

U-MATH: A UNIVERSITY-LEVEL BENCHMARK FOR EVALUATING MATHEMATICAL SKILLS IN LLMs

Anonymous authors

Paper under double-blind review

ABSTRACT

The current evaluation of mathematical skills in LLMs is limited, as existing benchmarks are relatively small, primarily focus on elementary and high-school problems, or lack diversity in topics. Additionally, the inclusion of visual elements in tasks remains largely under-explored.

To address these gaps, we introduce **U-MATH**, a novel benchmark of 1,125 unpublished open-ended university-level problems sourced from teaching materials. It is balanced across six core subjects, with 20% of problems requiring image understanding. Given the open-ended nature of U-MATH problems, we employ an LLM to judge the correctness of generated solutions. To this end, we release **μ -MATH**, an dataset to evaluate the LLMs’ capabilities in judging solutions.

The evaluation of general domain, math-specific, and multimodal LLMs highlights the challenges presented by U-MATH. Our findings reveal that LLMs achieve a maximum accuracy of only 53% on text-based tasks, with even lower 30% on visual problems. The solution assessment proves challenging for LLMs, with the best LLM judge having an F1-score of 76% on **μ -MATH**

During review, we publish the U-MATH and **μ -MATH** datasets on OSF.¹

Example: Differential Calculus.

U-MATH Problem:

The function $s(t) = 2 \cdot t^3 - 3 \cdot t^2 - 12 \cdot t + 8$ represents the position of a particle traveling along a horizontal line.

1. Find the velocity and acceleration functions.
2. Determine the time intervals when the object is slowing down or speeding up.

Reference Solution (shortened):

The velocity is $v(t) = s'(t) = 6 \cdot t^2 - 6 \cdot t - 12$, zeros of the $v(t)$ are $t = -1, 2$.

The acceleration is $a(t) = v'(t) = 12 \cdot t - 6$, zero of the $a(t)$ is $t = \frac{1}{2}$.

It speeds up when $v(t)$ and $a(t)$ have the same sign, and slows down when opposite.

Interval	$v(t)$	$a(t)$	Behavior
$(-\infty, -1)$	> 0	< 0	Slowing down
$(-1, \frac{1}{2})$	< 0	< 0	Speeding up
$(\frac{1}{2}, 2)$	< 0	> 0	Slowing down
$(2, \infty)$	> 0	> 0	Speeding up

Accounting for non-negative time, speed up on $(0, 1/2)$ and $(2, \infty)$, slow down on $(1/2, 2)$.

Figure 1: U-MATH covers university-level topics and require multiple steps to solve. A random sample is provided: reference solution is shortened to save space. In this example, common model errors is overlooking the non-negativity of time.

1 INTRODUCTION

Mathematical reasoning is a fundamental domain for assessing the true capabilities of Large Language Models (LLMs) to reason (Ahn et al., 2024). While existing benchmarks like GSM8K (Cobbe et al.,

¹https://osf.io/jpsa4/?view_only=d588b9fa862345cb98ccf7238a157cea

2021) and MATH (Hendrycks et al., 2021) provide valuable insights, they primarily focus on school-level mathematics. This leaves a significant gap in understanding how LLMs perform on more advanced, university-level problems. Moreover, these benchmarks are becoming saturated, as GPT-4, using advanced prompting techniques, has achieved over 92% success rate on GSM8K and 80% on MATH (Achiam et al., 2023).

Recent works, such as CHAMP (Mao et al., 2024) and MathOdyssey (Fang et al., 2024), aim to introduce more challenging problems but are limited in size (<400 samples) and lack comprehensive topic coverage. The most challenging problems stem from school-level competitions or olympiads, missing the crucial middle ground of university-level coursework that reflects academic demands.

Furthermore, there is a growing interest in assessing multi-modal LLMs’ abilities to perform mathematical reasoning involving visual elements (Ahn et al., 2024). Large datasets like MathVista (Lu et al., 2023), We-Math (Qiao et al., 2024), or MathVerse (Zhang et al., 2024) provide an extensive set of (mostly) visual tasks but may lack university-level problems and often rely on multiple-choice validation, leading to easier problems and faster saturation of benchmarks.

In turn, evaluating complex free-form answers remains a significant challenge for the field (Hendrycks et al., 2021). Current methods often rely on LLM judges to assess problems, which introduces potential biases and inconsistencies (Zheng et al., 2023). Errors introduced by automatic evaluators are often overlooked in popular benchmarks. This oversight makes it impossible to account for judge biases, which detracts from the reliability of the evaluation results.

Recent studies also indicate that evaluation of mathematical solutions is a demanding task (Zeng et al., 2023; Xia et al., 2024) and that an LLM’s ability to judge mathematical solutions is correlated with its problem-solving performance (Stephan et al., 2024), further signifying the importance of evaluations designed to assess the evaluators themselves — also called meta-evaluations.

Popular datasets for the task of mathematical meta-evaluation are PRM800K (Lightman et al., 2023), MR-GSM8K (Zeng et al., 2023) and MR-MATH (Xia et al., 2024). However, these are all based on the GSM8K and MATH datasets, still leaving a gap in meta-evaluations for university-level problems.

Aiming to bridge these gaps and provide a comprehensive evaluation of LLMs’ mathematical capabilities, we introduce **U-MATH** (*University Math*) and a supplementary meta-evaluation dataset, which we refer to as μ -**MATH** (*Meta U-MATH*). Our main contributions are:

1. **U-MATH Benchmark** (Section 3): We open-source a set of 1,125 of university-level problems collected from actual coursework with final answers and solutions. About 20% of problems require image *understanding* to be solved. The text-only part of the benchmark is balanced across 6 key subjects: Precalculus, Algebra, Differential Calculus, Integral Calculus, Multivariable Calculus, and Sequences&Series.
2. **μ -MATH Meta-Evaluation Benchmark** (Section 3.3): Additionally, we introduce a set of 340 meta-evaluation tasks sourced from U-MATH problems and designed to rigorously assess the quality of LLM judges. We manually select approximately 30% of the U-MATH problem statements and golden answers, supplying them with LLM-generated solutions, and label them based on whether the generated solutions are correct or not. The benchmark is designed to be challenging for LLM judges yet representative of the typical university-level math grading tasks.
3. **Comparison of Models** (Section 4): We conduct a comparative analysis of various open-source and proprietary LLMs on U-MATH. Our analysis highlights the high performance of specialized models in text-only problems and the superiority of proprietary models in visual tasks with the best U-MATH accuracy of 49%. Additionally, we examine several popular LLMs on μ -MATH to assess their ability to judge free-form mathematical problems. Our results show the best model achieving the macro F1-score of 76%.

We release the U-MATH and μ -MATH benchmarks under a permissive license to facilitate further research and ensure reproducibility.

2 BACKGROUND

Enhancing and evaluating the mathematical reasoning capabilities of LLMs is essential in AI research (Ahn et al., 2024). Studies show that finetuning with mathematical and code-related data enhances models’ general skills (Prakash et al., 2024). Mathematical tasks require logical thinking and multi-step problem-solving, thus improving overall reasoning abilities in LLMs (Chen et al., 2024).

This leads to the problem of evaluating LLM’s math abilities. Despite the significant progress, many existing benchmarks are limited in scope, focusing primarily on school-level mathematics or limited in size and topic coverage. Table 1 summarizes popular text-only and visual mathematical benchmarks.

Dataset	Levels	%Uni. Level	#Test	%Visual	%Free Form Answer
MMLU _{Math} (Hendrycks et al., 2020)	E H C	0	1.3k	0	0
GSM8k (Cobbe et al., 2021)	E	0	1k	0	0
MATH (Hendrycks et al., 2021)	H O	0	5k	0	100
MiniF2F (Zheng et al., 2021)	E H O	0	244	0	100
OCWCourses (Lewkowycz et al., 2022)	U	100	272	0	100
ProofNet (Azerbaiyev et al., 2023)	C U	≈50	371	0	100
CHAMP (Mao et al., 2024)	H	0	270	0	100
MathOdyssey (Fang et al., 2024)	H U O	26	387	0	100
MMMU _{Math} (Yue et al., 2023)	C	0	505	100	0
MathVista (Lu et al., 2023)	E H C	0	5k	100	46
MATH-V (Wang et al., 2024)	E H O	0	3k	100	50
We-Math (Qiao et al., 2024)	E H U	≈20	1.7k	100	0
MathVerse (Zhang et al., 2024)	H	0	4.7k	83.3	45
U-MATH (this work)	U	100	1.1k	20	100

Table 1: Existing Auto-evaluation Math benchmarks with corresponding test samples *published*, visual samples percent, and percent of multiple-choice questions. Level denotes E Elementary to Middle School, H High School, C College, U University, O Different Olympiads.

Textual Mathematical Benchmarks. Early efforts to assess LLMs’ mathematical abilities have emerged in datasets like MathQA (Amini et al., 2019) and the mathematics subset of MMLU (Hendrycks et al., 2020). These early benchmarks emphasized the importance of operation-based reasoning in solving mathematical word problems, typically in a multiple-choice format. Nowadays, even smaller models (e.g., 7B parameters) have achieved high scores on these tasks (Li et al., 2024b), suggesting that these benchmarks are becoming saturated. In response, more comprehensive datasets have emerged, such as GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021), or MGSM (Shi et al., 2022) (multilingual version of 250 GSM8K samples). These popular benchmarks are crucial for evaluating LLMs’ mathematical reasoning skills. However, they primarily focus on school-level problems, which may not fully assess the depth of mathematical reasoning.

Recent efforts attempt to address more advanced mathematical concepts. MathOdyssey (Fang et al., 2024) with competition problems, OCWCourses (Lewkowycz et al., 2022) from actual MIT courses, and ProofNet (Azerbaiyev et al., 2023) focusing on proofs aim to evaluate undergraduate-level or olympiad-level knowledge. However, these datasets are constrained by their small sizes (e.g., 387, 272, and 371 samples), limiting their statistical robustness and topic coverage. For example, MathOdyssey is limited to 101 samples in university-level topics (Calculus, Algebra, and Diff. Equations and Statistics). Other specialized datasets like MiniF2F (Zheng et al., 2021) provide valuable parallel corpora in formal languages, while CHAMP (Mao et al., 2024) offers helpful context and hints, but both are similarly limited in scale with 244 and 270 samples. Additionally, both heavily rely on already published resources: CHAMP sources material from a book, while MiniF2F re-uses international olympiads and MATH dataset problems. An attempt to provide a more robust evaluation, GHOSTS (Frieder et al., 2024) dataset, provides 728 problems (both from other datasets and new ones) but does not provide reference solutions and answers, focusing instead on human evaluation, making cheap automatic evaluation impossible.

