

# IS YOUR MODEL REALLY A GOOD MATH REASONER? EVALUATING MATHEMATICAL REASONING WITH CHECKLIST

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Exceptional mathematical reasoning ability is one of the key features that demonstrate the power of large language models (LLMs). How to comprehensively define and evaluate the mathematical abilities of LLMs, and even reflect the user experience in real-world scenarios, has emerged as a critical issue. Current benchmarks predominantly concentrate on problem-solving capabilities, presenting a substantial risk of model overfitting and fails to accurately measure the genuine mathematical reasoning abilities. In this paper, we argue that if a model really understands a problem, it should be robustly and readily applied across a diverse array of tasks. To this end, we introduce **MATHCHECK**, a well-designed checklist for testing task generalization and reasoning robustness, as well as an automatic tool to generate checklists efficiently. **MATHCHECK** includes multiple mathematical reasoning tasks and robustness tests to facilitate a comprehensive evaluation of both mathematical reasoning ability and behavior testing. Utilizing **MATHCHECK**, we develop **MATHCHECK-GSM** and **MATHCHECK-GEO** to assess mathematical textual reasoning and multi-modal reasoning capabilities, respectively, serving as upgraded versions of benchmarks including GSM8k, GeoQA, UniGeo, and Geometry3K. We adopt **MATHCHECK-GSM** and **MATHCHECK-GEO** to evaluate over 26 LLMs and 17 multi-modal LLMs, assessing their comprehensive mathematical reasoning abilities. Our results demonstrate that while frontier LLMs like GPT-4o continue to excel in various abilities on the checklist, many other model families exhibit a significant decline. Further experiments indicate that, compared to traditional math benchmarks, **MATHCHECK** better reflects true mathematical abilities and represents mathematical intelligence more linearly, thereby supporting our design. Using **MATHCHECK**, we can also efficiently conduct informative behavior analysis to deeply investigate models. Finally, we show that our proposed checklist paradigm can easily extend to other reasoning tasks for their comprehensive evaluation.<sup>1</sup>

## 1 INTRODUCTION

The AI community has been placing significant emphasis on mathematical reasoning as a means to explore the upper limits of intelligence in large language models (LLMs) (Achiam et al., 2023; Team et al., 2023; Meta, 2024; Jiang et al., 2024; Wei et al., 2022; Trinh et al., 2024; Romera-Paredes et al., 2024) and multi-modal large language models (MLLMs) (OpenAI, 2024c; Lu et al., 2023). A large number of efforts have been made on how to enhance (M)LLMs’ mathematical reasoning abilities. In pre-training, Wang et al. (2023d); Shao et al. (2024); Lin et al. (2024); Zhang et al. (2024c) studied the impact of the quality of mathematical corpus; in post-training, Yue et al. (2023); Yu et al. (2023); Li et al. (2024a) augmented a huge number of synthetic data, and then developed supervised fine-tuning (SFT) for math problem-solving. Recently, Luong et al. (2024) and Sun et al. (2024b) explored variants of reinforcement learning (RL) for further improvements.

To guarantee the high mathematical reasoning ability has been reached, it is crucial to fairly evaluate models’ performance. Current mainstream methods rely on the performance across math problem-solving tasks of varying difficulty levels, such as GSM8k (Cobbe et al., 2021) of elementary level,

<sup>1</sup>Data and code can be found here: <https://anonymous.4open.science/r/MathCheck>

	Problem Solving	Answerable Judging	Outcome Judging	Process Judging	
054					
055					
056	Original Problem	A robe takes 2 bolts of blue fiber and half that much white fiber. How many bolts in total does it take? "answer": 3.0	A robe <b>takes bolts</b> of blue fiber and half that much white fiber. How many bolts in total does it take? "answer": Unanswerable	A robe takes 2 bolts of blue fiber ... How many bolts in total does it take? "solution": Step 1: 2 bolts of blue fiber...The answer is <b>4 bolts</b> in total. "answer": Incorrect	A robe takes 2 bolts of blue fiber ... How many bolts in total does it take? "solution": Step 1: Identify the amount ... Step 3: <b>Multiply</b> the bolts of blue and white fiber together to find the total number of bolts. The answer is 2 bolts. "answer": Step 3
057					
058		★ Seed Data			
059					
060	Problem Understanding	To make a robe, you need 2 bolts of blue fiber and half as many bolts of white fiber compared to blue. What is the total number of bolts required for the robe? "answer": 3.0	To make a robe, you need <b>bolts</b> of blue fiber and half as many bolts of white fiber compared to blue. What is the total number of bolts required for the robe? "answer": Unanswerable	To make a robe, you need 2 bolts ... What is the total number of bolts required for the robe? "solution": Step 1: Calculate the number of blue bolts... So, 2 (blue)+ 1 (white) = 3.The answer is <b>3</b> . "answer": Correct	To make a robe, you need 2 bolts ... What is the total number of bolts required for the robe? "solution": Step 1: ... Step 2: Determine the number of white bolts, which <b>as many as blue bolts</b> . ... The answer is 4. "answer": Step 2
061					
062					
063					
064	Irrelevant Disturbance	A tailor is crafting a luxurious robe. The design requires 2 bolts of blue fiber and half that amount of white fiber. To add grandeur, the tailor also considered using 3 bolts of golden thread from the sun's rays, but eventually decided it would be too gaudy for the ceremony. How many bolts in total are needed for the robe, disregarding the golden thread? "answer": 3.0	A tailor is crafting a luxurious robe. The design requires 2 bolts of blue fiber and half that amount of white fiber. ... How many bolts in total are needed for the robe, disregarding the golden thread? "answer": Answerable	A tailor is crafting a luxurious robe. The design requires 2 bolts of blue fiber and half that amount of white fiber. ... How many bolts in total are needed for the robe, disregarding the golden thread? "solution": Step 1: Calculate the amount of blue fiber. The design requires ... The answer is: <b>300 yards</b> . "answer": Incorrect	A tailor is crafting a luxurious robe. The design requires 2 bolts of blue fiber and half that amount of white fiber. ... How many bolts in total are needed for the robe, disregarding the golden thread? "solution": Step 1: ... Step 2: Calculate the amount of white fiber required, which is double the blue fiber amount, so <b>2 bolts * 2 = 4 bolts</b> . Step 3: ... The answer is 6 bolts. "answer": Step 2
065					
066					
067					
068	Scenario Understanding	A robe takes x bolts of blue fiber and half that much white fiber. It takes 3 bolts in total. What is the value of unknown variable x? "answer": 2.0	A robe takes x bolts of blue fiber and <b>fewer</b> white fiber. It takes 3 bolts in total. What is the value of unknown variable x? "answer": Unanswerable	A robe takes x bolts of blue fiber and half that ... What is the value of unknown variable x? "solution": Step 1: Let's say the value of x is ... The answer is <b>2</b> . "answer": Correct	A robe takes x bolts of blue fiber and half that ... What is the value of unknown variable x? "solution": Step 1: Let's ... Step 3: To find out how many bolts of fiber are needed in total, the equation should be <b>x - 0.5x = 3</b> ... The answer is x equals 6. "answer": Step 3
069					
070					
071					
072					
073					
074					
075					
076					
077					
078					
079					

Figure 1: Overview of MATHCHECK design. The horizontal axis examines the task generalization of four math tasks while the vertical axis examines the reasoning robustness through four problem varieties. All data are generated from seed data, which is also from a mainstream benchmark dataset.

MATH (Hendrycks et al., 2021) of high school level, and TheromQA (Chen et al., 2023a) of university level. Recently, some mathematical datasets that are more challenging, diverse, and multi-modal have been proposed to enhance the mathematical evaluation (He et al., 2024; Liu et al., 2024c; Lu et al., 2023; Zhang et al., 2024b). However, these current evaluation methods focus on *individual* tasks (most of which are problem-solving) and robustness tests for each problem. In other words, they do not provide comprehensive guidance on whether LLMs really achieve mathematical reasoning ability. In this paper, we argue that: *if a model really understands a problem, it should work robustly across various tasks about this problem.* Therefore, it is necessary to evaluate models by multi-tasks with diverse robustness test. Through such investigation, the real reasoning ability of a model can be comprehensively evaluated. As a result, we can also perform detailed behavior tests on models (Ribeiro et al., 2020).

Drawing motivations from this insight, we introduce **MATHCHECK**, a well-designed checklist for testing task generalization and reasoning robustness. MATHCHECK includes general mathematical reasoning tasks and diverse robustness testing types to facilitate a comprehensive evaluation of mathematical reasoning ability and reasoning behavior testing. As shown in Figure 1, horizontally, we examine the task generalization including problem solving, answerable judging, outcome judging, and process judging. Vertically, we test the reasoning robustness through the original problem and its three robustness variants consisting of problem understanding, irrelevant disturbance, and scenario understanding. The data of each cell in the checklist corresponds to a specific type of robustness test and task form. To facilitate the construction of checklist, we propose an (M)LLMs-driven generation framework to automatically generate this data. Figure 2 illustrates the MATHCHECK data collection process, where the seed solving problem is firstly rewritten to its robustness problems, next all generated solving data are utilized to construct other task forms.

Utilizing MATHCHECK, we propose **MATHCHECK-GSM**, a MATHCHECK dataset generated from GSM8k (Cobbe et al., 2021). It contains a total of 3,096 high-quality samples consisting of 129

groups checklist matrix, which can be used to evaluate mathematical textual reasoning ability comprehensively. Besides, acknowledging the community’s focus on multi-modal reasoning capabilities, we further propose **MATHCHECK-GEO** to evaluate the multi-modal geometry reasoning ability. Generated from GeoQA (Chen et al., 2021), UniGeo (Chen et al., 2022), and Geometry3K (Lu et al., 2021), it contains a total of 1,440 samples with a checklist matrix of 60 groups. It is noteworthy that the construction pipeline of MATHCHECK can be applied to most mathematical datasets to dynamically establish a comprehensive and flexible evaluation benchmark, thereby mitigating data contamination (Zhou et al., 2023a; Zhu et al., 2024a;b).

We conduct extensive experiments on 26 LLMs and 17 MLLMs including different scales, API-base and open source, generalist and mathematical models. We find that frontier LLMs like GPT-4o continue to achieve superior performance in our MATHCHECK, but many other model families exhibit a significant decline. Further experiments indicate that compared to solving original problems which is the paradigm of mainstream benchmark, our MATHCHECK evaluation aligns more accurately with the genuine mathematical reasoning ability of the model. Utilizing MATHCHECK, we extensively analyze the models’ behaviors including training on massive solving data, reasoning consistency, performance on different complexity problems and applying different prompting technologies. Finally, we show the potential of applying MATHCHECK paradigm to other reasoning tasks such as commonsense reasoning and code generation, promoting more comprehensive evaluation of reasoning ability.

## 2 MATHCHECK

MATHCHECK is a well-designed checklist that includes general mathematical reasoning tasks and diverse robustness testing types for comprehensive evaluation, as well as a tool to automatically generate a large number of test cases in the manner of checklist. In our checklist, various mathematical tasks are arranged in rows to assess task generalization, whereas diverse variants of mathematical problems are placed in columns to evaluate reasoning robustness. We will elaborate on the task types in Section 2.1, problem variants in Section 2.2, and how we construct checklist data in Section 2.3.

### 2.1 TASK GENERALIZATION

Testing models across different tasks on the same domain not only offers a comprehensive and profound evaluation of their capabilities (Frank, 2023) but also caters to the practical demands and complexities of real-world applications (Ji et al., 2023). In MATHCHECK, we incorporate four math tasks including Problem Solving, Answerable Judging, Outcome Judging, and Process Judging.

**Problem Solving.** In this task, we ask the model to solve a given math problem. As the most widely used method to test mathematical reasoning ability in contemporary research (Cobbe et al., 2021; Hendrycks et al., 2021), it necessitates the model to analyze the problem, recall and apply appropriate math knowledge, and finally conclude reasoning results.

**Answerable Judging.** Given a math problem, models need to determine whether the problem provides sufficient information to answer the question. This task requires the model to analyze the question, then identify the essential conditions required for solving this question, subsequently verify whether these conditions are provided within the problem statement. Previous works utilized it to examine whether the model is a reasoner with critical thinking instead of a random parrot (Li et al., 2024b; Sun et al., 2024a; Ma et al., 2024).

**Outcome Judging.** Given a math problem and one of its solutions, let the model determine whether the final answer of the given solution is correct. Outcome-Judging is a coarse-grained judgment of solutions since the model only focuses on the correctness of the final answer. Researchers often apply the outcome-judging ability of models to verify the correctness of augmented data (Tang et al., 2024) and provide outcome rewards in reinforcement learning (Luong et al., 2024).

**Process Judging.** Given a math problem along with its wrong solution, the model is required to identify the step where the errors begin. Compared with the outcome-judging, the process-judging task is a more fine-grained judgment on the solution, which demands the model to judge step by step until the wrong step is located. It can help to debug the given wrong solution.

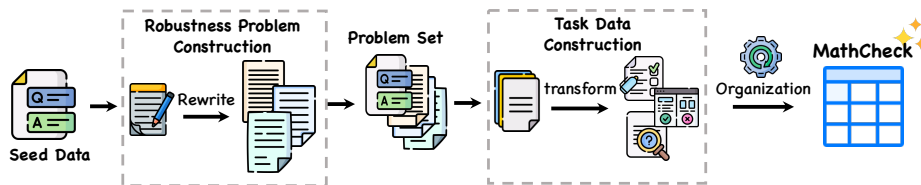


Figure 2: MATHCHECK generation pipeline.

## 2.2 REASONING ROBUSTNESS

A model that truly understands the inherent mathematical logic of a problem will exhibit reasoning robustness to diverse variations of this problem (Stolfo et al., 2023). Motivated by this, we utilize four problem forms including the original problem and its three rewritten variants to examine the reasoning robustness of models.

**Original Problem.** It is the seed problem of other reasoning robustness variants. At a minimum functionality test, it can check whether the model has the basic mathematical capabilities when no modifications have been made.

**Problem Understanding.** It refers to transforming the original problem into a new one that uses different wording or different sentence structures but does not change the mathematical logic of its original version (Patel et al., 2021; Zhou et al., 2024; Li et al., 2024b). It pays more attention to semantic robustness, and aims to examine whether models can correctly reason when dealing with different descriptions of the same mathematical logic.

**Irrelevant Disturbance.** It refers to inserting irrelevant conditions that are related to the topic of the original question, but have no impact on the final answer. Previous studies have disclosed that large language models are easily distracted by such perturbations (Shi et al., 2023). It needs the model to distinguish which conditions are necessary and which are irrelevant to the problem.

**Scenario Understanding.** When models comprehend the scenario of a math problem and its underlying logic, they should be able to solve other questions within that scenario (Liu et al., 2021; Yu et al., 2023; Zhou et al., 2023b). Therefore, we alter the original question to evaluate whether a model has a comprehensive understanding of the scenario. For example, as shown in Figure 1, we ask the question “the number of blue bolts” instead of “the number of total bolts”.

## 2.3 CHECKLIST CONSTRUCTION

Creating MATHCHECK data is a labor-intensive and time-consuming process. The advent of LLMs has introduced a new level of flexibility and quality to generate mathematical content (Norberg et al., 2023; Li et al., 2024b). Therefore, we employ (M)LLMs (e.g., GPT-4-Turbo in our experiments) as engines to automatically generate our MATHCHECK data. The data construction pipeline is shown in Figure 2. Users first assemble a collection of math problems with labels as seed data. Second, (M)LLMs initially rewrite these problems into their robustness varieties to make up the robustness problem set. Third, each problem in this set will be extended to construct multiple mathematical tasks about this problem. Finally, all data are manually checked to form MATHCHECK dataset correctly.

Based on the seed data, we automatically generate another three robustness problems as shown in the first column of Figure 1. *Problem Understanding* and *Irrelevant Disturbance* are the tasks of rewriting problems without altering the final answer. Hence, we prompt the model to rewrite our math problems while maintaining the original answer. For *Scenario Understanding*, we first extract a variable from the problem as a new answer, then prompt the model to change the question based on the extracted variable. Once we obtain the four robustness reasoning problems of the solving task, we rewrite them respectively to construct multiple tasks, including *Answerable Judging*, *Outcome Judging* and *Process Judging* as shown in the corresponding row of Figure 1. For the *Answerable Judging* task, we prompt the model to eliminate a condition from the original problem which is crucial for solving it to obtain an unanswerable problem. For *Outcome Judging* task, we ask the model to solve the problem and acquire candidate solutions, then these solutions are labeled (Correct

or Incorrect) according to the final answer. For *Process Judging* task, we apply the solution rewritten ability of (M)LLMs to construct process-judging data. Specifically, given a problem along with its correct solution, we prompt the model to make mistakes from the given steps and results in a wrong answer. In such a way, we can get a wrong solution while its mistake steps remain simultaneously. All of our prompts are listed in Appendix F.2.

### 3 EXPERIMENTS

#### 3.1 DATASETS

We use MATHCHECK to comprehensively measure the mathematical reasoning ability across textual and multi-modal settings. Consequently, two benchmarks MATHCHECK-GSM and MATHCHECK-GEO are introduced.

**MATHCHECK-GSM** is a MATHCHECK dataset generated from GSM8k (Cobbe et al., 2021). We choose GSM8k as the seed benchmark since (1) it is most widely used for evaluating mathematical textual reasoning capability. (2) we aim to determine whether advanced models are genuinely capable of reasoning at the grade school level. We first collect a test-mini set of GSM8k, which includes 129 problems sampled evenly according to the difficulty<sup>2</sup>. Subsequently, we generate 129 MATHCHECK style groups, totaling 3,096 high-quality samples by MATHCHECK. It can be used to evaluate the real mathematical reasoning ability of LLMs on GSM8k-level problems. A group of MATHCHECK-GSM case problems are listed in Appendix G.1.

**MATHCHECK-GEO** is a dataset for geometry problems, which is the representative task for evaluating multi-modal reasoning capability. First, we collect seed geometry problems from GeoQA (Chen et al., 2021), UniGeo (Chen et al., 2022), and Geometry3K (Lu et al., 2021), containing 60 problems in both English and Chinese. Subsequently, we generate 60 MATHCHECK style groups, totaling 1,440 high-quality samples. Notably, this is the first geometry problem dataset involving answerable, outcome, and process judgment tasks. MATHCHECK-GEO gives research community a harder and multi-modal MATHCHECK style dataset, as well as showing the extensibility of MATHCHECK. A group of MATHCHECK-GEO case problems are shown in Appendix G.2.

All datasets are checked with meticulous manual validation to ensure high quality and reliability. To this end, we recruited three graduate students who underwent training tailored to the requirements of our research. This rigorous verification process not only enhances the quality of our data but also reinforces the validity of our findings. Finally, our automatic data generation pipeline can achieve an average pass rate of 84.61% (Appendix C.2). The detailed data statistics and quality discussion of our checklist are reported in Appendix C.

