# Boosting Large Language Models with Socratic Method for Conversational Mathematics Teaching

### Yuyang Ding
East China Normal University
Shanghai, China
51265900017@stu.ecnu.edu.cn

### Hanglei Hu
East China Normal University
Shanghai, China
51254108023@stu.ecnu.edu.cn

### Jie Zhou*
East China Normal University
Shanghai, China
jzhou@cs.ecnu.edu.cn

### Qin Chen
East China Normal University
Shanghai, China
qchen@cs.ecnu.edu.cn

### Bo Jiang
East China Normal University
Shanghai, China
bjiang@deit.ecnu.edu.cn

### Liang He
East China Normal University
Shanghai, China
lhe@cs.ecnu.edu.cn

## ABSTRACT

With the introduction of large language models (LLMs), automatic math reasoning has seen tremendous success. However, current methods primarily focus on providing solutions or using techniques like Chain-of-Thought to enhance problem-solving accuracy. In this paper, we focus on improving the capability of mathematics teaching via a Socratic teaching-based LLM (SocraticLLM), which guides learners toward profound thinking with clarity and self-discovery via conversation. We collect and release a high-quality mathematical teaching dataset, named SocraticMATH, which provides Socratic-style conversations of problems with extra knowledge. Also, we propose a knowledge-enhanced LLM as a strong baseline to generate reliable responses with review, guidance/heuristic, rectification, and summarization. Experimental results show the great advantages of SocraticLLM by comparing it with several strong generative models. The codes and datasets are available on https://github.com/ECNU-ICALK/SocraticMath.

## CCS CONCEPTS

• **Computing methodologies** → **Natural language generation**;
• **Mathematics of computing** → *Mathematical software.*

## KEYWORDS

Socratic Teaching, LLMs, Mathematics, Conversation

---

*Corresponding author.

## 1 INTRODUCTION

Mathematics, viewed as a language that requires complex reasoning through structured symbols and systems, parallels the rules of spoken language, is a crucial aptitude for human intelligence. Recently, solving math problems autonomously via AI technology has attracted attention since as early as 1963 [8, 9, 16, 17].

The studies about mathematical AI (math word problems) are divided into three parts: statistical learning-based methodologies [34, 56], traditional machine learning techniques [23, 41, 42] and deep learning-based methods [14, 47]. Recently, large language models (LLMs) have achieved great successes in mathematics [33, 48], with various types of mathematical datasets [4, 13, 19, 26, 55] and mathematical LLMs [49, 52, 54] emerging.

However, previous research on mathematical LLMs mainly focused on improving mathematical ability while their application in teaching remains limited. LLMs tend to give the solution process directly when faced with mathematical problems, even when asked to play the role of a teacher [31]. There is a teaching method called Socratic questioning [10, 18], which ranges from providing direct hints to offering minimal guidance, allows instructors to customize their queries to appropriately challenge students while also aligning with their individual problem-solving abilities [2]. Taking Figure 1 as an example, the general LLMs tend to give the solutions with chain-of-though direct, whereas Socratic teaching-based LLMs guide the students using conversation.

Previous research has demonstrated the potential of utilizing LLMs for Socratic questioning. For instance, Shridhar et al. [45] utilized reinforcement learning to convert mathematical problem-solving steps into multiple sub-problems, thereby enabling the automatic generation of Socratic questioning. Similarly, Qi et al. [40] decomposed complex mathematical problems into a series of top-down questions, perceiving Socratic questioning as a recursive thinking process. This approach allows LLMs like ChatGPT to address fine-grained problems, ultimately resolving the main issue.