The current datasets are either too small, leading to higher measurement errors, or focus mainly on elementary and high school math, leaving a gap in evaluating LLMs’ proficiency in advanced university-level math topics.

Visual Mathematical Benchmarks. As multimodal LLMs gain prominence, there is a growing need for visual mathematical benchmarks (Zhang et al., 2024; Qiao et al., 2024). Early efforts in this domain focus primarily on geometric problems, as seen in datasets like GeoQA (Chen et al., 2022b), UniGeo (Chen et al., 2022a), and Geometry3K (Lu et al., 2021). These datasets have a narrow focus that does not encompass the breadth of mathematical visual reasoning required at advanced levels.

More recent benchmarks attempt to broaden the scope of visual mathematical evaluation. One of the first comprehensive attempts is the mathematical subset of MMMU (Yue et al., 2023), which offers 505 college-level multiple-choice questions, all with images. However, its multiple-choice format limits the complexity of problems that can be posed. MathVista (Lu et al., 2023) collects 28 existing datasets and introduces 3 new datasets with a total of 5k samples (1k testmini samples). However, as shown by Qiao et al. (2024), it faces challenges with data quality due to its compilation from older datasets.

The latest benchmarks, such as MATH-V (Vision) (Wang et al., 2024) and We-Math (Qiao et al., 2024), extend this approach to collect 3k and 1.7k visual samples, respectively. However, both datasets rely on multiple-choice questions in the test set, leading to faster saturation. MathVerse (Zhang et al., 2024) further extends this approach, relying on visual elements and providing some simple text problems with 1.2k brand-new samples. Among these, only the We-Math dataset includes university-level mathematical problems.

Our U-MATH dataset improves on existing benchmarks with 225 of 1,125 university-level problems that require visual elements (graph, table, diagram) to be solved. This balanced ratio ensures models are challenged to handle both traditional and visual problem-solving without over-relying on visuals, mirroring real-world scenarios.

Large Language Models for Mathematics. The application of LLMs to mathematical problem-solving shows promising results, particularly with models like GPT-3.5 and GPT-4 demonstrating strong reasoning abilities for complex tasks such as those in the MATH dataset (Achiam et al., 2023). While open-source models initially lagged in performance on advanced mathematical tasks, the Llama-3.1 (Dubey et al., 2024) is approaching parity with proprietary models. The most popular benchmarks, MATH and GSM8K, are nearing saturation, with Llama 3.1 405B achieving scores of 73.8% and 96.8%, respectively. Similarly, a Qwen2.5-Math-72B model (Yang et al., 2024b; Team, 2024) reach 85.9% on MATH while Qwen2-Math-72B (Yang et al., 2024a) reaches 96.7% on GSM8k.

To enhance LLMs’ mathematical capabilities, researchers develop various prompt-based methods (Liu et al., 2021). These include techniques for encouraging chain-of-thought generation (Wei et al., 2022), selecting final results from multiple sampled outputs (Wang et al., 2022), and using external tools such as calculators, WolframAlpha or Python interpreters (Gao et al., 2023) to reduce arithmetic errors. Additionally, instruction tuning during pre-training has been identified as a key factor in improving performance (Wang et al., 2017). While these approaches show promise, their effectiveness on university-level problems still needs to be explored due to the lack of suitable large-scale benchmarks.

Mathematical solution verification. Evaluating mathematical solutions is uniquely challenging due to the open-ended nature of answers and the inherent ambiguity in mathematical expressions. Consequently, many benchmarks opt for multiple-choice formats due to their grading simplicity. However, this approach often simplifies tasks, providing hints that models can exploit (Li et al., 2024c; Pezeshkpour and Hruschka, 2023).

While free-form evaluation using LLM judges is widespread (Zheng et al., 2023), it is known to introduce potential errors (Zheng et al., 2023), since evaluating mathematical solutions is a complex task in its own right (Zeng et al., 2023; Xia et al., 2024). These evaluation errors are largely overlooked and unaccounted for, limiting the reliability of inferences drawn from such evaluations.

Hence, it is important to be able to estimate the performance of automatic evaluators and to choose the most adequate among them. Recent studies show that evaluation performance is correlated with but does not equal problem-solving performance (Stephan et al., 2024). This underscores the importance of benchmarks designed specifically to assess the evaluators — also called meta-evaluations.

There are existing benchmarks that are well-suited for meta-evaluations. PRM800K (Lightman et al., 2023) contains 800K annotated steps from 75K solutions to 12K MATH dataset problems, designed

to confuse reward models. FELM (Zhao et al., 2024) provides GPT-3.5 annotations for solutions to 208 GSM8K and 194 MATH problems. MR-GSM8K (Zeng et al., 2023) and MR-MATH (Xia et al., 2024) introduce meta-evaluation datasets focused on the GSM8K and MATH datasets, respectively. However, these are either based on elementary to high-school level problems or feature specifically competition-style math, leaving a gap in meta-evaluations on complex and practical university tasks.

To address this, we introduce μ -MATH— a meta-evaluation dataset based on a subset of U-MATH problems. It provides LLM-generated solutions with verified labels, enabling precise and fine-grained assessment of LLMs’ evaluation abilities.

3 U-MATH

We present **U-MATH** (stands for University Math) — a benchmark designed to challenge LLMs with problems requiring deep understanding and advanced reasoning. The problems span 6 core topics and range in difficulty and number of questions. A subset of 20% of problems includes images to test the models’ ability to interpret and reason with graphical information. Reference solutions and answers accompany all problems.

Accuracy is the primary performance metric for **U-MATH**, its text-only problems (**U-MATH_T**) and problems that include a visual component (**U-MATH_V**). The main performance measure for μ -MATH is **macro-F1**.

We use an LLM as a judge (Zheng et al., 2023) to measure the accuracy of the free-form answers against the golden solutions. A problem is considered solved only if all required questions are answered and all requested items (e.g., all saddle points) are correctly identified.

3.1 DATASET COLLECTION

To create a benchmark that authentically reflects university-level mathematics, we collaborate with [ANONYMIZED], a platform providing learning content and software for top US universities specialized in mathematics. The problems are sourced from ongoing courses across various institutions currently run on the [ANONYMIZED] platform. Problems and solutions are crafted by subject matter experts and represent real-world academic standards. These samples are unpublished and have not been exposed to any external sources. Thus, the dataset could not be leaked to current LLMs.

We employ a multi-stage filtering process to select challenging problems from tens of thousands of available samples. First, we filter out problems with short solutions (< 100 characters) and problems in multiple-choice format. As LLMs are not designed to perform arithmetic calculations and are prone to errors (Hendrycks et al., 2021; Lewkowycz et al., 2022), we focus on testing mathematical reasoning rather than calculations. We filter out problems marked as allowing calculator usage. As for the visual problems selection, we chose to keep problems with a single image for convenience.

Next, we employ several small LLMs (LLaMA-3.1-8B (Dubey et al., 2024), Qwen2-7B (Yang et al., 2024a), Mistral-7B (Jiang et al., 2023), Mathstral-7B, NuminaMath-7B (Beeching et al., 2024)) to solve the problems. We select 175 most challenging problems for each subject based on the average problem solution rate. We randomly select 150 samples for the public test, keeping the rest for the private set. For this step, we use the same pipeline as described in Section 4. This way, we ensure that none of the individual models influence problem selection largely and that there is no overfitting to a specific LLM.

Next, we enlist a team of paid experts from the [ANONYMIZED], who actively teach various Calculus courses. The experts verify that each problem is suitable either for assessing the subject knowledge expected of college or university students or for testing prerequisite knowledge. The team thoroughly reviewed and affirmed that the selected problems meet these criteria.

3.2 DATASET STATISTICS

The U-MATH benchmark comprises **1,125** carefully curated and validated mathematical problems. These problems are distributed across **6 core subjects** with about 20% of the tasks incorporating visual

elements, such as graphs, tables, and geometric figures, mirroring the multi-modal nature of real-world mathematical problems: Precalculus (Review), Algebra, Differential Calculus (+Differential Equations), Integral Calculus, Multivariable Calculus, and Sequences & Series.

Math Subject	#Textual	#Visual	Avg. Questions	Avg. Answers
Algebra	150	30	1.93	1.28
Differential Calculus	150	68	2.37	1.15
Integral Calculus	150	80	1.09	1.01
Multivariable Calculus	150	28	1.74	1.09
Precalculus	150	13	1.51	1.23
Sequences and Series	150	6	1.36	1.00
All	900	225	1.66	1.12

Table 2: Average number of questions per problem and answers per question in U-MATH.

Table 2 summarizes the distribution of problems across different subjects. The average is **1.7** questions per problem (e.g., local minima, maxima, and increasing intervals are asked), and the average of **1.1** answers per question (for example, the number of saddle points in the correct answer).

3.3 META-EVALUATION FRAMEWORK (μ -MATH)

Mathematical problem evaluation is not straightforward. Even simple expressions like $x \cdot 0.5$ may have valid forms like $\frac{x}{2}$, $x \div 2$, $x/2$, or unsimplified variants like $9x/18$. In practice, evaluating free-form solutions requires testing expression equivalence in much less trivial cases, especially with more advanced problems (refer to Section A.3 in Appendix for an example).

To systematically study the ability of LLMs to evaluate free-form mathematical solutions on advanced, university-level problems, we introduce the μ -MATH (Meta U-MATH) benchmark. It consists of a curated subset of U-MATH samples, supplied with LLM-generated solutions — both correct and not. The solutions are labeled using a combination of manual inspection and automated verification via [ANONYMIZED]-API, which allows to test formal equivalence of mathematical expressions.

We selected 340 U-MATH problems (around 30%) based on their assessment difficulty to create a challenging meta-evaluation set. This benchmark does not aim to reflect the overall U-MATH distribution but rather provides a robust test for LLM judges. We focused on text-only problems, excluding those needing images, due to the limited size of the labeled U-MATH subset. The Qwen2.5-Math-7B model was used for solution generation, with about 15% of the solutions being intentionally incorrect. Ultimately, we have **340** samples in μ -MATH— one for each of the U-MATH-sourced problems.

A tested model is provided with a problem statement, a reference answer, and a solution to evaluate. We treat this as a binary classification task, using the macro-averaged **F1-score as the primary metric** to minimize the effect of class imbalance. Additionally, we report Positive Predictive Value (PPV or Precision) and True Positive Rate (TPR or Recall) for the positive class as well as Negative Predictive Value (NPV) and True Negative Rate (TNR) for the negative class, offering a finer-grained performance evaluation.

4 EXPERIMENTS AND RESULTS

4.1 EXPERIMENTAL SETUP

We select some top-performing recent LLMs to evaluate.

