

Uncovering the Intention Behind Equations in Mathematical Problems

Anonymous EMNLP submission

Abstract

Mathematical Equation Intent Recognition (MEIR) is a novel task aimed at identifying the intentions behind mathematical equations that people produce while solving **math world problems (MWP)**. We observe that, in previous research, researchers have often focused on how to let **large language models (LLMs)** correctly solve an MWP. However, focusing solely on the reasoning behind each step of a correct inference process is insufficient. We prefer that LLMs can provide guidance on the process of solving MWPs for students in educational settings. Therefore, they need to adjust the strategy based on the student's responses. We notice that, unlike existing mathematical datasets, students typically do not provide overly detailed descriptions of their steps in the real world. As a result, it is crucial for LLMs to possess the capability to understand the intention they produce those equations. We treat MEIR as a generation task, requiring models to summarize the intent in a single sentence. We also propose a data augmentation framework and utilized this framework to generate a benchmark called **Grade School Math Intention (GSMI)**. To evaluate MEIR task, we benchmark several LLMs on GSMI dataset. The results indicate that there is still significant room for improvement in the performance of general-purpose LLMs on the MEIR task. Conversely, capabilities acquired during pre-training and fine-tuning specifically in the field of mathematics significantly contribute to the model's ability to tackle those problems. Codes and datasets are available on <https://github.com/ch-666-six/MEIR>

1 Introduction

Recently, the capabilities of large language models (LLMs) (Minaee et al., 2024) have been extensively applied to tasks in the field of mathematics. Numerous researchers (Liu et al., 2023c; Wei et al., 2022; Kojima et al., 2022) have employed prompt-based methods or fine-tuning methods to further

Question:The price of a laptop is \$1000. If you get a 20% discount, how much do you have to pay?

Complete Answer (From GSM8k):
You will get a discount of $20/100 * \$1000 = \$\langle\langle 20/100 * 1000 = 200 \rangle\rangle 200$.
Therefore, you will have to pay $\$1000 - \$200 = \$\langle\langle 1000 - 200 = 800 \rangle\rangle 800$.
So the answer is: \$800.

Brief Answer (From Student):
According to the question, the solution process of this problem is as follows:
 $1000 - 20/100 * 100 = 800$.
As a result, we should pay 800 dollars.

Figure 1: An example of complete answer and brief answer. The standard answer is sourced from the GSM8K dataset (Cobbe et al., 2021), reflecting the ideal scenario of solving mathematical problems. However, in real-time scenarios, answers from students may resemble what is shown in the "brief answer". This answer may primarily consist of a series of mathematical equations.

enhance the ability of large language models to comprehend mathematical texts and solve mathematical problems.

However, we do not want LLMs simply become problem solvers. We hope to integrate the LLMs' mathematical capabilities closely with real-world educational scenarios. We find that, in real-time educational scenarios, particularly during homework or exams, students often arrive at their answers through a series of equations rather than a detailed step-by-step reasoning process. We show an example in Figure 1, citing a math question from GSM8k dataset (Cobbe et al., 2021).

As a result, to determine the correctness of the problem-solving process, we need to fully understand the intent behind these equations. Therefore, it is important to study the ability of LLMs to understand the intent behind the arithmetic equations.

We consider the task of **Mathematical Equation Intent Recognition (MEIR)** as a generation task. Specifically, we aim for the language model to produce concise descriptions for the equations within the solution steps of mathematical problems. To address this issue, relevant data is of necessity.

We use two different modules: Imitation-based Generator and Intention Extractor to generate data automatically, and propose a novel dataset called **Math World Problems Intention(MWPI)**. Details will be thoroughly explained in section 3.

MWPI is a benchmark to test whether language models can uncover the intention behind those equations appeared in the solutions of math world problems. The input context consists of a mathematical problem along with its solution steps expressed in the form of equations. The objective of the model is to produce, for each equation, a concise summary in the form of a sentence that encapsulates the intention behind the inclusion of that particular equation.

In our experimental evaluation, we observed that mainstream closed-source large language models (LLMs), such as GPT-4o, GPT-4 (OpenAI, 2024) and others, still exhibit potential for improvement in the MEIR task. This suggests that during the pre-training process, these models did not systematically acquire the ability to parse mathematical equations. In addition, we selected several open-source models and utilized instruction tuning to train them in the process of parsing mathematical expressions. We demonstrate that through specific instruction tuning, models with smaller parameter sizes can also achieve good performance. Meanwhile, through imitation-based generator, language model can improve themselves sustainably.

