong Niu, Xiangliang Zhang, Yingbo Zhou, and Semih Yavuz. 2024b. Improving llm reasoning through scaling inference computation with collaborative verification. arXiv preprint arXiv:2410.05318 .	947
894		948
895		949
896		950
897		951
898	Zhenwen Liang, Tianyu Yang, Jipeng Zhang, and Xiangliang Zhang. 2023a. Unimath: A foundational and multimodal mathematical reasoner. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 7126–7133.	952
899		953
900		954
901		955
902		956
		957
903	Zhenwen Liang, Dian Yu, Xiaoman Pan, Wenlin Yao, Qingkai Zeng, Xiangliang Zhang, and Dong Yu. 2023b. Mint: Boosting generalization in mathematical reasoning via multi-view fine-tuning. arXiv preprint arXiv:2307.07951 .	958
904		959
905		960
906		961
907		962
908	Zhenwen Liang, Dian Yu, Wenhao Yu, Wenlin Yao, Zhihan Zhang, Xiangliang Zhang, and Dong Yu. 2024c. Mathchat: Benchmarking mathematical reasoning and instruction following in multi-turn interactions. arXiv preprint arXiv:2405.19444 .	963
909		964
910		965
911		966
912		967
		968
913	Zhenwen Liang, Jipeng Zhang, Lei Wang, Wei Qin, Yunshi Lan, Jie Shao, and Xiangliang Zhang. 2021. Mwp-bert: Numeracy-augmented pre-training for math word problem solving. arXiv preprint arXiv:2107.13435 .	969
914		970
915		971
916		972
917		973
918	Zhenghao Lin, Zhibin Gou, Yeyun Gong, Xiao Liu, Yelong Shen, Ruochen Xu, Chen Lin, Yujiu Yang, Jian Jiao, Nan Duan, et al. 2024. Rho-1: Not all tokens are what you need. arXiv preprint arXiv:2404.07965 .	974
919		975
920		976
921		977
922	Haoxiong Liu, Yifan Zhang, Yifan Luo, and Andrew Chi-Chih Yao. 2024a. Augmenting math word problems via iterative question composing. arXiv preprint arXiv:2401.09003 .	978
923		979
924		980
925		981
		982
926	Hongwei Liu, Zilong Zheng, Yuxuan Qiao, Haodong Duan, Zhiwei Fei, Fengzhe Zhou, Wenwei Zhang, Songyang Zhang, Dahua Lin, and Kai Chen. 2024b. Mathbench: Evaluating the theory and application proficiency of llms with a hierarchical mathematics benchmark. arXiv preprint arXiv:2405.12209 .	983
927		
928		
929		
930		
931		
932	Junling Liu, Ziming Wang, Qichen Ye, Dading Chong, Peilin Zhou, and Yining Hua. 2023a. Qilin-med-vl: Towards chinese large vision-language model for general healthcare. arXiv preprint arXiv:2310.17956 .	984
933		985
934		986
935		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

990	Jingkun Ma, Runzhe Zhan, Derek F Wong, Yang Li, Di Sun, Hou Pong Chan, and Lidia S Chao. 2024. Vi-saidmath: Benchmarking visual-aided mathematical reasoning. arXiv preprint arXiv:2410.22995 .	1041
991		1042
992		1043
993		1044
994	Yujun Mao, Yoon Kim, and Yilun Zhou. 2024. Champ: A competition-level dataset for fine-grained analyses of llms’ mathematical reasoning capabilities. arXiv preprint arXiv:2401.06961 .	1045
995		1046
996		1047
997		1048
998	Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. 2024. Gsm-symbolic: Understanding the limitations of mathematical reasoning in large language models. arXiv preprint arXiv:2410.05229 .	1049
999		1050
1000		1051
1001		1052
1002		1053
1003	Swaroop Mishra, Matthew Finlayson, Pan Lu, Leonard Tang, Sean Welleck, Chitta Baral, Tanmay Rajpurohit, Oyvind Taffjord, Ashish Sabharwal, Peter Clark, et al. 2022. Lila: A unified benchmark for mathematical reasoning. arXiv preprint arXiv:2210.17517 .	1054
1004		1055
1005		1056
1006		1057
1007		1058
1008	MistralAI. 2024. Mathstral official platform .	1059
1009	MoonshotAI. 2024. k0-math official platform .	1060
1010	Harold Mouchere, Christian Viard-Gaudin, Dae Hwan Kim, Jin Hyung Kim, and Utpal Garain. 2011. Crohme2011: Competition on recognition of on-line handwritten mathematical expressions. In 2011 international conference on document analysis and recognition , pages 1497–1500. IEEE.	1061
1011		1062
1012		1063
1013		1064
1014		1065
1015		1066
1016	Zabir Al Nazi and Wei Peng. 2024. Large language models in healthcare and medical domain: A review. In Informatics , volume 11, page 57. MDPI.	1067
1017		1068
1018		1069
1019	Bolin Ni, JingCheng Hu, Yixuan Wei, Houwen Peng, Zheng Zhang, Gaofeng Meng, and Han Hu. 2024. Xwin-lm: Strong and scalable alignment practice for llms. arXiv preprint arXiv:2405.20335 .	1070
1020		1071
1021		1072
1022		1073
1023	OpenAI. 2024. Gpt-4o system card .	1074
1024	Jeff Z Pan, Simon Razniewski, Jan-Christoph Kalo, Sneha Singhania, Jiaoyan Chen, Stefan Dietze, Hajira Jabeen, Janna Omeliyanenko, Wen Zhang, Matteo Lissandrini, et al. 2023. Large language models and knowledge graphs: Opportunities and challenges. arXiv preprint arXiv:2308.06374 .	1075
1025		1076
1026		1077
1027		1078
1028		1079
1029		1080
1030	Shuai Peng, Di Fu, Liangcai Gao, Xiuqin Zhong, Hongguang Fu, and Zhi Tang. 2024. Multimath: Bridging visual and mathematical reasoning for large language models. arXiv preprint arXiv:2409.00147 .	1081
1031		1082
1032		1083
1033		1084
1034	Gabriel Poesia, David Broman, Nick Haber, and Noah D Goodman. 2024. Learning formal mathematics from intrinsic motivation. arXiv preprint arXiv:2407.00695 .	1085
1035		1086
1036		1087
1037		1088
1038	Stanislas Polu and Ilya Sutskever. 2021. Generative language modeling for automated theorem proving. arXiv preprint arXiv:2009.03393 .	1089
1039		1090
1040		1091
	Jinghui Qin, Zhicheng Yang, Jiaqi Chen, Xiaodan Liang, and Liang Lin. 2023. Template-based contrastive distillation pretraining for math word problem solving. IEEE Transactions on Neural Networks and Learning Systems .	1092
		1093
	Qwen. 2024. Qwen2-math technical report .	1094
		1095
	AM Rahman, Junyi Ye, Wei Yao, Wenpeng Yin, and Guiling Wang. 2024. From blind solvers to logical thinkers: Benchmarking llms’ logical integrity on faulty mathematical problems. arXiv preprint arXiv:2410.18921 .	1096
		1097
	Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soriccut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. arXiv preprint arXiv:2403.05530 .	1098
		1099
	Ankit Satpute, Noah Gießing, André Greiner-Petter, Moritz Schubotz, Olaf Teschke, Akiko Aizawa, and Bela Gipp. 2024. Can llms master math? investigating large language models on math stack exchange. In Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval , pages 2316–2320.	1100
		1101
	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. arXiv preprint arXiv:2402.03300 .	1102
		1103
	Wenhao Shi, Zhiqiang Hu, Yi Bin, Junhua Liu, Yang Yang, See-Kiong Ng, Lidong Bing, and Roy Ka-Wei Lee. 2024. Math-llava: Bootstrapping mathematical reasoning for multimodal large language models. arXiv preprint arXiv:2406.17294 .	1104
		1105
	Shiven Sinha, Ameya Prabhu, Ponnurangam Kumaraguru, Siddharth Bhat, and Matthias Bethge. 2024. Wu’s method can boost symbolic ai to rival silver medalists and alphageometry to outperform gold medalists at imo geometry. arXiv preprint arXiv:2404.06405 .	1106
		1107
	Shezheng Song, Xiaopeng Li, Shasha Li, Shan Zhao, Jie Yu, Jun Ma, Xiaoguang Mao, and Weimin Zhang. 2023. How to bridge the gap between modalities: A comprehensive survey on multimodal large language model. arXiv preprint arXiv:2311.07594 .	1108
		1109
	SquirrelAiLearning. 2024. Squirrel ai official platform .	1110
		1111
	Pragya Srivastava, Manuj Malik, Vivek Gupta, Tanuja Ganu, and Dan Roth. 2024. Evaluating llms’ mathematical reasoning in financial document question answering. In Findings of the Association for Computational Linguistics ACL 2024 , pages 3853–3878.	1112
		1113
	Kai Sun, Yushi Bai, Ji Qi, Lei Hou, and Juanzi Li. 2024a. Mm-math: Advancing multimodal math evaluation with process evaluation and fine-grained	1114