However, effectively integrating Socratic questioning methods into AI frameworks continues to be a significant challenge. First (**C1**), LLMs can not guarantee the reliability and quality of answers. The LLMs perform poorly on complex reasoning with hallucination problems. Second (**C2**), the strategies (when to ask, how to ask and ask what) to guide and heuristic students with just enough questions are unclear. Too many or too few queries that are too hard or too easy will influence students' earning process. Third
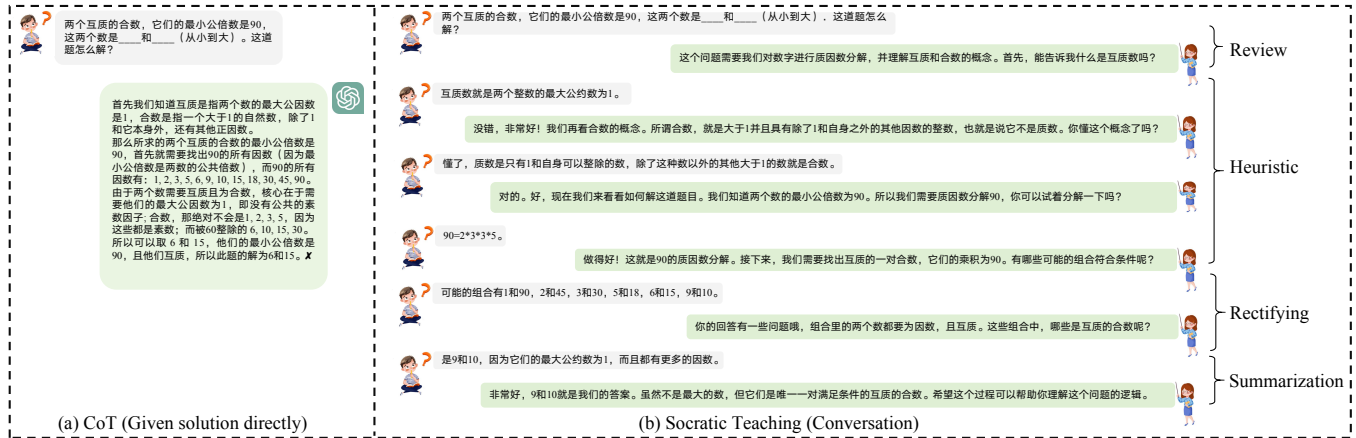
Figure 1: Examples of CoT and Socratic teaching.

(**C3**), there is a lack of relevant datasets for mathematics teaching [31]. Although classroom transcripts can provide a large amount of instructional tutoring data [15], there are issues like privacy security, crowdsourcing costs and annotation quality.

In this paper, we focus on integrating Socratic Questioning in mathematical education. For **C1**, we propose a knowledge-enhance Socratic teaching LLM (`SocraticLLM`) as a strong baseline to improve the reliability and quality of the generated response via extra knowledge. We design a strategy to tutor students step-by-step through the instructional structure of review, heuristic, rectify and summarize. Also, we create and release a high-quality Socratic-style mathematical dataset, `SocraticMATH`, which contains dialogue tutoring data with original questions, answers, and solutions and covers 513 knowledge points of primary school math. Extensive experiments show `SocraticLLM` outperforms strong baselines in terms of rich automatic and human/GPT-4 evaluation metrics.

The main contributions of this paper are summarized as follows. 1) We integrate the Socratic method with math teaching via a structured conversation with review, guidance/heuristic, rectification, and summarization. 2) We propose `SocraticLLM` as a strong baseline for mathematical teaching. To generate reliable responses, we design a Socratic-style prompt with extra knowledge to guide the teaching process. 3) We build and release a large-scale mathematical conversation dataset `SocraticMATH` with rich attributions over 500 knowledge points. A series of experiments on `SocraticMATH` indicate the effectiveness of our `SocraticLLM`.

## 2 RELATED WORK

**Mathematical Reasoning.** Mathematical reasoning refers to the process of using algorithms and computational models to solve complex mathematical problems [28, 29]. Recently, large language models have obtained great performance in mathematical reasoning using instruction tuning [27, 30, 48], in-context learning [11, 21] and tool-enhanced methods [38, 43]. Several studies solved the math problem step by step, such as Chain of Thought (CoT) [48], Tree of Thought (ToT) [51] and Graph of Thought (GoT) [7]. However, most of the previous studies focused on improving the accuracy of mathematical reasoning by giving the solutions to the problem. The goal of this paper is mathematical teaching using Socratic teaching.