All LLMs are tested using the same prompts and settings for fair comparison. The LLMs are restricted to a single generation of 4096 tokens with the temperature set to 0. We employ chain-of-thought (CoT) prompting (Wei et al., 2022) to encourage models to ‘think’ before providing an answer. Images are included directly in the prompts for multimodal LLMs. To text-only LLMs the problem description is provided as-is without visual elements. Refer to Appendix C.1 for the full prompts.

324
325
326
327
328
329
330
331
332
333
334
335

Model	Source	Size(s)	Visual	Open-weights
Mathstral-v0.1	(Mistral.ai, 2024)	7B	×	✓
NuminaMath-CoT	(Beeching et al., 2024)	7B	×	✓
LLaMA-3.1	(Dubey et al., 2024)	8B, 70B	×	✓
Qwen2-Math	(Yang et al., 2024a)	7B, 72B	×	✓
Qwen2.5-Math	(Yang et al., 2024b)	7B, 72B	×	✓
Qwen2.5	(Team, 2024)	7B, 72B	×	✓
Pixtral-12B-2409	(Mistral AI, 2024)	12B	✓	✓
LLAVA One Vision _(Qwen2-7B)	(Li et al., 2024a)	8B	✓	✓
Qwen2-VL	(Yang et al., 2024a)	7B, 72B	✓	✓
gpt-4o-2024-05-13	(OpenAI, 2024)	unknown	✓	×
Gemini-1.5-Pro-002	(Team et al., 2024)	unknown	✓	×

336
337
338

Table 3: LLMs name, version and sizes used for our experiments.

339
340
341
342
343
344
345
346

We report accuracy based on widely available gpt-4o-2024-05-13 as-a-judge for the final results, despite this not being the best model, yet conservative in false negative rate (as discussed in Section 4.2.2). The judge is presented with the problem statement, golden answer, and generated solutions. The temperature is set to 0. The judge is asked to extract the ‘student’s answer’, make derivations that may be necessary, and compare solutions. After this ‘reasoning phase’, we ask the judge to provide a Yes/No response in the same chat, which we interpret as a desired binary metric. Refer to Appendix C.2 for full judge prompt.

347
348

4.2 RESULTS

349

Table 4 summarizes the performance of text-only and multimodal LLMs on the U-MATH benchmark.

350
351

Model	U-MATH	U-MATH		Algebra		Diff. C.		Integral C.		Multivar C.		Precalculus		Seq.& Series	
		T	V	T	V	T	V	T	V	T	V*	T	V*		
Text-only models															
Mathstral-7B-v0.1	16.36	20.1	1.3	56.0	3.3	4.0	1.5	2.0	1.2	6.0	0.0	40.7	0.0	12.0	0.0
Llama-3.1-8B	15.47	19.0	1.3	52.0	3.3	4.7	2.9	1.3	0.0	10.0	0.0	36.7	0.0	9.3	0.0
Qwen2.5-7B	32.18	39.8	1.8	77.3	3.3	14.0	0.0	10.7	1.2	34.7	7.1	72.7	0.0	29.3	0.0
Qwen2-Math-7B	28.62	35.1	2.7	77.3	3.3	9.3	1.5	8.0	2.5	23.3	7.1	70.7	0.0	22.0	0.0
Qwen2.5-Math-7B	33.24	41.0	2.2	64.0	0.0	16.7	1.5	21.3	0.0	35.3	7.1	72.0	7.7	36.7	16.7
Llama-3.1-70B	26.13	32.0	2.7	70.0	10.0	10.0	0.0	5.3	2.5	25.3	3.6	55.3	0.0	26.0	0.0
Qwen2.5-72B	34.04	41.9	2.7	69.3	0.0	25.3	0.0	10.7	1.2	39.3	14.3	70.0	0.0	36.7	16.7
Qwen2-Math-72B	35.20	43.6	1.8	80.0	0.0	22.7	0.0	16.0	2.5	32.0	0.0	76.7	0.0	34.0	33.3
Qwen2.5-Math-72B	41.16	50.6	3.6	73.3	3.3	33.3	0.0	23.3	7.5	48.7	3.6	82.0	0.0	42.7	0.0
Multimodal models															
Qwen2-VL-7B	17.33	20.3	5.3	51.3	6.7	8.7	10.3	1.3	1.2	6.0	3.6	44.0	7.7	10.7	0.0
LLaVA-OV _(Qwen2-7B)	14.40	17.8	0.9	48.7	0.0	4.7	1.5	0.7	0.0	7.3	3.6	36.7	0.0	8.7	0.0
Pixtral-12B-2409	15.64	18.1	5.8	46.7	16.7	4.7	10.3	0.7	1.2	6.7	0.0	41.3	0.0	8.7	0.0
Qwen2-VL-72B	26.93	30.1	14.2	70.0	13.3	11.3	20.6	5.3	7.5	18.7	21.4	58.7	15.4	16.7	0.0
GPT-4o	36.53	41.8	15.6	80.0	16.7	22.0	14.7	11.3	13.8	38.7	25.0	68.0	15.4	30.7	0.0
Gemini-1.5-Pro	48.89	53.4	30.7	84.7	56.7	37.3	27.9	27.3	22.5	42.7	28.6	81.3	46.2	47.3	16.7

362
363
364
365
366

Table 4: Comparison of models’ accuracy on our U-MATH benchmark and its subjects. Scores for various mathematical categories, including text and visual analysis, are displayed. For each subject 2 numbers are provided - text-only (T) and visual (V) problems. Asterisk denotes a small number of samples (< 15). Free-form solutions judged by gpt-4o-2024-05-13. Images are not included in the prompt for text-only models, only the problem statement. **Bold** indicates the best result in each group.

373
374
375
376
377

Among text-only models, the math-specific model Qwen2.5-Math-72B achieves the highest overall accuracy at 41.2%, showcasing strong mathematical reasoning capabilities. In the multi-modal model group, **Gemini-1.5-pro-002** leads with an overall accuracy of **48.9%**, highlighting the advantages of integrating visual processing. In contrast, LLaVA-OV-Qwen2-7B lacks mathematical abilities in visual and textual tasks with **15.6%** on a U-MATH benchmark. Building on these results, several key trends emerge:

- **Model Size vs. Specialization:** Larger models expectedly outperform smaller ones. However, the small specialized model Qwen2.5-Math-7B surpasses or performs on par with 10 times larger models like Qwen2.5-72B or LLaMA-3.1-70B and the leading model Gemini-1.5-Pro in Multivariable Calculus. Similarly, Pixtral-12B performs consistently worse than minor Qwen2-VL-7B, indicating a lack of university-level data in training.
- **Textual vs. Visual Problem-Solving:** Across multimodal models, text-only problems’ accuracy vastly exceeds visual problems, highlighting areas for further improvement. The text-only models can solve a small percentage of visual problems, primarily due to guessing or judging errors discussed in Section 4.2.2.
- **Proprietary vs. Open-weights model:** Proprietary models like GPT-4o and Gemini still offer top or competitive performance but lack transparency and flexibility. At the moment, the gap is evident in visual comprehension. Open-weight models like Qwen-Math have taken a big step toward hitting top performance.

4.2.1 SUBJECT-SPECIFIC RESULTS

We analyze model performance across different mathematical subjects to uncover underlying trends. The models excel in Precalculus and Algebra, particularly on text-based problems, aligning with previous findings and benchmark saturations (Ahn et al., 2024). However, they struggle with visual problems in these areas, indicating a need for better visual-symbolic integration. In Sequences and Series, models demonstrate strong performance on abstract, formula-based text tasks, reflecting the logical structure of the subject. However, insufficient visual training data limits a complete evaluation of their capabilities in this domain.

In contrast, the models show moderate success in Differential and Multivariable Calculus text tasks, highlighting challenges in handling abstract, multi-dimensional concepts, especially in visual form. Integral Calculus poses the most significant difficulty, with lower performance in both text and visual tasks. Interpreting curves and areas proves particularly challenging, emphasizing the necessity for more focused multimodal training. Additionally, the extensive expressions common in Integral Calculus problems often confuse the models, diminishing their effectiveness.

4.2.2 META-EVALUATION (μ -MATH)

In this section, we present the results of our experiments on evaluating the performance of LLMs judging mathematical solution correctness. Using the μ -MATH benchmark, we compare several top-performing models, as displayed in Table 5.

Model	U-MATH _T	μ -MATH				
		F1 _{macro}	TPR	TNR	PPV	NPV
Mathstral-7B-v0.1	21.1	56.56	65.40	66.67	91.75	25.37
LLaMA-3.1-8B	19.0	58.29	68.86	64.71	91.71	26.83
Qwen2-Math-7B	35.1	65.01	87.20	45.10	90.00	38.33
Qwen2.5-Math-7B	41.0	66.90	85.81	52.94	91.18	39.71
Qwen2.5-7B	39.8	63.89	85.81	45.10	89.86	35.94
LLaMA-3.1-70B	32.0	66.27	74.05	80.39	95.54	35.35
Qwen2-Math-72B	43.6	73.03	92.39	52.94	91.75	55.10
Qwen2.5-Math-72B	50.6	69.95	86.85	58.82	92.28	44.12
Qwen2.5-72B	41.9	75.74	89.27	68.63	94.16	53.03
gpt-4o-2024-05-13	41.8	70.62	78.89	82.35	96.20	40.78
Gemini-1.5-pro-002	53.4	64.92	79.93	60.78	92.03	34.83

Table 5: Comparison of different models on the μ -MATH benchmark. Presented are standard binary classification metrics — Macro F1-score (F1), True Positive Rate (TPR), True Negative Rate (TNR), Positive Predictive Value (PPV), and Negative Predictive Value (NPV), with F1 as the primary one. U-MATH_{Text} accuracy score from Table 4 is added for comparison of model’s performance as a math solver vs as a math judge. **Bold** indicates the best result for each column.

The results of our meta-evaluation experiments highlight several challenges in using LLMs to judge mathematical solutions. No model consistently excels across all the tasks, with **the best-performing**

432 **model**, Qwen2.5-72B, **achieving an F1-score of 75.74**, which still leaves significant room for
433 improvement.

434 The Qwen2-Math specialist model often outperforms GPT-4o while displaying qualitatively different
435 behavior. Qwen2-Math has the highest true positive rate and the lowest true negative rate among
436 the similar-size models, while GPT-4o exhibits the opposite pattern. This aligns with manual
437 observations, which indicate that GPT-4o focuses more on matching final answers and often cannot
438 perform complex derivations, making it a more conservative judge and increasing the rate of false
439 negative judgments. In contrast, Qwen2-Math ‘follows the solution’ and excels at mathematical
440 transformations, but this also leads to higher hallucination risk and more false positives.

