# Supplementary Material - MATHVERSE: Does Your Multi-modal LLM Truly See the Diagrams in Visual Math Problems?

Renrui Zhang<sup>\*‡1,2</sup>, Dongzhi Jiang<sup>\*1</sup>, Yichi Zhang<sup>\*2</sup>, Haokun Lin<sup>2</sup>, Ziyu Guo<sup>3</sup> Pengshuo Qiu<sup>2</sup>, Aojun Zhou<sup>1</sup>, Pan Lu<sup>4</sup>, Kai-Wei Chang<sup>4</sup> Yu Qiao<sup>2</sup>, Peng Gao<sup>2</sup>, and Hongsheng Li<sup>1,5</sup>

<sup>1</sup> CUHK MMLab
<sup>2</sup> Shanghai AI Laboratory
<sup>3</sup> CUHK MiuLar Lab
<sup>4</sup> University of California, Los Angeles
<sup>5</sup> CPII under InnoHK
{renruizhang, dzjiang, ziyuguo}@link.cuhk.edu.hk
gaopeng@pjlab.org.cn hsli@ee.cuhk.edu.hk

# Overview

- Section A: Related work.
- Section B: Additional experimental details.
- Section C: Additional experiments and analysis.
- Section D: More dataset details.
- Section E: Comparison to current benchmarks.
- Section F: Limitation and future work.
- Section G: Error analysis.
- Section H: Qualitative examples.

# A Related Work

*Multi-modal Large Language Models* (MLLMs), building upon the prevalence of Large Language Models (LLMs) [4, 27, 44, 52, 53], have become increasingly prominent for incorporating multi-modalities [28,48,64,66,67]. They extend LLMs to tackle a diverse range of tasks and domains, including the mainstream 2D images [1,14,29,32] and other modalities, such as 3D point clouds [22,26,56], audio [23,51], video [6,63], and robotics [33,37,38,57]. Noteworthy examples like OpenAI's GPT-4V [46] and Google's Gemini [20] exhibit exceptional visual understanding and reasoning capabilities, setting new benchmarks in multi-modal performance. However, their closed-source nature poses a barrier to the broader application and development of MLLMs. Concurrently, another line of work is

<sup>\*</sup>Equal contribution <sup>‡</sup>Project lead <sup>†</sup>Corresponding author

dedicated to exploring advanced MLLMs open-source to the community. Prior efforts like LLaMA-Adapter [18,65], LLaVA [31,35,36], and MiniGPT-4 [11,70] leverage a frozen CLIP [48] model for image encoding, and inject the visual cues into LLaMA [52] for multi-modal instruction tuning. The subsequent mPLUG-Owl [58,59], Qwen-VL [3], InternLM-XComposer [15], and SPHINX [19,34] further push the frontier of MLLMs in understanding and generalizing across visual contexts. MAVIS [68] introduces the first mathematical visual instructiontuning paradigm with two newly curated datasets, MAVIS-Caption and MAVIS-Instruct, demonstrating improved reasoning skills. Despite comprehensive benchmarks [16,21,30,39,55,60] on general visual instruction-following scenarios, the specific potential of MLLMs for visual mathematical problem-solving remains under-explored. In this paper, we introduce the MATHVERSE benchmark to comprehensively evaluate the visual mathematical reasoning skills of MLLMs, providing unique perspectives for future research directions.

Mathematical Reasoning Benchmarks have emerged as a significant area of focus, posing considerable challenges for large foundational models, e.g., LLMs and MLLMs. Initially, datasets in this realm are designed to address basic algebraic [25] and arithmetic [49] word problems, which are relatively limited in scope and volume. Subsequent efforts, including MATH [25], GSM8K [13], and MMLU [24], expand the range and quality of textual mathematical problems. These datasets feature a broader spectrum of difficulties, establishing a robust benchmark for the evaluation of general and math-specific LLMs [17, 43, 54, 62, 69]. Besides the text-only assessment, there is a growing demand for comparable, high-quality benchmarks for evaluating mathematical problem-solving in visual contexts, with the rapid progress of MLLMs. There are prior attempts, such as GeoQA [9], UniGeo [7], and Geometry3K [41], which focused exclusively on geometric problems. The recently proposed MathVista [40] broadens the scope to incorporate a variety of multi-modal tasks involving mathematical reasoning, and MMMU [61] covers college-level questions demanding intricate, domain-specific knowledge. However, our analysis identifies three main shortcomings within the current visual math benchmarks, as elaborated in Section 1 of the main paper. Therefore, we propose MATHVERSE specialized in the multi-modal mathematical evaluation of MLLMs, comprising twelve subjects, six problem versions, and 20K test samples. Our objective is to thoroughly investigate whether and how much MLLMs genuinely interpret visual diagrams for mathematical reasoning.

# **B** Additional Experimental Details

More Implementation Details. We conduct all experiments on NVIDIA A100 GPUs. For different MLLMs, we select their latest models and best-performing configurations for evaluation. Table 3 presents the source of the models used in MATHVERSE. As the text-only LLMs can only take text questions as input, we evaluate them with the first three problem versions, i.e., Text Dominant, Text Lite, and Text Only. For the 'w/o' results, we utilize the template in

Question	Prompt
Free-form Question	Please first conduct reasoning, and then an- swer the question and provide the final value, e.g., 1, 2.5, 300, at the end. - Question: {question}
Multiple-choice Question	Please first conduct reasoning, and then an- swer the question and provide the correct op- tion letter, e.g., A, B, C, D, at the end. - Question: {question}

 

 Table 1: Input Prompt of MLLMs for Response Generation. We adopt two different prompts for the free-form and multiple-choice questions. Note that these prompts are used for five problem versions except for the Vision-only version.

Table 2: Input Prompt for Vision-only Problems. Especially for the Vision-only version without textual input, we add "According to the question shown in the image" at the beginning of the prompt, and remove the "Question:" at the end.

Question	Prompt
Free-form Question	According to the question shown in the image, please first conduct reasoning, and then answer the question and provide the final value, e.g., 1, 2.5, 300, at the end.
Multiple-choice Question	According to the question shown in the im- age, please first conduct reasoning, and then answer the question and provide the correct option letter, e.g., A, B, C, D, at the end.

MathVista [40] to prompt GPT-4 [45] for answer extraction, and directly score the final answer without the intermediate reasoning process.

**Prompt for Response Generation.** We adopt two types of prompts respectively for the free-form and multiple-choice questions, as shown in Table 1. We inspire the Chain-of-Thought (CoT) reasoning capabilities of MLLMs by using the phrase "first conduct reasoning". Especially for the Vision-only problem version in Table 2, we add "According to the question shown in the image" at the beginning to remind MLLMs to read the questions rendered within diagrams, where the textual input for MLLMs only contains the prompt itself.