**Theoretical Background of Socratic Teaching.** Carefully formulated questions can encourage students to self-explain [12],

enhance their understanding of the task, and facilitate effective planning of solutions [24]. Additionally, they can help identify significant gaps in student knowledge [35]. This spectrum enables educators to craft questions that are appropriately challenging yet within a student's capacity to respond [2]. This method, called Socratic questioning, is based on the philosophy that knowledge is not simply transferred but uncovered through a dynamic process of inquiry and dialogue. In summary, although there are established principles and guidelines, the practical application of Socratic questioning with AI presents challenges.

**Technologies of Socratic Teaching.** Research in the realm of automatic Socratic tutoring systems has shown progress, but the applicability of such systems is often constrained by the predefined and manually tailored nature of Socratic utterances for specific exercises [3]. Al-Hossami [2] presented a dataset comprising Socratic dialogues aimed at assisting novice programmers in rectifying errors in fundamental computational problems. Shridhar's study [45] explored the strategies involved in the automatic generation of math problem solutions for educational purposes. Qi et al. [40] perceived Socratic questioning as a recursive thought process, which breaks down complex problems into simpler, related sub-problems. However, newer technologies such as LLMs are still not widely implemented in tutoring systems, particularly in the field of mathematics education.

## 3 DATASET

### 3.1 Dataset Construction

For lack of a Socratic teaching-based mathematics dataset, we collect and annotate a diverse dataset, `SocraticMATH`, to promote the research of this domain. We construct the dataset with three phases: data collection, pre-annotation and human annotation.

**Data Collection.** The questions are mainly derived from the real primary school exams in China. To guarantee diversity, these problems cover the main maths knowledge points at the primary school level, ensuring that all the questions are manually labeled with solutions using markdown format. The questions consist of multiple-choice, fill-in-the-blank, and answer questions.

**Pre-Annotation.** To reduce the cost of human annotation, we pre-annotate the conversations for all 8935 questions using GPT-4. We design an in-context prompt using a manually authored

**Table 1: Comparison with existing datasets**

| Dataset | Size | Lang | Ans | Solution | Conv | Socratic | KG | Difficulty |
|---|---|---|---|---|---|---|---|---|
| SocraticMATH | 6,846 | CH | √ | Textual Steps | √ | √ | √ | √ |
| Math23k [47] | 23,162 | CH | √ | Equation | × | × | × | × |
| AQuA [26] | 97,975 | EN | √ | Textual Steps | × | × | × | × |
| MathQA [4] | 37,297 | EN | √ | Equation | × | × | × | × |
| GSM8K [13] | 8,792 | EN | √ | Textual Steps | × | × | × | × |
| SVAMP [39] | 1,000 | EN | √ | Equation | × | × | × | × |
| MATHDIAL [32] | 2,861 | EN | √ | Generation | Semi | × | × | × |

**Table 2: The statistical information of SocraticMATH.**

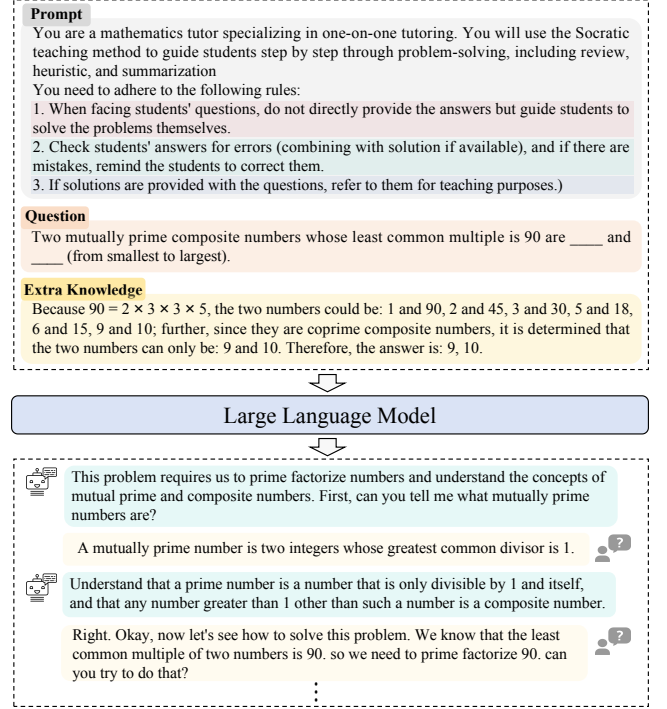|  | TRAIN | DEV | TEST | TOTAL |
|---|---|---|---|---|
| #CONV | 5,476 | 685 | 685 | 6,846 |
| #TURN/CONV | 4.95 | 4.95 | 5.02 | 4.96 |
| #WORD/SOLUTION | 73.29 | 73.17 | 72.70 | 73.21 |
| #WORD/UTTERANCE | 86.49 | 87.02 | 85.56 | 86.45 |
| #KG | 495 | 333 | 332 | 513 |
| #KG/CONV | 2.00 | 2.03 | 2.03 | 2.00 |