#### 3.2 EXPERIMENTAL SETUP

To systematically benchmark the mathematical reasoning capabilities of existing LLMs, we include a comprehensive evaluation of 43 models, comprising 26 LLMs and 17 MLLMs. These models are principally divided into two categories: generalist models encompassing both API-based commercial LLMs and open-sourced LLMs (large and small scale), and specialized mathematical models. We use the F1 metric for Outcome Judging and Answerable Judging tasks, and the Acc metric for the other two tasks. The list of selected models and details of evaluation setup can be found in Appendix D.

#### 3.3 MAIN RESULTS

Tables 1 and 2 illustrate the performance of various models on the MATHCHECK-GSM and MATHCHECK-GEO, respectively. The leftmost column represents the average performance across all tasks and all question variants. The middle four columns detail the performance on various mathematical reasoning tasks, while the right four columns display performance across different question variants. Consequently, each model is represented by a  $4 \times 4$  checklist table, which showcases the model’s performance in various dimensions. The details of all checklist tables are further elaborated in Appendix A and B.

<sup>2</sup>We define the difficulty according to the number of reasoning steps of its answers (2 steps to 8 steps)

Table 1: Model performance on MATHCHECK-GSM. **PS**: Problem Solving, **AJ**: Answerable Judging, **OJ**: Outcome Judging, **PJ**: Process Judging, **OP**: Original Problem, **PU**: Problem Understanding, **ID**: Irrelevant Disturbance, **SU**: Scenario Understanding. Each score is the average score of related units. For example, 'All' means all units, 'PS' includes solving units on four problem types, 'OP' includes original problems on four tasks units.

Models	All	PS	AJ	OJ	PJ	OP	PU	ID	SU
<i>Generalist Models</i>									
O1-preview	93.2	91.3	94.0	93.2	94.1	95.6	93.4	90.5	93.1
O1-mini	92.7	93.6	95.0	88.9	93.6	95.5	94.2	91.0	90.5
GPT-4o	92.0	95.0	95.0	90.1	87.8	94.6	91.6	92.0	89.6
GPT-4o-mini	87.2	90.1	89.6	88.6	80.4	88.9	89.4	85.6	85.1
GPT-4-Turbo-20240409	90.9	93.8	95.9	87.8	86.0	93.8	90.4	90.8	88.6
GPT-3.5-Turbo	61.4	73.5	64.3	48.3	59.5	65.4	64.6	60.1	55.4
Gemini-1.5-Pro	86.3	88.6	89.5	87.6	75.0	88.0	90.2	85.0	82.0
Claude-3.5-sonnet-20240620	90.2	94.8	95.3	90.9	79.9	92.5	92.1	89.9	86.3
Claude-3-opus-20240229	83.5	81.6	92.0	78.7	81.8	86.3	85.6	81.9	80.3
Claude-3-sonnet-20240229	75.0	77.9	88.9	65.1	68.0	76.5	77.8	73.7	71.9
Claude-3-haiku-20240229	57.5	79.7	49.9	44.3	56.0	61.9	62.4	55.9	49.6
Llama-3.1-70B-Instruct	90.5	95.2	95.3	89.4	82.2	93.3	91.2	89.8	87.7
Llama-3-70B-Instruct	84.7	90.1	87.5	84.6	76.7	87.7	86.7	84.7	79.9
DeepSeek V2	82.2	86.8	82.6	82.5	76.9	85.1	84.4	83.5	75.9
Mixtral 8 x 7B-Instruct	59.9	56.0	58.1	63.9	61.6	62.8	61.5	58.8	56.4
Mixtral 8 x 7B-Base	44.7	40.9	50.8	51.8	35.3	50.6	47.8	41.2	39.1
Qwen1.5-72B-Chat	50.6	71.1	64.2	31.9	35.1	57.0	51.1	43.6	50.6
Phi-3-Medium-4K-Instruct	72.0	89.7	70.8	63.2	64.1	77.6	78.7	71.1	60.4
Phi-3-Mini-4K-Instruct	64.1	71.3	64.5	62.9	57.6	68.5	66.6	61.2	60.0
Llama-3.1-8B-Instruct	71.0	76.9	65.8	77.2	64.0	74.6	73.6	66.0	69.6
Llama-3-8B-Instruct	64.2	68.6	61.4	64.9	61.8	67.8	68.8	62.9	57.1
ChatGLM3-6B	36.5	32.6	41.7	50.1	21.7	39.7	35.9	31.3	39.1
<i>Mathematical Models</i>									
DeepSeek-Math-7B-RL	50.7	79.5	50.0	45.1	28.1	53.3	51.2	47.5	50.6
DeepSeek-Math-7B-Instruct	50.2	70.0	64.8	40.4	25.8	51.6	54.4	45.8	49.2
DeepSeek-Math-7B-Base	44.0	49.8	51.5	44.0	30.8	49.0	46.0	37.0	44.1
MetaMath-LLama2-70B	45.7	70.0	35.7	45.3	31.6	49.9	51.5	43.4	37.8

Table 2: Model performance on MATHCHECK-GEO.

Models	All	PS	AJ	OJ	PJ	OP	PU	ID	SU
<i>Generalist Models</i>									
GPT-4o	65.3	57.5	75.5	69.5	58.8	65.2	67.0	64.3	64.8
GPT-4o-mini	59.0	50.8	69.8	61.4	53.8	61.9	62.0	54.1	57.8
GPT-4-Turbo-20240409	61.7	51.3	72.3	64.0	59.2	63.2	62.9	61.7	58.9
GPT-4-Vision-Preview	60.0	46.7	71.1	63.6	58.8	59.3	62.8	57.8	60.2
Gemini-1.5-Pro	58.7	47.5	67.4	55.0	64.6	62.3	58.6	57.1	56.9
Gemini-1.5-Flash	56.8	45.0	75.1	50.6	56.7	56.8	59.7	53.8	57.1
Claude-3.5-sonnet-20240620	58.7	54.2	71.0	53.0	56.7	59.9	63.8	54.3	56.8
Claude-3-opus-20240229	47.2	34.2	60.6	46.7	47.5	47.2	49.1	42.4	50.2
Claude-3-sonnet-20240229	49.9	35.8	59.0	51.6	52.9	51.2	53.0	44.7	50.4
Claude-3-haiku-20240307	36.7	27.9	41.3	41.7	35.8	39.2	38.8	33.3	35.4
QWen2-VL-72B-Instruct	61.4	60.0	53.1	61.3	71.3	69.0	62.4	58.0	56.4
QWen2-VL-7B-Instruct	42.1	35.8	49.4	46.4	36.7	40.9	45.6	41.7	40.0
InternVL-1.5-Chat	37.6	22.1	54.9	46.8	26.7	42.9	34.8	37.3	35.5
MiniCPM-LLama3-V-2.5	37.3	37.5	38.1	45.0	28.8	37.4	45.0	35.2	31.6
LLaVA-1.6-Mistral-7B-Instruct	31.8	10.0	38.8	51.2	27.1	33.8	35.5	28.4	29.2
Phi-3-Vision-128k-Instruct	29.6	12.9	35.0	48.6	22.9	32.6	31.8	28.2	26.0
CogVLM2-LLama3-Chat-19B	24.6	7.9	26.4	46.3	17.9	27.2	28.0	22.4	20.9

On MATHCHECK-GSM (Table 1), O1-preview and O1-mini exhibit outstanding performance with impressive overall score of 93.2 and 92.7, demonstrates strong effect of extending reasoning thought exploration. GPT-4o is closely followed with a score of 92.0 and demonstrates top performance

on the problem solving task and irrelevant disturbance variants. These results indicate that strong foundational models still possess formidable and robust performance across a variety of mathematical reasoning tasks. Among the open-source LLMs, LLaMa-3.1-70B-Instruct achieves the highest score of 90.5 and performs excellently across a range of tasks and problem variants. Its performance has significantly improved compared to LLaMA-3 version and surpasses that of GPT-4o-mini. Besides, Qwen1.5-72B-Chat underperforms in tasks other than problem solving, which we suspect is due to its special optimization of the solving task. This phenomenon is also observed across all math-customized models, which tend to be trained on similar mathematical problems and problem-solving processes, resulting in a relatively narrow scope of reasoning capabilities.

On MATHCHECK-GEO (Table 2), GPT-4o demonstrates the best performance, achieving a top score of 65.3 in the All category. The performance of GPT4-turbo-20240409 and GPT4-Vision-Preview is similar, reaching scores of 61.7 and 60.0, respectively. In particular, the performance of Claude-3-sonnet is slightly superior in visual contexts compared to that of its larger counterpart, Claude-3-opus. Among the open-source MLLMs, the large-size MLLMs demonstrate surprisingly strong performance, with Qwen-VL-70B attaining 60.4 over the GPT-4-Vision-Preview. However, the most of small-size MLLMs exhibited poor performance especially in problem solving, which suggests that the multi-modal reasoning capabilities of open-source small-size open-source MLLMs still have significant room for improvement.

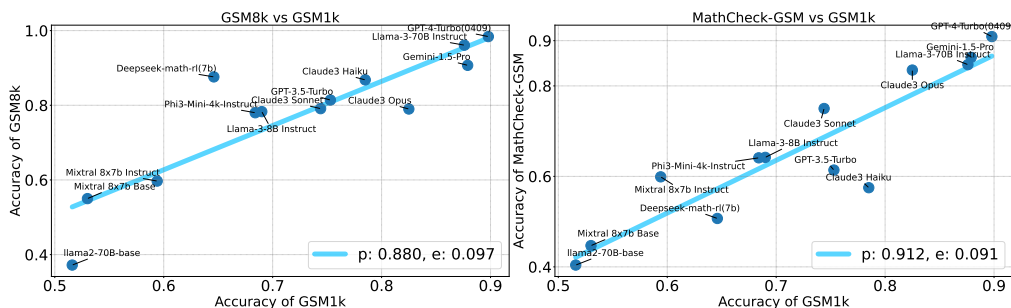


Figure 3: Correlation with GSM1k (Zhang et al., 2024a), a dataset that reflects real mathematical reasoning ability.  $p$  and  $e$  represent the Pearson Correlation Coefficient, and Root Mean Square Error.

### 3.4 MATHCHECK REPRESENTS MATHEMATICAL INTELLIGENCE MORE LINEARLY

One desiderata of a good mathematical benchmark is to reflect real mathematical intelligence perfectly. We follow previous works (Zhang et al., 2024a; Huang et al., 2024a) to assess “intelligence” from practical standpoints and use performance on private data (Zhang et al., 2024a) and compression efficiency (Du et al., 2024; Huang et al., 2024a) as surrogates to assess the genuine mathematical abilities of models. By examining the correlation between MATHCHECK and these surrogates, we can verify whether our design effectively reflects mathematical intelligence, and how it compares to traditional benchmarks.

**Correlation with Private Data.** Unlike traditional open-sourced benchmarks, private data is less likely to be contaminated or overfitted, making it an appropriate proxy of genuine mathematical intelligence. We adopt GSM1k (Zhang et al., 2024a), a new private GSM8k-level dataset, to measure the real mathematical reasoning of models. We compare the correlation of model performance between GSM1k and MATHCHECK-GSM/GSM8k. As shown in Figure 3, the left part illustrates the correlation between GSM8k and GSM1k. It reveals that most LLMs achieve scores up to 80% on GSM8k, with scores concentrated in the top half of the graph. However, on GSM1k, the scores are evenly distributed, indicating that some LLMs, such as deepseek-math-7B-RL, have inflated scores on GSM8k. This suggests that the GSM8k score is not a reliable benchmark for assessing the true mathematical reasoning ability of the models. In the right sub-figure, MATHCHECK-GSM and GSM1k display a good positive correlation, and some models that do not perform well on GSM1k can be detected by MATHCHECK-GSM. By comparing the Pearson correlation coefficient and the root mean square error, it shows that MATHCHECK has a higher correlation coefficient with GSM1k, mitigating bias evaluation caused by overfitting and data contamination.

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

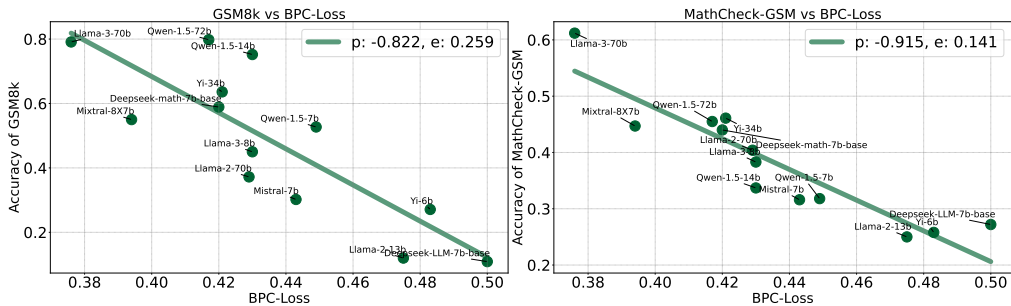


Figure 4: Performance correlation with BPC-loss, which reflects compression efficiency (Huang et al., 2024a). The lower BPC-loss represents the higher compression efficiency.

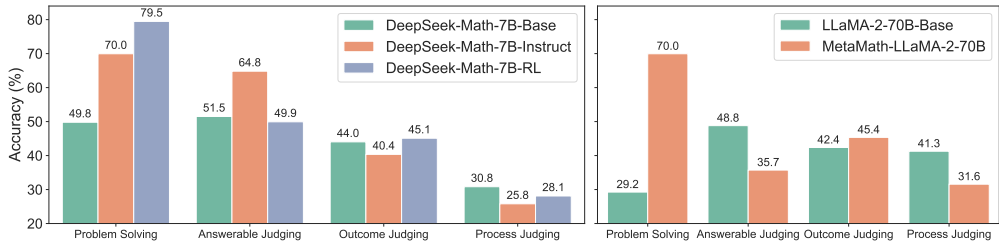


Figure 5: Behavior of mathematical models trained on massive solving data.

**Correlation with Compression Efficiency.** Compression efficiency has been empirically proven that represent intelligence well (Du et al., 2024) even linearly (Huang et al., 2024a), well aligned with the belief that compression is closely connected to intelligence (Deletang et al., 2024). Following Huang et al. (2024a), we use BPC-Loss in Arxiv papers tagged with “Math” to measure compression efficiency as a surrogate. Figure 4 shows the correlation between BPC-Loss and GSM8K/MathCheck-GSM. The left sub-figure reveals that a single traditional benchmark like GSM8K cannot adequately reflect genuine mathematical ability, as indicated by the low Pearson correlation coefficient ( $p = -0.822$ ). Many models, such as the Qwen series, deviate significantly from the regression line. In contrast, the right sub-figure displays the correlation with our MATHCHECK-GSM, demonstrating that MATHCHECK-GSM exhibits a significantly better correlation with genuine intelligence, with a Pearson correlation coefficient of  $p = -0.915$ . Our method shows that many models, such as the Qwen series, have scores on our benchmark that align more accurately with their true mathematical abilities. It shows that our design can represent mathematical intelligence more linearly.

#### 4 BEHAVIOR ANALYSIS

MATHCHECK contains multi-dimensional information for evaluation, therefore we can observe the behaviors of the models on it to help analyze the models.

**Behavior of Math Models.** Recently, some works claim that math reasoning ability is greatly improved by training on massive amounts of math solving data. To validate whether their mathematical reasoning ability really improves, we examine the behaviors of the math models and their base models on MATHCHECK. As shown in Figure 5, compared with the base model, the performance of DeepSeek-Math-7B-Instruct/RL on solving units is greatly improved. However, the performance improvement on other units is limited, or even downward. The same phenomenon can be observed on MetaMath. It implies that training solely on massive solving data (Yue et al., 2023; Li et al., 2024a; Tang et al., 2024) is not the right direction to improve mathematical reasoning ability. Instead, training models with diverse mathematical data, beyond just solving, should be considered.

**Reasoning Consistency.** We analyze the reasoning consistency of generalist models across each unit in MATHCHECK, and the detailed results are shown in Appendix A and B. We can see most of them show good reasoning consistency since they achieve similar scores on each unit, such as GPT



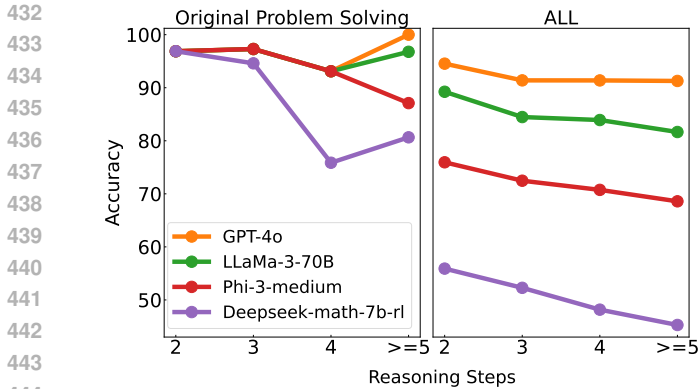


Figure 6: Performance on different complexity levels (i.e., reasoning steps) of MATHCHECK-GSM.

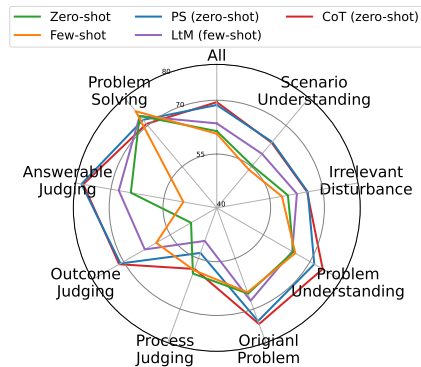


Figure 7: Different prompting technologies on MATHCHECK-GSM.

series, Llama-3 series and Mixtral series on MATHCHECK-GSM and GPT series on MATHCHECK-GEO. This is an interesting finding as it substantiates our assertion: *a model that really understands a problem can robustly work well on multiple related tasks*. Meanwhile, we also find that some models perform reasoning inconsistently. For example, Qwen1.5-72B-chat, Claude-3-Haiku and Phi-3-Medium show excellent performance on the solving task but much worse in other units of MATHCHECK-GSM. On MATHCHECK-GSM, Internet-VL achieves a high score of 40.0 on the original problem solving but decreases considerably when the problem switches to other robustness variants. These abnormal inconsistency behaviors of generalist models are highly similar to those mathematical models, revealing that they may conduct excessive decoration on original benchmarks.