441 Increased mathematical problem solving ability does not necessarily lead to lower judgment errors.
442 For example, Qwen2.5-Math outperforms Qwen2-Math on U-MATH tasks but does not show a
443 considerable advantage on the μ -MATH tasks. This too is in line with manual examinations revealing
444 that Qwen2.5-Math is overly specialized in problem-solving, which impairs its ability to follow
445 instructions properly and judge solutions. For instance, the model often starts solving the problem
446 from scratch instead of evaluating the solution. This indicates that finding a balance between
447 domain-specific skills and general capabilities is essential.

448 This point is further illustrated by the Qwen2.5 generalist model’s highest overall classification score:
449 it strikes a good balance between learning from general domain data and mathematical data produced
450 by specialist models.

452 5 CONCLUSION

453 In this paper, we introduced **U-MATH**, a novel benchmark designed to evaluate the mathematical
454 reasoning capabilities of LLMs at the university level. U-MATH comprises 1,125 unpublished
455 open-ended problems sourced from actual teaching materials, balanced across 6 core mathematical
456 subjects, with 20% of the problems requiring image understanding. Additionally, we presented
457 **μ -MATH**, a meta-evaluation dataset aimed at assessing the ability of LLMs to evaluate free-form
458 mathematical solutions.

459 Our experiments with general-domain, math-specific, and multimodal LLMs revealed significant
460 challenges in advanced mathematical reasoning and visual problem-solving. The best-performing
461 models achieved an accuracy of only 53% on text-based tasks and even lower, 30%, on visual
462 problems for Gemini-1.5-pro-002 model. Furthermore, the task of solution assessment has proven
463 to be challenging for LLMs, with the highest μ -MATH F1-score being 76% for Qwen2.5-72B,
464 indicating room for improvement in LLMs’ evaluation capabilities and highlighting that widely used
465 GPT-4o is not a silver bullet for judging.

466 **Limitations.** While U-MATH offers a substantial and diverse set of university-level problems, it
467 does not cover the full breadth and depth of advanced mathematical topics taught at universities.
468 The carefully curated selection may introduce biases, potentially favoring certain problem types or
469 difficulty levels (e.g., more accessible topics like Precalculus and Algebra). The inclusion of only
470 20% visual problems limits the assessment of LLMs’ capabilities in visual mathematical reasoning.
471 Additionally, relying on LLMs for problem-solving and solution evaluation introduces potential
472 biases and inaccuracies, as models may struggle with complex derivations or misinterpret instructions,
473 as evidenced by our findings with the μ -MATH dataset. The μ -MATH meta-evaluation dataset,
474 while valuable, encompasses only about 30% of the U-MATH samples and targets exclusively non-
475 visual problems by its design. Moreover, since we use a single model to generate solutions that are
476 then labeled, the benchmark may be less diverse and representative. These factors may limit the
477 exhaustiveness of our assessment of LLMs’ evaluation capabilities and the accuracy of judge’s error
478 rates for our primary benchmark.

479 **Future Work.** Future research can focus on enhancing LLM performance by integrating existing
480 tool-augmented models and exploring their effectiveness on U-MATH and μ -MATH tasks. For
481 instance, incorporating external tools, such as formal solvers, could improve complex textual and
482 multimodal reasoning capabilities. Additionally, our findings indicate that widely used models
483 like GPT-4o are not a silver bullet for solution evaluation; thus, developing specialized (finetuned)
484 models or techniques for more accurate and unbiased assessment is a promising direction. Expanding
485

486 the μ -MATH dataset and incorporating formal verification methods could further refine evaluation
487 processes. Moreover, exploring mathematical assessment as a skill in its own right by examining the
488 generalization of our findings to other benchmarks could lead to significant advancements.
489

490 By open-sourcing U-MATH, μ -MATH, and the evaluation code, we aim to facilitate further research
491 in advancing the mathematical reasoning capabilities of LLMs and encourage the development of
492 models better equipped to tackle complex, real-world mathematical problems.
493

494 ETHICS STATEMENT

495 We collected all data in U-MATH and μ -MATH with appropriate permissions, ensuring no personal
496 or proprietary information is included. The datasets consist solely of mathematical problems and
497 solutions, without any sensitive content. The annotators from [ANONYMIZED] are employed in
498 the partner laboratory with [ANONYMIZED]; their annotation time is fully compensated at a fair
499 hourly rate. We open-sourced the datasets and code under suitable licenses to support transparency
500 and research advancement. There are no conflicts of interest associated with this work.
501
502

503 REPRODUCIBILITY STATEMENT

504 All datasets and code will be available on GitHub. Detailed descriptions of dataset collection and
505 processing are in Section 3. The experimental setup, including model configurations and prompts, is
506 outlined in Section 4, with full prompts provided in Appendices C.1 and C.2. These resources enable
507 replication of our experiments.
508
509

510 REFERENCES

- 511
512
513 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
514 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical
515 report. *arXiv preprint arXiv:2303.08774*.
516
- 517 Janice Ahn, Rishu Verma, Renze Lou, Di Liu, Rui Zhang, and Wenpeng Yin. 2024. Large language
518 models for mathematical reasoning: Progresses and challenges.
519
- 520 Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh
521 Hajishirzi. 2019. MathQA: Towards interpretable math word problem solving with operation-
522 based formalisms. In *Proceedings of the 2019 Conference of the North American Chapter of*
523 *the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long*
524 *and Short Papers)*, pages 2357–2367, Minneapolis, Minnesota. Association for Computational
525 Linguistics.
- 526 Zhangir Azerbayev, Bartosz Piotrowski, Hailey Schoelkopf, Edward W Ayers, Dragomir Radev,
527 and Jeremy Avigad. 2023. Proofnet: Autoformalizing and formally proving undergraduate-level
528 mathematics. *arXiv preprint arXiv:2302.12433*.
529
- 530 Edward Beeching, Shengyi Costa Huang, Albert Jiang, Jia Li, Benjamin Lipkin, Zihan Qina, Kashif
531 Rasul, Ziju Shen, Roman Soletskyi, and Lewis Tunstall. 2024. Numinamath 7b cot. <https://huggingface.co/AI-M0/NuminaMath-7B-CoT>.
532
- 533 Jiaqi Chen, Tong Li, Jinghui Qin, Pan Lu, Liang Lin, Chongyu Chen, and Xiaodan Liang. 2022a.
534 UniGeo: Unifying geometry logical reasoning via reformulating mathematical expression. In
535 *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages
536 3313–3323, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
537
- 538 Jiaqi Chen, Jianheng Tang, Jinghui Qin, Xiaodan Liang, Lingbo Liu, Eric P. Xing, and Liang
539 Lin. 2022b. Geoqa: A geometric question answering benchmark towards multimodal numerical
reasoning.

- 540 Nuo Chen, Ning Wu, Jianhui Chang, and Jia Li. 2024. Controlmath: Controllable data generation
541 promotes math generalist models.
542
- 543 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
544 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to
545 solve math word problems. *arXiv preprint arXiv:2110.14168*.
- 546 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
547 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of
548 models. *arXiv preprint arXiv:2407.21783*.
- 549
- 550 Meng Fang, Xiangpeng Wan, Fei Lu, Fei Xing, and Kai Zou. 2024. Mathodyssey: Benchmarking
551 mathematical problem-solving skills in large language models using odyssey math data. *arXiv*
552 *preprint arXiv:2406.18321*.
- 553
- 554 Simon Frieder, Luca Pinchetti, Ryan-Rhys Griffiths, Tommaso Salvatori, Thomas Lukasiewicz,
555 Philipp Petersen, and Julius Berner. 2024. Mathematical capabilities of chatgpt. *Advances in*
556 *neural information processing systems*, 36.
- 557
- 558 Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and
559 Graham Neubig. 2023. Pal: Program-aided language models. In *International Conference on*
560 *Machine Learning*, pages 10764–10799. PMLR.
- 561 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Ja-
562 cob Steinhardt. 2020. Measuring massive multitask language understanding. *arXiv preprint*
563 *arXiv:2009.03300*.
- 564 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
565 and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *arXiv*
566 *preprint arXiv:2103.03874*.
- 567
- 568 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
569 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,
570 L el io Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas
571 Wang, Timoth ee Lacroix, and William El Sayed. 2023. Mistral 7b.
- 572
- 573 Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay
574 Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam
575 Neyshabur, Guy Gur-Ari, and Vedant Misra. 2022. Solving quantitative reasoning problems with
576 language models.
- 577
- 578 Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li,
579 Ziwei Liu, and Chunyuan Li. 2024a. Llava-onevision: Easy visual task transfer.
- 580
- 581 Chen Li, Weiqi Wang, Jingcheng Hu, Yixuan Wei, Nanning Zheng, Han Hu, Zheng Zhang, and
582 Houwen Peng. 2024b. Common 7b language models already possess strong math capabilities.
- 583
- 584 Wangyue Li, Liangzhi Li, Tong Xiang, Xiao Liu, Wei Deng, and Noa Garcia. 2024c. Can multiple-
585 choice questions really be useful in detecting the abilities of llms?
- 586
- 587 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike,
588 John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step. *arXiv preprint*
589 *arXiv:2305.20050*.
- 590
- 591 Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021.
592 Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language
593 processing.
- 594
- 595 Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng,
596 Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2023. Mathvista: Evaluating mathematical
597 reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*.