high-quality example to enhance the quality of the generated conversation. Additionally, we let GPT-4 act as a Socratic-style teacher with various student personality requirements (such as naughtiness, self-confidence, and carelessness) to ensure the richness of the generated dialogue.

**Human Annotation.** Though GPT-4 generates the conversation with a Socratic style for each math question, their quality is limited. First, LLMs are not good at math reasoning and their answers can be tainted with factual errors [28]. Second, the LLMs always give solutions to answer students' questions while lacking the teaching skills with inspiration and guidance. Thus, we clean and re-annotate the conversation to improve the data quality. Particularly, we first eliminate the data with an abnormal number of dialogue rounds. Then, we manually perform the annotation work on the data to optimize the logic and coherence of the conversation. Each sample is labeled by three experts, who are good at teaching and math. Due to the complexity of the conversation, the three experts label the conversation one by one to revise the errors iteratively. Particularly, we delete more than 23% dialogues and modify more than 18% dialogues, where over 5% utterances are revised.

## 3.2 Dataset Analysis

**Characters of SocraticMATH.** We present the statistical information of our SocraticMATH dataset in Table 2. Our dataset contains 513 knowledge points, almost all the knowledge points of math at the primary school level. To better guide the student to solve the math problem, the average number of turns for each conversation is about 5. The average length of utterances is about 86 words to provide detailed information patiently.

**Comparison with Exiting Datasets.** To show the advantages of SocraticMATH, we compare it with existing typical mathematics datasets (Table 1). Most datasets only provide the equation or textual steps directly to solve the math problem. MATHDIAL dataset contains the semi-annotated conversation where the students' questions are generated by LLMs. Furthermore, MATHDIAL mainly focuses on answering students' questions without the Socratic method, which requires teaching skills to inspire, guide, and inquire actively step by step. We also provide extensive attribution information, such as related knowledge points and difficulty levels,



**Figure 2: The framework of SocraticLLM.**

where #KG and #KG/Conv are the total number of knowledge points and the average number of knowledge points for each conversation.

## 4 OUR METHOD

We propose SocraticLLM as a strong and simple baseline for mathematics teaching (Figure 2). It generates responses with the teaching skills of review, guidance/heuristic, rectification, and summarization via LLMs. We design a Socratic-style prompt and integrate the original question with the extra knowledge (i.e., solution, answer) to improve the quality of the responses.

The Socratic-style prompt $P$ contains the task's definition and requirements. We ask the model to act as a one-on-one mathematical tutor using the Socratic teaching method. In particular, we require SocraticLLM to guide the student rather than answering questions directly. Then, we ask SocraticLLM to check and rectify the errors since the model tends to trust the users. To reduce the hallucination, we demand the model generate the response based on the extra knowledge by inputting the detailed solution and answer.