1096	classification. In <u>Findings of the Association for Computational Linguistics: EMNLP 2024</u> , pages 1358–1375.	enhanced mathematical reasoning. <u>arXiv preprint arXiv:2310.03731</u> .	1150
1097			1151
1098			
1099	Liangtai Sun, Yang Han, Zihan Zhao, Da Ma, Zhen-nan Shen, Baocai Chen, Lu Chen, and Kai Yu. 2024b. Scieval: A multi-level large language model evaluation benchmark for scientific research. In <u>Proceedings of the AAI Conference on Artificial Intelligence</u> , volume 38, pages 19053–19061.	Lei Wang, Shan Dong, Yuhui Xu, Hanze Dong, Yalu Wang, Amrita Saha, Ee-Peng Lim, Caiming Xiong, and Doyen Sahoo. 2024b. Mathhay: An automated benchmark for long-context mathematical reasoning in llms. <u>arXiv preprint arXiv:2410.04698</u> .	1152
1100			1153
1101			1154
1102			1155
1103			1156
1104			
1105	Lin Zhuang Sun, Hao Liang, Jingxuan Wei, Bihui Yu, Conghui He, Zenan Zhou, and Wentao Zhang. 2024c. Beats: Optimizing llm mathematical capabilities with backverify and adaptive disambiguate based efficient tree search. <u>arXiv preprint arXiv:2409.17972</u> .	Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang Sui. 2024c. Math-shepherd: Verify and reinforce llms step-by-step without human annotations. In <u>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</u> , pages 9426–9439.	1157
1106			1158
1107			1159
1108			1160
1109			1161
1110	Dídac Surís, Sachit Menon, and Carl Vondrick. 2023. Vipergpt: Visual inference via python execution for reasoning. In <u>Proceedings of the IEEE/CVF International Conference on Computer Vision</u> , pages 11888–11898.	Shen Wang, Tianlong Xu, Hang Li, Chaoli Zhang, Joleen Liang, Jiliang Tang, Philip S Yu, and Qingsong Wen. 2024d. Large language models for education: A survey and outlook. <u>arXiv preprint arXiv:2403.18105</u> .	1162
1111			1163
1112			1164
1113			1165
1114			1166
1115	TALEducation. 2023. Mathgpt official platform .		1167
1116	Jiamin Tang, Chao Zhang, Xudong Zhu, and Mengchi Liu. 2024. Tangram: A challenging benchmark for geometric element recognizing. <u>arXiv preprint arXiv:2408.13854</u> .	Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. 2023b. Scibench: Evaluating college-level scientific problem-solving abilities of large language models. <u>arXiv preprint arXiv:2307.10635</u> .	1168
1117			1169
1118			1170
1119			1171
1120	Amitayush Thakur, George Tsoukalas, Yeming Wen, Jimmy Xin, and Swarat Chaudhuri. 2024. An in-context learning agent for formal theorem-proving. In <u>First Conference on Language Modeling</u> .	Yiqi Wang, Wentao Chen, Xiaotian Han, Xudong Lin, Haiteng Zhao, Yongfei Liu, Bohan Zhai, Jianbo Yuan, Quanzeng You, and Hongxia Yang. 2024e. Exploring the reasoning abilities of multimodal large language models (mllms): A comprehensive survey on emerging trends in multimodal reasoning. <u>arXiv preprint arXiv:2401.06805</u> .	1172
1121			1173
1122			1174
1123			1175
1124	Yuxuan Tong, Xiwen Zhang, Rui Wang, Ruidong Wu, and Junxian He. 2024. Dart-math: Difficulty-aware rejection tuning for mathematical problem-solving. <u>arXiv preprint arXiv:2407.13690</u> .	Yixu Wang, Wenpin Qian, Hong Zhou, Jianfeng Chen, and Kai Tan. 2023c. Exploring new frontiers of deep learning in legal practice: A case study of large language models. <u>International Journal of Computer Science and Information Technology</u> , 1(1):131–138.	1176
1125			1177
1126			1178
1127			1179
1128	Shubham Toshniwal, Ivan Moshkov, Sean Narenthiran, Daria Gitman, Fei Jia, and Igor Gitman. 2024. Openmathinstruct-1: A 1.8 million math instruction tuning dataset. <u>arXiv preprint arXiv:2402.10176</u> .	Chenrui Wei, Mengzhou Sun, and Wei Wang. 2024. Proving olympiad algebraic inequalities without human demonstrations. <u>arXiv preprint arXiv:2406.14219</u> .	1180
1129			1181
1130			1182
1131			1183
1132	Trieu H Trinh, Yuhuai Wu, Quoc V Le, He He, and Thang Luong. 2024. Solving olympiad geometry without human demonstrations. <u>Nature</u> , 625(7995):476–482.	Chenrui Wei, Mengzhou Sun, and Wei Wang. 2024. Proving olympiad algebraic inequalities without human demonstrations. <u>arXiv preprint arXiv:2406.14219</u> .	1184
1133			1185
1134			1186
1135			1187
1136	George Tsoukalas, Jasper Lee, John Jennings, Jimmy Xin, Michelle Ding, Michael Jennings, Amitayush Thakur, and Swarat Chaudhuri. Putnambench: A multilingual competition-mathematics benchmark for formal theorem-proving. In <u>AI for Math Workshop@ICML 2024</u> .	Chenrui Wei, Mengzhou Sun, and Wei Wang. 2024. Proving olympiad algebraic inequalities without human demonstrations. <u>arXiv preprint arXiv:2406.14219</u> .	1188
1137			1189
1138			1190
1139			1191
1140			1192
1141			1193
1142	Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Mingjie Zhan, and Hongsheng Li. 2024a. Measuring multimodal mathematical reasoning with math-vision dataset. <u>arXiv preprint arXiv:2402.14804</u> .	Jiayang Wu, Wensheng Gan, Zefeng Chen, Shicheng Wan, and S Yu Philip. 2023a. Multimodal large language models: A survey. In <u>2023 IEEE International Conference on Big Data (BigData)</u> , pages 2247–2256. IEEE.	1194
1143			1195
1144			1196
1145			1197
1146	Ke Wang, Houxing Ren, Aojun Zhou, Zimu Lu, Sichun Luo, Weikang Shi, Renrui Zhang, Linqi Song, Mingjie Zhan, and Hongsheng Li. 2023a. Mathcoder: Seamless code integration in llms for	Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhanjan Kam-badur, David Rosenberg, and Gideon Mann. 2023b. Bloomberggpt: A large language model for finance. <u>arXiv preprint arXiv:2303.17564</u> .	1198
1147			1199
1148			1200
1149			1201
		Ting Wu, Xuefeng Li, and Pengfei Liu. 2024a. Progress or regress? self-improvement reversal in post-training. <u>arXiv preprint arXiv:2407.05013</u> .	1202
			1203