- 594 Pan Lu, Ran Gong, Shibiao Jiang, Liang Qiu, Siyuan Huang, Xiaodan Liang, and Song-Chun Zhu.
595 2021. Inter-gps: Interpretable geometry problem solving with formal language and symbolic
596 reasoning.
- 597 Yujun Mao, Yoon Kim, and Yilun Zhou. 2024. Champ: A competition-level dataset for fine-grained
598 analyses of llms’ mathematical reasoning capabilities. *arXiv preprint arXiv:2401.06961*.
- 600 Mistral AI. 2024. Announcing pixtral-12b. <https://mistral.ai/news/pixtral-12b/>. Accessed:
601 2024-10-01.
- 602 Mistral.ai. 2024. Mathstral. <https://mistral.ai/news/mathstral/>. Accessed: 2024-10-01.
- 603 OpenAI. 2024. Hello gpt-4o. <https://openai.com/index/hello-gpt-4o/>. Accessed: 2024-10-
604 01.
- 605
606 Pouya Pezeshkpour and Estevam Hruschka. 2023. Large language models sensitivity to the order of
607 options in multiple-choice questions.
- 608
609 Nikhil Prakash, Tamar Rott Shaham, Tal Haklay, Yonatan Belinkov, and David Bau. 2024. Fine-tuning
610 enhances existing mechanisms: A case study on entity tracking.
- 611
612 Runqi Qiao, Qiuna Tan, Guanting Dong, Minhui Wu, Chong Sun, Xiaoshuai Song, Zhuoma GongQue,
613 Shanglin Lei, Zhe Wei, Miaoxuan Zhang, et al. 2024. We-math: Does your large multimodal
614 model achieve human-like mathematical reasoning? *arXiv preprint arXiv:2407.01284*.
- 615
616 Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi,
617 Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, et al. 2022. Language models are
618 multilingual chain-of-thought reasoners. *arXiv preprint arXiv:2210.03057*.
- 619
620 Andreas Stephan, Dawei Zhu, Matthias Aßenmacher, Xiaoyu Shen, and Benjamin Roth. 2024. From
621 calculation to adjudication: Examining llm judges on mathematical reasoning tasks.
- 622
623 Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer,
624 Damien Vincent, Zhufeng Pan, Shibo Wang, Soroosh Mariooryad, Yifan Ding, Xinyang Geng, Fred
625 Alcober, Roy Frostig, Mark Omernick, Lexi Walker, Cosmin Paduraru, Christina Sorokin, Andrea
626 Tacchetti, Colin Gaffney, Samira Daruki, Olcan Sercinoglu, Zach Gleicher, Juliette Love, Paul
627 Voigtlaender, Rohan Jain, Gabriela Surita, Kareem Mohamed, Rory Blevins, Junwhan Ahn, Tao
628 Zhu, Kornrathop Kawintiranon, Orhan Firat, Yiming Gu, Yujing Zhang, Matthew Rahtz, Manaal
629 Faruqui, Natalie Clay, Justin Gilmer, JD Co-Reyes, Ivo Penchev, Rui Zhu, Nobuyuki Morioka,
630 Kevin Hui, Krishna Haridasan, Victor Campos, Mahdis Mahdieh, Mandy Guo, Samer Hassan,
631 Kevin Kilgour, Arpi Vezer, Heng-Tze Cheng, Raoul de Liedekerke, Siddharth Goyal, Paul Barham,
632 DJ Strouse, Seb Noury, Jonas Adler, Mukund Sundararajan, Sharad Vikram, Dmitry Lepikhin,
633 Michela Paganini, Xavier Garcia, Fan Yang, Dasha Valter, Maja Trebacz, Kiran Vodrahalli,
634 Chulayuth Asawaroengchai, Roman Ring, Norbert Kalb, Livio Baldini Soares, Siddhartha Brahma,
635 David Steiner, Tianhe Yu, Fabian Mentzer, Antoine He, Lucas Gonzalez, Bibo Xu, Raphael Lopez
636 Kaufman, Laurent El Shafey, Junhyuk Oh, Tom Hennigan, George van den Driessche, Seth Odoom,
637 Mario Lucic, Becca Roelofs, Sid Lall, Amit Marathe, Betty Chan, Santiago Ontanon, Luheng He,
638 Denis Teplyashin, Jonathan Lai, Phil Crone, Bogdan Damoc, Lewis Ho, Sebastian Riedel, Karel
639 Lenc, Chih-Kuan Yeh, Aakanksha Chowdhery, Yang Xu, Mehran Kazemi, Ehsan Amid, Anastasia
640 Petrushkina, Kevin Swersky, Ali Khodaei, Gowoon Chen, Chris Larkin, Mario Pinto, Geng Yan,
641 Adria Puigdomenech Badia, Piyush Patil, Steven Hansen, Dave Orr, Sebastien M. R. Arnold,
642 Jordan Grimstad, Andrew Dai, Sholto Douglas, Rishika Sinha, Vikas Yadav, Xi Chen, Elena
643 Gribovskaya, Jacob Austin, Jeffrey Zhao, Kaushal Patel, Paul Komarek, Sophia Austin, Sebastian
644 Borgeaud, Linda Friso, Abhimanyu Goyal, Ben Caine, Kris Cao, Da-Woon Chung, Matthew
645 Lamm, Gabe Barth-Maron, Thais Kagohara, Kate Olszewska, Mia Chen, Kaushik Shivakumar,
646 Rishabh Agarwal, Harshal Godhia, Ravi Rajwar, Javier Snaider, Xerxes Dotiwalla, Yuan Liu,
647 Aditya Barua, Victor Ungureanu, Yuan Zhang, Bat-Orgil Batsaikhan, Mateo Wirth, James Qin, Ivo
648 Danihelka, Tulsee Doshi, Martin Chadwick, Jilin Chen, Sanil Jain, Quoc Le, Arjun Kar, Madhu
649 Gurumurthy, Cheng Li, Ruoxin Sang, Fangyu Liu, Lampros Lamprou, Rich Munoz, Nathan Lintz,
650 Harsh Mehta, Heidi Howard, Malcolm Reynolds, Lora Aroyo, Quan Wang, Lorenzo Blanco, Albin
651 Cassirer, Jordan Griffith, Dipanjan Das, Stephan Lee, Jakub Sygnowski, Zach Fisher, James Besley,
652 Richard Powell, Zafarali Ahmed, Dominik Paulus, David Reitter, Zalan Borsos, Rishabh Joshi,

648 Aedan Pope, Steven Hand, Vittorio Selo, Vihan Jain, Nikhil Sethi, Megha Goel, Takaki Makino,
649 Rhys May, Zhen Yang, Johan Schalkwyk, Christina Butterfield, Anja Hauth, Alex Goldin, Will
650 Hawkins, Evan Senter, Sergey Brin, Oliver Woodman, Marvin Ritter, Eric Noland, Minh Giang,
651 Vijay Bolina, Lisa Lee, Tim Blyth, Ian Mackinnon, Machel Reid, Obaid Sarvana, David Silver,
652 Alexander Chen, Lily Wang, Loren Maggiore, Oscar Chang, Nithya Attaluri, Gregory Thornton,
653 Chung-Cheng Chiu, Oskar Bunyan, Nir Levine, Timothy Chung, Evgenii Eltyshev, Xiance Si,
654 Timothy Lillicrap, Demetra Brady, Vaibhav Aggarwal, Boxi Wu, Yuanzhong Xu, Ross McIlroy,
655 Kartikeya Badola, Paramjit Sandhu, Erica Moreira, Wojciech Stokowiec, Ross Hemsley, Dong
656 Li, Alex Tudor, Pranav Shyam, Elahe Rahimtoroghi, Salem Haykal, Pablo Sprechmann, Xiang
657 Zhou, Diana Mincu, Yujia Li, Ravi Addanki, Kalpesh Krishna, Xiao Wu, Alexandre Frechette,
658 Matan Eyal, Allan Dafoe, Dave Lacey, Jay Whang, Thi Avrahami, Ye Zhang, Emanuel Taropa,
659 Hanzhao Lin, Daniel Toyama, Eliza Rutherford, Motoki Sano, HyunJeong Choe, Alex Tomala,
660 Chalence Safranek-Shrader, Nora Kassner, Mantas Pajarskas, Matt Harvey, Sean Sechrist, Meire
661 Fortunato, Christina Lyu, Gamaleldin Elsayed, Chenkai Kuang, James Lottes, Eric Chu, Chao Jia,
662 Chih-Wei Chen, Peter Humphreys, Kate Baumli, Connie Tao, Rajkumar Samuel, Cicero Nogueira
663 dos Santos, Anders Andreassen, Nemanja Rakićević, Dominik Grewe, Aviral Kumar, Stephanie
664 Winkler, Jonathan Caton, Andrew Brock, Sid Dalmia, Hannah Sheahan, Iain Barr, Yingjie Miao,
665 Paul Natsev, Jacob Devlin, Feryal Behbahani, Flavien Prost, Yanhua Sun, Artiom Myaskovsky,
666 Thanumalayan Sankaranarayanan Pillai, Dan Hurt, Angeliki Lazaridou, Xi Xiong, Ce Zheng, Fabio
667 Pardo, Xiaowei Li, Dan Horgan, Joe Stanton, Moran Ambar, Fei Xia, Alejandro Lince, Mingqiu
668 Wang, Basil Mustafa, Albert Webson, Hyo Lee, Rohan Anil, Martin Wicke, Timothy Dozat,
669 Abhishek Sinha, Enrique Piqueras, Elahe Dabir, Shyam Upadhyay, Anudhyan Boral, Lisa Anne
670 Hendricks, Corey Fry, Josip Djolonga, Yi Su, Jake Walker, Jane Labanowski, Ronny Huang, Vedant
671 Misra, Jeremy Chen, RJ Skerry-Ryan, Avi Singh, Shruti Rijhwani, Dian Yu, Alex Castro-Ros,
672 Beer Changpinyo, Romina Datta, Sumit Bagri, Arnar Mar Hrafnkelsson, Marcello Maggioni,
673 Daniel Zheng, Yury Sulsky, Shaobo Hou, Tom Le Paine, Antoine Yang, Jason Riesa, Dominika
674 Rogozinska, Dror Marcus, Dalia El Badawy, Qiao Zhang, Luyu Wang, Helen Miller, Jeremy
675 Greer, Lars Lowe Sjos, Azade Nova, Heiga Zen, Rahma Chaabouni, Mihaela Rosca, Jiepu Jiang,
676 Charlie Chen, Ruibo Liu, Tara Sainath, Maxim Krikun, Alex Polozov, Jean-Baptiste Lespiau,
677 Josh Newlan, Zeynep Cankara, Soo Kwak, Yunhan Xu, Phil Chen, Andy Coenen, Clemens
678 Meyer, Katerina Tsihlias, Ada Ma, Juraj Gottweis, Jinwei Xing, Chenjie Gu, Jin Miao, Christian
679 Frank, Zeynep Cankara, Sanjay Ganapathy, Ishita Dasgupta, Steph Hughes-Fitt, Heng Chen,
680 David Reid, Keran Rong, Hongmin Fan, Joost van Amersfoort, Vincent Zhuang, Aaron Cohen,
681 Shixiang Shane Gu, Anhad Mohananey, Anastasija Ilic, Taylor Tobin, John Wieting, Anna Bortsova,
682 Phoebe Thacker, Emma Wang, Emily Caveness, Justin Chiu, Eren Sezener, Alex Kaskasoli,
683 Steven Baker, Katie Millican, Mohamed Elhawaty, Kostas Aisopos, Carl Lebsack, Nathan Byrd,
684 Hanjun Dai, Wenhao Jia, Matthew Wiethoff, Elnaz Davoodi, Albert Weston, Lakshman Yagati,
685 Arun Ahuja, Isabel Gao, Golan Pundak, Susan Zhang, Michael Azzam, Khe Chai Sim, Sergi
686 Caelles, James Keeling, Abhanshu Sharma, Andy Swing, YaGuang Li, Chenxi Liu, Carrie Grimes
687 Bostock, Yamini Bansal, Zachary Nado, Ankesh Anand, Josh Lipschultz, Abhijit Karmarkar,
688 Lev Proleev, Abe Ittycheriah, Soheil Hassas Yeganeh, George Polovets, Aleksandra Faust, Jiao
689 Sun, Alban Rustemi, Pen Li, Rakesh Shivanna, Jeremiah Liu, Chris Welty, Federico Lebron,
690 Anirudh Baddepudi, Sebastian Krause, Emilio Parisotto, Radu Soricut, Zheng Xu, Dawn Bloxwich,
691 Melvin Johnson, Behnam Neyshabur, Justin Mao-Jones, Renshen Wang, Vinay Ramasesh, Zaheer
692 Abbas, Arthur Guez, Constant Segal, Duc Dung Nguyen, James Svensson, Le Hou, Sarah York,
693 Kieran Milan, Sophie Bridgers, Wiktor Gworek, Marco Tagliasacchi, James Lee-Thorp, Michael
694 Chang, Alexey Guseynov, Ale Jakse Hartman, Michael Kwong, Ruizhe Zhao, Sheleem Kashem,
695 Elizabeth Cole, Antoine Miech, Richard Tanburn, Mary Phuong, Filip Pavetic, Sebastien Cevy,
696 Ramona Comanescu, Richard Ives, Sherry Yang, Cosmo Du, Bo Li, Zizhao Zhang, Mariko Iinuma,
697 Clara Huiyi Hu, Aurko Roy, Shaan Bijwadia, Zhenkai Zhu, Danilo Martins, Rachel Saputro, Anita
698 Gergely, Steven Zheng, Dawei Jia, Ioannis Antonoglou, Adam Sadovsky, Shane Gu, Yingying
699 Bi, Alek Andreev, Sina Samangooei, Mina Khan, Tomas Kocisky, Angelos Filos, Chintu Kumar,
700 Colton Bishop, Adams Yu, Sarah Hodkinson, Sid Mittal, Premal Shah, Alexandre Moufarek, Yong
701 Cheng, Adam Bloniarz, Jaehoon Lee, Pedram Pejman, Paul Michel, Stephen Spencer, Vladimir
Feinberg, Xuehan Xiong, Nikolay Savinov, Charlotte Smith, Siamak Shakeri, Dustin Tran, Mary
Chesus, Bernd Bohnet, George Tucker, Tamara von Glehn, Carrie Muir, Yiran Mao, Hideto Kazawa,
Ambrose Slone, Kedar Soparkar, Disha Shrivastava, James Cobon-Kerr, Michael Sharman, Jay
Pavagadhi, Carlos Araya, Karolis Misiunas, Nimesh Ghelani, Michael Laskin, David Barker,
Qiujia Li, Anton Briukhov, Neil Houlsby, Mia Glaese, Balaji Lakshminarayanan, Nathan Schucher,