Formally, given the question $Q$, the prompt $P$, and the extra knowledge $K$, we aim to generate the response $R_i$ based on the history conversation $H_{i-1} = \{R_1, U_1, R_2, U_2, ..., R_{i-1}, U_{i-1}\}$, where $U_i$ is the user's answer of $i$-th turn.

$$p(R_i|P, Q, K, H_{i-1}) = \prod_j^{|R_i|} p(R_i^j|P, Q, K, H_{i-1}, R_i^{1:j-1}) \quad (1)$$

where $R_i^j$ is the $j$-th word of response $R_i$, $|R_i|$ is the length of $R_i$. We use a language model $\mathcal{M}_\theta$ to model the generation probability, $p(R_i^j|P, Q, K, H_{i-1}, R_i^{1:j-1}) = \mathcal{M}_\theta(P, Q, K, H_{i-1}, R_i^{1:j-1})$, where $\theta$ is the learnable parameters of $\mathcal{M}_\theta$. Particularly, we use Low-Rank

**Table 3: Main results of automatic and human evaluation.**

| | Automatic Evaluation | | | | | | | | | Human Evaluation | | GPT-4 Evaluation | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | B-1 | B-2 | B-3 | B-4 | METEOR | R-1 | R-2 | R-L | BARTScore | Reliability | Socratic | Reliability | Socratic |
| mT5 | 0.303 | 0.225 | 0.174 | 0.135 | 0.320 | 0.439 | 0.227 | 0.323 | 0.724 | 5.714 | 6.405 | 6.80 | 5.62 |
| LLaMA2-7B | 0.301 | 0.213 | 0.157 | 0.116 | 0.335 | 0.454 | 0.216 | 0.308 | 0.710 | 5.310 | 6.333 | 5.76 | 4.87 |
| Qwen1.5-7B | 0.341 | 0.247 | 0.185 | 0.139 | 0.367 | 0.481 | 0.241 | 0.341 | 0.726 | 6.595 | 6.762 | 7.84 | 6.96 |
| ChatGPT | 0.273 | 0.194 | 0.147 | 0.111 | 0.398 | 0.431 | 0.197 | 0.257 | 0.695 | 6.500 | 6.405 | 8.07 | 6.70 |
| GPT-4 | 0.332 | 0.240 | 0.181 | 0.137 | **0.410** | 0.471 | 0.227 | 0.306 | 0.715 | 7.024 | 6.833 | - | - |
| SocraticLLM | **0.352** | **0.256** | **0.193** | **0.147** | 0.378 | **0.490** | **0.250** | **0.351** | **0.730** | **7.119** | **7.190** | **8.40** | **7.14** |
| - Prompt | 0.341 | 0.248 | 0.188 | 0.143 | 0.369 | 0.484 | 0.246 | 0.344 | 0.727 | 7.048 | 6.857 | 8.16 | 6.98 |
| - Knowledge | 0.347 | 0.253 | 0.191 | 0.145 | 0.374 | 0.488 | 0.247 | 0.350 | 0.729 | 6.643 | 6.692 | 7.83 | 6.91 |

Adaptation (LoRA) technology to improve the efficiency of training, where only a small number of extra parameters $\theta$ are trainable [20].

Then, the cross-entropy function is used to measure the generation losses of the response,

$$\mathcal{L} = \sum_i^N \sum_j^{|R_i|} log(p(R_i^j | P, Q, K, H_{i-1}, R_i^{1:j-1})) \qquad (2)$$

where $N$ is the number of turns in the conversation.

## 5 EXPERIMENTS

### 5.1 Experimental Setups.

**Metrics.** We adopt several typical automatic metrics for generation tasks, including BLEU [37] (marked as B-1/2/3/4), ROUGE [25] (marked as R-1/2/L), METEOR [6] and BARTScore [53], to evaluate the effectiveness of SocraticLLM turn by turn. We also conduct human evaluation and GPT-4 evaluation [22].

**Baselines.** We compare SocraticLLM with several typical and strong seq-to-seq models. **mT5** [50] trains a text-to-text transformer model on multilingual datasets via multi-task learning. **LLaMA2-7B** [46] and **Qwen1.5-7B** [5] are strong LLMs fine-tuned on our dataset. **ChatGPT** [44] and **GPT-4** [1] act as a mathematics tutor with the Socratic method, one of the SOTA conversation models.

**Implementation Details.** We select Qwen1.5-7B as the base LLMs and train it on A800 GPU with 80G. We set the rank of LoRA as 64. The learning rate is 3e-4 and the batch size is 64.