**Prompt for the CoT Evaluation.** Our proposed CoT evaluation contains two steps, i.e., key-step extraction and multi-step scoring, which prompt GPT-4 [45] and GPT-4V [46], respectively. The input configuration is listed in Table 4. We utilize the text-only GPT-4 in the first step to extract multiple key steps within the model's unstructured output, without feeding the question information. In

Model	Release Time	Source
ChatGPT [47]	2022-11	https://platform.openai.com/
GPT-4 [45]	2023-03	https://platform.openai.com/
Qwen-VL-Plus [3]	2023-11	https://help.aliyun.com/ zh/dashscope/developer- reference/vl-plus-quick-start
Gemini-Pro [20]	2023-12	https://ai.google.dev/
Qwen-VL-Max [3]	2024-01	https://help.aliyun.com/ zh/dashscope/developer- reference/vl-plus-quick-start
GPT-4V [46]	2023-09	https://platform.openai.com/
LLaMA-Adapter V2 [18]	2023-04	https://github.com/OpenGVLab/ LLaMA – Adapter / tree / main / llama_adapter_v2_multimodal7b
LLaVA-1.5 [35]	2023-10	https : / / huggingface . co / liuhaotian/llava-v1.5-13b
MiniGPT-v2 [11]	2023-10	https://github.com/Vision- CAIR/MiniGPT-4
mPLUG-Owl2 [59]	2023-11	https : / / huggingface . co / MAGAer13/mplug-ow12-11ama2-7b
G-LLaVA [17]	2023-12	https://github.com/pipilurj/ G-LLaVA/tree/main
ImageBind-LLM [23]	2023-05	https://github.com/OpenGVLab/ LLaMA - Adapter / tree / main / imagebind_LLM
ShareGPT4V [12]	2023-11	https://huggingface.co/Lin- Chen/ShareGPT4V-13B
SPHINX-Plus [34]	2023-11	https://huggingface.co/Alpha- VLLM/LLaMA2-Accessory/tree/ main / finetune / mm / SPHINX / SPHINX-v2-1k
LLaVA-NeXT [36]	2024-01	https : / / huggingface . co / liuhaotian/llava-v1.6-vicuna- 13b
SPHINX-MoE [19]	2024-01	https://huggingface.co/Alpha- VLLM/LLaMA2-Accessory/tree/ main / finetune / mm / SPHINX / SPHINX-MoE
InternLM-XComposer2 [15]	2024-01	<pre>https : / / huggingface . co / internlm/internlm-xcomposer2- v1-7b</pre>

Table 3: The Source of the Models Used in MATHVERSE.

#### MathVerse 5



Fig. 1: Manual Annotations for Function Problems in MATHVERSE. We provide detailed annotations, e.g., function expression and properties, for the diagrams of 534 function problems, which benefits the accuracy of GPT-4V [46] for CoT evaluation.

the second step, we input the extracted key-step reasoning and all the available content related to the problem into GPT-4V, allowing for a holistic assessment, including diagram interpretation, logical reasoning, and numerical computation. In Figure 1, we showcase the manual annotation for critical information within functional diagrams, e.g., function expression and properties. This assists GPT-4V in evaluating the visual perception accuracy of MLLMs for function graphs.

Human Performance Assessment. We recruit ten qualified college students specifically for the evaluation of human performance on MATHVERSE. These individuals are kept separate from the data curation stage, eliminating the possibility of them encountering the solutions beforehand. We allocate to each student the questions from a specific problem version. This strategy is to prevent them from gaining additional information from another version to answer questions, e.g., leveraging the textual *Implicit Property* from the Text-lite version to solve Text-intensive problems. They are asked to directly provide the final answer

**Table 4: Configuration for the CoT Evaluation Strategy.** We conduct two evaluation phases, respectively by prompting the text-only GPT-4 [45] and GPT-4V [46]. The symbol 'XXX' denotes the given one-shot sample, abbreviated for brevity. The 'Annotation' in the second phase is only required for function problems.

Phase	Input	Prompt
Key-step Extraction (GPT-4)	Model Output	I will give you a detailed solving procedure or a single answer for a math problem. If it is a procedure, you need to extract the key solution steps and list them accordingly in markdown syntax. If it is just a single answer, output the answer directly. Here are examples: - Model output: XXX - Extracted: 1. XXX 2. XXX 3. XXX - Model output: 2.2 - Extracted: The single answer is 2.2 Here is what you need to extract: - Model output: {model output} - Extracted:
Multi-step Scoring (GPT-4V)	Extracted Steps Question Diagram Answer Annotation	I will first give you a visual math problem, including the question, diagram, ground- truth answer, and detailed annotation of the diagram, and then give you a model out- put containing multiple key solution steps. Please think step by step and output the Average score, along with the Final answer score in the end, as described below: - Average score: Evaluate, based on the given question, answer, diagram, and di- agram annotation, whether each solution step is correct in logical reasoning, vi- sual perception, and numerical computa- tion, with an incorrect score of 0 and a cor- rect score of 1. Then, calculate the average score of multiple steps. - Final answer score: Match the model's fi- nal answer with the ground truth answer, scoring 1 if it matches and 0 if it doesn't. - If the model output only includes a single step or answer, the Average score and Final answer score are the same. - Question: {question} - Ground-truth answer: {answer} - Diagram annotation: {annotation} - Model output: {extracted steps} - Average score: - Final answer score:

without detailed reasoning. Therefore, we do not report the CoT evaluation results for human performance, alongside the 'Random Chance' baseline.

# C Additional Experimental Analysis

**Closed-source MLLMs are Better-performed.** From the performance in both accuracy and CoT evaluation, we observe a consistently better performance achieved by closed-source MLLMs than open-source ones. Despite the gap with humans, GPT-4V attains the leading position among MLLMs, showcasing superior mathematical capabilities over problem versions and subjects. MAVIS-7B is the best-performing open-source MLLM, surpassing Gemini-Pro and Qwen-VL-Max, while still lagging behind GPT-4V, suggesting the improvement space.

LLMs Achieve Competitive Results to MLLMs. Utilizing solely question texts as input, two LLMs, i.e., GPT-4 and ChatGPT, attain superior accuracy to most MLLMs in Text Dominant and Lite versions. Even in the absence of redundant *Descriptive Information* within Text-lite problems, GPT-4 outperforms LLaVA-NeXT and MAVIS-7B models. These findings not only indicate the strong mathematical reasoning skills of LLMs, but further emphasize the deficiencies in diagram interpretation of existing MLLMs. Importantly, the performance of GPT-4 is only exceeded by GPT-4V, which demonstrates that a satisfactory diagram perception capability can enhance the problem-solving for visual mathematics.

Mathematical Training Benefits the Performance. In addition to foundational visual instruction-following datasets, both InternLM-XComposer2 and LLaVA-NeXT extend their training regimes to include specialized mathematical problems that are either text-only or visual, such as MathQA [2], Geometry3K [41], and MathInstruct [62]. This approach of math-specific tuning contributes to their leading performance in MATHVERSE. Furthermore, G-LLaVA fine-tunes LLaVA-1.5 by a large-scale visual geometric dataset containing 170K enriched problems. This targeted refinement can improve several fields ('Len', 'Angle', and 'Apply') within the plane geometry subject. However, since G-LLaVA's fine-tuning data does not include problems of analytic geometry, solid geometry, and functions, it harms the related results of LLaVA-1.5 due to catastrophic forgetting. Moreover, due to a three-stage training strategy and largescale multi-modal mathematical data, MAVIS-7B exhibits superior performance over other MLLMs. This phenomenon underscores the critical role of developing extensive, high-quality visual math datasets for the effective training of MLLMs.

**Discrepancy Between CoT and Accuracy Evaluation Scores.** As illustrated by the two tables in the main paper, the CoT scores for MLLMs, in most cases, are much higher than accuracy scores. This observation demonstrates that our proposed CoT evaluation strategy identifies numerous correct intermediate reasoning steps, despite the final incorrect answer, highlighting the effectiveness



Fig. 2: Results with and without CoT Evaluation in MATHVERSE. Referring to Table ??, we denote the 'w/o' results in blue pillars, and highlight the increase and decrease magnitude with 'CoT-E' by green and red, respectively.

of fine-grained assessment. In Figure 2, we present the statistics of variance between the CoT and accuracy scores within different MLLMs. Although GPT-4V attains top-tier performance, it exhibits a pronounced gap concerning the evaluation of CoT reasoning quality. Conversely, SPHINX-MoE showcases favorable precision among open-source MLLMs, while preserving a relatively lower variance of two evaluation methods compared to InternLM-XComposer. This indicates its consistent step-by-step reasoning throughout the solving process.