**Behavior on Different Complexity Levels.** We categorize the complexity of problems based on the number of reasoning steps of the original problems, and select representative models of varying sizes for evaluation, as depicted in Figure 6. We can observe that the models’ accuracy on the original problem solving fluctuates and does not show an obvious downward trend as the problems are more difficult. While the score "ALL" shows a steady downward trend, it implies that MATHCHECK better demonstrates the reasoning skills and capabilities required when problems become difficult.

**Behavior on Different Prompting Technologies.** We evaluate five prompting techniques including Zero-shot, Few-shot (Brown et al., 2020), CoT (Wei et al., 2022), Least to Most prompting (Zhou et al., 2022), and Plan-and-Solve prompting (Wang et al., 2023b). The results of GPT-3.5-Turbo on MATHCHECK-GSM are illustrated in Figure 7. Overall, Chain of Thought (CoT) and Plan-and-Solve (PS) in the zero-shot setting demonstrate superior performance, though this is not consistently the case across all tasks and settings. In contrast, the Few-shot prompt generally yields worse results than the Zero-shot prompt. Through detailed analyses, we find that the math reasoning generalization of LLMs is sensitive to Few-shot samples, which inspires us that Zero-shot with advanced prompt techniques (e.g., CoT or PS) may be a better choice in mathematical reasoning tasks.

## 5 MATHCHECK APPLIED TO OTHER REASONING TASKS

MATHCHECK can be adapted to other reasoning tasks beyond mathematical problems. We attempt the migration of the MATHCHECK paradigm in both commonsense reasoning and code generation.

**Commonsense Reasoning:** It requires LLMs to apply parametric knowledge to reason and solve problems. In this paper, we choose the date understanding task in Big-bench (bench authors, 2023) as test-bed since it is widely used to measure commonsense reasoning ability (Wei et al., 2022). Appendix E.1 shows the case of applying MATHCHECK to date understanding. Similar to mathematical reasoning, date understanding is a numerical reasoning task, where it can easily utilize variants of each unit in MATHCHECK. With MATHCHECK, a simple raw data of date understanding have various corresponding test cases to examine the reasoning robustness and task generalization, helping us better evaluate model’s understanding of dates and avoiding hallucination.

**Code Generation:** We would like to show the possibility of transforming MATHCHECK in some real-world reasoning tasks such as code generation. Appendix E.2 demonstrates a case of applying

MATHCHECK to code generation. Unlike numerical reasoning, the adaptation of code generation should consider task relevance. For real-world tasks such as agents and robotics application, multiple variants reflects the diversity of environment and user requirements.

## 6 RELATED WORK

**Benchmarks of Textual Mathematical Reasoning.** Numerous benchmarks have been proposed to evaluate the mathematical reasoning capabilities including (Amini et al., 2019; Cobbe et al., 2021; Frieder et al., 2024). Some datasets, such as the elementary-level GSM8k (Cobbe et al., 2021). Consequently, more challenging datasets have been introduced, including those at the high-school level (Hendrycks et al., 2021), university level (Sawada et al., 2023; Zheng et al., 2021) and olympic level (Huang et al., 2024b). Additionally, to provide a more comprehensive evaluation of mathematical reasoning abilities, numerous benchmarks have been developed that measure the robustness of mathematical reasoning (Li et al., 2024b), including semantic perturbations (Wang et al., 2023a; Zhou et al., 2024), reverse problem-solving (Yu et al., 2023; Berglund et al., 2023), irrelevant distractions (Shi et al., 2023; Li et al., 2023) and functional variation questions (Srivastava et al., 2024; Gulati et al., 2024). Above benchmarks paradigm can not comprehensively reflect reasoning ability at a given level. Therefore, MATHCHECK tries to go for better reasoning benchmark paradigm.

**Benchmarks of Visual Mathematical Reasoning.** Recently, multi-modal large language models have demonstrated outstanding capabilities in visual-language reasoning tasks (Allaway et al., 2022; Chen et al., 2023b; Yang et al., 2023; Team et al., 2023). Several benchmarks (Lin et al., 2014; Antol et al., 2015; Hudson & Manning, 2019; Marino et al., 2019; Mobasher et al., 2022) have been introduced to assess the visual reasoning capabilities of multi-modal large language models across various modalities including abstract scenes, geometric diagrams, graphics, and charts (Lu et al., 2021; Chen et al., 2021; 2022; Masry et al., 2022; Kazemi et al., 2023; Lu et al., 2023). MATHCHECK-GEO offers a comprehensive evaluation and testing platform for the research on visual math reasoning.

**Benchmarks of Reasoning Consistency.** Prior studies have identified limitations in reasoning consistency. Wu et al. (2023) designed counterfactual tasks to demonstrate that LLMs often rely on memorization to address general reasoning tasks. Berglund et al. (2023) found that LLMs struggle to answer inverse questions such as “B is A” after training on “A is B”. In code reasoning, Gu et al. (2024) and Liu et al. (2024a) observed that LLMs successfully generate solution but fail to correct the wrong one. Similarly, Oh et al. (2024) found the gap between generation and evaluation in TriviaQA (Joshi et al., 2017). These findings inspire the design of MATHCHECK.

**Strategies of Improving Mathematical Reasoning.** Community has made significant efforts to enhance mathematical reasoning. In pre-training stage, previous works focus on collecting (Wang et al., 2023d; Paster et al., 2024; Shao et al., 2024) and synthesizing (Akter et al., 2024) math documents. In addition, Lin et al. (2024) selected key tokens in math data during pre-training. In post-training, numerous works generated massive problem-solving data for SFT (Yue et al., 2023; Li et al., 2024a; Tang et al., 2024). Besides, reinforcement learning such as GRPO (Shao et al., 2024) PRM (Lightman et al., 2024) can further improve reasoning ability. In inference, prompt and search strategies make LLMs reasoning better (Zhou et al., 2022; Wang et al., 2023b; Yao et al., 2024a).

## 7 CONCLUSION

In this paper, we argue that if a model really understands a problem, it should be able to successfully solve various tasks and variations of that problem. Based on this insight, we introduce **MATHCHECK**, a well-designed checklist for testing task generalization and reasoning robustness. To this end, we also propose an automatic tool for efficiently generating checklist for most of math reasoning datasets. Our proposed MATHCHECK allows the research community to clearly observe model performance across different dimensions, yielding more comprehensive and objective evaluation results. Using MATHCHECK, we develop **MATHCHECK-GSM** for textual reasoning and **MATHCHECK-GEO** for multi-modal reasoning. We evaluate massive (M)LLMs and conduct detailed analysis of model behaviors on MATHCHECK. Subsequently, we reveal that the evaluation on MATHCHECK is closer to the true reasoning abilities than previous benchmark paradigm. Finally, we show the potential of applying MATHCHECK paradigm to other reasoning tasks. We hope our practice and observation can constitute a significant stride towards better reasoning benchmark paradigm.

## REFERENCES

- 540  
541  
542 Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany  
543 Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report:  
544 A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.
- 545 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
546 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.  
547 *arXiv preprint arXiv:2303.08774*, 2023.
- 548 Syeda Nahida Akter, Shrimai Prabhumoye, John Kamalu, Sanjeev Satheesh, Eric Nyberg, Mostofa  
549 Patwary, Mohammad Shoeybi, and Bryan Catanzaro. Mind: Math informed synthetic dialogues  
550 for pretraining llms. *arXiv preprint arXiv:2410.12881*, 2024.
- 551  
552 Emily Allaway, Jena D Hwang, Chandra Bhagavatula, Kathleen McKeown, Doug Downey, and Yejin  
553 Choi. Penguins don’t fly: Reasoning about generics through instantiations and exceptions. *arXiv*  
554 *preprint arXiv:2205.11658*, 2022.
- 555 Aida Amini, Saadia Gabriel, Peter Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi.  
556 Mathqa: Towards interpretable math word problem solving with operation-based formalisms.  
557 *arXiv preprint arXiv:1905.13319*, 2019.
- 558 Anthropic. Claude 3, 2024a. URL <https://www.anthropic.com/index/claude-3>.
- 559  
560 Anthropic. Claude 3.5 sonnet, 2024b. URL [https://www.anthropic.com/news/](https://www.anthropic.com/news/claude-3-5-sonnet)  
561 [claude-3-5-sonnet](https://www.anthropic.com/news/claude-3-5-sonnet).
- 562 Stanislaw Antol, Aishwarya Agrawal, Jiaseen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick,  
563 and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international*  
564 *conference on computer vision*, pp. 2425–2433, 2015.
- 565  
566 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,  
567 Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- 568 BIG bench authors. Beyond the imitation game: Quantifying and extrapolating the capabilities of  
569 language models. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL  
570 <https://openreview.net/forum?id=uyTL5Bvosj>.
- 571  
572 Lukas Berglund, Meg Tong, Max Kaufmann, Mikita Balesni, Asa Cooper Stickland, Tomasz Korbak,  
573 and Owain Evans. The reversal curse: LLMs trained on "a is b" fail to learn "b is a". *arXiv preprint*  
574 *arXiv:2309.12288*, 2023.
- 575 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
576 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
577 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 578 Jiaqi Chen, Jianheng Tang, Jinghui Qin, Xiaodan Liang, Lingbo Liu, Eric P Xing, and Liang Lin.  
579 Geoqa: A geometric question answering benchmark towards multimodal numerical reasoning.  
580 *arXiv preprint arXiv:2105.14517*, 2021.
- 581  
582 Jiaqi Chen, Tong Li, Jinghui Qin, Pan Lu, Liang Lin, Chongyu Chen, and Xiaodan Liang. Unigeo:  
583 Unifying geometry logical reasoning via reformulating mathematical expression. *arXiv preprint*  
584 *arXiv:2212.02746*, 2022.
- 585 Wenhui Chen, Ming Yin, Max Ku, Pan Lu, Yixin Wan, Xueguang Ma, Jianyu Xu, Xinyi Wang, and  
586 Tony Xia. Theoremqa: A theorem-driven question answering dataset. In *The 2023 Conference on*  
587 *Empirical Methods in Natural Language Processing*, 2023a.
- 588  
589 Xi Chen, Josip Djolonga, Piotr Padlewski, Basil Mustafa, Soravit Changpinyo, Jialin Wu, Car-  
590 los Riquelme Ruiz, Sebastian Goodman, Xiao Wang, Yi Tay, et al. Pali-x: On scaling up a  
591 multilingual vision and language model. *arXiv preprint arXiv:2305.18565*, 2023b.
- 592 Zhe Chen, Jiannan Wu, Wenhui Wang, Weijie Su, Guo Chen, Sen Xing, Zhong Muyan, Qinglong  
593 Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning  
for generic visual-linguistic tasks. *arXiv preprint arXiv:2312.14238*, 2023c.

- 594 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
595 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve  
596 math word problems. *arXiv preprint arXiv:2110.14168*, 2021.  
597
- 598 Gregoire Deletang, Anian Ruoss, Paul-Ambroise Duquenne, Elliot Catt, Tim Genewein, Christopher  
599 Mattern, Jordi Grau-Moya, Li Kevin Wenliang, Matthew Aitchison, Laurent Orseau, Marcus  
600 Hutter, and Joel Veness. Language modeling is compression. In *The Twelfth International  
601 Conference on Learning Representations*, 2024. URL [https://openreview.net/forum?  
602 id=jznbgiynus](https://openreview.net/forum?id=jznbgiynus).
- 603 Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. Glm:  
604 General language model pretraining with autoregressive blank infilling. In *Proceedings of the 60th  
605 Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp.  
606 320–335, 2022.  
607
- 608 Zhengxiao Du, Aohan Zeng, Yuxiao Dong, and Jie Tang. Understanding emergent abilities of  
609 language models from the loss perspective, 2024.
- 610 Michael C Frank. Baby steps in evaluating the capacities of large language models. *Nature Reviews  
611 Psychology*, 2(8):451–452, 2023.  
612
- 613 Simon Frieder, Luca Pinchetti, Ryan-Rhys Griffiths, Tommaso Salvatori, Thomas Lukasiewicz,  
614 Philipp Petersen, and Julius Berner. Mathematical capabilities of chatgpt. *Advances in Neural  
615 Information Processing Systems*, 36, 2024.  
616
- 617 Alex Gu, Wen-Ding Li, Naman Jain, Theo X Olausson, Celine Lee, Koushik Sen, and Armando  
618 Solar-Lezama. The counterfeit conundrum: Can code language models grasp the nuances of their  
619 incorrect generations? *arXiv preprint arXiv:2402.19475*, 2024.  
620
- 621 Aryan Gulati, Brando Miranda, Eric Chen, Emily Xia, Kai Fronsdal, Bruno de Moraes Dumont, and  
622 Sanmi Koyejo. Putnam-axiom: A functional and static benchmark for measuring higher level  
623 mathematical reasoning. In *The 4th Workshop on Mathematical Reasoning and AI at NeurIPS'24*,  
624 2024.
- 625 Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu,  
626 Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for  
627 promoting agi with olympiad-level bilingual multimodal scientific problems. *arXiv preprint  
628 arXiv:2402.14008*, 2024.
- 629 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,  
630 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv  
631 preprint arXiv:2103.03874*, 2021.  
632
- 633 Yuzhen Huang, Jinghan Zhang, Zifei Shan, and Junxian He. Compression represents intelligence  
634 linearly. *arXiv preprint arXiv:2404.09937*, 2024a.  
635
- 636 Zhen Huang, Zengzhi Wang, Shijie Xia, and Pengfei Liu. Olympicarena medal ranks: Who is the  
637 most intelligent ai so far? *arXiv preprint arXiv:2406.16772*, 2024b.  
638
- 639 Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning  
640 and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer  
641 vision and pattern recognition*, pp. 6700–6709, 2019.
- 642 Hyangeun Ji, Insook Han, and Yujung Ko. A systematic review of conversational ai in language  
643 education: Focusing on the collaboration with human teachers. *Journal of Research on Technology  
644 in Education*, 55(1):48–63, 2023.  
645
- 646 Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris  
647 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al.  
Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.

- 648 Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly  
649 supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual*  
650 *Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1601–  
651 1611, 2017.
- 652 Mehran Kazemi, Hamidreza Alvari, Ankit Anand, Jialin Wu, Xi Chen, and Radu Soricut. Geomverse:  
653 A systematic evaluation of large models for geometric reasoning. *arXiv preprint arXiv:2312.12241*,  
654 2023.
- 655 Chen Li, Weiqi Wang, Jingcheng Hu, Yixuan Wei, Nanning Zheng, Han Hu, Zheng Zhang, and  
656 Houwen Peng. Common 7b language models already possess strong math capabilities. *arXiv*  
657 *preprint arXiv:2403.04706*, 2024a.
- 658 Qintong Li, Leyang Cui, Xueliang Zhao, Lingpeng Kong, and Wei Bi. Gsm-plus: A comprehensive  
659 benchmark for evaluating the robustness of llms as mathematical problem solvers. *arXiv preprint*  
660 *arXiv:2402.19255*, 2024b.
- 661 Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E Gonzalez,  
662 and Ion Stoica. From crowdsourced data to high-quality benchmarks: Arena-hard and benchbuilder  
663 pipeline. *arXiv preprint arXiv:2406.11939*, 2024c.
- 664 Zekun Li, Baolin Peng, Pengcheng He, and Xifeng Yan. Do you really follow me? adversarial instruc-  
665 tions for evaluating the robustness of large language models. *arXiv preprint arXiv:2308.10819*,  
666 2023.
- 667 Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan  
668 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *The Twelfth*  
669 *International Conference on Learning Representations*, 2024.
- 670 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr  
671 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–*  
672 *ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings,*  
673 *Part V 13*, pp. 740–755. Springer, 2014.
- 674 Zhenghao Lin, Zhibin Gou, Yeyun Gong, Xiao Liu, Yelong Shen, Ruochen Xu, Chen Lin, Yujiu  
675 Yang, Jian Jiao, Nan Duan, et al. Rho-1: Not all tokens are what you need. *arXiv preprint*  
676 *arXiv:2404.07965*, 2024.
- 677 Changshu Liu, Shizhuo Dylan Zhang, Ali Reza Ibrahimzada, and Reyhaneh Jabbarvand. Code-  
678 mind: A framework to challenge large language models for code reasoning. *arXiv preprint*  
679 *arXiv:2402.09664*, 2024a.
- 680 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in*  
681 *neural information processing systems*, 36, 2024b.
- 682 Hongwei Liu, Zilong Zheng, Yuxuan Qiao, Haodong Duan, Zhiwei Fei, Fengzhe Zhou, Wenwei  
683 Zhang, Songyang Zhang, Dahua Lin, and Kai Chen. Mathbench: Evaluating the theory and  
684 application proficiency of llms with a hierarchical mathematics benchmark. *arXiv preprint*  
685 *arXiv:2405.12209*, 2024c.
- 686 Qianying Liu, Wenyu Guan, Sujian Li, Fei Cheng, Daisuke Kawahara, and Sadao Kurohashi. Roda:  
687 Reverse operation based data augmentation for solving math word problems. *IEEE/ACM Transac-*  
688 *tions on Audio, Speech, and Language Processing*, 30:1–11, 2021.
- 689 Pan Lu, Ran Gong, Shibiao Jiang, Liang Qiu, Siyuan Huang, Xiaodan Liang, and Song-Chun Zhu.  
690 Inter-gps: Interpretable geometry problem solving with formal language and symbolic reasoning.  
691 *arXiv preprint arXiv:2105.04165*, 2021.
- 692 Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng,  
693 Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning  
694 of foundation models in visual contexts. In *The Twelfth International Conference on Learning*  
695 *Representations*, 2023.