### 5.2 Experimental Results

We report the main results of SocraticLLM and the selected baselines using automatic, human and GPT-4 evaluation (Table 3).

**Automatic Evaluation.** From the results, we observe that SocraticLLM achieves better results by comparing with the strong baselines over automatic metrics in most cases, indicating our model's effectiveness. Note that GPT-4 outperforms SocraticLLM in terms of METEOR without training because SocraticMATH is modified based on the dataset generated by GPT-4. Furthermore, these automatic metrics can not measure the quality of conversation in Socratic mathematical teaching. They mainly calculate the semantic information between the generated responses and the reference while ignoring the logic and fact errors in the output text.

**Human/GPT-4 Evaluation.** Single reference-based automatic metrics are not always reliable to reflect the real quality of the generated responses [36]. Therefore, we also conduct human and GPT-4 evaluations by crowd-sourcing and GPT-4. Particularly, we

ask the three experts and GPT-4 to label 150 samples randomly selected from the test set with guidelines. Based on a pre-determined scoring rubric, they annotate the generated response from reliability and Socratic strategy with scores 1-10. Reliability judges whether the model corrects the students' errors precisely and Socratic represents the guide and heuristic abilities of the model. We report the average scores here. From the results, we observe that SocraticLLM obtains the best results in both human and GPT-4 evaluations, showing that our model can reduce the hallucination with the Socratic method. Moreover, in the human evaluation, we find that LLMs like ChatGPT tend to believe the users' responses without a doubt. It is interesting to explore in further work.

**Ablation Studies.** We also conduct an ablation test to explore the effectiveness of the main parts consisting of SocraticLLM by removing Socratic-style prompt (- Prompt), extra knowledge (- Knowledge) and all of them (Qwen1.5-7B), respectively. We observe that both the Socratic-style prompt and extra knowledge are useful for SocraticLLM. The Socratic-style prompt enhances the model to learn the teaching skills based on structured conversation. Then, incorporating the extra knowledge into SocraticLLM reduces the hallucination problem by correcting the fact errors.

## 6 CONCLUSIONS AND FURTHER WORK

This paper presents SocraticLLM as a strong baseline to tutor students through structured conversation, encompassing review, heuristic, rectification, and summarization by integrating the Socratic method into mathematical education. By infusing extra knowledge into the LLM architecture, we ensure the reliability and quality of generated responses, thereby overcoming the issue of poor performance in complex reasoning tasks. To mitigate the scarcity of relevant datasets for mathematical teaching, we contribute the SocraticMATH dataset, comprising diverse dialogue data, enabling further advancements in this domain. Our experiments demonstrate the efficacy of SocraticLLM in enhancing mathematical education. In further work, we would like to incorporate personal information and knowledge graphs to pave the way for future developments in adaptive and interactive learning environments.

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).

[2] Erfan Al-Hossami, Razvan Bunescu, Ryan Teehan, Laurel Powell, Khyati Mahajan, and Mohsen Dorodchi. 2023. Socratic questioning of novice debuggers: A benchmark dataset and preliminary evaluations. In *Proceedings of BEA*. 709–726.

[3] Zeyad Alshaikh, Lasang Jimba Tamang, and Vasile Rus. 2020. Experiments with a socratic intelligent tutoring system for source code understanding. In *FLAIRS*.

[4] Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. MathQA: Towards Interpretable Math Word Problem Solving with Operation-Based Formalisms. In *Proceedings of ACL*. 2357–2367.

[5] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609* (2023).

[6] Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*. 65–72.

[7] Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Michal Podstawski, Hubert Niewiadomski, Piotr Nyczyk, et al. 2023. Graph of thoughts: Solving elaborate problems with large language models. *arXiv preprint arXiv:2308.09687* (2023).

[8] Daniel Bobrow et al. 1964. Natural language input for a computer problem solving system. (1964).

[9] Diane J Briars and Jill H Larkin. 1984. An integrated model of skill in solving elementary word problems. *Cognition and instruction* 1, 3 (1984), 245–296.

[10] Thomas C Brickhouse and Nicholas D Smith. 2009. Socratic teaching and Socratic method. (2009).