702 Yunhao Tang, Eli Collins, Hyeontaek Lim, Fangxiaoyu Feng, Adria Recasens, Guangda Lai,
703 Alberto Magni, Nicola De Cao, Aditya Siddhant, Zoe Ashwood, Jordi Orbay, Mostafa Dehghani,
704 Jenny Brennan, Yifan He, Kelvin Xu, Yang Gao, Carl Saroufim, James Molloy, Xinyi Wu, Seb
705 Arnold, Solomon Chang, Julian Schrittwieser, Elena Buchatskaya, Soroush Radpour, Martin
706 Polacek, Skye Giordano, Ankur Bapna, Simon Tokumine, Vincent Hellendoorn, Thibault Sottiaux,
707 Sarah Cogan, Aliaksei Severyn, Mohammad Saleh, Shantanu Thakoor, Laurent Shefey, Siyuan
708 Qiao, Meenu Gaba, Shuo yiin Chang, Craig Swanson, Biao Zhang, Benjamin Lee, Paul Kishan
709 Rubenstein, Gan Song, Tom Kwiatkowski, Anna Koop, Ajay Kannan, David Kao, Parker Schuh,
710 Axel Stjerngren, Golnaz Ghiasi, Gena Gibson, Luke Vilnis, Ye Yuan, Felipe Tiengo Ferreira,
711 Aishwarya Kamath, Ted Klimentko, Ken Franko, Kefan Xiao, Indro Bhattacharya, Miteyan Patel,
712 Rui Wang, Alex Morris, Robin Strudel, Vivek Sharma, Peter Choy, Sayed Hadi Hashemi, Jessica
713 Landon, Mara Finkelstein, Priya Jhakra, Justin Frye, Megan Barnes, Matthew Mauger, Dennis
714 Daun, Khuslen Baatarsukh, Matthew Tung, Wael Farhan, Henryk Michalewski, Fabio Viola, Felix
715 de Chaumont Quitry, Charline Le Lan, Tom Hudson, Qingze Wang, Felix Fischer, Ivy Zheng,
716 Elspeth White, Anca Dragan, Jean baptiste Alayrac, Eric Ni, Alexander Pritzel, Adam Iwanicki,
717 Michael Isard, Anna Bulanova, Lukas Zilka, Ethan Dyer, Devendra Sachan, Srivatsan Srinivasan,
718 Hannah Muckenhirn, Honglong Cai, Amol Mandhane, Mukarram Tariq, Jack W. Rae, Gary Wang,
719 Kareem Ayoub, Nicholas FitzGerald, Yao Zhao, Woohyun Han, Chris Alberti, Dan Garrette,
720 Kashyap Krishnakumar, Mai Gimenez, Anselm Levskaya, Daniel Sohn, Josip Matak, Inaki Iturrate,
721 Michael B. Chang, Jackie Xiang, Yuan Cao, Nishant Ranka, Geoff Brown, Adrian Hutter, Vahab
722 Mirrokni, Nanxin Chen, Kaisheng Yao, Zoltan Egyed, Francois Galilee, Tyler Liechty, Praveen
723 Kallakuri, Evan Palmer, Sanjay Ghemawat, Jasmine Liu, David Tao, Chloe Thornton, Tim Green,
724 Mimi Jasarevic, Sharon Lin, Victor Cotruta, Yi-Xuan Tan, Noah Fiedel, Hongkun Yu, Ed Chi,
725 Alexander Neitz, Jens Heitkaemper, Anu Sinha, Denny Zhou, Yi Sun, Charbel Kaed, Brice Hulse,
726 Swaroop Mishra, Maria Georgaki, Sneha Kudugunta, Clement Farabet, Izhak Shafran, Daniel
727 Vlasic, Anton Tsitsulin, Rajagopal Ananthanarayanan, Alen Carin, Guolong Su, Pei Sun, Shashank
728 V, Gabriel Carvajal, Josef Broder, Iulia Comsa, Alena Repina, William Wong, Warren Weilun Chen,
729 Peter Hawkins, Egor Filonov, Lucia Loher, Christoph Hirsenschall, Weiyi Wang, Jingchen Ye, Andrea
730 Burns, Hardie Cate, Diana Gage Wright, Federico Piccinini, Lei Zhang, Chu-Cheng Lin, Ionel
731 Gog, Yana Kulizhskaya, Ashwin Sreevatsa, Shuang Song, Luis C. Cobo, Anand Iyer, Chetan Tekur,
732 Guillermo Garrido, Zhu Yun Xiao, Rupert Kemp, Huaixiu Steven Zheng, Hui Li, Ananth Agarwal,
733 Christel Ngani, Kati Goshvadi, Rebeca Santamaria-Fernandez, Wojciech Fica, Xinyun Chen,
734 Chris Gorgolewski, Sean Sun, Roopal Garg, Xinyu Ye, S. M. Ali Eslami, Nan Hua, Jon Simon,
735 Pratik Joshi, Yelin Kim, Ian Tenney, Sahitya Potluri, Lam Nguyen Thiet, Quan Yuan, Florian
736 Luisier, Alexandra Chronopoulou, Salvatore Scellato, Praveen Srinivasan, Minmin Chen, Vinod
737 Koverkathu, Valentin Dalibard, Yaming Xu, Brennan Saeta, Keith Anderson, Thibault Sellam,
738 Nick Fernando, Fantine Huot, Junehyuk Jung, Mani Varadarajan, Michael Quinn, Amit Raul,
739 Maigo Le, Ruslan Habalov, Jon Clark, Komal Jalan, Kalesha Bullard, Achintya Singhal, Thang
740 Luong, Boyu Wang, Sujeevan Rajayogam, Julian Eisenschlos, Johnson Jia, Daniel Finchelstein,
741 Alex Yakubovich, Daniel Balle, Michael Fink, Sameer Agarwal, Jing Li, Dj Dvijotham, Shalini
742 Pal, Kai Kang, Jaclyn Konzelmann, Jennifer Beattie, Olivier Dousse, Diane Wu, Remi Crocker,
743 Chen Elkind, Siddhartha Reddy Jonnalagadda, Jong Lee, Dan Holtmann-Rice, Krystal Kallarackal,
744 Rosanne Liu, Denis Vnukov, Neera Vats, Luca Invernizzi, Mohsen Jafari, Huanjie Zhou, Lilly
745 Taylor, Jennifer Prendki, Marcus Wu, Tom Eccles, Tianqi Liu, Kavya Koppurapu, Francoise
746 Beaufays, Christof Angermueller, Andreea Marzoca, Shourya Sarcar, Hilal Dib, Jeff Stanway,
747 Frank Perbet, Nejc Trdin, Rachel Sterneck, Andrey Khorlin, Dinghua Li, Xihui Wu, Sonam
748 Goenka, David Madras, Sasha Goldshtein, Willi Gierke, Tong Zhou, Yaxin Liu, Yannie Liang,
749 Anais White, Yunjie Li, Shreya Singh, Sanaz Bahargam, Mark Epstein, Sujoy Basu, Li Lao,
750 Adnan Oztirel, Carl Crous, Alex Zhai, Han Lu, Zora Tung, Neeraj Gaur, Alanna Walton, Lucas
751 Dixon, Ming Zhang, Amir Globerson, Grant Uy, Andrew Bolt, Olivia Wiles, Milad Nasr, Iliia
752 Shumailov, Marco Selvi, Francesco Piccinno, Ricardo Aguilar, Sara McCarthy, Misha Khalman,
753 Mrinal Shukla, Vlado Galic, John Carpenter, Kevin Vellela, Haibin Zhang, Harry Richardson,
754 James Martens, Matko Bosnjak, Shreyas Rammohan Belle, Jeff Seibert, Mahmoud Alnahlawi,
755 Brian McWilliams, Sankalp Singh, Annie Louis, Wen Ding, Dan Popovici, Lenin Simicich, Laura
Knight, Pulkit Mehta, Nishesh Gupta, Chongyang Shi, Saaber Fatehi, Jovana Mitrovic, Alex
Grills, Joseph Pagadora, Dessie Petrova, Danielle Eisenbud, Zhishuai Zhang, Damion Yates,
Bhavishya Mittal, Nilesh Tripuraneni, Yannis Assael, Thomas Brovelli, Prateek Jain, Mihajlo
Velimirovic, Canfer Akbulut, Jiaqi Mu, Wolfgang Macherey, Ravin Kumar, Jun Xu, Haroon
Qureshi, Gheorghe Comanici, Jeremy Wiesner, Zhitao Gong, Anton Ruddock, Matthias Bauer,