# D More Dataset Details

# D.1 Data Curation

This paper engages twelve expert annotators for data curation, consisting of senior undergraduate and graduate students from across the globe with a strong background in science. In collaboration with the authors, they are required to mainly complete five tasks concerning data collection, categorization, quality review, problem version transformation, and function diagram annotation.

**Data Collection.** We comprehensively collect visual math problems from existing datasets [10,41,50] and public question repositories<sup>1,2,3</sup>. We specifically select high-quality plane geometric problems from current benchmarks, which showcase various question types, moderate question length, diverse diagram styles, and appropriate solving difficulty. For the manually collected problems of three subjects (plane geometry, solid geometry, and functions), we apply the Mathpix tool<sup>4</sup> to accurately extract the question texts, diagrams, explanations, and answers from the website. We strictly comply with copyright and licensing rules,

<sup>&</sup>lt;sup>1</sup>https://homework.study.com

<sup>&</sup>lt;sup>2</sup>https://www.ixl.com/math

<sup>&</sup>lt;sup>3</sup>https://mathspace.co/us

<sup>&</sup>lt;sup>4</sup>https://mathpix.com



Fig. 3: Manual Annotations for *Descriptive Information* in MATHVERSE. For some collected problems, we are required to supplement additional *Descriptive Information* (highlighted in red) to distinguish the Text-dominant version.

ensuring that we refrain from using data from sites that forbid copying and redistribution. After the initial collection, we obtain around 3.5K visual math problems, with 1.5K from existing datasets and 2K newly collected.

**Data Categorization and Review.** We first ask the human annotators to categorize the problems into three primary subjects, i.e., plane geometry, solid geometry, and functions. Within each subject, according to the definitions in Section D.2, the math problems are further divided into twelve fine-grained categories. At the same time, we meticulously review the collected dataset. We manually rectify the problems with incorrect answers and discard the problems with multiple diagrams, visual solutions, and too much similar content to others. Finally, 2,612 high-quality math problems with paired diagrams are preserved for MATHVERSE, spanning diverse subjects and subfields.

**Transformation of Problem Versions.** Given the three types of textual information within questions, human annotators rigorously transform each problem into six different versions as discussed in Section 2.2 of the main paper. We utilize Microsoft PowerPoint to annotate the diagrams in the Vision-dominant version, and employ Matplotlib to render the questions onto the diagrams in the Vision-only version. As illustrated in Figure 3, for problems with minimal *Descriptive Information*, we manually enhance the question text with additional contextual description about the diagram to differentiate the Text-dominant version.



Fig. 4: Manual Annotations for *Essential Condition* in MATHVERSE. For the original problems shown, we transfer some of the *Essential Condition* from diagrams to question texts (highlighted in green) to mark the Vision-dominant version.

sion. In the case of questions in Figure 4, where the *Essential Condition* has been fully depicted in the diagrams, we remove some of this content from the diagram and incorporate it into the text to mark the Vision-dominant version.

# D.2 Subject and Subfield Definition

The visual math problems within MATHVERSE encompass three primary subjects, plane geometry, solid geometry, and functions, alongside twelve finergrained subfields, which comprehensively evaluate the diagram understanding and mathematical reasoning capabilities of MLLMs.

**Plane Geometry** is a fundamental area that explores the properties and relations of points, lines, and surfaces in a two-dimensional plane. This subject delves into concepts such as angles, triangles, circles, and polygons, offering a rich context for assessing the spatial comprehension and logical deduction skills of MLLMs. We divide it into five subfields, as exemplified in Figure 5:

- Length focuses on the measurement and comparison of distances between points. This subfield includes understanding the properties of lines, segments, and their use in determining the perimeters of geometric shapes, which is foundational for MLLMs to solve plane geometry problems.

### A Plane Geometry:



Fig. 5: Examples of Five Subfields in Plane Geometry, spanning Length, Area, Angle, Analytic, and Applied Geometry problems. We showcase the Text-lite version.

- Area examines the size of two-dimensional surfaces. It encompasses calculating the areas of various shapes, such as triangles, rectangles, circles, and more complex polygons, by applying specific formulas and principles, which is crucial for comprehending the concept of space within geometry.
- Angle involves the study of angles and their properties, including different types of angles (acute, right, and obtuse), angle measurement, and the relationships between angles, particularly in polygons. This subfield demands the advanced spatial perception capacity of MLLMs.
- Analytic Geometry, also known as coordinate geometry, merges algebra and geometry to solve geometric problems using coordinate systems, exploring the calculation and reasoning of equations for geometric shapes. MLLMs are evaluated on their coordinate identification and algebraic capabilities.
- Applied Geometry relate to the application of geometric principles to solve real-world and theoretical problems. It challenges MLLMs to first understand the background information within questions, and apply their knowledge of lengths, areas, angles, and analytic geometry for problem-solving.

**Solid Geometry** focuses on the study of three-dimensional objects that have depth, length, and width, thereby offering a more complex and enriched exploration of spatial structures. This subject investigates a variety of shapes such as cubes, cylinders, spheres, and pyramids, and assesses MLLMs to tackle questions concerning the volume, surface area, and geometric properties of these solids. This subject contains three subfields, as exemplified in Figure 6:

- Length, extending from the 2D counterpart, focuses on measuring the edges and curves that define three-dimensional objects. It involves determining the linear distance between points in space, the perimeters of bases of solids, and the height or depth of objects. This measurement is a foundational element for MLLMs in analyzing geometric solids.



Fig. 6: Examples of Three Subfields in Solid Geometry, spanning Length, Area, and Volume problems. We showcase the Text-lite version.

- Area encompasses the calculation of the total area covered by the outer surfaces of solids. This normally requires MLLMs to break down complex shapes into several simpler components for area calculation in plane geometry, assessing their spatial and logical reasoning performance.
- **Volume** pertains to measuring the space enclosed within three-dimensional objects. This demands MLLMs to precisely identify the geometric solids and apply accurate formulas to calculate the volume, which evaluates their mathematical knowledge application and calculation skills.

*Functions* involve analyzing mathematical functions to understand the relationship between variables. These challenges range from simple tasks, like calculating a function value for a given input, to more complex scenarios, such as exploring the behavior and representation of various function types. We evaluate MLLMs by four types of function problems, exemplified in Figure 7:

- Function Coordinate focuses on interpreting and extracting coordinatelevel information from graphical representations of functions. It includes tasks such as identifying specific coordinate values of points on the graph and observing intersection points between functions and axes, which test the MLLM's basic proficiency in functional visual perception.
- Function Property emphasizes the model's capacity to discern and deduce the inherent properties of functions from their graphs, such as symmetry, asymptotes, extrema (maximum and minimum points), and intervals of increase or decrease. These problems can reveal the understanding of MLLMs for the deeper characteristics of functions.
- **Function Expression** refers to the direct analysis using the algebraic expressions of functions, widely including linear, quadratic, polynomial, exponential, logarithmic, and piece-wise functions. It challenges MLLMs to extract specific function expressions and apply transformations, bridging the gap between abstract mathematical reasoning and visual interpretation.



Fig. 7: Examples of Four Subfields in Functions, spanning Function Coordinate, Property, Expression, and Applied problems. We showcase the Text-lite version.

- Applied Function, similar to the applied geometry problems, requires MLLMs to leverage their functional knowledge and theorems in real-world scenarios, e.g., modeling economic data, predicting physical phenomena, and calculating probabilities. This assesses the MLLM's capabilities to understand functions in both theoretical contexts and practical situations.