[11] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *NeurIPS* 33 (2020), 1877–1901.

[12] Michelene TH Chi, Nicholas De Leeuw, Mei-Hung Chiu, and Christian LaVancher. 1994. Eliciting self-explanations improves understanding. *Cognitive science* 18, 3 (1994), 439–477.

[13] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training Verifiers to Solve Math Word Problems. arXiv:2110.14168 [cs.LG]

[14] Jelle Couperus. 2023. *Large Language Models and Mathematical Understanding*. Master's thesis.

[15] Dorottya Demszky and Heather Hill. 2022. The NCTE Transcripts: A dataset of elementary math classroom transcripts. *arXiv preprint arXiv:2211.11772* (2022).

[16] Edward A Feigenbaum, Julian Feldman, et al. 1963. *Computers and thought*. Vol. 7. New York McGraw-Hill.

[17] Charles R Fletcher. 1985. Understanding and solving arithmetic word problems: A computer simulation. *Behavior Research Methods, Instruments, & Computers* 17, 5 (1985), 565–571.

[18] Rick Garlikov. 2001. The Socratic method: Teaching by asking instead of by telling. *Website, http://www. garlikov. com/Soc_Meth. html* (2001).

[19] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874* (2021).

[20] Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-Rank Adaptation of Large Language Models. In *International Conference on Learning Representations*. https://openreview.net/forum?id=nZeVKeeFYf9

[21] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361* (2020).

[22] Tom Kocmi and Christian Federmann. 2023. Large language models are state-of-the-art evaluators of translation quality. *arXiv preprint arXiv:2302.14520* (2023).

[23] Nate Kushman, Yoav Artzi, Luke Zettlemoyer, and Regina Barzilay. 2014. Learning to Automatically Solve Algebra Word Problems. In *Proceedings of ACL*, Kristina Toutanova and Hua Wu (Eds.). Association for Computational Linguistics, Baltimore, Maryland, 271–281. https://doi.org/10.3115/v1/P14-1026

[24] H Chad Lane and Kurt VanLehn. 2005. Teaching the tacit knowledge of programming to noviceswith natural language tutoring. *Computer Science Education* 15, 3 (2005), 183–201.

[25] Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*. 74–81.

[26] Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017. Program Induction by Rationale Generation: Learning to Solve and Explain Algebraic Word Problems. In *Proceedings of ACL*. 158–167.

[27] Tiedong Liu and Bryan Kian Hsiang Low. 2023. Goat: Fine-tuned LLaMA Outperforms GPT-4 on Arithmetic Tasks. http://arxiv.org/abs/2305.14201 arXiv:2305.14201 [cs].

[28] Wentao Liu, Hanglei Hu, Jie Zhou, Yuyang Ding, Junsong Li, Jiayi Zeng, Mengliang He, Qin Chen, Bo Jiang, Aimin Zhou, et al. 2023. Mathematical Language Models: A Survey. *arXiv preprint arXiv:2312.07622* (2023).

[29] Pan Lu, Liang Qiu, Wenhao Yu, Sean Welleck, and Kai-Wei Chang. 2023. A Survey of Deep Learning for Mathematical Reasoning. In *Proceedings of ACL*, Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, Toronto, Canada, 14605–14631. https://doi.org/10.18653/v1/2023.acl-long.817

[30] Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qingwei Lin, Shifeng Chen, and Dongmei Zhang. 2023. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct. *arXiv preprint arXiv:2308.09583* (2023).

[31] Jakub Macina, Nico Daheim, Sankalan Pal Chowdhury, Tanmay Sinha, Manu Kapur, Iryna Gurevych, and Mrinmaya Sachan. 2023. MathDial: A Dialogue Tutoring Dataset with Rich Pedagogical Properties Grounded in Math Reasoning Problems. *arXiv preprint arXiv:2305.14536* (2023).