To conclude, the main contributions of this article are as follows:

1. To the best of our knowledge, we are the first to explore MEIR task.
2. We introduce a novel dataset called MWPI to evaluate the performance of models on MEIR task.
3. We employ an imitation-based generator to facilitate the generation of more diverse data under limited resources.

2 Related Work

2.1 Math World Problems Solving

Large language model have a strong ability to solve math world problems. Chain-of-thought prompting(Wei et al., 2022; Zhang et al., 2023b; Kojima et al., 2022) is a highly effective technique for eliciting detailed reasoning processes from LLMs to solve mathematical problems. Some researchers(Liu et al., 2023b; Gou et al., 2023; Imani

et al., 2023) also utilize external tools, like python executor and mathematical calculator, to enhance the calculate abilities of LLMs to solve mathematical problems. In addition, some researchers(Yu et al., 2023; Luo et al., 2023a; Ho et al., 2023; An et al., 2023) have organized mathematical corpora and fine-tuned open-source models using these corpora to enhance the mathematical reasoning capabilities. On several benchmark datasets (Cobbe et al., 2021; Hendrycks et al., 2021), LLMs have already demonstrated outstanding performance.

2.2 Instruction Tuning

Instruction tuning (Zhang et al., 2024) is an essential method for improving the capabilities and controllability of LLMs. This approach uses (INSTRUCTION, OUTPUT) pairs to train LLMs, where INSTRUCTION represents human instruction and OUTPUT denotes the target output that follows the instruction. LLMs like Instruct-GPT(Ouyang et al., 2022), Flan-T5(Chung et al., 2022), WizardLM(Xu et al., 2023), LLAVA(Liu et al., 2023a) and so on, are trained through instruction tuning. In domain-specific settings(Gupta et al., 2022; Zhang et al., 2023a; Luo et al., 2023b; Liu and Low, 2023), instruction tuning can also have a profound impact and contribute significantly to the performance. Compared to standard LLMs, instruction tuning enables more controllable and predictable model behavior. Due to the significant advantages of instruction fine-tuning, we also employed instruction tuning methods in our research.

2.3 Intent Understanding

Intent understanding(Louvan and Magnini, 2020) is one of the crucial tasks in artificial intelligence. In human-computer interaction(Jaimes and Sebe, 2007), accurately recognizing human intent facilitates machines in taking more appropriate actions to provide feedback. For example, in online shopping(Rahman et al., 2024; Yu et al., 2024), merchants always want to understand and accurately predict the buyer’s intention to promote consumption. Some researchers(Yin et al., 2024; Weld et al., 2022; Hariharan et al., 2022) treat intent understanding as intent classification and slot filling tasks. By contrast, in order to fully understand students’ intentions behind the equations, we view intent understanding as a text generation task(Li et al., 2021).

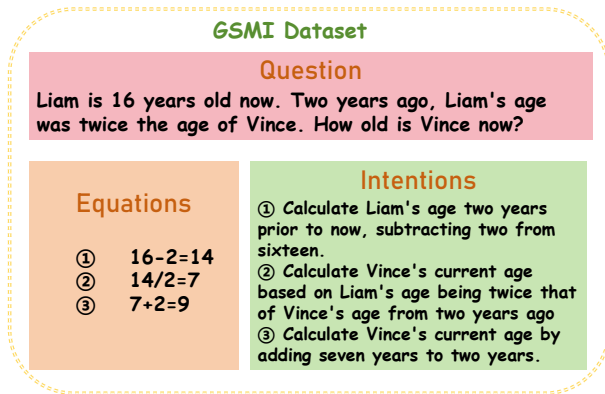


Figure 2: An example of GSMI dataset. Each data comprises a mathematical problem, a set of equations solving that problem and intention descriptions corresponding to each equation.

3 Dataset Construction

3.1 Overview

We propose a novel dataset called **Math World Problems Intention(MWPI)**. Each data in the dataset consists of 3 components: Question, Equations and Intentions. We provide an example in Figure 2 to illustrate the structure of this dataset.

For clarity in intent representation, we establish the rule that each sentence must begin with the word "calculate" and contain no more than 30 words.

Building on the existing dataset like GSM8k(Cobbe et al., 2021) and MATH(Hendrycks et al., 2021), we adopt the following two modules to generate the GSMI dataset: an imitation-based generator and an intention extractor. Figure 3 shows the following process.

Through the imitation-based generator the intention extractor, we can generate more valid instances to evaluate MEIR task. We utilize Chatgpt and GPT-4o(OpenAI, 2024) as LLMs to construct this two modules. The most labor-intensive step in this process is verifying the correctness of the expanded mathematical problems generated by LLMs. In this version, the GSMI dataset contains 8K training samples and 600 testing samples.