1204	Zhenyu Wu, Meng Jiang, and Chao Shen. 2024b. Get an a in math: Progressive rectification prompting. In <u>Proceedings of the AAAI Conference on Artificial Intelligence</u> , volume 38, pages 19288–19296.	Xin Xu, Tong Xiao, Zitong Chao, Zhenya Huang, Can Yang, and Yang Wang. 2024e. Can llms solve longer math word problems better? <u>arXiv preprint arXiv:2405.14804</u> .	1260
1205			1261
1206			1262
1207			1263
1208	Zijian Wu, Suozhi Huang, Zhejian Zhou, Huaiyuan Ying, Jiayu Wang, Dahua Lin, and Kai Chen. 2024c. Internlm2. 5-stepprover: Advancing automated theorem proving via expert iteration on large-scale lean problems. <u>arXiv preprint arXiv:2410.15700</u> .	Yifan Xu, Xiao Liu, Xinghan Liu, Zhenyu Hou, Yueyan Li, Xiaohan Zhang, Zihan Wang, Aohan Zeng, Zhengxiao Du, Wenyi Zhao, et al. 2024f. Chatglm-math: Improving math problem-solving in large language models with a self-critique pipeline. <u>arXiv preprint arXiv:2404.02893</u> .	1264
1209			1265
1210			1266
1211			1267
1212			1268
1213	Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, et al. 2023. The rise and potential of large language model based agents: A survey. <u>arXiv preprint arXiv:2309.07864</u> .	Yibo Yan and Joey Lee. 2024. Georeasoner: Reasoning on geospatially grounded context for natural language understanding. In <u>Proceedings of the 33rd ACM International Conference on Information and Knowledge Management</u> , pages 4163–4167.	1270
1214			1271
1215			1272
1216			1273
1217			1274
1218	Changrong Xiao, Sean Xin Xu, and Kunpeng Zhang. 2023. Multimodal data augmentation for image captioning using diffusion models. In <u>Proceedings of the 1st Workshop on Large Generative Models Meet Multimodal Applications</u> , pages 23–33.	Yibo Yan, Shen Wang, Jiahao Huo, Hang Li, Boyan Li, Jiamin Su, Xiong Gao, Yi-Fan Zhang, Tianlong Xu, Zhendong Chu, et al. 2024a. Errorradar: Benchmarking complex mathematical reasoning of multimodal large language models via error detection. <u>arXiv preprint arXiv:2410.04509</u> .	1275
1219			1276
1220			1277
1221			1278
1222			1279
1223	Junlin Xie, Zhihong Chen, Ruifei Zhang, Xiang Wan, and Guanbin Li. 2024. Large multimodal agents: A survey. <u>arXiv preprint arXiv:2402.15116</u> .	Yibo Yan, Haomin Wen, Siru Zhong, Wei Chen, Haodong Chen, Qingsong Wen, Roger Zimmermann, and Yuxuan Liang. 2024b. Urbanclip: Learning text-enhanced urban region profiling with contrastive language-image pretraining from the web. In <u>Proceedings of the ACM on Web Conference 2024</u> , pages 4006–4017.	1280
1224			1281
1225			1282
1226	Huajian Xin, Daya Guo, Zhihong Shao, Zhizhou Ren, Qihao Zhu, Bo Liu, Chong Ruan, Wenda Li, and Xiaodan Liang. 2024a. Deepseek-prover: Advancing theorem proving in llms through large-scale synthetic data. <u>arXiv preprint arXiv:2405.14333</u> .	Yuchen Yan, Jin Jiang, Yang Liu, Yixin Cao, Xin Xu, Xunliang Cai, Jian Shao, et al. 2024c. S3c-math: Spontaneous step-level self-correction makes large language models better mathematical reasoners. <u>arXiv preprint arXiv:2409.01524</u> .	1283
1227			1284
1228			1285
1229			1286
1230			1287
1231	Huajian Xin, ZZ Ren, Junxiao Song, Zhihong Shao, Wanxia Zhao, Haocheng Wang, Bo Liu, Liyue Zhang, Xuan Lu, Qiushi Du, et al. 2024b. Deepseek-prover-v1. 5: Harnessing proof assistant feedback for reinforcement learning and monte-carlo tree search. <u>arXiv preprint arXiv:2408.08152</u> .	An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, et al. 2024a. Qwen2. 5-math technical report: Toward mathematical expert model via self-improvement. <u>arXiv preprint arXiv:2409.12122</u> .	1288
1232			1289
1233			1290
1234			1291
1235			1292
1236			1293
1237	Wei Xiong, Chengshuai Shi, Jiaming Shen, Aviv Rosenberg, Zhen Qin, Daniele Calandriello, Misha Khalman, Rishabh Joshi, Bilal Piot, Mohammad Saleh, et al. 2024. Building math agents with multi-turn iterative preference learning. <u>arXiv preprint arXiv:2409.02392</u> .	Hongyang Yang, Xiao-Yang Liu, and Christina Dan Wang. 2023a. Fingpt: Open-source financial large language models. <u>arXiv preprint arXiv:2306.06031</u> .	1294
1238			1295
1239			1296
1240			1297
1241			1298
1242			1299
1243	Hanyi Xu, Wensheng Gan, Zhenlian Qi, Jiayang Wu, and Philip S Yu. 2024a. Large language models for education: A survey. <u>arXiv preprint arXiv:2405.13001</u> .	Zhen Yang, Jinhao Chen, Zhengxiao Du, Wenmeng Yu, Weihang Wang, Wenyi Hong, Zhihuan Jiang, Bin Xu, Yuxiao Dong, and Jie Tang. 2024b. Mathglm-vision: Solving mathematical problems with multimodal large language model. <u>arXiv preprint arXiv:2409.13729</u> .	1300
1244			1301
1245			1302
1246			1303
1247	Liang Xu, Hang Xue, Lei Zhu, and Kangkang Zhao. 2024b. Superclue-math6: Graded multi-step math reasoning benchmark for llms in chinese. <u>arXiv preprint arXiv:2401.11819</u> .	Zhen Yang, Ming Ding, Qingsong Lv, Zhihuan Jiang, Zehai He, Yuyi Guo, Jinfeng Bai, and Jie Tang. 2023b. Gpt can solve mathematical problems without a calculator. <u>arXiv preprint arXiv:2309.03241</u> .	1304
1248			1305
1249			1306
1250			1307
1251	Tianlong Xu, Richard Tong, Jing Liang, Xing Fan, Haoyang Li, and Qingsong Wen. 2024c. Foundation models for education: Promises and prospects. <u>arXiv preprint arXiv:2405.10959</u> .	Zhicheng Yang, Jinghui Qin, Jiaqi Chen, Liang Lin, and Xiaodan Liang. 2022. Logicsolver: Towards interpretable math word problem solving with logical prompt-enhanced learning. <u>arXiv preprint arXiv:2205.08232</u> .	1308
1252			1309
1253			1310
1254			1311
1255	Tianlong Xu, Yi-Fan Zhang, Zhendong Chu, Shen Wang, and Qingsong Wen. 2024d. Ai-driven virtual teacher for enhanced educational efficiency: Leveraging large pretrain models for autonomous error analysis and correction. <u>arXiv preprint arXiv:2409.09403</u> .		1312
1256			1313
1257			1314
1258			1315
1259			1316

