EMBEDDING SELF-CORRECTION AS AN INHERENT ABILITY IN LARGE LANGUAGE MODELS FOR EN HANCED MATHEMATICAL REASONING

Anonymous authors

006

007

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

032

034

037

Paper under double-blind review

ABSTRACT

Accurate mathematical reasoning with Large Language Models (LLMs) is crucial in revolutionizing domains that heavily rely on such reasoning. However, LLMs often encounter difficulties in certain aspects of mathematical reasoning, leading to flawed reasoning and erroneous results. To mitigate these issues, we introduce a novel mechanism, the Chain of Self-Correction (CoSC), specifically designed to embed self-correction as an inherent ability in LLMs, enabling them to validate and rectify their own results. The CoSC mechanism operates through a sequence of self-correction stages. In each stage, the LLMs generate a program to address a given problem, execute this program using program-based tools to obtain an output, subsequently verify this output. Based on the verification, the LLMs either proceed to the next correction stage or finalize the answer. This iterative self-correction process allows the LLMs to refine its reasoning steps and improve the accuracy of its mathematical reasoning. To enable the CoSC mechanism at a low cost, we employ a two-phase finetuning approach. In the first phase, the LLMs are trained with a relatively small volume of seeding data generated from GPT-4, establishing an initial CoSC capability. In the second phase, the CoSC capability is further enhanced by training with a larger volume of selfgenerated data using the trained model in the first phase, without relying on the paid GPT-4. Our comprehensive experiments demonstrate that CoSC significantly improves performance on traditional mathematical datasets among existing opensource LLMs. Notably, our CoSC-Code-34B model achieved a 53.5% score on MATH, the most challenging mathematical reasoning dataset in the public domain, surpassing the performance of well-established models such as ChatGPT, GPT-4, and even multi-modal LLMs like GPT-4V, Gemini-1.0 Pro, and Gemini-1.0 Ultra. It's important to note that, unlike these proprietary models, our CoSC performs inference in a zero-shot manner, without the need for demonstrations. The code and data for this work will be released once this paper is accepted.

⁰³⁹ 1 INTRODUCTION

040

041 Large Language Models (LLMs), such as GPT-4 (OpenAI, 2023a), have recently demonstrated 042 state-of-the-art performance across a variety of natural language processing (NLP) tasks, including 043 natural language generation and understanding (Chowdhery et al., 2023; Team et al., 2023; Anil 044 et al., 2023; Penedo et al., 2023). Despite their success, LLMs often struggle with mathematical reasoning tasks due to their lack of explicit logical reasoning and judgment, which are crucial for solving such problems. Moreover, there is a fundamental gap between natural language and the 046 language of mathematical formulas, which further complicates these tasks. As a result, accurate 047 mathematical reasoning remains an essential yet challenging capability for LLMs to develop, in order 048 to further advance various domains. Consequently, it is still an open challenge to tackle mathematical problems for existing open-source LLMs, such as LLaMA-2 (Touvron et al., 2023). 050

To improve the mathematical reasoning abilities, numerous approaches have been investigated in previous research, including prompting (Wei et al., 2022; Chen et al., 2023a; Wang et al., 2023b), pretraining (Azerbayev et al., 2023; Fu et al., 2023; Shao et al., 2024), and finetuning (Luo et al., 2023; Yu et al., 2023; Yu et al., 2023; Gou et al., 2023b). In particular, finetuning



Figure 1: Comparison of four reasoning frameworks for solving an example mathematical question. 074 (a) Chain of Thoughts (CoT) (Wei et al., 2022). (b) Program of Thoughts (PoT) (Chen et al., 2023a). 075 (c) ToRA (Gou et al., 2023b) that incorporates CoT, PoT, and the utilization of tools. (d) Our proposed 076 CoSC consists of a sequence of multiple self-correction stages (two stages are shown in this example). 077 Each stage has four sub-stages: (p1) LLMs generate program w.r.t. the question; (o1) execute the 078 program to obtain program output; (v1) perform two-step verification for consistency of the question with both the generated program and the program output; (c1) conclude a refined answer or continue 079 the next subsequent self-correction stage depending on the verification result. The final answer is extracted from the last conclusion sub-stage with regular expression matching. 081

082

has become a favored technique among them, which updates open-sourced LLMs based on previously
 generated high-quality question-response pair datasets. Compared to open-source LLMs, finetuning
 has demonstrated significant improvement, but there is still potential room for further enhancement.

Most current methods (Yu et al., 2023; Gou et al., 2023b) generate finetuning datasets by prompting 087 GPT-4 to rephrase mathematical questions from various perspectives or incorporate chain-of-thoughts (CoT) analysis (Wei et al., 2022) and program-of-thoughts (PoT) code (Chen et al., 2023a) to diversify mathematical responses. As a result, LLMs trained on these datasets can comprehend 090 different questions and learn to use code to solve mathematical problems. However, precise multi-091 round reasoning capabilities remain challenging for them. Once potential errors occur throughout 092 reasoning stages, it can lead to incorrect results. Consequently, it is essential to incorporate a selfcorrection mechanism into mathematical responses, which can enable LLMs to learn to correct themselves in multiple rounds. In addition, while supervised finetuning datasets can be developed 094 with GPT-4's assistance, it still requires human experiments with different prompts and the cost of 095 using the interface. This highlights the need for effective finetuning with unlabeled datasets. 096

To address the aforementioned challenges, our study introduces the Chain of Self-Correction (CoSC), 098 a novel mechanism designed to embed self-correction as an inherent capability in LLMs, enabling 099 them to validate and rectify their own results. The CoSC mechanism operates through a sequence of self-correction stages, where LLMs generate a program to solve a given problem, execute the program 100 using program-based tools to obtain an output, and subsequently verify this output. Depending on the 101 verification, the LLMs either advance to a subsequent stage of self-correction or conclude with the 102 refined solution. An example of our CoSC reasoning trajectory with multiple self-correction stages 103 (two stages are shown in this example) is shown in Fig. 1. This iterative self-correction process allows 104 the LLMs to refine their reasoning steps and improve the accuracy of its mathematical reasoning. 105