3.2 Imitation-based Generator

In the MEIR task, we focus on arithmetic problems of elementary school difficulty and in text modality. To enhance the model’s ability to learn the extraction of mathematical expression intentions, we implement data augmentation techniques(Li et al., 2022; Zhou et al., 2024).

Algorithm 1 Imitation-based Generator

Input: Original Question Dataset S
 Large Language Model $LLM()$
 Input Prompt $P()$
Output: Expanded Question Dataset S'

```

1:  $S' = []$ 
2: while Normal Execution do
3:    $Q = \text{Random\_sample}(S)$ 
4:    $Q' = LLM(P(Q))$ 
5:   if Grammar_Error( $Q'$ ) then
6:     CONTINUE
7:   end if
8:   if Answer_Error( $Q'$ ) then
9:     CONTINUE
10:  end if
11:   $S' = S' + [Q']$ 
12: end while
13: return  $S'$ 

```

Motivated by (Wei et al., 2022), we conclude that large language models possess significant in-context learning capabilities(Dong et al., 2023; Li, 2023). In the Chain-of-Thought(CoT) prompting method, researchers provide a step-by-step reasoning example within the input prompt. Guided by this example, LLMs like GPT-4(OpenAI, 2024) can mimic the provided instance from the prompt to perform structured reasoning process on a new mathematical problem.

Similarly, we innovatively propose the concept of an imitation-based generator. In our approach, we present a mathematical problem along with its corresponding solution process in the input prompt, instructing the LLMs to imitate the contextual information and generate a new problem that is structurally similar and of comparable difficulty. In this process, we utilize text-only ChatGPT to generate problems. The corresponding algorithm is shown in the Algorithm 1.

3.3 Intention Extractor

The purpose of this module is to extract the intent within mathematical equations. Firstly, for each step in the chain of thought, we extracted the mathematical equations representing that step. Next, we used the textual information of each step as the input prompt, allowing the large language model to summarize the intention within each step. The corresponding algorithm is shown in the Algorithm 2.

In short, after employing the aforementioned

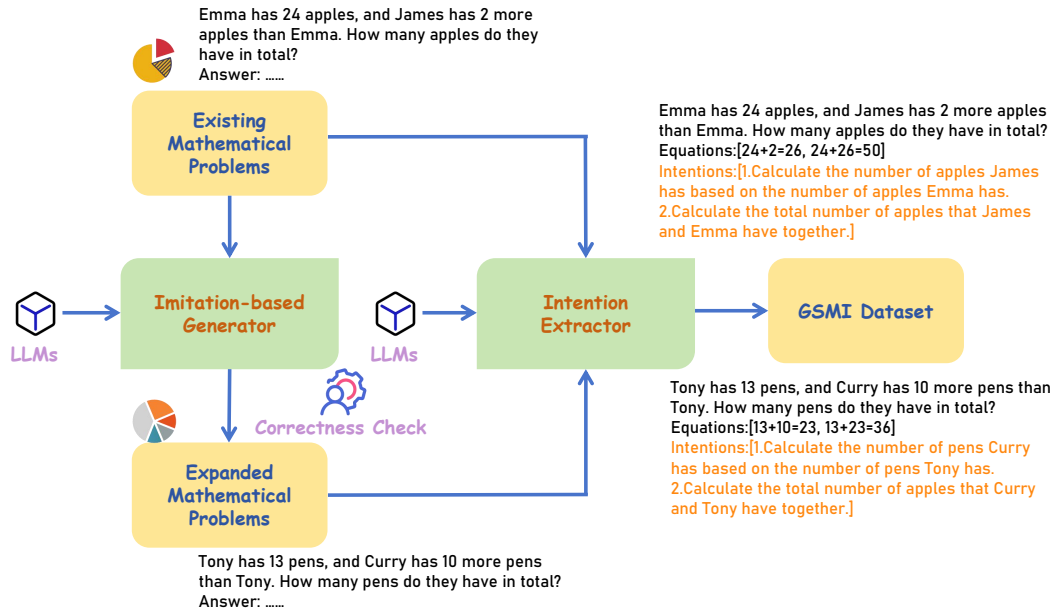


Figure 3: The Process of building GSMI dataset. The Imitation-based Generator is used to expand existing mathematical problems. The Intention Extractor is used to extract the intentions or objectives within each mathematical step. The figure illustrates a simple example from the GSMI dataset.