- 702 Trung Quoc Luong, Xinbo Zhang, Zhanming Jie, Peng Sun, Xiaoran Jin, and Hang Li. Reft:  
703 Reasoning with reinforced fine-tuning. *arXiv preprint arXiv:2401.08967*, 2024.  
704
- 705 Jingyuan Ma, Damai Dai, and Zhifang Sui. Large language models are unconscious of unreasonability  
706 in math problems. *arXiv preprint arXiv:2403.19346*, 2024.  
707
- 708 Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual  
709 question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf  
710 conference on computer vision and pattern recognition*, pp. 3195–3204, 2019.
- 711 Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A bench-  
712 mark for question answering about charts with visual and logical reasoning. *arXiv preprint  
713 arXiv:2203.10244*, 2022.
- 714 Meta. Introducing meta llama 3: The most capable openly available llm to date. <https://ai.meta.com/blog/meta-llama-3/>, 2024.  
715  
716
- 717 Shaghayegh Mobasher, Ghazal Zamaninejad, Maryam Hashemi, Melika Nobakhtian, and Sauleh  
718 Eetemadi. Parsvqa-caps: A benchmark for visual question answering and image captioning in  
719 persian. *people*, 101:404, 2022.
- 720 Kole Norberg, Husni Almoubayyed, Stephen E Fancsali, Logan De Ley, Kyle Weldon, April Murphy,  
721 and Steven Ritter. Rewriting math word problems with large language models. In *AIED23: artificial  
722 intelligence in education, empowering education with LLMs workshop*, 2023.  
723
- 724 Juhyun Oh, Eunsu Kim, Inha Cha, and Alice Oh. The generative ai paradox on evaluation: What it  
725 can solve, it may not evaluate. *arXiv preprint arXiv:2402.06204*, 2024.  
726
- 727 OpenAI. Gpt-3.5-turbo. 2022.
- 728 OpenAI. Gpt-4o, 2024a. URL <https://openai.com/index/hello-gpt-4o/>.
- 729 OpenAI. Gpt-4o mini, 2024b. URL [https://openai.com/index/  
730 gpt-4o-mini-advancing-cost-efficient-intelligence](https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence).
- 731 OpenAI. Gpt-4v, 2024c. URL [https://openai.com/research/  
732 gpt-4v-system-card](https://openai.com/research/gpt-4v-system-card).
- 733 OpenAI. O1-mini, 2024d. URL [https://openai.com/index/  
734 openai-o1-mini-advancing-cost-efficient-reasoning](https://openai.com/index/openai-o1-mini-advancing-cost-efficient-reasoning).
- 735 OpenAI. O1-preview, 2024e. URL [https://openai.com/index/  
736 introducing-openai-o1-preview](https://openai.com/index/introducing-openai-o1-preview).
- 737 Keiran Paster, Marco Dos Santos, Zhangir Azerbayev, and Jimmy Ba. Openwebmath: An open dataset  
738 of high-quality mathematical web text. In *The Twelfth International Conference on Learning  
739 Representations*, 2024.
- 740 Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are nlp models really able to solve simple math  
741 word problems? In *Proceedings of the 2021 Conference of the North American Chapter of the  
742 Association for Computational Linguistics: Human Language Technologies*, pp. 2080–2094, 2021.  
743
- 744 Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. Beyond accuracy:  
745 Behavioral testing of nlp models with checklist. *arXiv preprint arXiv:2005.04118*, 2020.  
746
- 747 Bernardino Romera-Paredes, Mohammadamin Barekatain, Alexander Novikov, Matej Balog,  
748 M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang,  
749 Omar Fawzi, et al. Mathematical discoveries from program search with large language models.  
750 *Nature*, 625(7995):468–475, 2024.  
751
- 752 Tomohiro Sawada, Daniel Paleka, Alexander Havrilla, Pranav Tadepalli, Paula Vidas, Alexander  
753 Krnias, John J Nay, Kshitij Gupta, and Aran Komatsuzaki. Arb: Advanced reasoning benchmark  
754 for large language models. *arXiv preprint arXiv:2307.13692*, 2023.  
755

- 756 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, YK Li, Y Wu,  
757 and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language  
758 models. *arXiv preprint arXiv:2402.03300*, 2024.
- 759  
760 Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed Chi, Nathanael Schärli,  
761 and Denny Zhou. Large language models can be easily distracted by irrelevant context. In  
762 *Proceedings of the 40th International Conference on Machine Learning*, pp. 31210–31227, 2023.
- 763  
764 Saurabh Srivastava, Anto PV, Shashank Menon, Ajay Sukumar, Alan Philipose, Stevin Prince, Sooraj  
765 Thomas, et al. Functional benchmarks for robust evaluation of reasoning performance, and the  
766 reasoning gap. *arXiv preprint arXiv:2402.19450*, 2024.
- 767  
768 Alessandro Stolfo, Zhijing Jin, Kumar Shridhar, Bernhard Schoelkopf, and Mrinmaya Sachan. A  
769 causal framework to quantify the robustness of mathematical reasoning with language models. In  
770 *The 61st Annual Meeting Of The Association For Computational Linguistics*, 2023.
- 771  
772 YuHong Sun, Zhangyue Yin, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Hui Zhao. Benchmarking  
773 hallucination in large language models based on unanswerable math word problem. In *Proceedings*  
774 *of the 2024 Joint International Conference on Computational Linguistics, Language Resources*  
775 *and Evaluation (LREC-COLING 2024)*, pp. 2178–2188, 2024a.
- 776  
777 Zhiqing Sun, Longhui Yu, Yikang Shen, Weiyang Liu, Yiming Yang, Sean Welleck, and Chuang  
778 Gan. Easy-to-hard generalization: Scalable alignment beyond human supervision. *arXiv preprint*  
779 *arXiv:2403.09472*, 2024b.
- 780  
781 Zhengyang Tang, Xingxing Zhang, Benyou Wan, and Furu Wei. Mathsacle: Scaling instruction  
782 tuning for mathematical reasoning. *arXiv preprint arXiv:2403.02884*, 2024.
- 783  
784 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu  
785 Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable  
786 multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- 787  
788 Trieu H Trinh, Yuhuai Wu, Quoc V Le, He He, and Thang Luong. Solving olympiad geometry  
789 without human demonstrations. *Nature*, 625(7995):476–482, 2024.
- 790  
791 Haoyu Wang, Guozheng Ma, Cong Yu, Ning Gui, Linrui Zhang, Zhiqi Huang, Suwei Ma, Yongzhe  
792 Chang, Sen Zhang, Li Shen, et al. Are large language models really robust to word-level perturba-  
793 tions? *arXiv preprint arXiv:2309.11166*, 2023a.
- 794  
795 Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim.  
796 Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language  
797 models. *arXiv preprint arXiv:2305.04091*, 2023b.
- 798  
799 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu,  
800 Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the  
801 world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.
- 802  
803 Weihang Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang,  
804 Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang.  
805 Cogvlm: Visual expert for pretrained language models, 2023c.
- 806  
807 Zengzhi Wang, Rui Xia, and Pengfei Liu. Generative ai for math: Part i–mathpile: A billion-token-  
808 scale pretraining corpus for math. *arXiv preprint arXiv:2312.17120*, 2023d.
- 809  
810 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny  
811 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*  
812 *neural information processing systems*, 35:24824–24837, 2022.
- 813  
814 Zhaofeng Wu, Linlu Qiu, Alexis Ross, Ekin Akyürek, Boyuan Chen, Bailin Wang, Najoung Kim,  
815 Jacob Andreas, and Yoon Kim. Reasoning or reciting? exploring the capabilities and limitations of  
816 language models through counterfactual tasks. *arXiv preprint arXiv:2307.02477*, 2023.

- 810 Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Li-  
811 juan Wang. The dawn of llms: Preliminary explorations with gpt-4v (ision). *arXiv preprint*  
812 *arXiv:2309.17421*, 9(1):1, 2023.
- 813  
814 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan.  
815 Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural*  
816 *Information Processing Systems*, 36, 2024a.
- 817 Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li,  
818 Weilin Zhao, Zhihui He, Qianyu Chen, Huarong Zhou, Zhensheng Zou, Haoye Zhang, Shengding  
819 Hu, Zhi Zheng, Jie Zhou, Jie Cai, Xu Han, Guoyang Zeng, Dahai Li, Zhiyuan Liu, and Maosong  
820 Sun. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv preprint 2408.01800*, 2024b.
- 821 Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo  
822 Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for  
823 large language models. *arXiv preprint arXiv:2309.12284*, 2023.
- 824  
825 Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhao Chen.  
826 Mammoth: Building math generalist models through hybrid instruction tuning. *arXiv preprint*  
827 *arXiv:2309.05653*, 2023.
- 828 Hugh Zhang, Jeff Da, Dean Lee, Vaughn Robinson, Catherine Wu, Will Song, Tiffany Zhao, Pranav  
829 Raja, Dylan Slack, Qin Lyu, et al. A careful examination of large language model performance on  
830 grade school arithmetic. *arXiv preprint arXiv:2405.00332*, 2024a.
- 831  
832 Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou,  
833 Pan Lu, Kai-Wei Chang, Peng Gao, et al. Mathverse: Does your multi-modal llm truly see the  
834 diagrams in visual math problems? *arXiv preprint arXiv:2403.14624*, 2024b.
- 835 Yifan Zhang, Yifan Luo, Yang Yuan, and Andrew C Yao. Autonomous data selection with language  
836 models for mathematical texts. In *ICLR 2024 Workshop on Navigating and Addressing Data*  
837 *Problems for Foundation Models*, 2024c.
- 838  
839 Kunhao Zheng, Jesse Michael Han, and Stanislas Polu. Minif2f: a cross-system benchmark for  
840 formal olympiad-level mathematics. *arXiv preprint arXiv:2109.00110*, 2021.
- 841 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
842 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and  
843 chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023.
- 844 Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans,  
845 Claire Cui, Olivier Bousquet, Quoc Le, et al. Least-to-most prompting enables complex reasoning  
846 in large language models. *arXiv preprint arXiv:2205.10625*, 2022.
- 847  
848 Kun Zhou, Yutao Zhu, Zhipeng Chen, Wentong Chen, Wayne Xin Zhao, Xu Chen, Yankai Lin,  
849 Ji-Rong Wen, and Jiawei Han. Don’t make your llm an evaluation benchmark cheater. *arXiv*  
850 *preprint arXiv:2311.01964*, 2023a.
- 851 Zihao Zhou, Maizhen Ning, Qiufeng Wang, Jie Yao, Wei Wang, Xiaowei Huang, and Kaizhu Huang.  
852 Learning by analogy: Diverse questions generation in math word problem. In *Findings of the Association*  
853 *for Computational Linguistics: ACL 2023*, pp. 11091–11104. Association for Computational  
854 Linguistics, 2023b. URL <https://aclanthology.org/2023.findings-acl.705>.
- 855 Zihao Zhou, Qiufeng Wang, Mingyu Jin, Jie Yao, Jianan Ye, Wei Liu, Wei Wang, Xiaowei Huang,  
856 and Kaizhu Huang. Mathattack: Attacking large language models towards math solving ability. In  
857 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 19750–19758, 2024.
- 858  
859 Kaijie Zhu, Jiaao Chen, Jindong Wang, Neil Zhenqiang Gong, Diyi Yang, and Xing Xie. Dyval:  
860 Dynamic evaluation of large language models for reasoning tasks. In *The Twelfth International*  
861 *Conference on Learning Representations*, 2024a. URL [https://openreview.net/forum?](https://openreview.net/forum?id=gjfoL9z5Xr)  
862 [id=gjfoL9z5Xr](https://openreview.net/forum?id=gjfoL9z5Xr).
- 863  
864 Kaijie Zhu, Jindong Wang, Qinlin Zhao, Ruochen Xu, and Xing Xie. Dyval 2: Dynamic evaluation  
865 of large language models by meta probing agents. *arXiv preprint arXiv:2402.14865*, 2024b.



864	APPENDIX	
865		
866	<b>A Heatmap of MATHCHECK-GSM</b>	<b>18</b>
867		
868	<b>B Heatmap of MATHCHECK-GEO</b>	<b>20</b>
869		
870		
871	<b>C Data Statistics and Quality</b>	<b>21</b>
872	C.1 Overview of Data . . . . .	21
873	C.2 Effectiveness of GPT-4-turbo Rewriting . . . . .	22
874	C.3 Discussion of data bias generated by GPT . . . . .	22
875		
876		
877	<b>D Evaluation Setup</b>	<b>23</b>
878		
879		
880	<b>E MATHCHECK Applied to Other Reasoning Tasks</b>	<b>24</b>
881	E.1 Date Understanding . . . . .	24
882	E.2 Code Generation . . . . .	25
883		
884		
885	<b>F Prompt List</b>	<b>26</b>
886	F.1 Evaluation Prompt . . . . .	26
887	F.2 Data Generation Prompt . . . . .	28
888		
889		
890	<b>G Case Problems</b>	<b>31</b>
891	G.1 Case Problems in MATHCHECK-GSM. Problem Group ID: GSM-54 . . . . .	31
892	G.2 Case Problems in MATHCHECK-GEO. Problem Group ID: GEO-15 . . . . .	37
893		
894		
895		
896		
897		
898		
899		
900		
901		
902		
903		
904		
905		
906		
907		
908		
909		
910		
911		
912		
913		
914		
915		
916		
917		

## A HEATMAP OF MATHCHECK-GSM

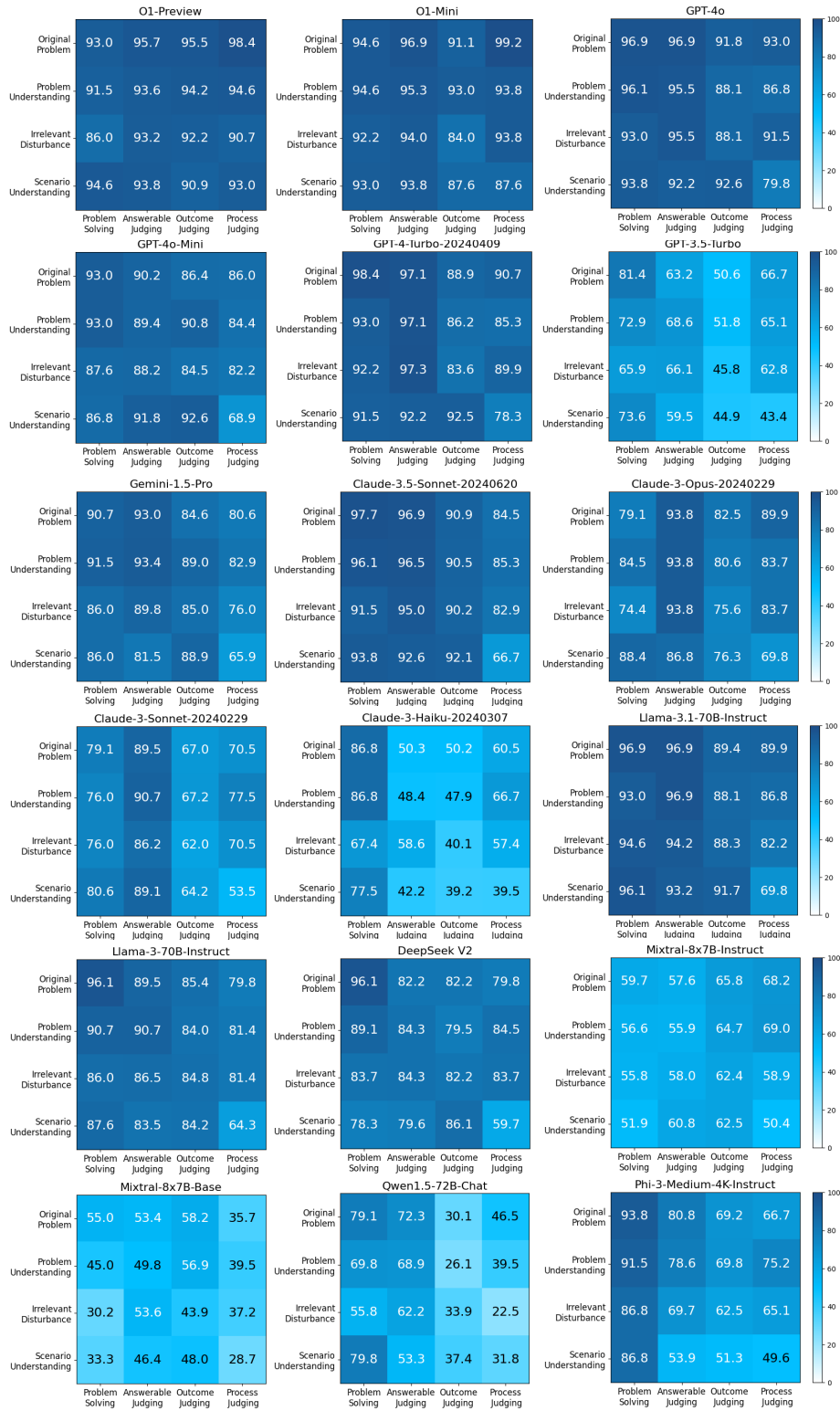


Figure 8: Visualized heatmap of MATHCHECK-GSM - Part 1.

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

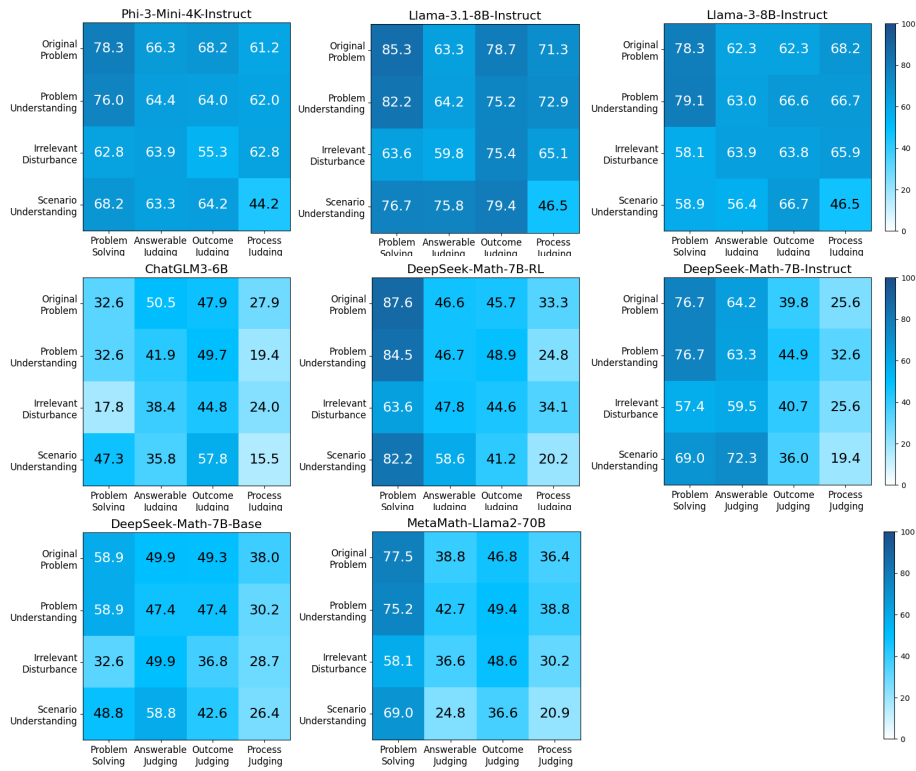


Figure 9: Visualized heatmap of MATHCHECK-GSM - Part 2.