1317	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2022.	Beichen Zhang, Kun Zhou, Xilin Wei, Xin Zhao, Jing Sha, Shijin Wang, and Ji-Rong Wen. 2024a.	1371
1318	React: Synergizing reasoning and acting in language models. arXiv preprint arXiv:2210.03629 .	Evaluating and improving tool-augmented computation-intensive math reasoning. Advances in Neural Information Processing Systems , 36.	1372
1319			1373
1320			1374
1321	Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. 2024.	Di Zhang, Jianbo Wu, Jingdi Lei, Tong Che, Jia-tong Li, Tong Xie, Xiaoshui Huang, Shufei Zhang, Marco Pavone, Yuqiang Li, et al. 2024b.	1376
1322	A survey on multimodal large language models. National Science Review , page nwae403.	Llama-berry: Pairwise optimization for o1-like olympiad-level mathematical reasoning. arXiv preprint arXiv:2410.02884 .	1377
1323			1378
1324			1379
1325	Huaiyuan Ying, Shuo Zhang, Linyang Li, Zhejian Zhou, Yunfan Shao, Zhaoye Fei, Yichuan Ma, Jiawei Hong, Kuikun Liu, Ziyi Wang, et al. 2024.	Duzhen Zhang, Yahan Yu, Jiahua Dong, Chenxing Li, Dan Su, Chenhui Chu, and Dong Yu. 2024c.	1380
1326	Internlm-math: Open math large language models toward verifiable reasoning. arXiv preprint arXiv:2402.06332 .	Mm-llms: Recent advances in multimodal large language models. arXiv preprint arXiv:2401.13601 .	1381
1327			1382
1328			1383
1329			1384
1330	Dian Yu, Baolin Peng, Ye Tian, Linfeng Song, Haitao Mi, and Dong Yu. 2024a.	Jiaxin Zhang, Zhongzhi Li, Mingliang Zhang, Fei Yin, Chenglin Liu, and Yashar Moshfeghi. 2024d.	1385
1331	Siam: Self-improving code-assisted mathematical reasoning of large language models. arXiv preprint arXiv:2408.15565 .	Geoeval: benchmark for evaluating llms and multi-modal models on geometry problem-solving. arXiv preprint arXiv:2402.10104 .	1386
1332			1387
1333			1388
1334	Fei Yu, Hongbo Zhang, Prayag Tiwari, and Benyou Wang. 2024b.	Mengxue Zhang, Zichao Wang, Zhichao Yang, Weiqi Feng, and Andrew Lan. 2023.	1389
1335	Natural language reasoning, a survey. ACM Computing Surveys , 56(12):1–39.	Interpretable math word problem solution generation via step-by-step planning. arXiv preprint arXiv:2306.00784 .	1390
1336			1391
1337	Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhen-guo Li, Adrian Weller, and Weiyang Liu. 2023.	Ming-Liang Zhang, Zhong-Zhi Li, Fei Yin, Liang Lin, and Cheng-Lin Liu. 2024e.	1392
1338	Metamath: Bootstrap your own mathematical questions for large language models. arXiv preprint arXiv:2309.12284 .	Fuse, reason and verify: Geometry problem solving with parsed clauses from diagram. arXiv preprint arXiv:2407.07327 .	1393
1339			1394
1340			1395
1341			1396
1342			1397
1343	Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhao Chen. 2024a.	Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Yu Qiao, et al. 2024f.	1398
1344	Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In Proceedings of CVPR .	Does your multi-modal llm truly see the diagrams in visual math problems? In European Conference on Computer Vision , pages 169–186. Springer.	1400
1345			1401
1346			1402
1347			1403
1348			1404
1349			1405
1350			1406
1351			1407
1352	Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhao Chen. 2023.	Xuanyu Zhang and Qing Yang. 2023.	1408
1353	Mammoth: Building math generalist models through hybrid instruction tuning. arXiv preprint arXiv:2309.05653 .	Xuanyuan 2.0: A large chinese financial chat model with hundreds of billions parameters. In Proceedings of the 32nd ACM international conference on information and knowledge management , pages 4435–4439.	1409
1354			1410
1355			1411
1356			1412
1357	Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang, Kai Zhang, Shengbang Tong, Yuxuan Sun, Botao Yu, Ge Zhang, Huan Sun, et al. 2024b.	Yiming Zhang, Baoyi He, Shengyu Zhang, Yuhao Fu, Qi Zhou, Zhijie Sang, Zijin Hong, Kejing Yang, Wenjun Wang, Jianbo Yuan, et al. 2024g.	1413
1358	Mmmu-pro: A more robust multi-discipline multimodal understanding benchmark. arXiv preprint arXiv:2409.02813 .	Unconstrained model merging for enhanced llm reasoning. arXiv preprint arXiv:2410.13699 .	1414
1359			1415
1360			1416
1361			1417
1362	Xiang Yue, Tuney Zheng, Ge Zhang, and Wenhao Chen. 2024c.	Zeren Zhang, Jo-Ku Cheng, Jingyang Deng, Lu Tian, Jinwen Ma, Ziran Qin, Xiaokai Zhang, Na Zhu, and Tuo Leng. 2024h.	1418
1363	Mammoth2: Scaling instructions from the web. arXiv preprint arXiv:2405.03548 .	Diagram formalization enhanced multi-modal geometry problem solver. arXiv preprint arXiv:2409.04214 .	1419
1364			1420
1365	Liang Zeng, Liangjun Zhong, Liang Zhao, Tianwen Wei, Liu Yang, Jujie He, Cheng Cheng, Rui Hu, Yang Liu, Shuicheng Yan, et al. 2024.	Zhihan Zhang, Tao Ge, Zhenwen Liang, Wenhao Yu, Dian Yu, Mengzhao Jia, Dong Yu, and Meng Jiang. 2024i.	1421
1366	Skywork-math: Data scaling laws for mathematical reasoning in large language models—the story goes on. arXiv preprint arXiv:2407.08348 .	Learn beyond the answer: Training language models with reflection for mathematical reasoning. arXiv preprint arXiv:2406.12050 .	1422
1367			1423
1368			1424
1369			
1370			