Algorithm 2 Intention Extractor

Input: Reasoning Step Set R [
 Large Language Model $LLM()$
 Input Prompt $P()$

Output: Intention Set I [
 1: $I = []$

- 2: **for** each item r in R **do**
 - 3: $i = LLM(P(r))$
 - 4: $I = I + [i]$
 - 5: **end for**
 - 6: **return** I
-

two modules, we can continuously generate data in GSMI, facilitating the model’s improved learning of the recognition of mathematical expression intents.

4 Experiments

4.1 Experimental Setup

We use the GSMI dataset to evaluate the model’s capability in recognizing the intent of mathematical expressions. For this purpose, we selected the following candidate models.

- **GPT-4o**(OpenAI, 2024) is a multilingual, multimodal generative pre-trained transformer designed by OpenAI. It was announced on 13 May, 2024, and released in the same day.

- **GPT-4** (OpenAI, 2024) is also a generative model designed by OpenAI, and it was announced in March, 2023.
- **GPT-3.5** (OpenAI, 2024), also known as ChatGpt, is a powerful large-scale language model. It was announced by OpenAI in March, 2022.
- **LLaMA** (Touvron et al., 2023) is a family of autoregressive large language models released by Meta AI, and we use the LLaMA-2 and LLaMA-3 models.
- **MetaMath** (Yu et al., 2023) is a fine-tuned model specifically for the field of mathematics. Researchers fine tune LLaMA (Touvron et al., 2023) on MetaMathQA dataset(Yu et al., 2023) and obtain MetaMath.
- **WizardMath** (Luo et al., 2023a) is a fine-tuned model for mathematics. Researchers train WizardMath model using reinforcement learning methods.

For open-source models, we performed instruction fine-tuning using the training dataset and then evaluated the models using the testing dataset. Due to the constraints on computational resources, we adopted the LoRA (Low-Rank Adaptation) parameter-efficient fine-tuning approach(Hu et al., 2021). By default, the open-source model is trained

MODEL	BLEU-1	ROUGE-L			BERT-SCORE		
		P	R	F1	P	R	F1
Prompting Closed-source Models							
GPT-4o	0.2278	0.5691	0.4533	0.4944	0.8121	0.7844	0.7971
GPT-4	0.2194	0.5826	0.3879	0.4570	0.8098	0.7601	0.7834
GPT-3.5	0.2304	0.5210	0.4819	0.4910	0.7972	0.7942	0.7949
Tuning Open-source Models							
LLAMA-2-7b	0.2335	0.5206	0.4935	0.5001	0.8015	0.7919	0.7961
LLAMA-2-13b	0.2386	0.5450	0.5306	0.5316	0.8129	0.8095	0.8107
LLAMA-3-8b	0.2373	0.5436	0.5202	0.5252	0.8087	0.8061	0.8068
LLAMA-3-8b-instruct	0.2376	0.5416	0.5235	0.5261	0.8097	0.8069	0.8077
MetaMath-7b	0.2377	0.5392	0.5233	0.5255	0.8111	0.8079	0.8089
MetaMath-13b	0.2369	0.5602	0.5308	0.5386	0.8186	0.8105	0.8139
WizardMath-7b	0.2372	0.5549	0.5400	0.5412	0.8152	0.8134	0.8138

Table 1: Results on MEIR task. P means Precision. R means Recall. And F1 means F1 score. In the closed-source models, the best-performing value in each row is highlighted in yellow. In the open-source models, the best-performing value in each row is highlighted in green. The best value in each row is highlighted in bold.

on the training set for 3 epochs with a learning rate of $2e-4$.

For closed-source models, we directly evaluated them using the testing dataset.

4.2 Evaluation metrics

We consider MEIR to be a text generation task. For the results generated by our candidate models for each equation, we need to evaluate their similarity to the ground truth. To this end, we selected the following evaluation metrics.

- **BLEU** (Papineni et al., 2002) is a metric for evaluating the quality of machine-generated text, which calculates precision for various n-gram lengths and combines these using a weighted geometric mean.
- **ROUGE** (Lin, 2004) is a set of metrics used to evaluate the quality machine-generated text. We use ROUGE-L, which captures the longest sequence of words that appear in both the candidate and reference summaries in the same order.
- **BERT-SCORE** (Zhang et al., 2020) leverages contextual embeddings from pre-trained transformer models, specifically BERT (Devlin et al., 2019) or RoBERTa (Liu et al., 2019), to capture semantic similarity between the candidate and reference texts. We use bert-large-

uncased as our base model, which contains 24 layers.

4.3 Experiment Result

The experimental results of the MEIR task are shown in the Table 1. We meticulously recorded the performance of all candidate models, retaining four decimal places, and documented the results in the table.