## B HEATMAP OF MATHCHECK-GEO

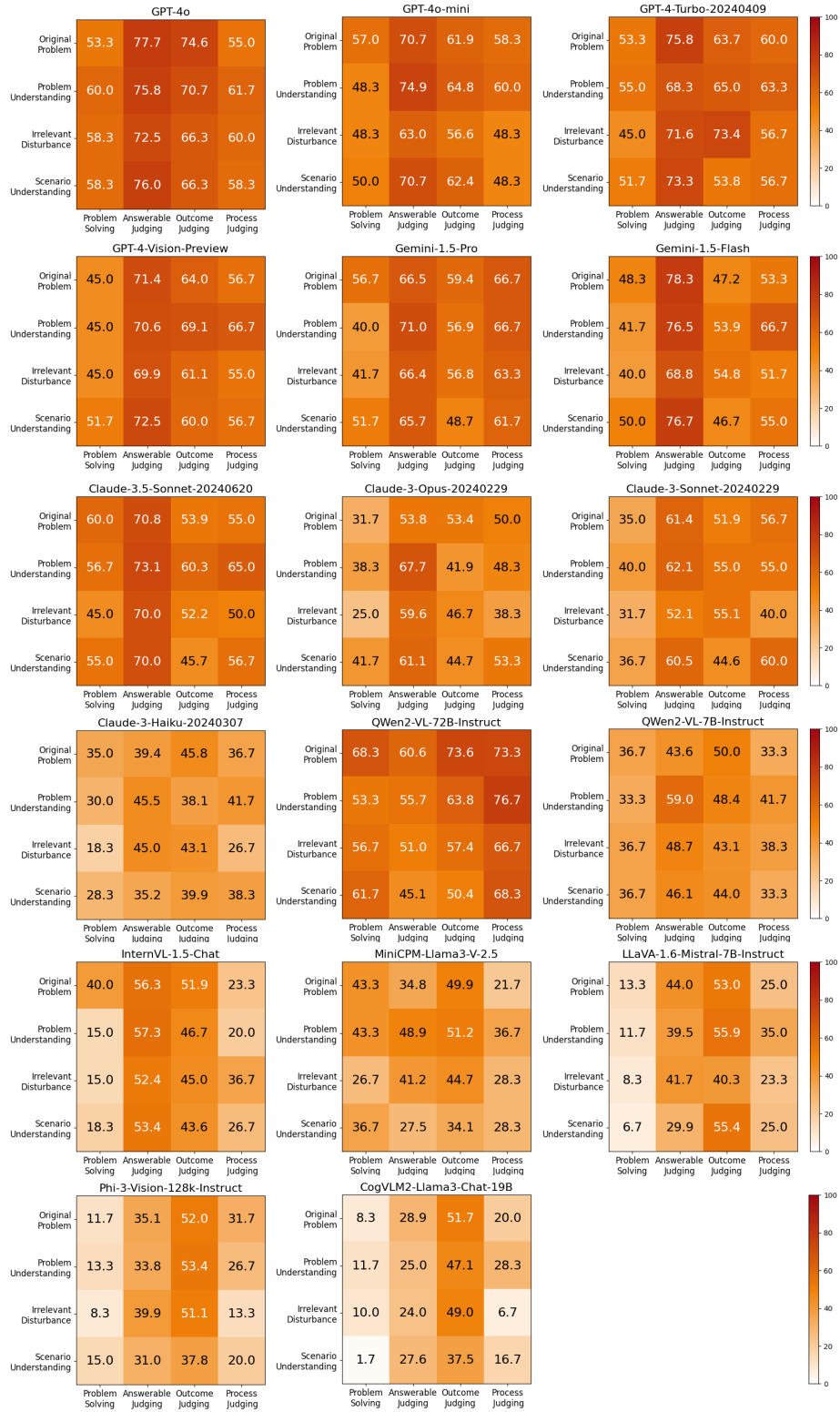


Figure 10: The visualized heatmap of MATHCHECK-GEO.

## C DATA STATISTICS AND QUALITY

### C.1 OVERVIEW OF DATA

Table 3 and Table 4 show the data statistics of MATHCHECK-GSM and MATHCHECK-GEO. Table 5 shows the data statistics of each group in MATHCHECK-GSM and MATHCHECK-GEO. In each group, since answerable judging and outcome judging are binary-classification tasks, we try our best to include two different labels in these units for fair evaluation.

Table 3: Data statistics of MATHCHECK-GSM

	<b>Problem Solving</b>	<b>Answerable Judging</b>	<b>Outcome Judging</b>	<b>Process Judging</b>
<b>Original Problem</b>	129	258	258	129
<b>Problem Understanding</b>	129	258	258	129
<b>Irrelevant Disturbance</b>	129	258	258	129
<b>Scenario Understanding</b>	129	258	258	129

Table 4: Data statistics of MATHCHECK-GEO

	<b>Problem Solving</b>	<b>Answerable Judging</b>	<b>Outcome Judging</b>	<b>Process Judging</b>
<b>Original Problem</b>	60	120	120	60
<b>Problem Understanding</b>	60	120	120	60
<b>Irrelevant Disturbance</b>	60	120	120	60
<b>Scenario Understanding</b>	60	120	120	60

Table 5: Data statistics of each group in MATHCHECK-GSM and MATHCHECK-GEO

	<b>Problem Solving</b>	<b>Answerable Judging</b>	<b>Outcome Judging</b>	<b>Process Judging</b>
<b>Original Problem</b>	1	2	2	1
<b>Problem Understanding</b>	1	2	2	1
<b>Irrelevant Disturbance</b>	1	2	2	1
<b>Scenario Understanding</b>	1	2	2	1

## C.2 EFFECTIVENESS OF GPT-4-TURBO REWRITING

In the process of human evaluation, we selected three graduate students as human annotators, all of them possess the mathematical skills required for evaluating the generated data. Our human evaluation principle is that the generated mathematical problems should maintain the correctness of mathematical logic. For example, in the “Problem Understanding”, the generated question should not alter the logical structure of original question, which ensures the consistency between rewritten question and answer. The generated data will be marked as a failure if any of annotators determines that the generation failed. Furthermore, annotators corrected each failed data instead of discarding them. This approach ensures our dataset is entirely accurate and the evaluation results are reliable.

We conduct statistics on the pass rate of MATHCHECK-GSM rewritten by GPT4-turbo, as shown in Table 6. It can be seen that the rewriting pass rate is high, which reflects the effectiveness of our generation method. The success rate of Problem Understanding and Scenario Understanding is higher than 90%. There is a pass rate of 86.82% in the Irrelevant Disturbance and 81.40% in Wrong Step Rewriting. It provides references when we use MATHCHECK generation.

Table 6: Pass rate (%) checked by human annotators for the data generated by GPT4-turbo.

Rewriting Type	Problem Understanding	Irrelevant Disturbance	Scenario Understanding	Unanswerable Question Rewriting	Wrong Step Rewriting
Human Pass Rate	93.02	86.82	91.47	85.38	81.40

## C.3 DISCUSSION OF DATA BIAS GENERATED BY GPT

While we acknowledge there are possible self-bias in LLM-rewritten questions, we assert that this bias is acceptable and does not undermine the conclusions or rationality of MATHCHECK. This is supported by considerations across several dimensions.

**Motivations.** The motivation behind MATHCHECK is to establish a paradigm that mitigates benchmark hacking in the evaluation of mathematical reasoning, thereby revealing the genuine mathematical reasoning abilities of language models more comprehensively. Rewriting is an integral part of the MATHCHECK pipeline, which can naturally be performed by either humans or LLMs. While we acknowledge that involving experts in the rewriting process might be the fairest approach, the scalability of this method is a significant concern, as noted in several of today’s LLM benchmarks, such as Arena Hard (Li et al., 2024c) and MT-Bench (Zheng et al., 2023), due to the high associated costs. To enhance scalability and practicality, we opted to use LLMs as the rewriters. Given that GPT-4 is widely recognized as the most advanced model accessible to the public, we believe that choosing GPT-4 as the rewriter is the closest approximation to the quality of expert human rewriting.

**Human-Checked Questions.** In fact, for the data construction which the LLM participates in, we mainly utilize the powerful rewriting ability of LLMs to edit the seed math problem instead of generating a new one from scratch. Moreover, we manually check the generated text to avoid some unnatural generated text.

**Experimental Results and Analysis.** On one hand, although the data are generated by GPT-4-Turbo in our experiments, they do not bring extra benefits to GPT-Family models to make them obviously outperform others. As shown in Table 1, the performance of Claude-3.5-sonnet is similar with GPT-4-Turbo, and even much better than GPT-4o-mini, which follows the commonsense on these LLMs. On the other hand, we compare the experimental results on Non-GPT-Rewritten and GPT-Rewritten Questions. In some data constructions where the LLM is not involved, GPT4-family exhibits the same performance ranking as the score “All”. Specifically, the samples in Original Problem&Outcome Judging (OP-OJ) belong to Non-GPT-Rewritten Questions, which are generated based on the rules. Table 7 shows that the performance ranking on non-LLM-generated data is close to the score “All”, where GPT-series continues to perform better than other advanced models. All of these results verify that the possible bias to GPT models is acceptable in our MATHCHECK.

Table 7: Model performance on Non-GPT-Rewritten Questions of MATHCHECK-GSM

Models	All	OP-OJ
GPT-4o	92.0	<b>91.8</b>
GPT-4-Turbo-20240409	90.9	<b>88.9</b>
Gemini-1.5-Pro	86.3	84.6
Claude-3-Opus-20240229	83.5	82.5
Llama-3-70B-Instruct	84.7	85.4

## D EVALUATION SETUP

We conduct evaluations of multiple representative generalist and mathematical models on our MATHCHECK benchmark. For MATHCHECK-GSM, the evaluation models encompass: (a) Generalist models, including proprietary models such as O1-Preview (OpenAI, 2024e), O1-Mini (OpenAI, 2024d), GPT-4o (OpenAI, 2024a), GPT-4o-mini (OpenAI, 2024b), GPT-4-Turbo (Achiam et al., 2023), GPT-3.5-Turbo (OpenAI, 2022), Gemini-1.5-Pro (Team et al., 2023), Claude-3 (Anthropic, 2024a), Claude-3.5-Sonnet Anthropic (2024b), Llama-3<sup>3</sup>, Llama-3.1<sup>4</sup>, DeepSeek V2 (Shao et al., 2024), Mixtral 8 x 7B (Jiang et al., 2024), Qwen1.5 (Bai et al., 2023), Phi-3 (Abdin et al., 2024), and ChatGLM3 (Du et al., 2022); (b) Mathematical models, including DeepSeek-Math (Shao et al., 2024) and MetaMath (Yu et al., 2023). For MATHCHECK-GEO, we conduct evaluations on generalist models: (a) proprietary models such as GPT-4o (OpenAI, 2024a), GPT-4o-mini (OpenAI, 2024b), GPT-4-Turbo (Achiam et al., 2023), GPT-4-vision (OpenAI, 2024c), Gemini-1.5-Pro (Team et al., 2023), Claude-3.5-Sonnet Anthropic (2024b) and Claude-3 (Anthropic, 2024a); (b) open-source models including Qwen2-VL (Wang et al., 2024), InternVL-1.5 (Chen et al., 2023c), Phi-3-Vision (Abdin et al., 2024), LLaVA-1.6-Mistral-7B-Instruct (Liu et al., 2024b), MiniCPM-Llama3-V-2.5 (Yao et al., 2024b) and CogVLM2-Llama3 (Wang et al., 2023c).

For Problem Solving and Process Judging tasks, we employ accuracy as the evaluation measure. For Outcome Judging and Answerable Judging tasks, we utilize Macro-F1 as the metric. We employ a zero-shot setting for generalist models and a few-shot setting (two-shot) for base models and mathematical models to enhance their ability to follow specific instructions and tasks. All the prompts used for evaluating (M)LLMs are provided in Appendix F.1.

For all the close-resourced models, we utilize the default hyper-parameters, setting the temperature to 0 and the max tokens to 1,024. Similarly, for all open-source models, the parameters are uniformly configured as follows: *do\_sample* is set to False, *max\_gen\_len* is set to 512, and the temperature is set to 0.1.

<sup>3</sup><https://ai.meta.com/blog/meta-llama-3>

<sup>4</sup><https://ai.meta.com/blog/meta-llama-3-1>

## E MATHCHECK APPLIED TO OTHER REASONING TASKS

	Solving	Answerable Judging	Outcome Judging	Process Judging	
Original Problem	Yesterday's date was 4/30/2021. What is the date tomorrow in MM/DD/YYYY? "answer": 5/2/2021 ★ Seed Data	Yesterday's date was 4/30. What is the date tomorrow in MM/DD/YYYY? "answer": Unanswerable	Yesterday's date was 4/30/2021. What is the date tomorrow in MM/DD/YYYY? "solution": If yesterday was 4/30/2021, then tomorrow would be 5/02/2021. "answer": Correct	Yesterday's date was 4/30/2021. What is the date tomorrow in MM/DD/YYYY? "solution": Step 1: ... Step 2: Determine the next day. Since the current date is the last day of April ... The answer is 05/01/2021. "answer": Step 2	Reasoning Robustness
Problem Understanding	Yesterday was April 30, 2021. What is the date tomorrow in MM/DD/YYYY? "answer": 5/2/2021	Yesterday was April, 2021. What is the date tomorrow in MM/DD/YYYY? "answer": Unanswerable	Yesterday was April 30, 2021. What is the date tomorrow in MM/DD/YYYY? "solution": If ... then tomorrow would be May 1, 2021, so the date in MM/DD/YYYY format is 05/01/2021. "answer": Incorrect	Yesterday was April 30, 2021. What is the date tomorrow in MM/DD/YYYY? "solution": Step 1: If yesterday was 4/30/2021, then tomorrow is May 1, 2021. Step 2: ... The answer is 05/01/2021. "answer": Step 1	
Irrelevant Disturbance	Yesterday was April 30, 2021. A week ago it was 4/23/2021. What is the date tomorrow in MM/DD/YYYY? "answer": 5/2/2021	Yesterday was April 30, 2021. A week ago it was 4/23/2021. What is the date tomorrow in MM/DD/YYYY? "answer": Answerable	Yesterday was April 30, 2021. A week ago it was 4/23/2021. What is the date tomorrow in MM/DD/YYYY? "solution": The date tomorrow will be 05/01/2021. "answer": Incorrect	Yesterday was April 30, 2021. A week ago it was 4/23/2021. What is the date tomorrow in MM/DD/YYYY? "solution": Step 1: ... Step 2: ... Step 3: Since we are moving forward by one day from April 30th, we add one day to the date. Step 4 :... The answer is 05/01/2021. "answer": Step 3	
Scenario Understanding	Yesterday was April x, 2021. The date tomorrow is 5/2/2021. What's the value of x? "answer": 30	Yesterday was April x, 2021. The date is 5/2/2021. What's the value of x? "answer": Unanswerable	Yesterday was April x, 2021. The date tomorrow is 5/2/2021. What's the value of x? "solution": If tomorrow is May 2, 2021, then ...So, the value of x is 30 "answer": Correct	Yesterday was April x, 2021. The date tomorrow is 5/2/2021. What's the value of x? "solution": Step1: ... Step 2: ... Therefore, yesterday would be May 1st, 2021. The value of x is 1. "answer": Step 2	

Task Generalization

Figure 11: Case of MATHCHECK in Date Understanding.

## E.1 DATE UNDERSTANDING

To show that our proposed benchmark paradigm MATHCHECK can be adapted to other reasoning tasks beyond mathematical problems, we try to transform some representative reasoning task into MATHCHECK paradigm. We firstly apply it in commonsense reasoning, which requires LLMs to apply world knowledge to reason and solve problems. Specifically, we choose the date understanding task in Big-bench (bench authors, 2023) since it is a wildly used task to measure commonsense reasoning ability (Wei et al., 2022).

Figure 11 shows the case of applying MATHCHECK to date understanding. Similar to mathematical reasoning, date understanding is a numerical reasoning task, therefore it can easily utilize variants of each unit in MATHCHECK. For example, in Irrelevant Disturbance, we can add some irrelevant date conditions to cause disturbance. In scenario understanding, we can ask for other variables in order to examine whether models have a comprehensive understanding of this date knowledge. This case demonstrates the high adaptability of MATHCHECK to commonsense reasoning task especially numerical reasoning.



	Solving	Answerable Judging	Outcome Judging	Process Judging					
1296									
1297									
1298	Original Problem Write a function in python that takes string and returns string without numbers. "answer": <pre>def remove_num(text):     text_without_nums = ""     for char in text:         if not char.isdigit():             text_without_nums += char     return text_without_nums</pre> Seed Data	Answerable Judging Write a function in python that takes string and returns string without some specific chars. "answer": Unanswerable	Outcome Judging Write a function in python that takes string and returns string without numbers. "solution": <pre>def remove_num(text):     text_without_nums = ""     for char in text:         if not char.isdigit():             text_without_nums += char     return text_without_nums</pre> "answer": Correct	Process Judging Write a function in python that that takes string and returns string without numbers. "solution": <pre>1 def remove_num(text): 2     text_without_nums = "" 3     for char in text: 4         if not char.isdigit(): 5             text_without_nums = char 6     return text_without_nums</pre> "answer": Step 5	Reasoning Robustness				
1302									
1303						Problem Understanding Write a python function that takes a string and returns it without a number "answer": <pre>def remove_num(text):     text_without_nums = ""     for char in text:         if not char.isdigit():             text_without_nums += char     return text_without_nums</pre>	Answerable Judging Write a python function that takes a string and returns it without a number. "answer": Answerable	Outcome Judging Write a python function that takes a string and returns it without a number. "solution": <pre>def remove_num(text):     text_without_nums = ""     for char in text:         if not char.isdigit():             text_without_nums += char     return text</pre> "answer": Incorrect	Process Judging Write a python function that takes a string and returns it without a number. "solution": <pre>1 def remove_num(text): 2     text_without_nums = "" 3     for char in text: 4         if not char.isdigit(): 5             text_without_nums += char 6     return text</pre> "answer": Step 6
1304									
1305									
1306									
1307	Irrelevant Disturbance Write a python function that takes a string containing letters, numbers, symbols, etc. and returns the string without the numbers. "answer": <pre>def remove_num(text):     text_without_nums = ""     for char in text:         if not char.isdigit():             text_without_nums += char     return text_without_nums</pre>	Answerable Judging Write a python function that takes a string containing letters, numbers, symbols, etc. and returns the string without the numbers. "answer": Answerable	Outcome Judging Write a python function that takes a string containing letters, numbers, symbols, etc. and returns the string without the numbers. "solution": <pre>def remove_num(text):     text_without_nums = ""     for char in text:         if char.isdigit():             text_without_nums += char     return text_without_nums</pre> "answer": Incorrect	Process Judging Write a python function that takes a string containing letters, numbers, symbols, etc. and returns the string without the numbers. "solution": <pre>1 def remove_num(text): 2     text_without_nums = "" 3     for char in text_without_nums: 4         if not char.isdigit(): 5             text_without_nums += char 6     return text_without_nums</pre> "answer": Step 3					
1308									
1309									
1310									
1311	Scenario Understanding Write a function in java that takes string and returns string without numbers. "answer": <pre>public static String removeNums(String input) {     String text = input.replaceAll("\\d", "");     return text }</pre>	Answerable Judging Write a function in java that takes string and returns string without "answer": Unanswerable	Outcome Judging Write a function in java that takes string and returns string without numbers. "solution": <pre>public static String removeNums(String input) {     String text = input.replaceAll("\\d", "");     return text; }</pre> "answer": Correct	Process Judging Write a function in java that takes string and returns string without numbers. "solution": <pre>1 public static String removeNums(String input) 2 { 3     String text = input.replaceAll("", "\\d"); 4     return text; 5 }</pre> "answer": Step 3					
1312									
1313									
1314									
1315	Task Generalization								
1316									
1317									
1318									
1319									
1320									
1321									
1322									
1323									
1324									
1325									
1326									
1327									

Figure 12: Case of MATHCHECK in Code Generation.