### D.3 Detailed Statistics of MATHVERSE

More Data Statistics. In Table 5, we provide a more detailed data statistics of MATHVERSE. Therein, the 534 newly annotated questions refer to all the function problems, for which we meticulously annotate critical functional information, as depicted in Figure 1. The number of newly annotated diagrams represents the 5,224 math problems in the Vision-dominant and Vision-only versions. For these problems, we respectively integrate the *Essential Condition* and all textual content with the diagrams. We also list the numbers of multiple-choice answers, where A, B, C, and D are almost uniformly distributed.

**Problem Length Variance.** In Table 6, we highlight the variance in question and answer lengths across the five problem versions in MATHVERSE, excluding the Vision-only category due to its absence of text. For both word and character levels, as we remove the pre-defined textual elements (*Descriptive Information*, *Implicit Property*, and *Essential Condition*), the maximum and average lengths of questions decrease accordingly, while the answer lengths remain the same. In Figure 8, we visualize the word-level variation of question length for the three problem versions: Text Dominant, Text Lite, and Vision Dominant. By progressively omitting *Descriptive Information* and *Essential Condition* from the Text-dominant version, we observe a clear downward trajectory for the question length distribution and average values.

Table 5: Statistics of MATHVERSE.

Statistic	Number
Total questions	2,612
- Subjects/subfields	3/12
- Multiple-choice questions	1,631 (62.4%)
- Free-form questions	981 (37.6%)
- Newly collected questions	1,236 (47.3%)
- Existing-dataset questions	1,376(52.7%)
- Questions with explanations	1,236 (47.3%)
- Newly annotated questions	534(20.4%)
Multiple-choice question	
- Proportion of answer A	585 (22.4%)
- Proportion of answer B	828 (31.7%)
- Proportion of answer C	703 (26.9%)
- Proportion of answer D	444 (17.0%)
- Proportion of answer E&F	52 (2.0%)
Total test samples	15,672
- Newly annotated samples	10,448 (66.7%)
- Newly annotated diagrams	5,224 ( $33.3%$ )
- Samples of each version	2,612 (16.7%)
Number of unique images	2,420 (92.6%)
Number of unique questions	2,573 (98.5%)
Number of unique answers	847 (32.4%)

Table 6: Length of Different Problem Versions in MATHVERSE.

Problem Version	Word	Character
Text Dominant & Text Only		
- Maximum question length	203	1,311
- Maximum answer length	17	102
- Average question length	35.7	204.8
- Average answer length	1.4	6.3
Text Lite		
- Maximum question length	179	1,173
- Maximum answer length	17	102
- Average question length	22	133.8
- Average answer length	1.4	6.3
Vision Intensive		
- Maximum question length	171	1,126
- Maximum answer length	17	102
- Average question length	18.8	116.8
- Average answer length	1.4	6.3
Vision Dominant		
- Maximum question length	176	1,132
- Maximum answer length	17	102
- Average question length	17.6	123.5
- Average answer length	1.4	6.3

# E Comparison to Current Benchmarks

In this section, we offer a detailed comparison between MATHVERSE and existing multi-modal mathematical benchmarks, i.e., geometry-specific benchmarks [5,8, 10,42,50], MathVista [40], and MMMU [61], from the following four aspects:

The Investigation of Diagram Interpretation Capacity. As discussed in Figure 1 of the main paper, the math problems in most existing datasets contain excessive redundant information in textual content, which is repetitive to the visual elements in diagrams. This issue enables MLLMs to potentially bypass the process of visual understanding, and thereby cannot determine whether and how much MLLMs truly interpret the math diagram. In contrast, our MATHVERSE includes six problem versions with different information content across text and vision. By comparing the performance variance between different problem versions, we can thoroughly investigate the mathematical diagram interpretation capabilities of MLLMs for the first time.

**Evaluation Approach.** Previous benchmarks adopt a simple True or False metric to score the response from MLLMs, which lacks fine-grained information and intermediate reasoning assessment, as analyzed in Figure 2 of the main paper. In contrast, MATHVERSE adopts a unique CoT evaluation strategy by examining each crucial solution step within the model output. This approach not only unveils the CoT reasoning quality of MLLMs, but also provides detailed error analysis, serving as valuable guidance for future enhancement.



Fig. 8: Distribution of Question Length for Three Problem Versions. We exclude the *Descriptive Information* and *Essential Condition* from the Text-dominant problems, respectively creating the Text-lite and Vision-dominant versions.

The Depth and Width in Math Problems. The geometry-specific benchmarks evaluate only a limited dimension of mathematical skills in MLLMs. Math-Vista instead incorporates a variety of math-related question-answering tasks, e.g., textbook figures, tables, plots, charts, puzzles, and synthetic scenes, as exemplified in Figure 9. However, the integration of these peripheral tasks (covering more than 70%) might divert the focus from the specialized mathematical evaluation of MLLMs. In addition, MMMU focuses on college-level complexity, requiring advanced domain-specific knowledge, as depicted in Figure 10. Given this, the lack of profound mathematical theorems would restrict the performance of MLLMs, biasing the evaluation of logical reasoning and visual perception proficiency. Therefore, our MATHVERSE concentrates on specialized visual math problems (plane geometry, solid geometry, and functions) with a moderate difficulty (high-school level), aiming to fully exert the capabilities of MLLMs.

**Total Volume of Test Samples.** We summarize the size of test instances for different datasets in Table 7. As demonstrated, our MATHVERSE offers a considerably larger number of samples than others, nearly three times to MathVista and twenty times to GeoQA+, including meticulously annotated six versions of visual math problems. This contributes to a comprehensive and robust evaluation of visual mathematical reasoning capabilities.

Table 7: Number of Test Samples in Different Benchmarks.

Benchmark	GEOS	Geo3K	$\mathrm{GeoQA}+$	MathVista	MMMU- Math	MATHVERSE
Test Samples	119	601	755	6,141	540	$15,\!672$

### 🚺 Table QA:

Items sold last week		Making leaf rubbings		Rounds in the spelling bee		Athletes per country		neanut hutter cookie dough	\$3 per lh
ltam	Erencency	Leaf rubbings made	Frequency	Voor	Pounds	Stem	Leaf		
item	riequoney	0	14	reta	rtounus	1	122	double chocolate cookie dough	\$3 per lb
TV	16	1	5	2008	17	1	133	chocolate chin cookie dough	\$5 per lb
tablet computer	7	2	15	2009	10	2	0125799	chocolate chip cookie dough	op per in
sneaker	32	3	11	2010	9	3	4578	oatmeal raisin cookie dough	\$7 per lb
operater	02	4	3	2010		4	1	dingerspan cookie dough	\$3 per lh
cell phone	23	5	10	2011	15	5	2455668	gingerondy coonie dough	to per in
video game console	19	6	17	2012	2 7 6 16		16	snickerdoodle cookie dough	\$8 per lb

## C Textbook and Science QA:



### C Plot and Chart QA:



IQ Test and Synthetic QA:



C General and Icon QA:



Fig. 9: Diagram Examples of Math-related Tasks in MathVista [40]. These tasks are not strongly correlated to the mathematical reasoning skills of MLLMs, probably skewing the assessment emphasis towards visual math problems.



Fig. 10: Diagram Examples with Required Knowledge in MMMU [61]. These math problems demand MLLMs to comprehend college-level domain knowledge, potentially hindering them from fully exerting mathematical reasoning skills.

# F Limitation and Future Work

While our MATHVERSE takes a step forward in the field of visual mathematical evaluation for MLLMs, it is important to recognize several limitations as follows.

We have categorized the math problems in MATHVERSE by various criteria, including subjects, subfields, and versions featuring differing degrees of multi-modal content. These categorization approaches evaluate the capabilities of MLLMs from multiple dimensions. Nevertheless, it is also meaningful to further divide the problems based on their difficulty levels, akin to MATH [25], a text-only benchmark defining five levels of difficulty. This additional layer of differentiation can provide deeper insights into the problem-solving abilities of MLLMs across a spectrum of challenges, which we leave as future work.