We observed that in the closed-source models, GPT-4o (OpenAI, 2024), as the latest proposed model, performs the best on the ROUGE-L and BERT-SCORE metrics. This indicates that GPT-4o surpasses the previously proposed GPT-4 and GPT-3.5 models in executing the MEIR task.

However, it is important to note that as a general-purpose large language model, GPT-4o, along with other closed-source models, lacks sufficient pre-training in the field of mathematics. As seen in the table, models with relatively smaller parameter sizes can outperform general-purpose large language models on the MEIR task after undergoing instruction tuning.

5 Analysis and Discussion

5.1 Mathematical Fine-tuning

MetaMath (Yu et al., 2023) and WizardMath (Luo et al., 2023a) are both fine-tuned versions of the LLAMA model. They were fine-tuned using extensive mathematical data, for example, the Meta-

Question	Chenny is 10 years old. Alyana is 4 years younger than Chenny. How old is Anne if she is 2 years older than Alyana?
Equations	["10-4=6", "6+2=8"]
GPT-4o Results	(1) Calculate how much younger Alyana is than Chenny. (2) Calculate how much older Anne is than Alyana.
MetaMath-13b Results	(1) Calculate Alyana’s age by subtracting four from ten. (2) Calculate Anne’s age by adding two to six.
WizardMath-7b Results	(1) Calculate the age of Alyana by subtracting her age from Chenny’s age.(2) Calculate Anne’s age by adding six and two.

Table 2: An simple example in GSMI testing set. In this example, GPT-4o clearly misunderstood the intent of the intermediate steps and provided an incorrect answer. MetaMath-13b and WizardMath-7b accurately grasped the intent of the intermediate steps.

MathQA dataset reached a size of 395K. As shown in Table 1, compared to LLAMA, MetaMath and WizardMath exhibit a significant advantage in handling the MEIR task. Notably, on the BERT-SCORE metric, which closely aligns with human evaluation, both models demonstrate remarkable capability.

Through controlled experiments, we have concluded that: **Mathematical Fine-tuning is highly effective and necessary for downstream mathematical tasks.**

5.2 Model Size

With LLMs demonstrating powerful capabilities across various domains, many people have begun to believe that there is a positive correlation between the parameter size of a model and its ability to handle complex problems.

However, as shown in the Table 1, on the MEIR task, the performance of smaller open-source models surpasses that of larger closed-source models. This indicates that for the MEIR task, **high-quality data refinement is more crucial than larger model sizes.** We require models to acquire knowledge and capability within a specific domain.

Table 2 presents an example from the evaluation set. In this instance, GPT-4o made evident errors in summarizing the intent, whereas MetaMath and WizardMath accurately summarized the intent.

5.3 Data Augmentation

In machine learning, richer datasets often yield better results, while a lack of data can easily lead to overfitting on the training data.

As shown in Figure 3, through Imitation-based Generator and Intention Extractor, we can continuously generate new data to further train the model

on the MEIR task. Compared to collecting mathematical problems and answers from the real world, the method illustrated in Figure 3 clearly requires significantly less human effort and time.

However, we need to investigate the effectiveness of this data augmentation method. We raise a question that does generating more examples through Imitation-based Generator and Intention Extractor on the existing datasets enable the model to perform better on the MEIR task?

In this regard, we introduce a variable K. K represents the total number of examples involved in instruction tuning. We select K values of 500, 1000, 2000, and 5000 for experimentation on MetaMath and WizardMath models. The experimental results are shown in Table 3.

It is evident that as K increases, both ROUGE-L and BERT-SCORE metrics show an overall increase trend. When all data generated through Imitation-based Generator and Intention Extractor modules is used in the instruction tuning process, the performance also improves significantly. These two modules can continuously generate new data. This indicates that we can leverage the imitation generation capability of LLMs to produce richer training data with limited resources. This part of training data truly helps language models better acquire the ability to uncover the intentions behind mathematical equations.

In summary, we state that **appropriate data augmentation strategies contributes to enhancing the language models’ performance on the MEIR task.**

6 Conclusions

In this artical, we introduce the research efforts on uncovering the intention behind equations in

ROUGE-L F1 SCORE					
MODEL	K=500	K=1000	K=2000	K=5000	All Training Data
MetaMath-7B	0.5036	0.5197	0.5189	0.5246	0.5255
WizardMath-7B	0.5241	0.5322	0.5408	0.5316	0.5412
MetaMath-13B	0.5110	0.5201	0.5337	0.5373	0.5386
BERT-SCORE F1 SCORE					
MODEL	K=500	K=1000	K=2000	K=5000	All Training Data
MetaMath-7B	0.7979	0.8051	0.8053	0.8084	0.8089
WizardMath-7B	0.8087	0.8117	0.8126	0.8116	0.8138
MetaMath-13B	0.8008	0.8076	0.8129	0.8132	0.8139

Table 3: Results of the impact of generated data. The table records the performance of the model for different values of K. The maximum value in each row is highlighted with a pink shade, and the maximum value in each column is indicated in bold.

mathematical problems.