1425	Wayne Xin Zhao, Kun Zhou, Zheng Gong, Beichen Zhang, Yuanhang Zhou, Jing Sha, Zhigang Chen, Shijin Wang, Cong Liu, and Ji-Rong Wen. 2022. Jiuzhang: A chinese pre-trained language model for mathematical problem understanding. In <u>Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining</u> , pages 4571–4581.	Solving math word problems via cooperative reasoning induced language models. <u>arXiv preprint arXiv:2210.16257</u> .	1481
1426			1482
1427			1483
1428			
1429		Wenwen Zhuang, Xin Huang, Xiantao Zhang, and Jin Zeng. 2024. Math-puma: Progressive upward multimodal alignment to enhance mathematical reasoning. <u>arXiv preprint arXiv:2408.08640</u> .	1484
1430			1485
1431			1486
1432			1487
1433	Xin Zhao, Kun Zhou, Beichen Zhang, Zheng Gong, Zhipeng Chen, Yuanhang Zhou, Ji-Rong Wen, Jing Sha, Shijin Wang, Cong Liu, et al. 2023. Jiuzhang 2.0: A unified chinese pre-trained language model for multi-task mathematical problem solving. In <u>Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining</u> , pages 5660–5672.	Chengke Zou, Xingang Guo, Rui Yang, Junyu Zhang, Bin Hu, and Huan Zhang. 2024. Dynamath: A dynamic visual benchmark for evaluating mathematical reasoning robustness of vision language models. <u>arXiv preprint arXiv:2411.00836</u> .	1488
1434			1489
1435			1490
1436			1491
1437			1492
1438		Xingchen Zou, Yibo Yan, Xixuan Hao, Yuehong Hu, Haomin Wen, Erdong Liu, Junbo Zhang, Yong Li, Tianrui Li, Yu Zheng, et al. 2025. Deep learning for cross-domain data fusion in urban computing: Taxonomy, advances, and outlook. <u>Information Fusion</u> , 113:102606.	1493
1439			1494
1440			1495
1441	Xueliang Zhao, Xinting Huang, Wei Bi, and Lingpeng Kong. 2024. Sego: Sequential subgoal optimization for mathematical problem-solving. <u>arXiv preprint arXiv:2310.12960</u> .		1496
1442			1497
1443			1498
1444			
1445	Tianyang Zhong, Zhengliang Liu, Yi Pan, Yutong Zhang, Yifan Zhou, Shizhe Liang, Zihao Wu, Yanjun Lyu, Peng Shu, Xiaowei Yu, et al. 2024. Evaluation of openai o1: Opportunities and challenges of agi. <u>arXiv preprint arXiv:2409.18486</u> .		
1446			
1447			
1448			
1449			
1450	Kun Zhou, Beichen Zhang, Jiapeng Wang, Zhipeng Chen, Wayne Xin Zhao, Jing Sha, Zhichao Sheng, Shijin Wang, and Ji-Rong Wen. 2024a. Jiuzhang3.0: Efficiently improving mathematical reasoning by training small data synthesis models. <u>arXiv preprint arXiv:2405.14365</u> .		
1451			
1452			
1453			
1454			
1455			
1456	Minxuan Zhou, Hao Liang, Tianpeng Li, Zhiyu Wu, Mingan Lin, Linzhuang Sun, Yaqi Zhou, Yan Zhang, Xiaoqin Huang, Yicong Chen, et al. 2024b. Mathscape: Evaluating mllms in multimodal math scenarios through a hierarchical benchmark. <u>arXiv preprint arXiv:2408.07543</u> .		
1457			
1458			
1459			
1460			
1461			
1462	Zhi Zhou, Jiang-Xin Shi, Peng-Xiao Song, Xiao-Wen Yang, Yi-Xuan Jin, Lan-Zhe Guo, and Yu-Feng Li. 2024c. Lawgpt: A chinese legal knowledge-enhanced large language model. <u>arXiv preprint arXiv:2406.04614</u> .		
1463			
1464			
1465			
1466			
1467	Zihao Zhou, Shudong Liu, Maizhen Ning, Wei Liu, Jindong Wang, Derek F Wong, Xiaowei Huang, Qifeng Wang, and Kaizhu Huang. 2024d. Is your model really a good math reasoner? evaluating mathematical reasoning with checklist. <u>arXiv preprint arXiv:2407.08733</u> .		
1468			
1469			
1470			
1471			
1472			
1473	Zihao Zhou, Qifeng Wang, Mingyu Jin, Jie Yao, Jianan Ye, Wei Liu, Wei Wang, Xiaowei Huang, and Kaizhu Huang. 2024e. Mathattack: Attacking large language models towards math solving ability. In <u>Proceedings of the AAAI Conference on Artificial Intelligence</u> , volume 38, pages 19750–19758.		
1474			
1475			
1476			
1477			
1478			
1479	Xinyu Zhu, Junjie Wang, Lin Zhang, Yuxiang Zhang, Ruyi Gan, Jiaying Zhang, and Yujiu Yang. 2022.		
1480			

A Details of Math-LLMs' Progress

The rapid development of general-purpose LLMs has made significant advancements in natural language processing tasks. However, the development of domain-specific models remains a core requirement, as they are better equipped to handle specialized tasks that general models may not address effectively. This is particularly true in fields such as healthcare (Liu et al., 2023a; Nazi and Peng, 2024), law (Cui et al., 2023; Zhou et al., 2024c; Wang et al., 2023c), finance (Wu et al., 2023b; Yang et al., 2023a; Zhang and Yang, 2023), and urban science (Yan et al., 2024b; Zou et al., 2025; Yan and Lee, 2024), where domain-specific knowledge is critical for high accuracy and performance.

In the case of mathematical reasoning, general models may struggle with tasks that require deep understanding of complex mathematical concepts, structures, and problem-solving steps. Therefore, the development of math-specific LLMs is of paramount importance, as these models are designed to enhance performance in mathematical reasoning, theorem proving, equation solving, and other math-intensive tasks.

Therefore, Table 4 provides a detailed overview of various math-specific LLMs (*i.e.*, Math-(LLMs)), sorted by their release date. It includes information about the organization behind each model, the release date, publication details, language(s) supported, parameter size, evaluation benchmarks, and whether the model is open source.

Key findings are summarized as follows:

- 1. Release Trends:** The models started emerging in 2020, with a significant increase in the number of releases from 2022 onward, indicating a growing interest in developing math-specific LLMs.
- 2. Parameter Sizes:** There is a noticeable trend towards larger parameter sizes, with some models offering up to 130B parameters, reflecting the increasing computational capacity for handling complex mathematical tasks.
- 3. Evaluation Benchmarks:** Many models are evaluated on popular benchmarks like GSM8K, MATH, and MMLU, highlighting the focus on improving performance across well-established mathematical reasoning datasets.

4. Multilingual Support: While most models are focused on English, a few (*e.g.*, MathGPT & Math-LLM) also support Chinese, showing a trend towards multilingual capabilities.

5. Open Source: A significant number of models are open-source, allowing broader access and fostering further research and development in the field.