[32] Jakub Macina, Nico Daheim, Sankalan Pal Chowdhury, Tanmay Sinha, Manu Kapur, Iryna Gurevych, and Mrinmaya Sachan. 2023. MathDial: A Dialogue Tutoring Dataset with Rich Pedagogical Properties Grounded in Math Reasoning Problems. In *EMNLP*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, 5602–5621. https://aclanthology.org/2023.findings-emnlp.372

[33] Nikolaos Matzakos, Spyridon Doukakis, and Maria Moundridou. 2023. Learning Mathematics with Large Language Models: A Comparative Study with Computer Algebra Systems and Other Tools. *International Journal of Emerging Technologies in Learning (Online)* 18, 20 (2023), 51.

[34] Arindam Mitra and Chitta Baral. 2016. Learning to use formulas to solve simple arithmetic problems. In *Proceedings of ACL*. 2144–2153.

[35] Laurie Murphy and Josh Tenenberg. 2005. Do computer science students know what they know? A calibration study of data structure knowledge. In *Proceedings of SIGCSE*. 148–152.

[36] Jekaterina Novikova, Ondřej Dušek, Amanda Cercas Curry, and Verena Rieser. 2017. Why We Need New Evaluation Metrics for NLG. In *Proceedings of EMNLP*, Martha Palmer, Rebecca Hwa, and Sebastian Riedel (Eds.). 2241–2252.

[37] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*. 311–318.

[38] Aaron Parisi, Yao Zhao, and Noah Fiedel. 2022. TALM: Tool Augmented Language Models. arXiv:2205.12255 [cs.CL]

[39] Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are NLP Models really able to Solve Simple Math Word Problems?. In *Proceedings of NAACL*. 2080–2094.

[40] Jingyuan Qi, Zhiyang Xu, Ying Shen, Minqian Liu, Di Jin, Qifan Wang, and Lifu Huang. 2023. The Art of SOCRATIC QUESTIONING: Recursive Thinking with Large Language Models. In *Proceedings of EMNLP*. 4177–4199.

[41] Subhro Roy and Dan Roth. 2015. Solving General Arithmetic Word Problems. In *Proceedings of EMNLP*. 1743–1752.

[42] Subhro Roy, Tim Vieira, and Dan Roth. 2015. Reasoning about quantities in natural language. *TACL* 3 (2015), 1–13.

[43] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language Models Can Teach Themselves to Use Tools. arXiv:2302.04761 [cs.CL]

[44] John Schulman, B Zoph, C Kim, J Hilton, J Menick, J Weng, JFC Uribe, L Fedus, L Metz, M Pokorny, et al. 2022. ChatGPT: Optimizing language models for dialogue. In *OpenAI blog*.

[45] Kumar Shridhar, Alessandro Stolfo, and Mrinmaya Sachan. 2022. Distilling multi-step reasoning capabilities of large language models into smaller models via semantic decompositions. *arXiv preprint arXiv:2212.00193* (2022).

[46] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971* (2023).

[47] Yan Wang, Xiaojiang Liu, and Shuming Shi. 2017. Deep neural solver for math word problems. In *Proceedings of EMNLP*. 845–854.

[48] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *NeurIPS* 35 (2022), 24824–24837.

[49] Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. Wizardlm: Empowering large language models to follow complex instructions. *arXiv preprint arXiv:2304.12244* (2023).

[50] Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A Massively Multilingual Pre-trained Text-to-Text Transformer. In *Proceedings of NAACL*. 483–498.

[51] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. *arXiv preprint arXiv:2305.10601* (2023).

[52] Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2023. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284* (2023).

[53] Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021. Bartscore: Evaluating generated text as text generation. *Advances in Neural Information Processing Systems* 34 (2021), 27263–27277.

[54] Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. 2023. Mammoth: Building math generalist models through hybrid instruction tuning. *arXiv preprint arXiv:2309.05653* (2023).

[55] Wei Zhao, Mingyue Shang, Yang Liu, Liang Wang, and Jingming Liu. 2020. Ape210k: A large-scale and template-rich dataset of math word problems. *arXiv preprint arXiv:2009.11506* (2020).

[56] Lipu Zhou, Shuaixiang Dai, and Liwei Chen. 2015. Learn to solve algebra word problems using quadratic programming. In *Proceedings of EMNLP*. 817–822.