Firstly, we stated the importance of MEIR task. In real life, when handling mathematical problems, students might not provide very detailed descriptions for each step. However, the mathematical equations at each step are essential. Therefore, understanding the intention behind those listing equations at each step means comprehending the student’s problem-solving approach. This is highly beneficial in the field of education.

Next, we introduced two modules: Imitation-based Generator and Intention Extractor. The Imitation-based Generator is used to increase data diversity. and the Intention Extractor is used to extract the intention behind each step. Through these two modules, we constructed the GSMI dataset. With minimal human resource consumption, these two modules can be used to generate more varied data. Experimental evidence has shown that the data generated by this structure is highly beneficial for improving model performance on MEIR tasks.

Subsequently, we selected a subset of candidate models and evaluated their ability to solve MEIR tasks on the GSMI dataset. The experimental results indicate that powerful general LLMs like GPT-4o still have shortcomings in understanding equations. Conversely, following a series of instruction tuning processes, small-scale open-source models demonstrate outstanding performance in understanding equations. Those models that have undergone mathematical fine-tuning, like MetaMath and WizardMath, excel in MEIR tasks.

In conclusion, we pioneered the study of equation intention analysis. We are the first to propose the MEIR task and have conducted thorough ex-

periments to explore the capability of LLMs in addressing this task. Exploring equation intention is an interesting and important topic, and needs further attention and in-depth research.

7 Limitations and Future Works

In this experimental work, we have exposed certain limitations. Due to computational constraints, the maximum model parameter size we used for fine-tuning open-source models was 13 billion. In the future, we will run the MEIR task on larger-scale open-source models to explore their capabilities in understanding mathematical equations.

In our experiments, we have demonstrated that the data generated through these Imitation-based Generator and Intention Extractor modules helps improve the model’s ability to understand equation intentions. In future work, we will propose more refined data augmentation mechanisms and introduce a larger-scale GSMI dataset.

Furthermore, for the generated data from Imitation-based Generator and Intention Extractor modules, we did not conduct comprehensive comparative analyses with existing datasets. In future work, a thorough comparative analysis is of necessity to make sure that our training data is of high quality.

Finally, as we have stated, the MEIR task closely aligns with educational settings. It is not sufficient to merely identify the intended meaning of correct equations. In the future, we aim to intelligently identify errors students make when producing equations in mathematical education scenarios. This places higher demands on language models, that they not only need to recognize and generalize the