The curated dataset in MATHVERSE focuses on math problems in the high school level with moderate difficulty, which aims to fully demonstrate the mathematical reasoning skills within current MLLMs. However, with the advancement of architecture and training methodologies, future MLLMs have the potential to grasp more complex knowledge and theorems across a variety of domains. Therefore, there is significant value in further augmenting MATHVERSE with problems spanning broader complexity and disciplines, including those at the college level and within scientific fields. By transforming the expanded problems into different versions, we can facilitate a more comprehensive and robust evaluation of MLLMs for their diagram interpretation and reasoning capabilities.

Moreover, the problems in MATHVERSE and other current mathematical benchmarks are mainly in English. Given that some multilingual MLLMs [3,23]



Fig. 11: Distribution of GPT-4V's [46] Errors in Reasoning and Answers. For the six problem versions in MATHVERSE, we provide the statistics of errors made by GPT-4V based on their occurrence in answers ('Ans.') and reasoning processes ('Rea.').

have been developed, existing evaluation cannot reveal their full capabilities when confined to a single language. The incorporation of multilingual visual math problems would not only extend the dataset's global applicability, but also enhance the assessment of MLLMs for linguistic diversity and understanding.

# G Error Analysis

To delve into the fine-grained predictions, we select the best-performing MLLM, GPT-4V [46], to understand its modes of success and failure. Our proposed CoT evaluation strategy has produced a detailed assessment of model output, including step-wise scores and explanation, reducing extensive manual effort in identifying and analyzing errors. We conduct our analysis on the two-step output from the CoT evaluation across the entire dataset, focusing on two key dimensions.

*Errors in Reasoning or Answer?* In Figure 11, we showcase the statistics of different error distributions in six problem versions of MATHVERSE. We define the following six error categories: correct final answer with correct/partially correct/incorrect CoT reasoning and incorrect final answer with correct/partially



Fig. 12: Distribution of GPT-4V's [46] Errors within Different Types. We present the statistics of four error types by GPT-4V in the six problem versions, i.e., Visual Perception Error, Reasoning Error, Calculation Error, and Knowledge Error.

correct/incorrect CoT reasoning. For all six versions, the incorrect final answers are mostly caused by the partially incorrect reasoning process. In addition, a number of problems with correct answers are accompanied by partially or entirely incorrect reasoning, e.g., 15.3% in Text Dominant, which cannot be detected by the traditional True or False evaluation. As we remove the content within textual questions and enrich the visual diagram, e.g., from Text Dominant and Lite to Vision Dominant and Only, we observe a progressive increase in the error rate of 'incorrect final answer with incorrect CoT reasoning', indicating that MLLMs are challenged to conduct high-quality intermediate reasoning by capturing more information from the visual input.

What Types of Errors? To further investigate the specific error types, we survey the problems with errors that occur either within the reasoning process or the final answer. As depicted in Figure 12, we divide the errors of GPT-4V into four distinct types: visual perception error, reasoning error, knowledge error, and calculation error. Consistent with our findings in the main paper, the primary source of errors in problem-solving attributes to the inaccurate interpretation of mathematical diagrams, which significantly impedes the performance of MLLMs. For the problem versions that demand advanced diagram interpretation, e.g., Vision Dominant and Only, we observe a notable increase in the rate of visual perception errors, demonstrating an urgent need for stronger visual encoders in MLLMs. Moreover, reasoning errors also account for a considerable percentage, indicating that the logical deduction skills of MLLMs still

require improvement. As expected, knowledge errors do not significantly hinder the mathematical reasoning capabilities of MLLMs in MATHVERSE.

# **H** Qualitative Examples

To ease the understanding, we offer a variety of qualitative examples in MATH-VERSE. In Section H.1, we showcase the meticulously transformed six versions of visual math problems. In Section H.2, we compare the response of different MLLMs on Text-lite problems, including GPT-4V [46], LLaVA-NeXT [36], and SPHINX-MoE [19]. Specifically, we present the key-step extraction output by the CoT evaluation, and mark the multi-step scoring results aside. In Section H.3, we provide the response comparison of GPT-4V for three problem versions in MATHVERSE, i.e., Text Dominant, Text Lite, and Vision Dominant.

# H.1 Comparison of Six Problem Versions

Please refer to Figures  $13 \sim 15$ .

# H.2 Response of Different MLLMs

Please refer to Figures  $16 \sim 21$ .

# H.3 Response of Different Problem Versions

Please refer to Figures  $22 \sim 27$ .



Fig. 13: Comparison of Six Problem Versions in MATHVERSE.

	Descriptive Information	Implicit Property	Essential Condition	
	🕮 Text Input	🔍 Vision Input	🛄 Text Input	🔍 Vision Input
Text Dominant	A soft drink can has a height of 13 cm and a radius of 3 cm. Find L, the length of the longest straw that can fit into the can (so that the straw is not bent and fits entirely inside the can).		A square pyramid is shown left with the following dimensions. The height of this square pyramid is 5x+7. Length L is the side of the base of the pyramid. Write down an expression for L, in terms of the variable x.	
Text Lite	The radius is 3 cm. Find L, the length of the longest straw that can fit into the can.		A square pyramid is shown left with the following dimensions. The height of this square pyramid is 5x+7. Write down an expression for L, in terms of the variable x.	
Text Only	A soft drink can has a height of 13 cm and a radius of 3 cm. Find L, the length of the longest straw that can fit into the can (so that the straw is not bent and fits entirely inside the can).		A square pyramid is shown left with the following dimensions. The height of this square pyramid is 5x+7. Length L is the side of the base of the pyramid. Write down an expression for L, in terms of the variable x.	
Vision Intensive	The radius is 3 cm. Find L, the length of the longest straw that can fit into the can.	13 13	The height of this solid is $5x+7$ . Write down an expression for L, in terms of the variable x.	
Vision Dominant	Find L, the length of the longest straw that can fit into the can.		A square pyramid is shown left with the following dimensions. Find an expression for L, in terms of the variable x.	
Vision Only				Approximate of the solution of

Fig. 14: Comparison of Six Problem Versions in MATHVERSE.



Fig. 15: Comparison of Six Problem Versions in MATHVERSE.



Fig. 16: Response Comparison of GPT-4V [46], LLaVA-NeXT [36], and SPHINX-MoE [19]. We adopt the Text-lite version of the problem, and highlight the key-step extraction and scoring by the CoT evaluation strategy.



Fig. 17: Response Comparison of GPT-4V [46], LLaVA-NeXT [36], and SPHINX-MoE [19]. We adopt the Text-lite version of the problem, and highlight the key-step extraction and scoring by the CoT evaluation strategy.



Fig. 18: Response Comparison of GPT-4V [46], LLaVA-NeXT [36], and SPHINX-MoE [19]. We adopt the Text-lite version of the problem, and highlight the key-step extraction and scoring by the CoT evaluation strategy.



Fig. 19: Response Comparison of GPT-4V [46], LLaVA-NeXT [36], and SPHINX-MoE [19]. We adopt the Text-lite version of the problem, and highlight the key-step extraction and scoring by the CoT evaluation strategy.



Fig. 20: Response Comparison of GPT-4V [46], LLaVA-NeXT [36], and SPHINX-MoE [19]. We adopt the Text-lite version of the problem, and highlight the key-step extraction and scoring by the CoT evaluation strategy.



Fig. 21: Response Comparison of GPT-4V [46], LLaVA-NeXT [36], and SPHINX-MoE [19]. We adopt the Text-lite version of the problem, and highlight the key-step extraction and scoring by the CoT evaluation strategy.