## E.2 CODE GENERATION

In addition to commonsense reasoning task, we would like to show the possibility of transforming MATHCHECK in some real-world reasoning tasks. Specifically, we choose the code generation task due to its high relevance to Text2Sql, agents and robotics. Figure 12 demonstrates a case of applying MATHCHECK to code generation. Unlike numerical reasoning tasks, the adaptation of code generation needs to consider task relevance. For example, in Scenario Understanding, we can ask models to write the same function in other program languages (Python to Java in our case) in order to examine whether models have a comprehensive understanding of this function requirements. It shows that MATHCHECK have potential for real-world tasks such as agents and robotics application. Meanwhile, we encourage researchers to design more specific variants towards their reasoning task on MATHCHECK framework to test reasoning robustness and task generalization.

## 1350 F PROMPT LIST

1351

## 1352 F.1 EVALUATION PROMPT

1353

1354 You are an AI assistant that determines whether math problems are solved  
 1355 correctly. Answer the question. Finally give the answer in the format:  
 1356 The answer is: ...

1357

1358 Question: [QUESTION]

1359 Answer:

1360

## 1361 1: Zero-shot Prompt of Problem Solving

1362

1362 You are an AI assistant that determines whether math problems are solved  
 1363 correctly. I will first give you a math problem and its solution, help me  
 1364 judge whether the final answer is correct or incorrect. Give your  
 1365 judgment between Correct or Incorrect. Finally summarize your answer in  
 1366 the format:  
 1367 The answer is: ...

1368

1369 Question: [QUESTION]

1370 Solution: [SOLUTION]

1371 Judgement:

1372

## 1373 2: Zero-shot Prompt of Outcome Judging

1374

1374 You are an AI assistant that identify which step begins the error in  
 1375 solution. I will give you a math problem along with a wrong solution.  
 1376 Please help me identify the step where the errors begin. Finally give the  
 1377 wrong step in the format:  
 1378 The answer is: Step i

1379

1380 Question: [QUESTION]

1381 Solution: [SOLUTION]

1382 Judgement:

1383

## 1384 3: Zero-shot Prompt of Process Judging

1385

1384 You are an AI assistant that determines whether math problems are  
 1385 answerable or unanswerable. Please analyze whether the question provides  
 1386 sufficient information to obtain an answer. Give your judgment between  
 1387 Answerable or Unanswerable. Finally summarize your answer in the format:  
 1388 The answer is: ...

1389

1390 Question: [QUESTION]

1391 Judgement:

1392

## 1393 4: Zero-shot Prompt of Answerable Judging

1394

1394 You are an AI assistant to help me solve math problems. Answer the  
 1395 question. Finally give the answer in the format: The answer is: ...  
 1396 Follow the given examples and answer the question.

1397

1397 Question: Leah had 32 chocolates and her sister had 42. If they ate 35,  
 1398 how many pieces do they have left in total?

1399

1399 Answer: Step 1: Originally, Leah had 32 chocolates.

1400

1400 Step 2: Her sister had 42. So in total they had  $32 + 42 = 74$ .

1401

1401 Step 3: After eating 35, they had  $74 - 35 = 39$ .

1402

1402 The answer is 39.

1403

1403 Question: Jason had 20 lollipops. He gave Denny some lollipops. Now Jason  
 has 12 lollipops. How many lollipops did Jason give to Denny?

1404 Answer: Step 1: Jason started with 20 lollipops.  
1405 Step 2: Then he had 12 after giving some to Denny.  
1406 Step 3: So he gave Denny  $20 - 12 = 8$ .  
1407 The answer is 8.

1408  
1409 Question: [QUESTION]  
1410 Answer:  
1411

### 5: Few-shot Prompt of Problem Solving

1412  
1413  
1414 You are an AI assistant that determines whether math problems are solved  
1415 correctly. I will first give you a math problem and its solution, help me  
1416 judge whether the final answer is correct or incorrect.

1417  
1418 Give your judgment between Correct or Incorrect. Finally summarize your  
1419 answer in the format: The answer is: ...  
1420 Follow the given examples and give your judgment.

1421 Question: Leah had 32 chocolates and her sister had 42. If they ate 35,  
1422 how many pieces do they have left in total?

1423 Solution: Step 1: Originally, Leah had 32 chocolates.  
1424 Step 2: Her sister had 42. So in total they had  $32 + 42 = 74$ .  
1425 Step 3: After eating 35, they had  $74 - 35 = 39$ .  
1426 The answer is 39.

1427 Judgment: Step 1 and Step 2 accurately calculate the total number of  
1428 chocolates they both had originally.  
1429 Step 3 correctly calculates how many they have left after eating 35  
1430 chocolates.  
1431 The answer is: Correct.

1432 Question: Jason had 20 lollipops. He gave Denny some lollipops. Now Jason  
1433 has 12 lollipops. How many lollipops did Jason give to Denny?

1434 Solution: Step 1: Jason started with 20 lollipops.  
1435 Step2: Then he had 12 after giving some to Denny.  
1436 Step3: So he gave Denny  $20 - 12 = 8$ .  
1437 The answer is 32.

1438 Judgment: Jason ended up with 12 lollipops after giving some to Denny,  
1439 having started with 20. Therefore, the calculation to find out how many  
1440 lollipops Jason gave to Denny should be:  $20 - 12 = 8$ .  
1441 The answer is: Incorrect.

1442 Question: [QUESTION]  
1443 Solution: [SOLUTION]  
1444 Judgement:

### 6: Few-shot Prompt of Outcome Judging

1445  
1446  
1447 You are an AI assistant that identify which step begins the error in  
1448 solution. I will give you a math problem along with a wrong solution.  
1449 Please help me identify the step where the errors begin.

1450 Finally give the wrong step in the format: The answer is: Step I  
1451 Follow the given examples and give your judgment.

1452  
1453 Question: Leah had 32 chocolates and her sister had 42. If they ate 35,  
1454 how many pieces do they have left in total?

1455 Solution: Step 1: Originally, Leah had 32 chocolates.  
1456 Step 2: Her sister had 42. So in total they had  $32 + 42 = 84$ .  
1457 Step 3: After eating 35, they had  $84 - 35 = 49$ .  
The answer is 49.

1458 Judgment: The judgment of the given steps is as follows:  
1459 Step 1: Correctly states Leah's initial amount of chocolates.

1458 Step 2: Incorrectly calculates the total number of chocolates both Leah  
 1459 and her sister had originally.  
 1460 The answer is: Step 2.  
 1461  
 1462 Question: Jason had 20 lollipops. He gave Denny some lollipops. Now Jason  
 1463 has 12 lollipops. How many lollipops did Jason give to Denny?  
 1464 Solution: Step 1: Jason started with 20 lollipops.  
 1465 Step 2: Then he had 12 after giving some to Denny.  
 1466 Step 3: So he gave Denny  $20 + 12 = 8$ .  
 1467 The answer is 32.  
 1468 Judgment: The correct method to find out how many lollipops Jason gave to  
 1469 Denny would be to subtract the amount he had left from the amount he  
 1470 started with:  $20 - 12 = 8$ . Thus, The reasoning error begins at Step 3.  
 1471 The answer is: Step 3.  
 1472  
 1473 Question: [QUESTION]  
 1474 Solution: [SOLUTION]  
 1475 Judgement:

### 7: Few-shot Prompt of Process Judging

1476  
 1477  
 1478 You are an AI assistant that determines whether math problems are  
 1479 answerable or unanswerable. Please analyze whether the question provides  
 1480 sufficient information to obtain an answer.  
 1481  
 1482 Give your judgment between Answerable or Unanswerable. Finally summarize  
 1483 your answer in the format: The answer is: ...  
 1484 Follow the given examples and give your judgment.  
 1485  
 1486 Question: Leah had 32 chocolates and her sister had 42. If they ate 35,  
 1487 how many pieces do they have left in total?  
 1488 Judgment: The question provides all necessary information to perform the  
 1489 calculation.  
 1490 The answer is: Answerable.  
 1491  
 1492 Question: Jason had 20 lollipops. He gave Denny some lollipops. How many  
 1493 lollipops did Jason give to Denny?  
 1494 Judgment: The question is not answerable as given. The reason is that  
 1495 there is insufficient information to determine the exact number of  
 1496 lollipops Jason gave to Denny.  
 1497 The answer is: Unanswerable.  
 1498  
 1499 Question: [QUESTION]  
 1500 Judgement:

### 8: Few-shot Prompt of Answerable Judging

1501  
 1502  
 1503  
 1504  
 1505  
 1506  
 1507  
 1508  
 1509  
 1510  
 1511

## F.2 DATA GENERATION PROMPT

Your objective is to rewrite a given math question using the following  
 perturbation strategy. The rewritten question should be reasonable,  
 understandable, and able to be responded to by humans.  
 Perturbation strategy: Problem Understanding: It refers to transforming  
 the original problem into a new problem that uses different wording or  
 different sentence structures but does not change the solution of the  
 original problem.  
 The given question: {QUESTION}  
 Answer of the given question: {ANSWER}

1512 Please rewrite the question using the specified perturbation strategy  
1513 while minimizing edits to avoid significant deviation in the question  
1514 content.  
1515 It is important to ensure that the rewritten question has only one  
1516 required numerical answer. You just need to print the rewritten question  
1517 without answer.  
1518 The rewritten question:  
1519 Question: {QUESTION}  
1520 Answer: {ANSWER}  
1521 Given step: {STEP}  
1522 The rewritten answer:

### 9: Prompt of Problem Understanding Rewriting

1524 Your objective is to rewrite a given math question using the following  
1525 perturbation strategy. The rewritten question should be reasonable,  
1526 understandable, and able to be responded to by humans.  
1527  
1528 Perturbation strategy: Irrelevant Disturbance: It involves introducing  
1529 distracting conditions that have no impact on the final answer. These  
1530 introduced conditions should be relevant to the topic of the original  
1531 question and preferably include numerical values. However, the rewritten  
1532 problem must maintain an identical solution to that of the original  
1533 problem.  
1534 The given question: {QUESTION}  
1535 Answer of the given question: {ANSWER}  
1536  
1537 Please rewrite the question using the specified perturbation strategy  
1538 while minimizing edits to avoid significant deviation in the question  
1539 content.  
1540 It is important to ensure that the rewritten question has only one  
1541 required numerical answer. You just need to print the rewritten question  
1542 without answer.  
1543 The rewritten question:  
1544 Question: {QUESTION}  
1545 Answer: {ANSWER}  
1546 Given step: {STEP}  
1547 The rewritten answer:

### 10: Prompt of Irrelevant Disturbance Rewriting

1548 Your objective is to rewrite a given math question using the following  
1549 perturbation strategy. The rewritten question should be reasonable,  
1550 understandable, and able to be responded to by humans.  
1551  
1552 Perturbation strategy: Unanswerable question: It refers to eliminating a  
1553 condition from the original question that is crucial for solving it while  
1554 keeping the rest of the content unchanged. The rewritten problem should  
1555 no longer have a valid answer, as it lacks the constraint that was  
1556 removed.  
1557 The given question: {QUESTION}  
1558 Answer of the given question: {ANSWER}  
1559  
1560 Please rewrite the question using the specified perturbation strategy  
1561 while minimizing edits to avoid significant deviation in the question  
1562 content.  
1563 It is important to ensure that the rewritten question has only one  
1564 required numerical answer. You just need to print the rewritten question  
1565 without answer.  
1566 The rewritten question:  
1567 Question: {QUESTION}  
1568 Answer: {ANSWER}

1566 Given step: {STEP}  
1567 The rewritten answer:

### 11: Prompt of Unanswerable Question Rewriting

1571 You are an AI assistant to help me rewrite question into a declarative  
1572 statement when its answer is provided.  
1573 Follow the given examples and rewrite the question.

1574 Question: How many cars are in the parking lot? The answer is 5.  
1575 Result: There are 5 cars in the parking lot.

1576

1577 Question: How many trees did the grove workers plant today? The answer is  
1578 6.  
1579 Result: The grove workers planted 6 trees today.

1580

1581 Question: If they ate 35, how many pieces do they have left in total? The  
1582 answer is 39.  
1583 Result: They have 39 pieces left in total if they ate 35.

1584

1585 Question: How many lollipops did Jason give to Denny? The answer is 8.  
1586 Result: Jason gave 8 lollipops to Denny.

1587

1588 Question: How many toys does he have now? The answer is 9.  
1589 Result: He now has 9 toys.

1590

1591 Question: How many computers are now in the server room? The answer is  
1592 29.  
1593 Result: There are 29 computers now in the server room.

1594

1595 Question: How many golf balls did he have at the end of wednesday? The  
1596 answer is 33.  
1597 Result: He had 33 golf balls at the end of Wednesday.

1598

1599 Question: How much money does she have left? The answer is 8.  
1600 Result: She has 8 money left.

1601

1602 Question: {QUESTION} The answer is {ANSWER}.  
1603 Result:

### 12: Prompt to Rewrite Question and Answer into a Declarative Statement

1604 Following is a question and its correct solution. Rewrite the solution  
1605 according to following requirements: (1) Do not change the format (2)  
1606 Keep those steps before the given step unchanged (3) Make minor changes  
1607 to the given step so that the reasoning of this step and subsequent steps  
1608 are incorrect, resulting in an incorrect answer.

1609

1610 Question: {QUESTION}  
1611 Answer: {ANSWER}  
1612 Given step: {STEP}  
1613 The rewritten answer:

### 13: Prompt to Generate the Wrong Step

1614  
1615  
1616  
1617  
1618  
1619

1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673

## G CASE PROBLEMS

### G.1 CASE PROBLEMS IN MATHCHECK-GSM. PROBLEM GROUP ID: GSM-54

[Question]: Mike plays ping pong for 40 minutes. In the first 20 minutes, he scores 4 points. In the second 20 minutes, he scores 25% more points. How many total points did he score?  
[Answer]: 9.0

#### 14: Problem Solving - Original Problem

[Question]: During a 40-minute ping pong session, Mike scores 4 points in the initial half. In the latter half, he manages to increase his score by 25% compared to the first half. What is the total score Mike achieved in this session?  
[Answer]: 9.0

#### 15: Problem Solving - Problem Understanding

[Question]: Mike plays ping pong in a local tournament and decides to practice for 40 minutes before the first match. During his practice session, in the first 20 minutes, while intermittently checking his phone and hydrating, he manages to score 4 points. In the following 20 minutes, feeling more warmed up and despite a short break to adjust his paddle's grip tape, he scores 25% more points than in the first session. Considering these distractions, how many total points did Mike score in his 40-minute practice session?  
[Answer]: 9.0

#### 16: Problem Solving - Irrelevant Disturbance

[Question]: Mike plays ping pong for 40 minutes. In the first 20 minutes, he scores  $x$  points. In the second 20 minutes, he scores 25% more points. He scored 9 total points. What is the value of unknown variable  $x$ ?  
[Answer]: 4.0

#### 17: Problem Solving - Scenario Understanding

[Question]: Mike plays ping pong for 40 minutes. In the first 20 minutes, he scores 4 points. In the second 20 minutes, he scores 25% more points. How many total points did he score?  
[Answer]: Answerable

#### 18: Answerable Judging (*Answerable*) - Original Problem

[Question]: Mike plays ping pong for minutes. In the first 20 minutes, he scores 4 points. In the second 20 minutes, his performance increases by 25%. How many total points did he score?  
[Answer]: Unanswerable

#### 19: Answerable Judging (*Unanswerable*) - Original Problem

[Question]: During a 40-minute ping pong session, Mike scores 4 points in the initial half. In the latter half, he manages to increase his score by 25% compared to the first half. What is the total score Mike achieved in this session?  
[Answer]: Answerable

#### 20: Answerable Judging (*Answerable*) - Problem Understanding

1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727

[Question]: During a 40-minute ping pong session, Mike scores points in the initial half. In the latter half, he manages to increase his score by 25% compared to the first half. What is the total score Mike achieved in this session?

[Answer]: Unanswerable

**21: Answerable Judging (*Unanswerable*) - Problem Understanding**

[Question]: Mike plays ping pong in a local tournament and decides to practice for 40 minutes before the first match. During his practice session, in the first 20 minutes, while intermittently checking his phone and hydrating, he manages to score 4 points. In the following 20 minutes, feeling more warmed up and despite a short break to adjust his paddle's grip tape, he scores 25% more points than in the first session.

Considering these distractions, how many total points did Mike score in his 40-minute practice session?

[Answer]: Answerable

**22: Answerable Judging (*Answerable*) - Irrelevant Disturbance**

[Question]: Mike plays ping pong in a local tournament and decides to practice for 40 minutes before the first match. During his practice session, in the first 20 minutes, while intermittently checking his phone and hydrating, he manages to score points. In the following 20 minutes, feeling more warmed up and despite a short break to adjust his paddle's grip tape, he scores 25% more points than in the first session.

Considering these distractions, how many total points did Mike score in his 40-minute practice session?

[Answer]: Unanswerable

**23: Answerable Judging (*Unanswerable*) - Irrelevant Disturbance**

[Question]: Mike plays ping pong for 40 minutes. In the first 20 minutes, he scores  $x$  points. In the second 20 minutes, he scores 25% more points. He scored 9 total points. What is the value of unknown variable  $x$ ?

[Answer]: Answerable

**24: Answerable Judging (*Answerable*) - Scenario Understanding**

[Question]: Mike plays ping pong for 40 minutes. In the first 20 minutes, he scores  $x$  points. In the second 20 minutes, he scores 25% more points. What is the value of unknown variable  $x$ ?

[Answer]: Unanswerable

**25: Answerable Judging (*Unanswerable*) - Scenario Understanding**

[Question]: Mike plays ping pong for 40 minutes. In the first 20 minutes, he scores 4 points. In the second 20 minutes, he scores 25% more points. How many total points did he score?

[Solution]:

Step 1: Find the number of points Mike scored in the second 20 minutes. To do this, take the number of points he scored in the first 20 minutes, which is 4 points, and increase this by 25%.