- 756 Nick Felt, Anirudh GP, Anurag Arnab, Dustin Zelle, Jonas Rothfuss, Bill Rosgen, Ashish Shenoy,
757 Bryan Seybold, Xinjian Li, Jayaram Mudigonda, Goker Erdogan, Jiawei Xia, Jiri Simsa, Andrea
758 Michi, Yi Yao, Christopher Yew, Steven Kan, Isaac Caswell, Carey Radebaugh, Andre Elisseeff,
759 Pedro Valenzuela, Kay McKinney, Kim Paterson, Albert Cui, Eri Latorre-Chimoto, Solomon Kim,
760 William Zeng, Ken Durden, Priya Ponnappalli, Tiberiu Sosea, Christopher A. Choquette-Choo,
761 James Manyika, Brona Robenek, Harsha Vashisht, Sebastien Pereira, Hoi Lam, Marko Velic,
762 Denese Owusu-Afriyie, Katherine Lee, Tolga Bolukbasi, Alicia Parrish, Shawn Lu, Jane Park,
763 Balaji Venkatraman, Alice Talbert, Lambert Rosique, Yuchung Cheng, Andrei Sozanschi, Adam
764 Paszke, Praveen Kumar, Jessica Austin, Lu Li, Khalid Salama, Wooyeol Kim, Nandita Dukkupati,
765 Anthony Baryshnikov, Christos Kaplanis, XiangHai Sheng, Yuri Chervonyi, Caglar Unlu, Diego
766 de Las Casas, Harry Askham, Kathryn Tunyasuvunakool, Felix Gimeno, Siim Poder, Chester
767 Kwak, Matt Miecznikowski, Vahab Mirrokni, Alek Dimitriev, Aaron Parisi, Dangyi Liu, Tomy Tsai,
768 Toby Shevlane, Christina Kouridi, Drew Garmon, Adrian Goedeckemeyer, Adam R. Brown, Anitha
769 Vijayakumar, Ali Elqursh, Sadegh Jazayeri, Jin Huang, Sara Mc Carthy, Jay Hoover, Lucy Kim,
770 Sandeep Kumar, Wei Chen, Courtney Biles, Garrett Bingham, Evan Rosen, Lisa Wang, Qijun Tan,
771 David Engel, Francesco Pongetti, Dario de Cesare, Dongseong Hwang, Lily Yu, Jennifer Pullman,
772 Srini Narayanan, Kyle Levin, Siddharth Gopal, Megan Li, Asaf Aharoni, Trieu Trinh, Jessica
773 Lo, Norman Casagrande, Roopali Vij, Loic Matthey, Bramandia Ramadhana, Austin Matthews,
774 CJ Carey, Matthew Johnson, Kremena Goranova, Rohin Shah, Shereen Ashraf, Kingshuk Dasgupta,
775 Rasmus Larsen, Yicheng Wang, Manish Reddy Vuyyuru, Chong Jiang, Joana Ijazi, Kazuki Osawa,
776 Celine Smith, Ramya Sree Boppana, Taylan Bilal, Yuma Koizumi, Ying Xu, Yasemin Altun, Nir
777 Shabat, Ben Bariach, Alex Korchemniy, Kiam Choo, Olaf Ronneberger, Chimezie Iwuanyanwu,
778 Shubin Zhao, David Soergel, Cho-Jui Hsieh, Irene Cai, Shariq Iqbal, Martin Sundermeyer, Zhe
779 Chen, Elie Bursztein, Chaitanya Malaviya, Fadi Biadisy, Prakash Shroff, Inderjit Dhillon, Tejasi
780 Latkar, Chris Dyer, Hannah Forbes, Massimo Nicosia, Vitaly Nikolaev, Somer Greene, Marin
781 Georgiev, Pidong Wang, Nina Martin, Hanie Sedghi, John Zhang, Praseem Banzal, Doug Fritz,
782 Vikram Rao, Xuezhi Wang, Jiageng Zhang, Viorica Patraucean, Dayou Du, Igor Mordatch, Ivan
783 Jurin, Lewis Liu, Ayush Dubey, Abhi Mohan, Janek Nowakowski, Vlad-Doru Ion, Nan Wei, Reiko
784 Tojo, Maria Abi Raad, Drew A. Hudson, Vaishakh Keshava, Shubham Agrawal, Kevin Ramirez,
785 Zhichun Wu, Hoang Nguyen, Ji Liu, Madhavi Sewak, Bryce Pettrini, DongHyun Choi, Ivan Philips,
786 Ziyue Wang, Ioana Bica, Ankush Garg, Jarek Wilkiewicz, Priyanka Agrawal, Xiaowei Li, Danhao
787 Guo, Emily Xue, Naseer Shaik, Andrew Leach, Sadh MNM Khan, Julia Wiesinger, Sammy Jerome,
788 Abhishek Chakladar, Alek Wenjiao Wang, Tina Ornduff, Folake Abu, Alireza Ghaffarkhah, Marcus
789 Wainwright, Mario Cortes, Frederick Liu, Joshua Maynez, Andreas Terzis, Pouya Samangouei,
790 Riham Mansour, Tomasz Kępa, François-Xavier Aubet, Anton Algymer, Dan Banica, Agoston
791 Weisz, Andras Orban, Alexandre Senges, Ewa Andrejczuk, Mark Geller, Niccolo Dal Santo,
792 Valentin Anklin, Majd Al Meray, Martin Baeuml, Trevor Strohman, Junwen Bai, Slav Petrov,
793 Yonghui Wu, Demis Hassabis, Koray Kavukcuoglu, Jeffrey Dean, and Oriol Vinyals. 2024. Gemini
794 1.5: Unlocking multimodal understanding across millions of tokens of context.
- 795 Qwen Team. 2024. Qwen2.5: A party of foundation models.
- 796 Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Mingjie Zhan, and Hongsheng Li. 2024. Measuring
797 multimodal mathematical reasoning with math-vision dataset. *arXiv preprint arXiv:2402.14804*.
- 798 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-
799 ery, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language
800 models. *arXiv preprint arXiv:2203.11171*.
- 801 Yan Wang, Xiaojiang Liu, and Shuming Shi. 2017. Deep neural solver for math word problems.
802 In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*,
803 pages 845–854, Copenhagen, Denmark. Association for Computational Linguistics.
- 804 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
805 Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances
806 in neural information processing systems*, 35:24824–24837.
- 807 Shijie Xia, Xuefeng Li, Yixin Liu, Tongshuang Wu, and Pengfei Liu. 2024. Evaluating mathematical
808 reasoning beyond accuracy. *arXiv preprint arXiv:2404.05692*.

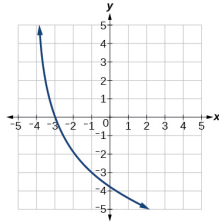
- 810 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,
811 Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang,
812 Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren
813 Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang,
814 Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin,
815 Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong
816 Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu,
817 Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang,
818 Zhifang Guo, and Zhihao Fan. 2024a. Qwen2 technical report.
- 819 An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu,
820 Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu,
821 Xingzhang Ren, and Zhenru Zhang. 2024b. Qwen2.5-math technical report: Toward mathematical
822 expert model via self-improvement.
- 823 Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu
824 Jiang, Weiming Ren, Yuxuan Sun, et al. 2023. Mmmu: A massive multi-discipline multimodal un-
825 derstanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference*
826 *on Computer Vision and Pattern Recognition*, pages 9556–9567.
- 827 Zhongshen Zeng, Pengguang Chen, Shu Liu, Haiyun Jiang, and Jiaya Jia. 2023. Mr-gsm8k: A
828 meta-reasoning benchmark for large language model evaluation. *CoRR*, abs/2312.17080.
- 829 Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan
830 Lu, Kai-Wei Chang, Peng Gao, et al. 2024. Mathverse: Does your multi-modal llm truly see the
831 diagrams in visual math problems? *arXiv preprint arXiv:2403.14624*.
- 832 Yiran Zhao, Jinghan Zhang, I Chern, Siyang Gao, Pengfei Liu, Junxian He, et al. 2024. Felm:
833 Benchmarking factuality evaluation of large language models. *Advances in Neural Information*
834 *Processing Systems*, 36.
- 835 Kunhao Zheng, Jesse Michael Han, and Stanislas Polu. 2021. Minif2f: a cross-system benchmark for
836 formal olympiad-level mathematics. *arXiv preprint arXiv:2109.00110*.
- 837 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
838 Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.
839 2023. Judging llm-as-a-judge with mt-bench and chatbot arena.
- 840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

A PROBLEM EXAMPLES

A.1 U-MATH PROBLEMS

Example 1: Algebra.

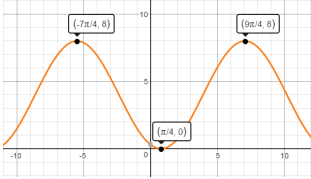
Write a logarithmic equation corresponding to the graph shown. Use $\log_3(x)$ as a parent function:



The final answer: $-3 \cdot \log_3(x + 4)$

Example 3: Precalculus Review.

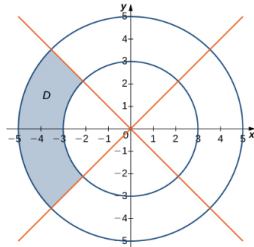
Find a formula for $f(x)$, the sinusoidal function whose graph is shown below:



The final answer:
 $f(x) = -4 \cdot \cos\left(\frac{1}{2} \cdot \left(x - \frac{\pi}{4}\right)\right) + 4$

Example 5: Multivariable Calculus.

The graph of the polar rectangular region D is given. Express the region D in polar coordinates:



1. The interval of r is $[3, 5]$
2. The interval of θ is $\left[\frac{3}{4} \cdot \pi, \frac{5}{4} \cdot \pi\right]$

Example 2: Integral Calculus.