Fig. 22: Response of Different Problem Versions by GPT-4V [46]. By comparing the output from three problem versions, we observe increasing reasoning errors when more textual content is visually depicted in diagrams. Note that we only showcase the Text-dominant version of the problem for briefness.



Fig. 23: Response of Different Problem Versions by GPT-4V [46]. By comparing the output from three problem versions, we observe increasing reasoning errors when more textual content is visually depicted in diagrams. Note that we only showcase the Text-dominant version of the problem for briefness.



Fig. 24: Response of Different Problem Versions by GPT-4V [46]. By comparing the output from three problem versions, we observe increasing reasoning errors when more textual content is visually depicted in diagrams. Note that we only showcase the Text-dominant version of the problem for briefness.

#### MathVerse 33



1	1/1	<ol> <li>Identify that OC is perpendicular to OB, making triangle OCB a right triangle with OC and OB as its legs. Since OC and OB are radii of the circle and the radius is given as 1.0, OC = OB = 1.0. Therefore, triangle OCB is an isosceles right triangle.</li> </ol>					
1	1/1	2. Identify that OD is perpendicular to AB, which means that OD is the altitude of triangle AOB from point O to side AB. Since triangle AOB is also isosceles ( $OA = OB = 1.0$ ), OD bisects AB. This means that $AD = DB$ .					
1	1/1	3. Consider triangle AOD and triangle BOD. Since OD is perpendicular to AB, both triangles AOD and BOD are right triangles. Also, since AD = DB, triangle AOD is congruent to triangle BOD. Therefore, AE = EB.					
×	0/1	4. Since AE = EB and AC = AB (because triangle ABC is inscribed in the circle and AC and AB are both radii), we can say that CE = EB.					
×	0/1	5. To find AE^2 + CE^2, since AE = CE, we can say that AE^2 + CE^2 = 2 * AE^2.					
×	0/1	6. To find AE, use the Pythagorean theorem in triangle AOD (or BOD, since they are congruent). We have:					
		AO^2 = AD^2 + OD^2.					
		Since AO = 1.0 (radius of the circle) and AD = $1/2$ * AB (because AD = DB and AB is the diameter of the circle, which is twice the radius), we have:					
		1.0^2 = (1/2 * AB)^2 + OD^2.					
		AB is the diameter of the circle, so AB = 2 * radius = 2 * 1.0 = 2.0. Therefore: 1.0 ^ 2 = (1/2 * 2.0)^2 + OD^2 OD^2 = 0 Since OD^2 = 0, OD must be 0, which means that point D coincides with point O, and AE = EB = 1.0 (since AE = EB and AB = 2 * radius = 2.0. Therefore, AE ^ 2 + C ^ 2 - 2 * AE ^ 2 + 2 * 10 ^ 2 = 2.					
- √	1/1	7. The correct answer is: B: 2					
		Multi-step Score: 0.57 Final CoT Evaluation Score: 0.7					
$\geq$							

Text Lite:



- 1/1 2. Since OC is perpendicular to OB and both are radii of the circle, triangle OCB is an isosceles right triangle. Therefore, the length of BC can be found using the Pythagorean theorem: BC = J2
- ✓ 1/1 3. Similarly, since OD is perpendicular to AB, triangle OAD is also an isosceles right triangle. Therefore, AD = OD = 1.
- × 0/1 4. Using the Pythagorean theorem again, AO = J2. Since AO is the diameter of the circle, AC is also the diameter, and thus AC = AO =
- × 0/1 5. We have AC = J2 and BC = J2. Since E is a point on AC, AE + EC = AC.
- X 0/1 6. We want to find AE^2 + CE^2. Let's denote AE as x and CE as (J2 x). Then we can express AE^2 + CE^2 as: 2x^2 2J2x + 2.
- X 0/1 7. To minimize this expression, we can take the derivative with respect to x and set it to zero to find the minimum value. However, since we are given that AE and CE are segments of a line, and we know that the sum of their squares in a right triangle is constant, we can infer that the sum of the squares of AE can CE will be equal to the squares of AC.
- ✓ 1/1 8. Therefore, AE<sup>2</sup> + CE<sup>2</sup> = AC<sup>2</sup> = (√2)<sup>2</sup> = 2.
- ✓ 1/1 9. The correct answer is: B: 2
  - Multi-step Score: 0.56 Final CoT Evaluation Score: 0.69

#### Vision Dominant:

- ✓ 1/1 1. Given that the radius of circle O is 1, this means that OA, OE, and OC are all radii of the circle, and thus they all have a length of 1. ✓ 1/1 2. Use the Pythagorean theorem (a^2 + b^2 = c^2) on triangle AOE, where AE is the hypotenuse, and OA and OE are the two legs: AE^2 = OA^2 + OE^2, AE^2 = 2.  $\checkmark$  1/1 3. By the same logic, apply the Pythagorean theorem to triangle COE, where CE is the hypotenuse: CE^2 = OC^2 + OE^2, CE^2 = 2.  $\times$  0/1 4. Now, according to the question, find the value of AE^2 + CE^2: AE^2 + CE^2 = 2 + 2 AE^2 + CE^2 = 4. X 0/1 5. Therefore, the value of AE<sup>2</sup> + CE<sup>2</sup> is 4. The correct option is D: 4.
  - Multi-step Score: 0.6 Final CoT Evaluation Score: 0.42

Fig. 25: Response of Different Problem Versions by GPT-4V [46]. By com-

paring the output from three problem versions, we observe increasing reasoning errors when more textual content is visually depicted in diagrams. Note that we only showcase the Text-dominant version of the problem for briefness.



Fig. 26: Response of Different Problem Versions by GPT-4V [46]. By comparing the output from three problem versions, we observe increasing reasoning errors when more textual content is visually depicted in diagrams. Note that we only showcase the Text-dominant version of the problem for briefness.



Fig. 27: Response of Different Problem Versions by GPT-4V [46]. By comparing the output from three problem versions, we observe increasing reasoning errors when more textual content is visually depicted in diagrams. Note that we only showcase the Text-dominant version of the problem for briefness.