464	intended meaning of correct equations, but also	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul	518
465	need to uncover the underlying reasons for errors	Arora, Steven Basart, Eric Tang, Dawn Song, and	519
466	in incorrect equations.	Jacob Steinhardt. 2021. Measuring mathematical	520
		problem solving with the math dataset .	521
		Namgyu Ho, Laura Schmid, and Se-Young Yun. 2023.	522
467	References	Large language models are reasoning teachers . In	523
		<i>Proceedings of the 61st Annual Meeting of the As-</i>	524
468	Shengnan An, Zexiong Ma, Zeqi Lin, Nanning Zheng,	<i>sociation for Computational Linguistics (Volume 1:</i>	525
469	Jian-Guang Lou, and Weizhu Chen. 2023. Learning	<i>Long Papers)</i> , ACL 2023, Toronto, Canada, July 9-14,	526
470	from mistakes makes llm better reasoner .	2023, pages 14852–14882. Association for Computa-	527
		tional Linguistics.	528
471	Hyung Won Chung, Le Hou, Shayne Longpre, Barret	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan	529
472	Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and	530
473	Wang, Mostafa Dehghani, Siddhartha Brahma, Al-	Weizhu Chen. 2021. Lora: Low-rank adaptation of	531
474	bert Webson, Shixiang Shane Gu, Zhuyun Dai,	large language models .	532
475	Mirac Suzgun, Xinyun Chen, Aakanksha Chowdh-	Shima Imani, Liang Du, and Harsh Shrivastava. 2023.	533
476	ery, Alex Castro-Ros, Marie Pellat, Kevin Robinson,	Mathprompter: Mathematical reasoning using large	534
477	Dasha Valter, Sharan Narang, Gaurav Mishra, Adams	language models . In <i>Proceedings of the The 61st An-</i>	535
478	Yu, Vincent Zhao, Yanping Huang, Andrew Dai,	<i>Annual Meeting of the Association for Computational</i>	536
479	Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Ja-	<i>Linguistics: Industry Track</i> , ACL 2023, Toronto,	537
480	cob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le,	Canada, July 9-14, 2023, pages 37–42. Association	538
481	and Jason Wei. 2022. Scaling instruction-finetuned	for Computational Linguistics.	539
482	language models .		
483	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,	Alejandro Jaimes and Nicu Sebe. 2007. Multimodal	540
484	Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias	human–computer interaction: A survey . <i>Computer</i>	541
485	Plappert, Jerry Tworek, Jacob Hilton, Reiichiro	<i>Vision and Image Understanding</i> , 108(1):116–134.	542
486	Nakano, Christopher Hesse, and John Schulman.	Special Issue on Vision for Human-Computer Inter-	543
487	2021. Training verifiers to solve math word prob-	action.	544
488	lems . <i>CoRR</i> , abs/2110.14168.		
489	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	Takeshi Kojima, Shixiang (Shane) Gu, Machel Reid, Yu-	545
490	Kristina Toutanova. 2019. BERT: Pre-training of	taka Matsuo, and Yusuke Iwasawa. 2022. Large lan-	546
491	deep bidirectional transformers for language under-	guage models are zero-shot reasoners . In <i>Advances in</i>	547
492	standing . In <i>Proceedings of the 2019 Conference of</i>	<i>Neural Information Processing Systems</i> , volume 35,	548
493	<i>the North American Chapter of the Association for</i>	pages 22199–22213. Curran Associates, Inc.	549
494	<i>Computational Linguistics: Human Language Tech-</i>		
495	<i>nologies, Volume 1 (Long and Short Papers)</i> , pages	Bohan Li, Yutai Hou, and Wanxiang Che. 2022. Data	550
496	4171–4186, Minneapolis, Minnesota. Association for	augmentation approaches in natural language pro-	551
497	Computational Linguistics.	cessing: A survey . <i>AI Open</i> , 3:71–90.	552
498	Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong	Junyi Li, Tianyi Tang, Wayne Xin Zhao, and Ji-Rong	553
499	Wu, Baobao Chang, Xu Sun, Jingjing Xu, Lei Li, and	Wen. 2021. Pretrained language model for text gener-	554
500	Zhifang Sui. 2023. A survey on in-context learning .	ation: A survey . In <i>Proceedings of the Thirtieth Inter-</i>	555
		<i>national Joint Conference on Artificial Intelligence,</i>	556
501	Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen,	<i>IJCAI-21</i> , pages 4492–4499. International Joint Con-	557
502	Yujia Yang, Minlie Huang, Nan Duan, and Weizhu	ferences on Artificial Intelligence Organization. Sur-	558
503	Chen. 2023. Tora: A tool-integrated reasoning	vey Track.	559
504	agent for mathematical problem solving . <i>CoRR</i> ,	Yinheng Li. 2023. A practical survey on zero-shot	560
505	abs/2309.17452.	prompt design for in-context learning . In <i>Proceed-</i>	561
		<i>ings of the Conference Recent Advances in Natural</i>	562
506	Prakhar Gupta, Cathy Jiao, Yi-Ting Yeh, Shikib Mehri,	<i>Language Processing - Large Language Models for</i>	563
507	Maxine Eskenazi, and Jeffrey Bigham. 2022. In-	<i>Natural Language Processings</i> , RANLP. INCOMA	564
508	structDial: Improving zero and few-shot general-	Ltd., Shoumen, BULGARIA.	565
509	ization in dialogue through instruction tuning . In	Chin-Yew Lin. 2004. ROUGE: A package for auto-	566
510	<i>Proceedings of the 2022 Conference on Empirical</i>	matic evaluation of summaries . In <i>Text Summariza-</i>	567
511	<i>Methods in Natural Language Processing</i> , pages 505–	<i>tion Branches Out</i> , pages 74–81, Barcelona, Spain.	568
512	525, Abu Dhabi, United Arab Emirates. Association	Association for Computational Linguistics.	569
513	for Computational Linguistics.		
514	Shruthi Hariharan, Vignesh Kumar Krishnamurthy,	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae	570
515	Utkarsh, and Jayantha Gowda Sarapanahalli. 2022.	Lee. 2023a. Visual instruction tuning . In <i>Ad-</i>	571
516	Enhancing slot tagging with intent features for task	<i>vances in Neural Information Processing Systems</i> ,	572
517	oriented natural language understanding using bert .	volume 36, pages 34892–34916. Curran Associates,	573
		Inc.	574