Solve the integral:

$$\int \frac{-9 \cdot \sqrt[3]{x}}{9 \cdot \sqrt[3]{x^2} + 3 \cdot \sqrt{x}} dx$$

$$\begin{aligned} & -\frac{2}{27} \cdot \ln\left(\frac{1}{3} \cdot |1 + 3 \cdot \sqrt[6]{x}|\right) \\ & -\frac{1}{3} \cdot \sqrt[6]{x^2} - \frac{3}{2} \cdot \sqrt[6]{x^4} + \frac{2}{3} \cdot \sqrt[6]{x^3} \\ & + \frac{2}{9} \cdot \sqrt[6]{x} + C \end{aligned}$$

Example 4: Multivariable Calculus.

E is located inside the cylinder $x^2 + y^2 = 1$ and between the circular paraboloids $z = 1 - x^2 - y^2$ and $z = x^2 + y^2$. Find the volume of E .

$$\text{Volume} = \frac{\pi}{4}$$

Example 6: Differential Calculus.

Sketch the curve:

$$y = \frac{x^3}{6 \cdot (x + 3)^2}$$

Provide the following:

1. The domain (in interval notation)
2. Vertical asymptotes
3. Horizontal asymptotes
4. Slant asymptotes
5. Intervals where the function is increasing
6. Intervals where the function is decreasing
7. Intervals where the function is concave up
8. Intervals where the function is concave down
9. Points of inflection

1. The domain: $(-\infty, -3) \cup (-3, \infty)$

Figure 2: Example text-only and visual problems from the U-MATH benchmark, illustrating the topic, problem, and golden answer.

A.2 U-MATH PROBLEM AND SOLUTION

Example: Differential Calculus.

U-MATH Problem:

The function $s(t) = 2 \cdot t^3 - 3 \cdot t^2 - 12 \cdot t + 8$ represents the position of a particle traveling along a horizontal line.

1. Find the velocity and acceleration functions.
2. Determine the time intervals when the object is slowing down or speeding up.

Golden answer:

1. The velocity function $v(t) = 6 \cdot t^2 - 6 \cdot t - 12$ and acceleration function $a(t) = 12 \cdot t - 6$
2. The time intervals when the object speeds up $(0, \frac{1}{2})$, $(2, \infty)$ and slows down $(\frac{1}{2}, 2)$

Reference solution:

The velocity is the derivative of the position function:

$$v(t) = s'(t) = 6t^2 - 6t - 12$$

The acceleration is the derivative of the velocity function:

$$a(t) = v'(t) = 12t - 6$$

To determine when the object is speeding up or slowing down, we compare the signs of $v(t)$ and $a(t)$.

Step 1: Find the Zeros of $v(t)$ and $a(t)$

First, solve for $v(t) = 0$:

$$6t^2 - 6t - 12 = 0 \Rightarrow t^2 - t - 2 = 0 \Rightarrow (t - 2)(t + 1) = 0$$

Thus, $t = 2$ and $t = -1$.

Next, solve for $a(t) = 0$:

$$12t - 6 = 0 \Rightarrow t = \frac{1}{2}$$

Step 2: Analyze the Signs of $v(t)$ and $a(t)$

We analyze the signs of $v(t)$ and $a(t)$ on the intervals determined by $t = -1$, $t = \frac{1}{2}$, and $t = 2$.

Interval	$v(t)$	$a(t)$	Behavior
$(-\infty, -1)$	> 0	< 0	Slowing down
$(-1, \frac{1}{2})$	< 0	< 0	Speeding up
$(\frac{1}{2}, 2)$	< 0	> 0	Slowing down
$(2, \infty)$	> 0	> 0	Speeding up

Step 3: Account for non-negative time

The object is speeding up on $(0, \frac{1}{2})$ and $(2, \infty)$ and slowing down on $(\frac{1}{2}, 2)$.

Figure 3: An example problem from the U-MATH benchmark, illustrating the problem, reference solution and golden answer.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

A.3 μ -MATH META-EVALUATION

Example: Integral Calculus.

U-MATH Problem:
Solve the integral:

$$\int \frac{20 \cdot \cos(-10 \cdot x)^3}{21 \cdot \sin(-10 \cdot x)^7} dx$$

Golden answer:

$$C + \frac{1}{21} \cdot \left(\frac{1}{2} \cdot (\cot(10 \cdot x))^4 + \frac{1}{3} \cdot (\cot(10 \cdot x))^6 \right)$$

LLM-generated answer:

$$-\frac{3 \sin(10x)^2 - 2}{126 \sin(10x)^6} + C$$

Golden judge label: correct

Comment:
The reference answer and the submitted one can be simplified, respectively, to

$$C + \frac{\cot^4(10x)}{42} + \frac{\cot^6(10x)}{63} \quad \text{and} \quad C + \frac{\cot^6(10x)}{63} + \frac{\cot^4(10x)}{42} + \frac{1}{126},$$

which differ by a constant term of $1/126$.

Figure 4: An example problem from the μ -MATH meta-evaluation benchmark, illustrating the comparison between the golden (reference) answer and the answer generated by an LLM.

1026 **B SUB-TOPICS DISTRIBUTION**

1027

1028 The U-MATH dataset cover variety of topics across 6 core subjects. Below is the count of unique
1029 topics per subject:

1030

- 1031 • Differential Calculus: 51 unique topics
- 1032 • Sequences and Series: 28 unique topics
- 1033 • Integral Calculus: 35 unique topics
- 1034 • Precalculus Review: 19 unique topics
- 1035 • Algebra: 74 unique topics
- 1036 • Multivariable Calculus: 53 unique topics

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

	Subject	Topic Count	Topic Name
1080			
1081	Differential Calculus	29	Curve Sketching
1082		13	Limits
1083		12	One-Sided Limits
1084		12	L'Hospital's Rule
1084		11	Increasing and Decreasing Functions
1085		11	Higher Derivatives
1085		10	Applications of Derivatives (Local Extrema)
1086		10	Concavity
1087		9	Product Rule
1087		7	Critical Numbers
1088	Sequences and Series	40	Taylor Series
1089		30	Fourier Series
1090		18	Maclaurin Series
1090		12	Approximating Constants Using Power Series
1091		6	Radius of Convergence (Center of Convergence)
1092		5	Differentiate Power Series
1092		4	Error in Approximation
1093		4	Approximating Integrals Using Power Series
1094		3	Series
1094		3	Sum of Numerical Series
1095	Integral Calculus	83	The Substitution Rule
1096		24	Antiderivatives
1096		10	Volumes of Solids of Revolution About the X-Axis
1097		9	Trigonometric Substitutions and Inverse Substitutions
1098		9	Integrate Respect Independent Variable
1099		7	Applications of Integrals
1099		7	Single Variable Surface Area Integrals
1100		6	Volume of Solids of Revolution About the Y-Axis
1101		5	Integration by Parts
1101		4	The Definite Integral Definition
1102	Precalculus Review	55	Trigonometric Functions
1103		24	Zeros
1103		11	Inverses of Functions
1104		8	Inequalities
1104		7	Equations with Exponents and Logarithms
1105		7	Properties of Functions
1106		6	Exponential Functions
1107		6	Logarithmic Functions
1107		6	Linear Modeling
1108		5	Complex Numbers
1109	Algebra	18	Equations and Inequalities
1110		13	Polynomial Equations
1110		8	Find Composition of Two Functions
1111		7	Polynomials
1111		6	Find Slope Line
1112		6	Applications of Exponential Function
1113		6	Quadratic Equations
1113		6	Divide Rational Expressions
1114		6	Solve Linear Equation
1115		5	Zeros of Polynomials
1116	Multivariable Calculus	13	Triple Integrals
1117		11	Lagrange Multipliers
1117		9	Double Integrals in Polar Coordinates
1118		8	Derivatives of Parametric Equations
1118		8	Integrals of Multivariable Functions
1119		8	Double Integral Over General Region
1120		6	Classification of Critical Points
1120		6	Limit of 2 Variable Function
1121		6	Plane Parametric Polar and Conic Equations
1122		5	Applications of Double Integrals

Table 6: Top 8 Topics for Each Subject.

C PROMPTS

C.1 PREDICTION PROMPT

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

Solution CoT Prompt.

System:
You are an expert in mathematics.

User:
Solve the following problem. Make sure to show your work before giving the final answer.

Problem:
{{problem text}}

Comment:
Images (if present) are passed with native for provider API schema. For OpenAI-compatible endpoints it is `image_url` field.^a

^a<https://platform.openai.com/docs/guides/vision>

Figure 5: Prediction for comparing student’s answer and reference answer

C.2 JUDGMENT PROMPT

1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241

Judgment CoT Prompt.

System:
 You are a math teacher responsible for grading student’s homework. You need to compare the student’s final answer with the provided reference answer to determine if the student’s answer is correct. You should not consider any additional reasoning or explanations given.

Instructions:

1. Equivalence: If the student’s answer is mathematically equivalent to the reference answer (e.g., 1.5 vs. $3/2$ or $a^2 - b^2$ vs. $(a + b)(a - b)$ or $(1, 2)$, $(4, 5)$ vs $(1, 2) \cup (4, 5)$), the answer should be judged as correct.
2. Approximations: Do not accept approximate answers unless the problem explicitly allows them, however more precise answer should be accepted (e.g., $2 \cdot \pi$ when the reference answer is 6.2832).
3. Multiple Parts: For problems with multiple questions or parts, ALL parts must be correct for the overall answer to be correct.
4. Missing answer: Missing answers and general form answers should be considered incorrect unless the problem explicitly allows them. (e.g. all points should be found, all values asked should be provided etc.).
5. All the problem constraints have to be fulfilled (e.g., if the problem asks for 3 Taylor series terms, the student needs to provide exactly 3; or if the point(s) are asked the student need to provide point with all coordinates, e.g. (x, y)).
6. Reference Answer: Assume the provided reference answer is always correct. Do not attempt to solve the problem yourself.

Provide a clear and concise evaluation without adding extra information or solving the problem yourself.

User: {{image, if any}}
 Please compare the student’s answer with the provided reference answer.
Problem:
 {{problem text}}

Reference Answer:
 {{golden answer}}

Student’s Answer:
 {{generated answer}}

—

For each question or part:

1. Write the reference answer and the student’s final answer.
2. Make any derivations or transformations that may be necessary to compare the reference answer and the student’s answer.
3. Only then perform the comparison.

After comparing all parts, provide a final judgment is the student’s answer correct or incorrect.

Assistant:
 {{CoT solution}}

User:
 Now, summarize the judgment above in single word. Is the student answer fully correct? Please answer with a SINGLE word — either Yes or No

Assistant:
 {{extracted answer}}

Figure 6: Judgment for comparing student’s answer and reference answer