# References

- Alayrac, J.B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., Lenc, K., Mensch, A., Millican, K., Reynolds, M., et al.: Flamingo: a visual language model for few-shot learning. Advances in Neural Information Processing Systems 35, 23716– 23736 (2022)
- Amini, A., Gabriel, S., Lin, P., Koncel-Kedziorski, R., Choi, Y., Hajishirzi, H.: Mathqa: Towards interpretable math word problem solving with operation-based formalisms. arXiv preprint arXiv:1905.13319 (2019)
- Bai, J., Bai, S., Yang, S., Wang, S., Tan, S., Wang, P., Lin, J., Zhou, C., Zhou, J.: Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. arXiv preprint arXiv:2308.12966 (2023)
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language models are few-shot learners. In: Advances in neural information processing systems. pp. 1877–1901 (2020)
- Cao, J., Xiao, J.: An augmented benchmark dataset for geometric question answering through dual parallel text encoding. In: Proceedings of the 29th International Conference on Computational Linguistics. pp. 1511–1520 (2022)
- Chen, G., Zheng, Y.D., Wang, J., Xu, J., Huang, Y., Pan, J., Wang, Y., Wang, Y., Qiao, Y., Lu, T., et al.: Videollm: Modeling video sequence with large language models. arXiv preprint arXiv:2305.13292 (2023)
- Chen, J., Li, T., Qin, J., Lu, P., Lin, L., Chen, C., Liang, X.: Unigeo: Unifying geometry logical reasoning via reformulating mathematical expression. ArXiv abs/2212.02746 (2022)
- Chen, J., Li, T., Qin, J., Lu, P., Lin, L., Chen, C., Liang, X.: Unigeo: Unifying geometry logical reasoning via reformulating mathematical expression. arXiv preprint arXiv:2212.02746 (2022)
- Chen, J., Tang, J., Qin, J., Liang, X., Liu, L., Xing, E.P., Lin, L.: Geoqa: A geometric question answering benchmark towards multimodal numerical reasoning. ArXiv abs/2105.14517 (2021), https://api.semanticscholar.org/CorpusID: 235253782
- Chen, J., Tang, J., Qin, J., Liang, X., Liu, L., Xing, E.P., Lin, L.: Geoqa: A geometric question answering benchmark towards multimodal numerical reasoning. arXiv preprint arXiv:2105.14517 (2021)
- Chen, J., Li, D.Z.X.S.X., Zhang, Z.L.P., Xiong, R.K.V.C.Y., Elhoseiny, M.: Minigpt-v2: Large language model as a unified interface for vision-language multitask learning. arXiv preprint arXiv:2310.09478 (2023)
- Chen, L., Li, J., wen Dong, X., Zhang, P., He, C., Wang, J., Zhao, F., Lin, D.: Sharegpt4v: Improving large multi-modal models with better captions. ArXiv abs/2311.12793 (2023), https://api.semanticscholar.org/CorpusID: 265308687
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., et al.: Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168 (2021)
- Dai, W., Li, J., Li, D., Tiong, A.M.H., Zhao, J., Wang, W., Li, B., Fung, P., Hoi, S.: Instructblip: Towards general-purpose vision-language models with instruction tuning (2023)
- Dong, X., Zhang, P., Zang, Y., Cao, Y., Wang, B., Ouyang, L., Wei, X., Zhang, S., Duan, H., Cao, M., et al.: Internlm-xcomposer2: Mastering free-form text-image

composition and comprehension in vision-language large model. arXiv preprint arXiv:2401.16420 (2024)

- Fu, C., Chen, P., Shen, Y., Qin, Y., Zhang, M., Lin, X., Yang, J., Zheng, X., Li, K., Sun, X., Wu, Y., Ji, R.: Mme: A comprehensive evaluation benchmark for multimodal large language models. arXiv preprint arXiv:2306.13394 (2023)
- Gao, J., Pi, R., Zhang, J., Ye, J., Zhong, W., Wang, Y., Hong, L., Han, J., Xu, H., Li, Z., et al.: G-llava: Solving geometric problem with multi-modal large language model. arXiv preprint arXiv:2312.11370 (2023)
- Gao, P., Han, J., Zhang, R., Lin, Z., Geng, S., Zhou, A., Zhang, W., Lu, P., He, C., Yue, X., Li, H., Qiao, Y.: Llama-adapter v2: Parameter-efficient visual instruction model. arXiv preprint arXiv:2304.15010 (2023)
- Gao, P., Zhang, R., Liu, C., Qiu, L., Huang, S., Lin, W., Zhao, S., Geng, S., Lin, Z., Jin, P., et al.: Sphinx-x: Scaling data and parameters for a family of multi-modal large language models. arXiv preprint arXiv:2402.05935 (2024)
- Gemini Team, G.: Gemini: a family of highly capable multimodal models. arXiv preprint arXiv:2312.11805 (2023)
- Guo, Z., Zhang, R., Chen, H., Gao, J., Gao, P., Li, H., Heng, P.A.: Sciverse. https://sciverse-cuhk.github.io (2024), https://sciverse-cuhk.github.io/
- 22. Guo, Z., Zhang, R., Zhu, X., Tang, Y., Ma, X., Han, J., Chen, K., Gao, P., Li, X., Li, H., et al.: Point-bind & point-llm: Aligning point cloud with multi-modality for 3d understanding, generation, and instruction following. arXiv preprint arXiv:2309.00615 (2023)
- Han, J., Zhang, R., Shao, W., Gao, P., Xu, P., Xiao, H., Zhang, K., Liu, C., Wen, S., Guo, Z., et al.: Imagebind-llm: Multi-modality instruction tuning. arXiv preprint arXiv:2309.03905 (2023)
- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., Steinhardt, J.: Measuring massive multitask language understanding. Proceedings of the International Conference on Learning Representations (ICLR) (2021)
- Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart, S., Tang, E., Song, D., Steinhardt, J.: Measuring mathematical problem solving with the math dataset. NeurIPS (2021)
- Hong, Y., Zhen, H., Chen, P., Zheng, S., Du, Y., Chen, Z., Gan, C.: 3d-llm: Injecting the 3d world into large language models. Advances in Neural Information Processing Systems 36 (2024)
- 27. Jiang, A.Q., Sablayrolles, A., Roux, A., Mensch, A., Savary, B., Bamford, C., Chaplot, D.S., de Las Casas, D., Hanna, E.B., Bressand, F., Lengyel, G., Bour, G., Lample, G., Lavaud, L.R., Saulnier, L., Lachaux, M., Stock, P., Subramanian, S., Yang, S., Antoniak, S., Scao, T.L., Gervet, T., Lavril, T., Wang, T., Lacroix, T., Sayed, W.E.: Mixtral of experts. Arxiv 2401.04088 (2024)
- Jiang, D., Song, G., Wu, X., Zhang, R., Shen, D., Zong, Z., Liu, Y., Li, H.: Comat: Aligning text-to-image diffusion model with image-to-text concept matching. arXiv preprint arXiv:2404.03653 (2024)
- Li, B., Zhang, Y., Chen, L., Wang, J., Pu, F., Yang, J., Li, C., Liu, Z.: Mimic-it: Multi-modal in-context instruction tuning. arXiv preprint arXiv:2306.05425 (2023)
- Li, B., Wang, R., Wang, G., Ge, Y., Ge, Y., Shan, Y.: Seed-bench: Benchmarking multimodal llms with generative comprehension. ArXiv abs/2307.16125 (2023)
- Li, F., Zhang, R., Zhang, H., Zhang, Y., Li, B., Li, W., Ma, Z., Li, C.: Llava-nextinterleave: Tackling multi-image, video, and 3d in large multimodal models. arXiv preprint arXiv:2407.07895 (2024)