In summary, the table reflects the rapid development of specialized Math-LLMs, with an increasing trend towards larger models, comprehensive evaluation benchmarks, and support for multilingual applications.

B Details of Metrics

B.1 Discriminative Metrics

Discriminative tasks refer to evaluation processes where the outputs are typically binary, such as "Yes" or "No". These tasks often include multiple-choice questions, fill-in-the-blank problems, or judgment assessments. The evaluation metrics focus on LLM's accuracy in specific task types and its ability to control biases.

Accuracy (ACC): It measures the proportion of correctly predicted outcomes. The value should be as high as possible.

$$ACC = \frac{\sum_{1,m} x_i}{\sum_{1,n} y_j}$$

Where x_i represents the correct output for the i -th instance, y_j represents the j -th instance, m is the number of the correct instances and n is the number of the total instances.

Exact match: It evaluates the congruence between the answers generated by LLM and the correct ones. Specifically, in cases where the answer produced LLM coincides with the reference answer, a score of 1 point will be assigned. Conversely, if there is any discrepancy between them except for bias, a score of 0 point will be given.

F_1 score: It combines two crucial aspects, namely precision and recall, in order to comprehensively assess the accuracy of LLM. It is calculated as :

$$F_1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

The value of the F_1 score ranges from 0 to 1. A higher value of the F_1 score indicates better overall performance of LLM in terms of both precision and recall.

Math (LLMs)	Organization	Release Date	Publication	Language	Parameter Size	Evaluation Benchmarks	Open Source
GPT-4 (Pola and Sutskever, 2021)	OpenAI	Sep 2020	-	English	160M/400M/700M	-	✓
Hypertree Proof Search (Lample et al., 2022)	Meta	Nov 2022	NeurIPS'22	English	-	miniF2F/Metamath	-
Minerva (Lewkowycz et al., 2022)	Google	Jun 2022	NeurIPS'22	English	8B/6.2B/5.40B	MATH/MMLU-STEM/GSM8k	-
JiuZhang 1.0 (Zhao et al., 2022)	RUC & IFLYTEK	Jun 2022	KDD'22	English	145M	-	✓
GAIRMath-Abel (Chern et al., 2023)	Shanghai Jiaotong University	2023	-	English	7B/13B/70B	GSM8k/MATH/MMLU/SVAMP/SCQ5k-English/MathQA	✓
JiuZhang 2.0 (Zhao et al., 2023)	RUC & IFLYTEK	Jan 2023	KDD ADS'23	English	-	JCAG/BAG (MathBERT/DART/JiuZhang)	✓
KwaiTiMath (Fu et al., 2023)	Kuaishou	Jan 2023	-	English/Chinese	13B	GSM8k/MathBench	✓
MathCoder (Wang et al., 2023a)	CUHK	Jan 2023	ICLR'24	English	7B/13B	GSM8k/MATH	✓
Llemma (Azerbayev et al., 2023)	Princeton University & Eleuther AI	Jan 2023	-	English	7B/34B	MATHGSM8k/MMLU-STEM/SAT/OCW/Course	✓
Skywork-13B-Math (Zeng et al., 2024) ★	SkyworkAI	Jan 2023	-	English	7B/13B	GSM8k/MATH/MATH	✓
MathGPT (TAL Education, 2023) ★	TAL Education Group	Aug 2023	-	English/Chinese	130B	C-Eval-Math/AGIEval-Math/AMESK/CMMLU-Math/GAOKAO-Math/Math01	✓
WizardMath (Luo et al., 2023)	Microsoft	Aug 2023	-	English	7B/70B	GSM8k/MATH	✓
MAMmoTH (Yue et al., 2023)	UWaterloo	Sep 2023	ICLR'24	English	7B/13B/70B	GSM/MATH/MMLU-STEM/AQuA/NumGLUE	✓
MathGLM (Yang et al., 2023b)	Tsinghua & Zhipu AI	Sep 2023	-	English	10M/100M/500M/2B/Arith. & 335M/6B/10B (MWP)	BiG-bench/ Ape210K	✓
MetaMath (Yu et al., 2023)	Cambridge & Huawei	Sep 2023	-	English	7B/13B/70B	GSM8k/MATH	✓
DeepSeekMath (Shao et al., 2024)	DeepSeek AI	Jan 2024	-	English	7B	GSM8k/MATH/OCW/SAT/MMLU-STEM/CMATH/Gaokao-Math/Ctoze/Gaokao-MathQA	✓
InternLM2.5-StepProver (Wu et al., 2024c)	Shanghai AI Lab	Jan 2024	-	English/Chinese	7B	MathUserEval/Ape210k/CMATH/GSM8k/MATH/Hungarian	✓
ChatGLM-Math (Xu et al., 2024f)	Zhipu AI	Apr 2024	-	English/Chinese	32B	MathUserEval/Ape210k/CMATH/GSM8k/MATH/Hungarian	✓
Rho-Math (Lin et al., 2024)	Microsoft	Apr 2024	-	English	18/7B	GSM8k/MATH/MMLU-STEM/SAT/SVAMP/ASDiv/MAWPS/TAB/MQA	✓
DeepSeekProver-V1 (Xin et al., 2024b)	DeepSeek AI	May 2024	-	English	7B	miniF2F/IMO	✓
InternLM2-Math (Wu et al., 2024c)	Shanghai AI Lab	May 2024	-	English/Chinese	1.8B/7B/20B/8x22B	MiniF2F-test/MATH/MATH-Python/GSM8k/MathBench-A/Hungary/	✓
JiuZhang 3.0 (Zhao et al., 2024a)	RUC & IFLYTEK	May 2024	-	English	7B/8B	GSM8k/MATH/G-Head/SVAMP/MAWPS/ASDiv/TAB/MWP	✓
MAMmoTH2 (Yue et al., 2024c)	UWaterloo	May 2024	-	English	7B/8B	TheoremQA/MATH/GSM8k/GQA/MMLU-STEM/BIBH	✓
Math-LLaVA (Shi et al., 2024)	NUS	Jun 2024	EMNLP Finding'24	English	13B	MMLU-MATH-V/MathVista	✓
Mathstral (Mistral AI, 2024)	Mistral AI	Jul 2024	-	English	7B	MATHGSM8k/GREMath/AMC2023/AIME2024/MathOdysey	✓
DeepSeek-Prover-V1.5 (Xin et al., 2024a)	DeepSeek AI	Aug 2024	-	English	7B	miniF2F/IMO/Net	✓
Qwen2-Math (Qwen, 2024)	Alibaba	Aug 2024	-	English/Chinese	1.5B/7B/72B	GSM8k/Math/MMLU-STEM/CMATH/GaokaoMath Ctoze/Gaokao Math QA	✓
Qwen2-Math-Instruct (Qwen, 2024)	Alibaba	Aug 2024	-	English/Chinese	1.5B/7B/72B	GSM8k/MATH/Marwa Math/GaoKao2023 En/Olympiad/Bench/College Math/MMLU STEM/Gaokao/CMATH/CNMiddle School 24/AIME24/AMC23	✓
MathGLM-Vision (Yang et al., 2024b) ★	Tsinghua & Zhipu AI	Sep 2024	-	English	9B/19B/32B	MathVista/MathVista(GPS)/MathVerse/Math-Vision/MMLU/MathVL	-
Math-LLM (Liu et al., 2024c) ★	East China Normal University	Sep 2024	-	Chinese	-	CMM-Math/MathVista/Math-V	-
Qwen2.5-Math (Yang et al., 2024a)	Alibaba	Sep 2024	-	English/Chinese	8.26B/7B/72B	GSM8k/MATH/MMLU-STEM/CMATH/Gaokao Math	✓
Xwin-LM (Ni et al., 2024)	Microsoft	May 2024	-	English	1.5B/7B/70B	GSM8k/MATH	✓
math-specialized Gemini 1.5 Pro ★	Google	Not launched yet	-	English	-	MATH/AIME2024/MathOdysey/HiddenMath/IMO Bench	✓
10-math (Moonshot AI, 2024)	Moonshot AI	Nov 2024	-	English/Chinese	-	KAONAN/MATH/AIME/OMNI-MATH/GAOKAO/ZHONGKAO	✓
DuoLingo Math (DuoLingo, 2024)	DuoLingo	2024	-	English	-	-	-
Khanmigo (Khan Academy, 2024)	Khan Academy	2024	-	English	-	-	-
Squirrel LAM (SquirrelAI Learning, 2024) ★	Squirrel AI Learning	2024	-	Chinese	-	-	-