575	Tengxiao Liu, Qipeng Guo, Yuqing Yang, Xiangkun Hu, Yue Zhang, Xipeng Qiu, and Zheng Zhang. 2023b. Plan, verify and switch: Integrated reasoning with diverse x-of-thoughts . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 2807–2822. Association for Computational Linguistics.	631
576		632
577		633
578		634
579		635
580		636
581		637
582		638
583	Tiedong Liu and Bryan Kian Hsiang Low. 2023. Goat: Fine-tuned llama outperforms gpt-4 on arithmetic tasks .	639
584		640
585		641
586	Wentao Liu, Hanglei Hu, Jie Zhou, Yuyang Ding, Junsong Li, Jiayi Zeng, Mengliang He, Qin Chen, Bo Jiang, Aimin Zhou, and Liang He. 2023c. Mathematical language models: A survey .	642
587		643
588		644
589		645
590	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach .	646
591		647
592		648
593		649
594		650
595	Samuel Louvan and Bernardo Magnini. 2020. Recent neural methods on slot filling and intent classification for task-oriented dialogue systems: A survey . In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 480–496, Barcelona, Spain (Online). International Committee on Computational Linguistics.	651
596		652
597		653
598		654
599		655
600		656
601		657
602	Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qingwei Lin, Shifeng Chen, and Dongmei Zhang. 2023a. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct . <i>CoRR</i> , abs/2308.09583.	658
603		659
604		660
605		661
606		662
607		663
608	Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2023b. Wizardcoder: Empowering code large language models with evol-instruct .	664
609		665
610		666
611		667
612		668
613	Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. 2024. Large language models: A survey .	669
614		670
615		671
616		672
617	OpenAI. 2024. Gpt-4 technical report .	673
618	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback . In <i>Advances in Neural Information Processing Systems</i> , volume 35, pages 27730–27744. Curran Associates, Inc.	674
619		675
620		676
621		677
622		678
623		679
624		680
625		681
626		682
627		683
628	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation . In <i>Proceedings of the</i>	684
629		685
630		686
		687
		688
		689
		690
		691
		692
		693
		694
		695
		696
		697
		698
		699
		700
		701
		702
		703
		704
		705
		706
		707
		708
		709
		710
		711
		712
		713
		714
		715
		716
		717
		718
		719
		720
		721
		722
		723
		724
		725
		726
		727
		728
		729
		730
		731
		732
		733
		734
		735
		736
		737
		738
		739
		740
		741
		742
		743
		744
		745
		746
		747
		748
		749
		750
		751
		752
		753
		754
		755
		756
		757
		758
		759
		760
		761
		762
		763
		764
		765
		766
		767
		768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780
		781
		782
		783
		784
		785
		786
		787
		788
		789
		790
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800
		801
		802
		803
		804
		805
		806
		807
		808
		809
		810
		811
		812
		813
		814
		815
		816
		817
		818
		819
		820
		821
		822
		823
		824
		825
		826
		827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		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
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		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

690 Wen-Ju Yu, Shin-Yuan Hung, Annie Pei-I Yu, and Yu-Li
691 Hung. 2024. [Understanding consumers' continuance](#)
692 [intention of social shopping and social media partic-](#)
693 [ipation: The perspective of friends on social media.](#)
694 *Information & Management*, 61(4):103808.

695 Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang,
696 Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tian-
697 wei Zhang, Fei Wu, and Guoyin Wang. 2024. [Instruc-](#)
698 [tion tuning for large language models: A survey.](#)

699 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q.
700 Weinberger, and Yoav Artzi. 2020. [Bertscore: Evalu-](#)
701 [ating text generation with bert.](#)

702 Yue Zhang, Leyang Cui, Deng Cai, Xinting Huang,
703 Tao Fang, and Wei Bi. 2023a. [Multi-task instruction](#)
704 [tuning of llama for specific scenarios: A preliminary](#)
705 [study on writing assistance.](#)

706 Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex
707 Smola. 2023b. [Automatic chain of thought prompt-](#)
708 [ing in large language models.](#) In *The Eleventh In-*
709 *ternational Conference on Learning Representations,*
710 *ICLR 2023, Kigali, Rwanda, May 1-5, 2023.* Open-
711 Review.net.

712 Yue Zhou, Chenlu Guo, Xu Wang, Yi Chang, and Yuan
713 Wu. 2024. [A survey on data augmentation in large](#)
714 [model era.](#)