Step 2: To calculate 25% of 4, divide 4 by 4 which gives 1 (since 25% is the same as 1/4th).

Step 3: Add this 25% (which is 1 point) to the original 4 points to find out how many points he scored in the second 20 minutes. So, 4 points + 1 point = 5 points.

Step 4: Now, to find the total number of points Mike scored in the 40 minutes, add the points scored in the first 20 minutes to the points scored in the second 20 minutes.



1728 So, 4 points (first 20 minutes) + 5 points (second 20 minutes) = 9 points  
 1729 .  
 1730 The answer is 9 points.  
 1731 [Answer]: Correct

### 26: Outcome Judging (*Correct*) - Original Problem

1734 [Question]: Mike plays ping pong for 40 minutes. In the first 20 minutes  
 1735 , he scores 4 points. In the second 20 minutes, he scores 25% more points  
 1736 . How many total points did he score?  
 1737 [Solution]:  
 1738 Step 1: Mike scores 4 points in the first 20 minutes.  
 1739 Step 2: Mike scores 25% more points in the second 20 minutes.  
 1740 Step 3: Mike's total points are the sum of his points in the first and  
 1741 second 20 minutes.  
 1742 The answer is 28 points.  
 1743 [Answer]: Incorrect

### 27: Outcome Judging (*Incorrect*) - Original Problem

1744 [Question]: During a 40-minute ping pong session, Mike scores 4 points in  
 1745 the initial half. In the latter half, he manages to increase his score  
 1746 by 25% compared to the first half. What is the total score Mike achieved  
 1747 in this session?  
 1748 [Solution]:  
 1749 Step 1: Calculate Mike's score in the first half, which is given as 4  
 1750 points.  
 1751 Step 2: Calculate the 25% increase in score for the second half.  $(25\% \text{ of } 4 = \frac{25}{100} \times 4 = 1)$ .  
 1752 Step 3: Add this increase to the initial score to find the score for the  
 1753 second half.  $(4 + 1 = 5)$ .  
 1754 Step 4: Combine the scores from both halves to get the total score.  $(4 \text{ (first half)} + 5 \text{ (second half)} = 9)$ .  
 1755 The answer is 9.  
 1756 [Answer]: Correct

### 28: Outcome Judging (*Correct*) - Problem Understanding

1759 [Question]: During a 40-minute ping pong session, Mike scores 4 points in  
 1760 the initial half. In the latter half, he manages to increase his score  
 1761 by 25% compared to the first half. What is the total score Mike achieved  
 1762 in this session?  
 1763 [Solution]:  
 1764 Step 1: Mike scores 4 points in the initial half, so his score after 20  
 1765 minutes is 20 minutes  $\times 4 \text{ points}/20 \text{ minutes} = 8 \text{ points}$ .  
 1766 Step 2: In the latter half, Mike's score increases by 25% compared to the  
 1767 first half, so his score after 20 minutes is 8 points + (8 points  $\times$   
 1768 25% increase) = 8 points + 2 points = 10 points.  
 1769 Step 3: Therefore, Mike's total score after 40 minutes is 10 points + 4  
 1770 points (from the initial half) = 14 points.  
 1771 The answer is 14 points.  
 1772 [Answer]: Incorrect

### 29: Outcome Judging (*Incorrect*) - Problem Understanding

1773 [Question]: Mike plays ping pong in a local tournament and decides to  
 1774 practice for 40 minutes before the first match. During his practice  
 1775 session, in the first 20 minutes, while intermittently checking his phone  
 1776 and hydrating, he manages to score 4 points. In the following 20 minutes  
 1777 , feeling more warmed up and despite a short break to adjust his paddle's  
 1778 grip tape, he scores 25% more points than in the first session.  
 1779 Considering these distractions, how many total points did Mike score in  
 1780 his 40-minute practice session?  
 1781

1782 [Solution]:  
 1783 Step 1: Identify the points scored by Mike in the first 20 minutes, which  
 1784 is 4 points.  
 1785 Step 2: Calculate the number of points scored in the next 20 minutes by  
 1786 finding 25% more than the points scored in the first session. 25% of 4  
 1787 points is 1 (which is  $0.25 * 4$ ).  
 1788 Step 3: Add the additional points to the first session's points to find  
 1789 the total points scored in the second session. So, 4 points + 1 point = 5  
 1790 points in the second session.  
 1791 Step 4: To find the total points scored in the 40-minute practice session  
 1792 , add the points from the first and second sessions together.  
 1793 So, 4 points + 5 points = 9 points.\n\nThe answer is 9 points.  
 [Answer]: Correct

### 1794 30: Outcome Judging (*Correct*) - Irrelevant Disturbance

1795  
 1796 [Question]: Mike plays ping pong in a local tournament and decides to  
 1797 practice for 40 minutes before the first match. During his practice  
 1798 session, in the first 20 minutes, while intermittently checking his phone  
 1799 and hydrating, he manages to score 4 points. In the following 20 minutes  
 1800 , feeling more warmed up and despite a short break to adjust his paddle's  
 1801 grip tape, he scores 25% more points than in the first session.  
 1802 Considering these distractions, how many total points did Mike score in  
 1803 his 40-minute practice session?  
 1804 [Solution]:  
 1805 Step 1: Calculate the points scored in the first 20 minutes.\n4 points in  
 1806 20 minutes = 2 points per minute  
 1807 Step 2: Calculate the points scored in the second 20 minutes.\n25% more  
 1808 points than in the first session =  $1.25 * 2 \text{ points/minute} = 2.5 \text{ points/}$   
 1809 minute  
 1810 Step 3: Calculate the total points scored in the 40-minute practice  
 1811 session.\n2 points/minute x 20 minutes = 40 points\n2.5 points/minute x  
 1812 20 minutes = 50 points  
 1813 Step 4: Add the points scored in the first and second sessions.\n40  
 1814 points + 50 points = 90 points  
 1815 The answer is 90 points.  
 1816 [Answer]: Incorrect

### 1817 31: Outcome Judging (*Incorrect*) - Irrelevant Disturbance

1818 [Question]: Mike plays ping pong for 40 minutes. In the first 20 minutes,  
 1819 he scores  $x$  points. In the second 20 minutes, he scores 25% more points.  
 1820 He scored 9 total points. What is the value of unknown variable  $x$ ?  
 1821 [Solution]:  
 1822 Step 1: Define the given information: Mike plays ping pong for 40 minutes  
 1823 , scores  $x$  points in the first 20 minutes, and scores 25% more points in  
 1824 the second 20 minutes.  
 1825 Step 2: Since he scores 25% more points in the second 20 minutes, we can  
 1826 express the points scored in the second 20 minutes as  $(x + 0.25x)$ .  
 1827 Step 3: This can be simplified as  $(1.25x)$ , representing the points  
 1828 scored in the second 20 minutes.  
 1829 Step 4: Since he scored 9 points in total, we can write the equation:\n  
 1830  $x + 1.25x = 9$   
 1831 Step 5: By simplifying the equation, we combine like terms:\n  
 1832  $2.25x = 9$   
 1833 Step 6: To find the value of  $x$ , we divide both sides of the equation by  
 1834 2.25:\n  
 1835  $x = \frac{9}{2.25}$   
 Step 7: Simplify the division:\n  
 The answer is  $x = 4$ .  
 [Answer]: Correct

### 1836 32: Outcome Judging (*Correct*) - Scenario Understanding

1836

1837 [Question]: Mike plays ping pong for 40 minutes. In the first 20 minutes,  
 1838 he scores  $x$  points. In the second 20 minutes, he scores 25% more points.  
 1839 He scored 9 total points. What is the value of unknown variable  $x$ ?

1840 [Solution]:

1841 Step 1: Let's first find out how many points Mike scored in the first 20  
 1842 minutes.

1843 Step 2: Since Mike scored  $x$  points in the first 20 minutes, we know that  
 1844  $x/2 = 9$  (since he scored 9 points in total).

1845 Step 3: Now, we know that  $x = 2(9) = 18$ .

1846 The answer is  $x = 18$ .

1847 [Answer]: Incorrect

1848

### 33: Outcome Judging (*Incorrect*) - Scenario Understanding

1849

1850 [Question]: Mike plays ping pong for 40 minutes. In the first 20 minutes,  
 1851 he scores 4 points. In the second 20 minutes, he scores 25% more points.  
 1852 How many total points did he score?

1853 [Solution]:

1854 Step 1: Find the number of points Mike scored in the second 20 minutes.  
 1855 To do this, take the number of points he scored in the first 20 minutes,  
 1856 which is 4 points, and increase this by 25%.

1857 Step 2: To calculate 25% of 4, divide 4 by 4 which gives 1 (since 25% is  
 1858 the same as 1/4th).

1859 Step 3: Add this 25% (which is 2 points) to the original 4 points to find  
 1860 out how many points he scored in the second 20 minutes. So, 4 points + 2  
 1861 points = 6 points.

1862 Step 4: Now, to find the total number of points Mike scored in the 40  
 1863 minutes, add the points scored in the first 20 minutes to the points  
 1864 scored in the second 20 minutes. So, 4 points (first 20 minutes) + 6  
 1865 points (second 20 minutes) = 10 points.

1866 The answer is 10 points.

1867 [Answer]: Step 3

1868

### 34: Process Judging - Original Problem

1869

1870 [Question]: During a 40-minute ping pong session, Mike scores 4 points in  
 1871 the initial half. In the latter half, he manages to increase his score  
 1872 by 25% compared to the first half. What is the total score Mike achieved  
 1873 in this session?

1874 [Solution]:

1875 Step 1: Calculate Mike's score in the first half, which is given as 4  
 1876 points.

1877 Step 2: Calculate the 25% increase in score for the second half.  $\frac{4}{4} \times 25 = 0.16$ .

1878 Step 3: Add this increase to the initial score to find the score for the  
 1879 second half.  $4 + 0.16 = 4.16$ .

1880 Step 4: Combine the scores from both halves to get the total score.  $4 + 4.16 = 8.16$ .

1881 The answer is 8.16.

1882 [Answer]: Step 2

1883

### 35: Process Judging - Problem Understanding

1884

1885 [Question]: Mike plays ping pong in a local tournament and decides to  
 1886 practice for 40 minutes before the first match. During his practice  
 1887 session, in the first 20 minutes, while intermittently checking his phone  
 1888 and hydrating, he manages to score 4 points. In the following 20 minutes  
 1889 , feeling more warmed up and despite a short break to adjust his paddle's  
 grip tape, he scores 25% more points than in the first session.  
 Considering these distractions, how many total points did Mike score in  
 his 40-minute practice session?

[Solution]:

1890 Step 1: Identify the points scored by Mike in the first 20 minutes, which  
 1891 is 4 points.  
 1892 Step 2: Calculate the number of points scored in the next 20 minutes by  
 1893 finding 25% more than the points scored in the first session. 25% of 4  
 1894 points is 1 (which is  $0.25 * 4$ ).  
 1895 Step 3: Add the additional points to the first session's points to find  
 1896 the total points scored in the second session. So, 4 points + 1 point = 6  
 1897 points in the second session.  
 1898 Step 4: To find the total points scored in the 40-minute practice session  
 1899 , add the points from the first and second sessions together. So, 4  
 1900 points + 6 points = 10 points.  
 1901 The answer is 10 points.  
 1902 [Answer]: Step 3

### 36: Process Judging - Irrelevant Disturbance

1904 [Question]: Mike plays ping pong for 40 minutes. In the first 20 minutes,  
 1905 he scores  $x$  points. In the second 20 minutes, he scores 25% more points.  
 1906 He scored 9 total points. What is the value of unknown variable  $x$ ?  
 1907 [Solution]:  
 1908 Step 1: Define the given information: Mike plays ping pong for 40 minutes  
 1909 , scores  $x$  points in the first 20 minutes, and scores 25% more points in  
 1910 the second 20 minutes.  
 1911 Step 2: Since he scores 25% more points in the second 20 minutes, we can  
 1912 express the points scored in the second 20 minutes as  $(x + 0.25x)$ .  
 1913 Step 3: This can be simplified as  $(1.25x)$ , representing the points  
 1914 scored in the second 20 minutes.  
 1915 Step 4: Since he scored 9 points in total, we can write the equation:  
 1916  $x + 1.25x = 9$   
 1917 Step 5: By simplifying the equation, we combine like terms:  
 1918  $2.25x = 9$   
 1919 Step 6: To find the value of  $x$ , we divide both sides of the equation by  
 1920 2.25:  
 1921  $x = \frac{9}{2.25}$   
 1922 Step 7: Simplify the division:  
 1923  $x = 5$   
 1924 The answer is  $x = 5$ .  
 1925 [Answer]: Step 7

### 37: Process Judging - Scenario Understanding

1922  
 1923  
 1924  
 1925  
 1926  
 1927  
 1928  
 1929  
 1930  
 1931  
 1932  
 1933  
 1934  
 1935  
 1936  
 1937  
 1938  
 1939  
 1940  
 1941  
 1942  
 1943

## G.2 CASE PROBLEMS IN MATHCHECK-GEO. PROBLEM GROUP ID: GEO-15

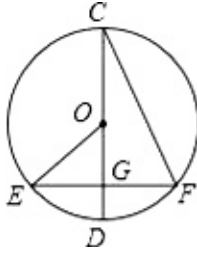


Figure 13: Geometry diagram for geometry problems in group 15.

[Question]: As shown in the figure, the diameter  $CD$  of  $\odot O$  crosses the midpoint  $G$  of chord  $EF$ ,  $\angle DCF = 20.0$ , then  $\angle EOD$  is equal to  $( )^\circ$   
 [Answer]: 40.0

## 38: Problem Solving - Original Problem

[Question]: In the circle with center  $O$ , diameter  $CD$  intersects the midpoint  $G$  of the chord  $EF$ , and the measure of angle  $DCF$  is 20 degrees. Determine the measurement of angle  $EOD$  in degrees.  
 [Answer]: 40.0

## 39: Problem Solving - Problem Understanding

[Question]: In the figure of circle  $O$ , the diameter  $CD$  intersects the midpoint  $G$  of the chord  $EF$ . The length of the chord  $EF$  is 7.5 cm, which is irrelevant to our angle measurements. The angle  $\angle DCF$  is given to be 20.0 degrees. We need to calculate the angle  $\angle EOD$ . What is the measure of this angle in degrees?  
 [Answer]: 40.0

## 40: Problem Solving - Irrelevant Disturbance

[Question]: As shown in the figure, the diameter  $CD$  of  $\odot O$  crosses the midpoint  $G$  of chord  $EF$ ,  $\angle DCF = x$ ,  $\angle EOD$  is equal to  $40^\circ$ . What is the value of unknown variable  $x$ ?  
 [Answer]: 20.0

## 41: Problem Solving - Scenario Understanding

[Question]: As shown in the figure, the diameter  $CD$  of  $\odot O$  crosses the midpoint  $G$  of chord  $EF$ ,  $\angle DCF = 20.0$ , then  $\angle EOD$  is equal to  $( )^\circ$   
 [Answer]: Answerable

## 42: Answerable Judging (Answerable) - Original Problem

[Question]: As shown in the figure, the diameter  $CD$  of  $\odot O$  crosses chord  $EF$ ,  $\angle DCF = 20.0$ , then  $\angle EOD$  is equal to  $( )^\circ$   
 [Answer]: Unanswerable

## 43: Answerable Judging (Unanswerable) - Original Problem

[Question]: In the circle with center  $O$ , diameter  $CD$  intersects the midpoint  $G$  of the chord  $EF$ , and the measure of angle  $DCF$  is 20 degrees. Determine the measurement of angle  $EOD$  in degrees.  
 [Answer]: Answerable

## 44: Answerable Judging (Answerable) - Problem Understanding

1998

1999

2000

2001

2002

[Question]: In the circle with center  $O$ , diameter  $CD$  intersects the midpoint  $G$  of the chord  $EF$ . Determine the measurement of angle  $EOD$  in degrees.

[Answer]: Unanswerable

2003

#### 45: Answerable Judging (*Unanswerable*) - Problem Understanding

2004

2005

2006

2007

2008

2009

2010

2011

[Question]: In the figure of circle  $O$ , the diameter  $CD$  intersects the midpoint  $G$  of the chord  $EF$ . The length of the chord  $EF$  is 7.5 cm, which is irrelevant to our angle measurements. The angle  $\angle DCF$  is given to be 20.0 degrees. We need to calculate the angle  $\angle EOD$ . What is the measure of this angle in degrees?

[Answer]: Answerable

#### 46: Answerable Judging (*Answerable*) - Irrelevant Disturbance

2012

2013

2014

2015

2016

2017

2018

[Question]: In the figure of circle  $O$ , the diameter  $CD$  intersects the midpoint  $G$  of the chord  $EF$ . The length of the chord  $EF$  is 7.5 cm, which is irrelevant to our angle measurements. We need to calculate the angle  $\angle EOD$ . What is the measure of this angle in degrees?

[Answer]: Unanswerable

#### 47: Answerable Judging (*Unanswerable*) - Irrelevant Disturbance

2019

2020

2021

2022

2023

2024

[Question]: As shown in the figure, the diameter  $CD$  of  $\odot O$  crosses the midpoint  $G$  of chord  $EF$ ,  $\angle DCF = x$ ,  $\angle EOD$  is equal to  $40^\circ$ . What is the value of unknown variable  $x$ ?

[Answer]: Answerable

#### 48: Answerable Judging (*Answerable*) - Scenario Understanding

2025

2026

2027

2028

2029

2030

[Question]: As shown in the figure, the diameter  $CD$  of  $\odot O$  crosses the midpoint  $G$  of chord  $EF$ ,  $\angle DCF = x$ . What is the value of unknown variable  $x$ ?

[Answer]: Unanswerable

#### 49: Answerable Judging (*Unanswerable*) - Scenario Understanding

2031

2032

2033

2034

2035

2036

2037

2038

2039

2040

2041

2042

2043

2044

2045

2046

2047

2048

2049

2050

2051

[Question]: As shown in the figure, the diameter  $CD$  of  $\odot O$  crosses the midpoint  $G$  of chord  $EF$ ,  $\angle DCF = 20.0$ , then  $\angle EOD$  is equal to  $(\ )^\circ$

[Solution]:

Step 1: Identify that the angle  $DCF$  is an angle formed by a chord ( $CF$ ) and a tangent ( $CD$ ) at the point  $C$ .

Step 2: Recognize that the angle between a chord and a tangent at the point of contact is equal to the angle in the alternate segment, which is angle  $ECF$  in this case.

Step 3: Thus, the measure of angle  $ECF$  is also 20 degrees.

Step 4: Observe that the diameter  $CD$  passes through the center  $O$  of the circle and cuts the chord  $EF$  at its midpoint  $G$ .

Step 5: Recall that a line drawn from the center of the circle to the midpoint of a chord is perpendicular to the chord.