Table 4: Overview of math-specific LLMs (sort by release date). ★ refers to those designed to support the multimodal mathematical setting.

Macro- F_1 score: It calculates the F_1 score for each category separately and then takes the average of the F_1 scores of all categories, so as to obtain the overall performance of LLM on all categories.

Round-r accuracy: It is the proportion of correct answers given by a model on the question set Q_r in round r . It is calculated as follows:

$$ACC_r(M) = \frac{\sum_{q \in Q_r} I[M(q) = g_t(q)]}{|Q_r|}$$

Here, $ACC_r(M)$ represents the accuracy of LLM M on question set Q_r in round r . I is an indicator function. When the answer $M(q)$ given by M for question q is consistent with the true answer $g_t(q)$ of the question, the value of I is 1; otherwise, it is 0. The symbol $\sum_{q \in Q_r}$ means summing over all questions in question set Q_r . $|Q_r|$ indicates the number of questions in question set Q_r .

ACC_{step} : It is used to evaluate LLM's ability to identify the first step where an error occurs. The accuracy for identifying the first erroneous step is calculated as follows:

$$ACC_{step} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(S_{step,i} = G_{step,i})$$

Here, N is the total number of samples. For the i -th sample, $S_{step,i}$ is the predicted step where the error occurs, and $G_{step,i}$ is the ground truth label for the first erroneous step. The indicator function $\mathbb{I}(\cdot)$ returns 1 if the predicted step matches the ground truth and 0 otherwise.

ACC_{cate} : It is for assessing LLM's performance in categorizing the type of error. The accuracy for error categorization is defined by

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

Here, N is the total number of samples. For the i -th sample, $C_{error,i}$ is the predicted error category, and $G_{error,i}$ is the ground truth label for the error category. The indicator function $\mathbb{I}(\cdot)$ has the same meaning as in the previous metric, returning 1 if the predicted error category matches the ground truth and 0 otherwise.

The skill success rate: It measures the proportion of a model correctly applying major skills in problem-solving. It's calculated by analyzing test questions and determining correct use of major skills, then finding the ratio to total questions. For example, in triangle area calculation, checking use of the area formula. Similarly, **the secondary skill success rate** focuses on the proportion of correct application of secondary skills like understanding graphic properties and unit conversion, calculated by analyzing problem-solving and finding the ratio to total questions.

The False Positive Rate (FPR): It is the proportion of cases where the evaluation LLM misjudges an incorrect answer as a correct one. A low FPR indicates that LLM rarely misjudges incorrect student answers as correct.

The False Negative Rate (FNR): It is the proportion of cases where the evaluation LLM mis-

judges a correct answer as an incorrect one. A low FNR indicates that LLM is relatively accurate in correctly determining whether a student’s answer is correct.

Mean Squared Error (MSE): It is a metric that measures the average of the squares of the differences between the LLM’s predicted values and the actual true values. It is calculated as:

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Here, n represents the number of samples. For the i -th sample, y_i is the true value and \hat{y}_i is the predicted value by LLM. The summation symbol $\sum_{i=1}^n$ means summing up the squared differences for all n samples. Dividing by n gives the average squared difference, which is the MSE. MSE should be as low as possible.

Average-Case Accuracy (A_{avg}): This metric evaluates the average accuracy of LLM across all variants of a seed question. It is calculated as the proportion of correct answers across all variants and seed questions. The formula is:

$$A_{avg} = \frac{1}{N} \sum_{i=1}^N \frac{1}{M} \sum_{j=1}^M I[\text{Ans}(i, j) = \text{GT}(i, j)]$$

where N is the total number of seed questions, M is the number of variants per seed question, and $I[\text{Ans}(i, j) = \text{GT}(i, j)]$ checks if the answer matches the ground truth.

Worst-Case Accuracy (A_{wst}): This evaluates the worst-case performance by considering the minimum accuracy across all variants of a seed question. It reflects the robustness of LLM against challenging variations. The formula is:

$$A_{wst} = \frac{1}{N} \sum_{i=1}^N \min_{j \in [1, M]} I[\text{Ans}(i, j) = \text{GT}(i, j)]$$

B.2 Generative Metrics

Generative tasks involve evaluating the content generated by LLM, typically encompassing free-form answers and responses to open-ended questions. These tasks focus primarily on assessing the extent of hallucinations in the generated content, especially when the content is not faithful to the given images. Evaluating generative tasks often requires more complex metrics, such as CHAIR and Faithscore, which measure hallucinations across different categories, including objects, attributes, and relationships within the generated content. These

metrics provide a nuanced understanding of the fidelity and reliability of MLLMs in producing content aligned with the visual and textual inputs.

Reasoning Robustness (RR): This metric measures the relative robustness of LLM by comparing the worst-case performance to the average-case performance. The formula is:

$$RR = \frac{A_{wst}}{A_{avg}}$$

Repetition Consistency (RC): This evaluates the consistency of LLM’s responses across repeated queries for the same question variant. It helps distinguish between variability due to randomness and systematic errors. The formula is:

$$RC(i, j) = \frac{1}{K} \sum_{k=1}^K I[\text{Ans}_k(i, j) = \text{Ans}(i, j)]$$

where K is the number of repetitions.

OpenCompass Scoring: It is a comprehensive evaluation framework that leverages the OpenCompass platform to assess the generative capabilities of LLM across multiple dimensions. Perplexity (PPL) evaluates the naturalness and fluency of generated text, with lower scores indicating greater model confidence and the ability to produce contextually coherent sequences. Simultaneously, CircularEval assesses the robustness and consistency of LLM in multiple-choice scenarios by evaluating its performance across N random permutations of the options in an N -option question. A question is deemed correctly answered only if LLM provides the correct response for all permutations, highlighting its ability to handle randomized inputs reliably.