- 38 R. Zhang and D. Jiang et al.
- Li, J., Li, D., Xiong, C., Hoi, S.: Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In: International Conference on Machine Learning. pp. 12888–12900. PMLR (2022)
- 33. Li, X., Zhang, M., Geng, Y., Geng, H., Long, Y., Shen, Y., Zhang, R., Liu, J., Dong, H.: Manipllm: Embodied multimodal large language model for object-centric robotic manipulation. CVPR 2024 (2023)
- 34. Lin, Z., Liu, C., Zhang, R., Gao, P., Qiu, L., Xiao, H., Qiu, H., Lin, C., Shao, W., Chen, K., et al.: Sphinx: The joint mixing of weights, tasks, and visual embeddings for multi-modal large language models. arXiv preprint arXiv:2311.07575 (2023)
- Liu, H., Li, C., Li, Y., Lee, Y.J.: Improved baselines with visual instruction tuning (2023)
- Liu, H., Li, C., Li, Y., Li, B., Zhang, Y., Shen, S., Lee, Y.J.: Llava-next: Improved reasoning, ocr, and world knowledge (January 2024), https://llava-vl.github. io/blog/2024-01-30-llava-next/
- 37. Liu, J., Li, C., Wang, G., Lee, L., Zhou, K., Chen, S., Xiong, C., Ge, J., Zhang, R., Zhang, S.: Self-corrected multimodal large language model for end-to-end robot manipulation. arXiv preprint arXiv:2405.17418 (2024)
- Liu, J., Liu, M., Wang, Z., Lee, L., Zhou, K., An, P., Yang, S., Zhang, R., Guo, Y., Zhang, S.: Robomamba: Multimodal state space model for efficient robot reasoning and manipulation. arXiv preprint arXiv:2406.04339 (2024)
- 39. Liu, Y., Duan, H., Zhang, Y., Li, B., Zhang, S., Zhao, W., Yuan, Y., Wang, J., He, C., Liu, Z., et al.: Mmbench: Is your multi-modal model an all-around player? arXiv preprint arXiv:2307.06281 (2023)
- Lu, P., Bansal, H., Xia, T., Liu, J., yue Li, C., Hajishirzi, H., Cheng, H., Chang, K.W., Galley, M., Gao, J.: Mathvista: Evaluating math reasoning in visual contexts with gpt-4v, bard, and other large multimodal models. ArXiv abs/2310.02255 (2023)
- 41. Lu, P., Gong, R., Jiang, S., Qiu, L., Huang, S., Liang, X., Zhu, S.C.: Inter-gps: Interpretable geometry problem solving with formal language and symbolic reasoning. In: Annual Meeting of the Association for Computational Linguistics (2021), https://api.semanticscholar.org/CorpusID:234337054
- 42. Lu, P., Gong, R., Jiang, S., Qiu, L., Huang, S., Liang, X., Zhu, S.C.: Inter-gps: Interpretable geometry problem solving with formal language and symbolic reasoning. arXiv preprint arXiv:2105.04165 (2021)
- 43. Luo, H., Sun, Q., Xu, C., Zhao, P., Lou, J., Tao, C., Geng, X., Lin, Q., Chen, S., Zhang, D.: Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct. arXiv preprint arXiv:2308.09583 (2023)
- 44. OpenAI: Chatgpt. https://chat.openai.com (2023)
- 45. OpenAI: Gpt-4 technical report. ArXiv abs/2303.08774 (2023)
- 46. OpenAI: GPT-4V(ision) system card (2023), https://openai.com/research/ gpt-4v-system-card
- 47. Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Gray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Askell, A., Welinder, P., Christiano, P., Leike, J., Lowe, R.: Training language models to follow instructions with human feedback. In: Advances in Neural Information Processing Systems (2022)
- Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., Sutskever, I.: Learning transferable visual models from natural language supervision. In: International Conference on Machine Learning (2021), https://api.semanticscholar.org/CorpusID: 231591445

- 49. Roy, S., Roth, D.: Solving general arithmetic word problems. ArXiv abs/1608.01413 (2016), https://api.semanticscholar.org/CorpusID:560565
- Seo, M., Hajishirzi, H., Farhadi, A., Etzioni, O., Malcolm, C.: Solving geometry problems: Combining text and diagram interpretation. In: Proceedings of the 2015 conference on empirical methods in natural language processing. pp. 1466–1476 (2015)
- 51. Su, Y., Lan, T., Li, H., Xu, J., Wang, Y., Cai, D.: Pandagpt: One model to instruction-follow them all. arXiv preprint arXiv:2305.16355 (2023)
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al.: Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971 (2023)
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al.: Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288 (2023)
- 54. Wang, K., Ren, H., Zhou, A., Lu, Z., Luo, S., Shi, W., Zhang, R., Song, L., Zhan, M., Li, H.: Mathcoder: Seamless code integration in LLMs for enhanced mathematical reasoning. In: The Twelfth International Conference on Learning Representations (2024), https://openreview.net/forum?id=z8TW0ttBPp
- 55. Xu, P., Shao, W., Zhang, K., Gao, P., Liu, S., Lei, M., Meng, F., Huang, S., Qiao, Y., Luo, P.: Lvlm-ehub: A comprehensive evaluation benchmark for large visionlanguage models. arXiv preprint arXiv:2306.09265 (2023)
- Xu, R., Wang, X., Wang, T., Chen, Y., Pang, J., Lin, D.: Pointllm: Empowering large language models to understand point clouds. arXiv preprint arXiv:2308.16911 (2023)
- 57. Yang, S., Liu, J., Zhang, R., Pan, M., Guo, Z., Li, X., Chen, Z., Gao, P., Guo, Y., Zhang, S.: Lidar-llm: Exploring the potential of large language models for 3d lidar understanding. arXiv preprint arXiv:2312.14074 (2023)
- 58. Ye, Q., Xu, H., Xu, G., Ye, J., Yan, M., Zhou, Y., Wang, J., Hu, A., Shi, P., Shi, Y., Jiang, C., Li, C., Xu, Y., Chen, H., Tian, J., Qian, Q., Zhang, J., Huang, F.: mplug-owl: Modularization empowers large language models with multimodality (2023)
- Ye, Q., Xu, H., Ye, J., Yan, M., Hu, A., Liu, H., Qian, Q., Zhang, J., Huang, F., Zhou, J.: mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration (2023)
- 60. Ying, K., Meng, F., Wang, J., Li, Z., Lin, H., Yang, Y., Zhang, H., Zhang, W., Lin, Y., Liu, S., et al.: Mmt-bench: A comprehensive multimodal benchmark for evaluating large vision-language models towards multitask agi. ICML 2024 (2024)
- 61. Yue, X., Ni, Y., Zhang, K., Zheng, T., Liu, R., Zhang, G., Stevens, S., Jiang, D., Ren, W., Sun, Y., Wei, C., Yu, B., Yuan, R., Sun, R., Yin, M., Zheng, B., Yang, Z., Liu, Y., Huang, W., Sun, H., Su, Y., Chen, W.: Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. arXiv preprint arXiv:2311.16502 (2023)
- Yue, X., Qu, X., Zhang, G., Fu, Y., Huang, W., Sun, H., Su, Y., Chen, W.: Mammoth: Building math generalist models through hybrid instruction tuning. arXiv preprint arXiv:2309.05653 (2023)
- Zhang, H., Li, X., Bing, L.: Video-llama: An instruction-tuned audio-visual language model for video understanding. arXiv preprint arXiv:2306.02858 (2023)
- Zhang, R., Guo, Z., Zhang, W., Li, K., Miao, X., Cui, B., Qiao, Y., Gao, P., Li, H.: Pointclip: Point cloud understanding by clip. In: CVPR 2022 (2022)

- 40 R. Zhang and D. Jiang et al.
- 65. Zhang, R., Han, J., Zhou, A., Hu, X., Yan, S., Lu, P., Li, H., Gao, P., Qiao, Y.: LLaMA-adapter: Efficient fine-tuning of large language models with zero-initialized attention. In: The Twelfth International Conference on Learning Representations (2024), https://openreview.net/forum?id=d4UiXAHN2W
- 66. Zhang, R., Jiang, Z., Guo, Z., Yan, S., Pan, J., Dong, H., Gao, P., Li, H.: Personalize segment anything model with one shot. ICLR 2024 (2023)
- Zhang, R., Wang, L., Qiao, Y., Gao, P., Li, H.: Learning 3d representations from 2d pre-trained models via image-to-point masked autoencoders. CVPR 2023 (2023)
- Zhang, R., Wei, X., Jiang, D., Zhang, Y., Guo, Z., Tong, C., Liu, J., Zhou, A., Wei, B., Zhang, S., et al.: Mavis: Mathematical visual instruction tuning. arXiv preprint arXiv:2407.08739 (2024)
- 69. Zhou, A., Wang, K., Lu, Z., Shi, W., Luo, S., Qin, Z., Lu, S., Jia, A., Song, L., Zhan, M., et al.: Solving challenging math word problems using gpt-4 code interpreter with code-based self-verification. arXiv preprint arXiv:2308.07921 (2023)
- Zhu, D., Chen, J., Shen, X., Li, X., Elhoseiny, M.: Minigpt-4: Enhancing visionlanguage understanding with advanced large language models. arXiv preprint arXiv:2304.10592 (2023)