Step 6: Therefore, triangle  $EOG$  is a right-angled triangle at  $G$ .

Step 7: Understand that angle  $EOD$  is double the angle  $ECF$  because angle  $EOD$  is the angle at the center whereas angle  $ECF$  is at the circumference based on the Circle Theorem.

Step 8: Calculate angle  $EOD$  as 2 times angle  $ECF$ , which is  $2 * 20$  degrees.

Step 9: Perform the multiplication:  $2 * 20$  degrees = 40 degrees.

The answer is: Angle  $EOD$  measures 40 degrees.

[Answer]: Correct

#### 50: Outcome Judging (*Correct*) - Original Problem

2052

2053 [Question]: As shown in the figure, the diameter CD of  $\odot O$  crosses  
 2054 the midpoint G of chord EF,  $\angle DCF = 20.0$ , then  $\angle EOD$  is equal  
 2055 to  $( )^\circ$

2056 [Solution]:

2057 Step 1: Notice that chord EF is bisected by diameter CD at point G.

2058 Step 2: Recognize that  $\angle DCF$  is an angle formed by the diameter CD  
 and chord EF, so it is inscribed in the semicircle of  $\odot O$ .

2059 Step 3: Recall that an angle inscribed in a semicircle is a right angle;  
 2060 hence  $\angle DCF = 90^\circ$  (not the given  $20^\circ$ , which seems to  
 be a mistake).

2061 Step 4: Realize that  $\angle EOD$  is the angle at the center of the circle  
 2062 , and it stands on the arc EF.

2063 Step 5: Understand that the angle at the center of a circle is twice the  
 2064 angle at the circumference when both stand on the same arc.

2065 Step 6: Since  $\angle DCF$  is a right angle ( $90^\circ$ ) and  $\angle EOD$   
 2066 is twice that angle at the center,  $\angle EOD$  is  $2 * 90^\circ$ .

2067 Step 7: Calculate  $2 * 90^\circ$  to get  $180^\circ$ .

2068 The answer is  $\angle EOD$  is equal to  $180^\circ$ .

2069 [Answer]: Incorrect

### 2070 51: Outcome Judging (*Incorrect*) - Original Problem

2071

2072 [Question]: In the circle with center O, diameter CD intersects the  
 2073 midpoint G of the chord EF, and the measure of angle DCF is 20 degrees.  
 2074 Determine the measurement of angle EOD in degrees.

2075 [Solution]:

2076 Step 1: Identify that the angle DCF is an angle formed by a chord (CF)  
 and a tangent (CD) at the point C.

2077 Step 2: Recognize that the angle between a chord and a tangent at the  
 2078 point of contact is equal to the angle in the alternate segment, which is  
 angle ECF in this case.

2079 Step 3: Thus, the measure of angle ECF is also 20 degrees.

2080 Step 4: Observe that the diameter CD passes through the center O of the  
 2081 circle and cuts the chord EF at its midpoint G.

2082 Step 5: Recall that a line drawn from the center of the circle to the  
 2083 midpoint of a chord is perpendicular to the chord.

2084 Step 6: Therefore, triangle EOG is a right-angled triangle at G.

2085 Step 7: Understand that angle EOD is double the angle ECF because angle  
 2086 EOD is the angle at the center whereas angle ECF is at the circumference  
 based on the Circle Theorem.

2087 Step 8: Calculate angle EOD as 2 times angle ECF, which is  $2 * 20$  degrees  
 2088 .

2089 Step 9: Perform the multiplication:  $2 * 20$  degrees = 40 degrees.

2090 The answer is: Angle EOD measures 40 degrees.

2091 [Answer]: Correct

### 2092 52: Outcome Judging (*Correct*) - Problem Understanding

2093

2094 [Question]: In the circle with center O, diameter CD intersects the  
 2095 midpoint G of the chord EF, and the measure of angle DCF is 20 degrees.  
 2096 Determine the measurement of angle EOD in degrees.

2097 [Solution]:

2098 Step 1: Identify that chord EF is intersected at its midpoint G by the  
 diameter CD; this implies that triangle CGF is isosceles with  $CG = GF$ .

2099 Step 2: Recognize that angle DCF is given as 20 degrees and that angle  
 2100 DGF is also 20 degrees because CGF is an isosceles triangle and angle CGF  
 2101 equals angle DGF.

2102 Step 3: Realize that angle EOF is twice angle DGF since angle EOF is the  
 2103 angle subtended at the center by chord EF, and angle DGF is the angle  
 2104 subtended at the circumference by the same chord (the Angle at the Center  
 Theorem).

2105 Step 4: Calculate angle EOF as  $2 * \text{angle DGF}$  which is  $2 * 20$  degrees = 40  
 degrees.

2106 Step 5: Note that angle EOD is twice angle EOF because EO is a diameter  
 2107 of the circle, and therefore angle EOD subtends a semicircle, making it a  
 2108 straight angle or 180 degrees.  
 2109 Step 6: Divide 180 degrees by the angle EOF to find angle EOD. Since EOF  
 2110 is 40 degrees, angle EOD is 180 degrees / 40 degrees = 4.5 times EOF.  
 2111 Step 7: Determine that there is a mistake in the previous step: angle EOF  
 2112 is part of angle EOD and not the other way around. Therefore, angle EOD  
 2113 consists of two EOF angles plus the straight angle DOF, which is 180  
 2114 degrees.  
 2115 Step 8: Sum up the angles to find angle EOD: angle EOF (40 degrees) +  
 2116 angle DOF (180 degrees) + angle EOF (40 degrees) which is 40 + 180 + 40.  
 2117 Step 9: Add the angles to get the measurement of angle EOD: 40 + 180 + 40  
 2118 = 260 degrees.  
 2119 The answer is 260 degrees.  
 2120 [Answer]: Incorrect

### 53: Outcome Judging (*Incorrect*) - Problem Understanding

2121 [Question]: In the figure of circle O, the diameter CD intersects the  
 2122 midpoint G of the chord EF. The length of the chord EF is 7.5 cm, which  
 2123 is irrelevant to our angle measurements. The angle  $\angle DCF$  is given  
 2124 to be 20.0 degrees. We need to calculate the angle  $\angle EOD$ . What is  
 2125 the measure of this angle in degrees?  
 2126 [Solution]:  
 2127 Step 1: Recognize that  $\angle DCF$  is an angle formed by intersecting a  
 2128 chord (EF) and a diameter (CD) of the circle.  
 2129 Step 2: Acknowledge that an angle formed by a chord and a tangent (or in  
 2130 this case a diameter) at the point of contact is equal to the angle in  
 2131 the alternate segment, which means  $\angle DCF$  is equal to angle  $\angle EOF$ .  
 2132 Step 3: Consider that CD is a diameter, so triangle COD is an isosceles  
 2133 triangle with CO = OD (both are radii of the circle).  
 2134 Step 4: Recall that the base angles in an isosceles triangle are equal,  
 2135 which means  $\angle CDO = \angle DCO$ .  
 2136 Step 5: Since  $\angle CDF$  is given as 20.0 degrees and  $\angle CDO = \angle DCO$ , each of these angles  
 2137 is also 20.0 degrees.  
 2138 Step 6: Understand that  $\angle EOD$  is the sum of  $\angle EOF$  and  $\angle FOD$   
 2139 (since  $\angle EOF$  and  $\angle FOD$  are adjacent angles).  
 2140 Step 7: Since  $\angle EOF$  is equal to  $\angle DCF$ , which is 20.0 degrees,  
 2141 then  $\angle EOF$  is also 20.0 degrees.  
 2142 Step 8: Calculate the measure of  $\angle EOD$  by adding the measures of  $\angle$   
 2143 angle EOF and  $\angle FOD$ , which are both 20.0 degrees.  
 2144 Step 9: Perform the calculation:  $\angle EOD = \angle EOF + \angle FOD$   
 2145 = 20.0 degrees + 20.0 degrees.  
 2146 Step 10: Conclude that  $\angle EOD$  measures 40.0 degrees.  
 2147 The answer is 40.0 degrees.  
 2148 [Answer]: Correct

### 54: Outcome Judging (*Correct*) - Irrelevant Disturbance

2149 [Question]: In the figure of circle O, the diameter CD intersects the  
 2150 midpoint G of the chord EF. The length of the chord EF is 7.5 cm, which  
 2151 is irrelevant to our angle measurements. The angle  $\angle DCF$  is given  
 2152 to be 20.0 degrees. We need to calculate the angle  $\angle EOD$ . What is  
 2153 the measure of this angle in degrees?  
 2154 [Solution]:  
 2155 Step 1: Calculate the points scored in the first 20 minutes. 4 points in  
 2156 20 minutes = 2 points per minute  
 2157 Step 2: Calculate the points scored in the second 20 minutes. 25% more  
 2158 points than in the first session = 1.25 x 2 points/minute = 2.5 points/  
 2159 minute  
 2160 Step 3: Calculate the total points scored in the 40-minute practice  
 2161 session. 2 points/minute x 20 minutes = 40 points 2.5 points/minute x  
 2162 20 minutes = 50 points



2160 Step 4: Add the points scored in the first and second sessions.\n40  
 2161 points + 50 points = 90 points  
 2162 The answer is 90 points.  
 2163 [Answer]: Incorrect

#### 2164 55: Outcome Judging (*Incorrect*) - Irrelevant Disturbance

2165  
 2166  
 2167 [Question]: As shown in the figure, the diameter CD of  $\odot O$  crosses  
 2168 the midpoint G of chord EF,  $\angle DCF = x$ ,  $\angle EOD$  is equal to  
 2169  $40^\circ$ . What is the value of unknown variable x?  
 2170 [Solution]:  
 2171 Step 1: Recognize that diameter CD of circle O passes through the center  
 2172 of the circle (point O) and bisects chord EF at point G.  
 2173 Step 2: Since CD is a diameter, it creates a semicircle from E to F  
 2174 passing through C.  
 2175 Step 3: Recognize that  $\angle EOD$  is the angle subtended by the diameter  
 2176 and thus is a right angle ( $90^\circ$ ) because it is an angle in a  
 2177 semicircle.  
 2178 Step 4: Understand that triangle EOD is a right triangle with  $\angle EOD$   
 2179 being  $90^\circ$  and  $\angle EOF$  being half of  $\angle EOD$  because O is  
 2180 the center of the circle and EOF is half of the circle (an arc of  $180^\circ$   
 2181 degree).  
 2182 Step 5: Since  $\angle EOD$  is  $40^\circ$  as given,  $\angle EOF$  must be  
 2183 half of that, so  $\angle EOF = 40^\circ / 2 = 20^\circ$ .  
 2184 Step 6: Notice that  $\angle DCF$  and  $\angle EOF$  are vertical angles and  
 2185 hence are equal. Therefore,  $x = \angle DCF = \angle EOF$ .  
 2186 Step 7: Substitute the value from Step 5 into the equation for x,  
 2187 therefore  $x = 20^\circ$ .  
 2188 The answer is  $x = 20^\circ$ .  
 2189 [Answer]: Correct

#### 2187 56: Outcome Judging (*Correct*) - Scenario Understanding

2188  
 2189  
 2190 [Question]: As shown in the figure, the diameter CD of  $\odot O$  crosses  
 2191 the midpoint G of chord EF,  $\angle DCF = x$ ,  $\angle EOD$  is equal to  
 2192  $40^\circ$ . What is the value of unknown variable x?  
 2193 [Solution]:  
 2194 Step 1: Recognize that angle  $\angle EOD$  lies on the diameter of the  
 2195 circle, which means it's a central angle and the arc ED it intercepts is  
 2196  $2x \angle EOD = 2x40^\circ = 80^\circ$ .  
 2197 Step 2: Since CD is a diameter, any angle formed on the circumference by  
 2198 diameter CD is a right angle ( $90^\circ$ ). Hence, angles  $\angle CFE$  and  
 2199  $\angle CED$  are 90 degrees each.  
 2200 Step 3: The arc EG is half the circumference of the circle intercepted by  
 2201 angle  $\angle CFE$ , so the measure of arc EG is  $180^\circ$  since  $\angle$   
 2202  $\angle CFE$  is a right angle.  
 2203 Step 4: Arc GF is the same as arc EG because G is the midpoint of EF,  
 2204 hence arc GF is also  $180^\circ$ .  
 2205 Step 5: The entire arc EF, which is the sum of arcs EG and GF, is  $180^\circ$   
 2206  $+ 180^\circ = 360^\circ - 80^\circ = 280^\circ$  since the  
 2207 circle is  $360^\circ$  and we must subtract arc ED ( $80^\circ$ ) which is  
 2208 not part of arc EF.  
 2209 Step 6: The angle  $\angle DCF$  is an inscribed angle that intercepts arc  
 2210 EF. The measure of an inscribed angle is half the measure of the arc it  
 2211 intercepts.  
 2212 Step 7: Calculate the measure of angle  $\angle DCF$  which is half of the  
 2213 measure of arc EF:  $x = 1/2 * 280^\circ = 140^\circ$ .  
 The answer is  $140^\circ$ .  
 [Answer]: Incorrect

#### 2212 57: Outcome Judging (*Incorrect*) - Scenario Understanding

2213

2214 [Question]: As shown in the figure, the diameter CD of  $\odot O$  crosses  
 2215 the midpoint G of chord EF,  $\angle DCF = 20.0$ , then  $\angle EOD$  is equal  
 2216 to  $( )^\circ$   
 2217 [Solution]:  
 2218 Step 1: OE and OF are the radius of the circle, so  $OE=OF$ .  
 2219 Step 2: G is the midpoint of EF, so  $EG=FG$ .  
 2220 Step 3: Because triangle OGE and triangle OGF are congruent triangles.  
 2221 Step 4: so the angle EOD is equal to the angle FOG.  
 2222 Step 5: because the angle FOG is the outer angle of the triangle COF, the  
 2223 angle  $FOG = \angle DCF + \angle CFO$ .  
 2224 Step 6: because  $OC=OF$  (they are all radius of circles),  $\angle DCF = \angle$   
 2225  $CFO$ .  
 2226 Step 7: therefore the angle  $FOG = 2 * \angle DCF$ .  
 2227 Step 8: so the angle  $EOD = 3 * \angle DCF$ .  
 2228 Step 9: Given the angle  $DCF = 20.0$ , the angle  $EOD = 60.0$ .  
 2229 The answer is: Angle EOD measures 60.0 degrees.  
 2230 [Answer]: Step 8

### 58: Process Judging - Original Problem

2231  
 2232 [Question]: In the circle with center O, diameter CD intersects the  
 2233 midpoint G of the chord EF, and the measure of angle DCF is 20 degrees.  
 2234 Determine the measurement of angle EOD in degrees.  
 2235 [Solution]:  
 2236 Step 1: OE and OF are the radius of the circle, so  $OE=OF$ .  
 2237 Step 2: G is the midpoint of EF, so  $EG=FG$ .  
 2238 Step 3: Because triangle OGE and triangle OGF are congruent triangles.  
 2239 Step 4: so the angle EOD is equal to the angle FOG.  
 2240 Step 5: because the angle FOG is the outer angle of the triangle COF, the  
 2241 angle  $FOG = \angle DCF + \angle CFO$ .  
 2242 Step 6: because  $OC=OF$  (they are all radius of circles),  $\angle DCF = 2 * \angle$   
 2243  $CFO$ .  
 2244 Step 7: therefore the angle  $FOG = 1.5 * \angle DCF$ .  
 2245 Step 8: so the angle  $EOD = 1.5 * \angle DCF$ .  
 2246 Step 9: Given the angle  $DCF = 20.0$ , the angle  $EOD = 30.0$ .  
 2247 The answer is: Angle EOD measures 30.0 degrees.  
 2248 [Answer]: Step 6

### 59: Process Judging - Problem Understanding

2249  
 2250 [Question]: In the figure of circle O, the diameter CD intersects the  
 2251 midpoint G of the chord EF. The length of the chord EF is 7.5 cm, which  
 2252 is irrelevant to our angle measurements. The angle  $\angle DCF$  is given to  
 2253 be 20.0 degrees. We need to calculate the angle  $\angle EOD$ . What is the  
 2254 measure of this angle in degrees?  
 2255 [Solution]:  
 2256 Step 1: OE and OF are the radius of the circle, so  $OE=OF$ .  
 2257 Step 2: G is the midpoint of EF, so  $EG=FG$ .  
 2258 Step 3: Because triangle OGE and triangle OGF are congruent triangles.  
 2259 Step 4: so the angle EOD is equal to the angle FOG.  
 2260 Step 5: because the angle FOG is the outer angle of the triangle COF, the  
 2261 angle  $FOG = \angle DCF + \angle CFO$ .  
 2262 Step 6: because  $OC=OF$  (they are all radius of circles),  $\angle DCF = \angle$   
 2263  $CFO$ .  
 2264 Step 7: therefore the angle  $FOG = \angle DCF$ .  
 2265 Step 8: so the angle  $EOD = \angle DCF$ .  
 2266 Step 9: Given the angle  $DCF = 20.0$ , the angle  $EOD = 20.0$ .  
 2267 The answer is: Angle EOD measures 20.0 degrees.  
 2268 [Answer]: Step 7

### 60: Process Judging - Irrelevant Disturbance

2268 [Question]: As shown in the figure, the diameter CD of  $\odot O$  crosses  
2269 the midpoint G of chord EF,  $\angle DCF = x$ ,  $\angle EOD$  is equal to  $40^\circ$   
2270 degree. What is the value of unknown variable x?  
2271 [Solution]:  
2272 Step 1: OE and OF are the radius of the circle, so OE=OF.  
2273 Step 2: G is the midpoint of EF, so EG=FG.  
2274 Step 3: Because triangle OGE and triangle OGF are congruent triangles.  
2275 Step 4: so the angle EOD is equal to the angle FOG.  
2276 Step 5: because the angle FOG is the outer angle of the triangle COF, the  
2277 angle FOG=angle DCF+angle CFO.  
2278 Step 6: because OC=OF (they are all radius of circles), angle DCF=angle  
2279 CFO.  
2280 Step 7: therefore the angle FOG=4\* angle DCF.  
2281 Step 8: so the angle EOD = 4 \* angle DCF.  
2282 Step 9: Given the angle EOD=40.0, the angle DCF = x = 10.0.  
2283 The answer is x = 10 degrees.  
2284 [Answer]: Step 7

61: Process Judging - Scenario Understanding

2283  
2284  
2285  
2286  
2287  
2288  
2289  
2290  
2291  
2292  
2293  
2294  
2295  
2296  
2297  
2298  
2299  
2300  
2301  
2302  
2303  
2304  
2305  
2306  
2307  
2308  
2309  
2310  
2311  
2312  
2313  
2314  
2315  
2316  
2317  
2318  
2319  
2320  
2321