Bilingual Evaluation Understudy (BLEU): It evaluates the quality of text generation by measuring n-gram overlap between generated and reference texts, focusing on precision and brevity. Its formula is:

$$\text{BLEU} = \text{BP} \cdot \exp \left(\sum_{n=1}^N w_n \log p_n \right)$$

where BP is the brevity penalty, calculated as 1 if $c > r$, or $\exp(1 - r/c)$ if $c \leq r$, with c and r representing the lengths of the generated and reference texts, respectively. w_n denotes n-gram weights (typically uniform), and p_n is the precision of n-grams of size n . BLEU scores range from 0 to 1 (often expressed as percentages, 0-100%), with higher scores indicating greater similarity between the generated and reference texts.

Recall-Oriented Understudy for Gisting Evaluation-L (ROUGE-L): It evaluates the quality of generated text by measuring its similarity to reference text, focusing on sequence alignment and structural consistency through the Longest Common Subsequence (LCS). It calculates recall as the proportion of the LCS length relative to the reference text length. The formula of recall is:

$$R = \frac{\text{LCS}(\text{Generated}, \text{Reference})}{\text{Length}(\text{Reference})}$$

It also calculates precision as the proportion of the LCS length relative to the generated text length. The formula is:

$$P = \frac{\text{LCS}(\text{Generated}, \text{Reference})}{\text{Length}(\text{Generated})}$$

The F_1 score is a harmonic mean of precision and recall, expressed as:

$$F_1 = \frac{(1 + \beta^2) \cdot P \cdot R}{\beta^2 \cdot P + R}$$

where β (commonly set to 1) controls the weighting of recall and precision. ROUGE-L scores range from 0 to 1, with higher scores indicating greater similarity between the generated and reference texts.

Consensus-based Image Description Evaluation (CIDEr): It is designed for image description tasks, measuring the semantic relevance of generated descriptions by calculating the TF-IDF weighted n-gram similarity with reference descriptions. The formula is:

$$\text{CIDEr}_n(c_i, S_i) = \frac{1}{m} \sum_{j=1}^m \frac{g^n(c_i) \cdot g^n(s_{ij})}{\|g^n(c_i)\| \cdot \|g^n(s_{ij})\|}$$

$$\text{CIDEr}(c_i, S_i) = \sum_{n=1}^N w_n \text{CIDEr}_n(c_i, S_i)$$

Here, c_i is the candidate description, $S_i = \{s_{i1}, s_{i2}, \dots, s_{im}\}$ is the set of reference descriptions, and m is the number of references. $g^n(c_i)$ and $g^n(s_{ij})$ are the TF-IDF weighted n-gram vectors for the candidate and reference descriptions, with $\|g^n(c_i)\|$ and $\|g^n(s_{ij})\|$ being their magnitudes. w_n is the weight for n-grams of different lengths, usually $w_n = 1/N$, where N is the maximum n-gram length. Scores range from 0 to 10, with higher scores indicating stronger alignment between candidate and reference descriptions.

Mathematical Symbol Similarity: This metric measures the similarity between the correct steps in a reasoning process and the steps generated by LLM, using symbolic computation software to perform the evaluation.

GPT Scoring: This metric evaluates the generated content based on scores assigned by GPT or other language models, focusing on the linguistic coherence and logical consistency of the text.

Context Length Generalization Efficacy (CoLeG-E): It is a metric used to measure LLM's consistency in answering variations of the same question across different context lengths. It is defined as:

$$\text{CoLeG-E}(M) = \frac{\sum_{q \in Q_R} [\bigwedge_{r=1}^R I[M(q^r) = gt(q^r)]]}{|Q_R|}$$

where Q_R represents the set of all questions under evaluation, and q^r refers to the r -th variation of a question q , corresponding to a specific context length. $M(q^r)$ is LLM's predicted answer for the r -th variation, while $gt(q^r)$ denotes the ground truth answer. The indicator function $I[\cdot]$ equals 1 if LLM's answer matches the ground truth, and 0 otherwise. The logical AND operator $\bigwedge_{r=1}^R$ ensures that the model must answer all variations of a question correctly for that question to be considered correctly answered.

Context Length Generalization Robustness (CoLeG-R): It measures LLM's robustness to context length expansion by quantifying the relative drop in accuracy from initial to extended questions. It is defined as:

$$\text{CoLeG-R}(M) = 1 - \frac{\text{ACC}_0(M) - \text{ACC}_R(M)}{\text{ACC}_0(M)}$$

Here, $\text{ACC}_0(M)$ is the LLM's accuracy on the initial set of shorter-context questions Q_0 , and $\text{ACC}_R(M)$ is its accuracy on the extended longer-context questions Q_R . Higher CoLeG-R values indicate better robustness, with less performance degradation across context lengths.

Performance Drop Rate (PDR): This metric measures the relative decline in model performance when transitioning from the original dataset to the perturbed dataset. It is defined as:

$$\text{PDR} = 1 - \frac{\sum_{(x,y) \in D_a} I[\text{LLM}(x), y] / |D_a|}{\sum_{(x,y) \in D} I[\text{LLM}(x), y] / |D|}$$

where D is the original dataset and D_a is the perturbed dataset. $I[\text{LLM}(x), y]$ is an indicator function that checks if the LLM's output matches the ground truth y .

Accurately Solved Pairs (ASP): ASP measures the percentage of seed questions and their perturbed variations that are both correctly answered by LLM. It is defined as:

$$ASP = \frac{\sum_{x,y;x',y'} I[\text{LLM}(x), y] \cdot I[\text{LLM}(x', y')]}{N \cdot |D|}$$

where x and x' are a seed question and its variation, respectively. N is the number of perturbations per question. $|D|$ is the total number of seed questions.

Mean Average Precision (mAP): It is a metric that evaluates LLM’s ability to rank relevant answers higher in its output list for a given query. It is defined as:

$$mAP = \frac{1}{|Q|} \sum_{q \in Q} AP(q)$$

$$AP(q) = \frac{1}{m} \sum_{k=1}^m P(k)$$

$$P(k) = \frac{\# \text{ relevant ans retrieved up to position } k}{k}$$

Here, Q represents the set of all queries in the dataset. $AP(q)$ is the Average Precision for query q , calculated as the mean of the precision values $P(k)$ at ranks where relevant answers appear. $P(k)$ is the precision at rank k , representing the proportion of relevant answers retrieved up to position k . m is the total number of relevant answers for query q .

Training Set Coverage (TSC): It measures how effectively LLM has learned to generate correct solutions for tasks similar to those in its training set. TSC is particularly useful in cross-domain or cross-modal tasks, where it assesses LLM’s ability to generalize learned patterns to problems aligned with its training data. Higher TSC scores indicate better learning and consistency, while lower scores suggest insufficient training or overfitting.

Pass@N: This metric measures the likelihood of LLM generating at least one correct solution within N attempts for a given problem. Formally:

$$Pass@N = \mathbb{E}_{\text{Problems}}[\min(c, 1)]$$

where c represents the number of correct answers out of N responses. A higher Pass@N indicates a greater chance of producing a correct answer in multiple attempts, reflecting LLM’s potential capability.

PassRatio@N: This metric calculates the proportion of correct answers among N generated responses for a given problem. It is defined as:

$$PassRatio@N = \mathbb{E}_{\text{Problems}} \left[\frac{c}{N} \right]$$

where c is the count of correct answers. This metric reflects LLM’s stability in consistently generating correct answers. It can be considered analogous to Pass@1 but offers reduced variance.