To implement the CoSC mechanism at a low cost, we adopt a two-phase finetuning approach. In the
 first phase, termed the CoSC foundational learning, LLMs are trained with a relatively small volume
 of seeding data generated from GPT-4, equipping them with a baseline proficiency in the CoSC

methodology. In particular, we prompt GPT-4 with training questions from MATH (Hendrycks et al., 2021) and GSM8K (Cobbe et al., 2021) datasets to generate mathematical reasoning trajectories that adhere to the CoSC protocol. Specifically, each generated trajectory consists of program-of-thoughts code, program output, a two-step verification process that ensures the alignment of the question with both the generated program and the resulting output, and a conclusion to determine whether the trajectory should be refined or if the final answer can be provided.

114 Utilizing GPT-4 can be expensive, especially when dealing with large volumes of training data. 115 Alternatively, we propose a cost-free method to further boost performance through self-enhancement 116 in the second phase. This second phase, referred to as **CoSC self-enhancement**, builds upon 117 the foundational learning by further adapting the LLMs obtained from the first phase with self-118 generated trajectories. These trajectories are produced by the models trained in the foundational phase, allowing for the generation of a substantial volume of data without additional GPT-4 intervention. 119 In both phases, we retain only the trajectories whose answers match the ground-truth labels of the 120 corresponding questions. 121

¹²² In summary, our study makes the following contributions:

123 124

125

127

128

129

130

136

137

138

140

141

142

143

144

145

146

147 148

149 150

- We propose the Chain of Self-Correction (CoSC) mechanism, which effectively embeds self-correction as an inherent ability in Large Language Models (LLMs). Once LLMs learn this ability during training, they can self-correct in a zero-shot setting during inference without the need for external feedback or few-shot demonstrations. With the inherent self-correction ability, even an originally weak LLM is able to achieve excellent performance in mathematical reasoning, as strongly supported by our experimental results. This unique contribution distinguishes our work from related works discussed in Section 2.2 of this field.
- To implement the CoSC mechanism at a low cost, we introduce a two-phase finetuning approach. The first phase involves CoSC foundational learning, where we use a relatively small volume of seeding data generated by GPT-4. In the second phase, CoSC self-enhancement occurs using a larger volume of self-generated data with the model obtained from the first phase, without relying on paid GPT-4.
 - Our comprehensive experiments demonstrate that the CoSC mechanism provides a new benchmark for performance on established mathematical datasets when compared to existing open-source LLMs. Notably, our CoSC-Code-34B model achieves superior performance over both closed-source non-multi-modal LLMs and some of multi-modal LLMs, particularly on the challenging MATH dataset. These findings underscore the potential of our CoSC to improve the accuracy and reliability of LLMs in complex problem-solving domains.
 - The proposed CoSC mechanism, by embedding self-correction as an inherent capability in LLMs, enables them to think before responding to a question, creating an internal chain of self-correction to progressively verify and rectify their original answers. It is more akin to the slow thinking process of humans, which is particularly helpful in solving difficult mathematical reasoning problems. This approach can provide valuable insights for future research and contribute to the ongoing advancement of LLMs across various domains.

2 RELATED WORK

151 2.1 LLMs FOR MATHEMATICAL REASONING

152 Mathematical reasoning (Liu et al., 2023a; Wang et al., 2023a; Huang et al., 2024b; Toshniwal 153 et al., 2024; Chen et al., 2024a; Zhang et al., 2024) is a challenging reasoning task for LLMs, which 154 requires the ability to understand mathematical concepts, computation and multi-round reasoning. Existing mathematical reasoning approaches can be broadly classified into three categories: (1) 156 Prompting methods (Wei et al., 2022; Chen et al., 2023a; Wang et al., 2023b) focus on extracting 157 the inherent mathematical reasoning skills of LLMs by utilizing well-crafted prompting strategies 158 during inference. Notably, they leverage the existing knowledge in LLMs without the need for parameter updates. (2) Pretraining methods (Azerbayev et al., 2023; Fu et al., 2023; Yang et al., 2024) 159 pre-train LLMs on large-scale corpora containing mathematical problems and related content with 160 language modeling objectives. The goal is to train a base foundation language model as a platform 161 for mathematical domain. (3) Finetuning methods (Luo et al., 2023; Yu et al., 2023; Yu et al., 2023;

Gou et al., 2023b; Liu et al., 2023b; Gao et al., 2024; Mitra et al., 2024) refine the mathematical reasoning capabilities of LLMs by offering more targeted training, which is inspired by distillation techniques (Furlanello et al., 2018; Hsieh et al., 2023). Central to this approach is the generation of high-quality question-response pair datasets, which needs the assistance of complexity-based chain-of-thoughts prompting or tools-based augmentation. In this paper, we propose a Chain of Self-Correction (CoSC) along this direction that incorporates an iterative self-correction process into datasets generation.

169 170

182 183

184

2.2 EXISTING METHODS RELATED TO SELF-CORRECTION IN LLMS

171 Mathematical reasoning poses a significant challenge due to its demand for precise multi-round 172 logical reasoning to solve problems. The potential for errors increases with each reasoning step, 173 making it crucial for LLMs to have the ability to self-correct in order to produce accurate results. 174 There are some recent studies (Chen et al., 2023b; Gou et al., 2023a; Lightman et al., 2023; Huang 175 et al., 2024a; Chen et al., 2024b) attempt to enable large language models to perform self-correction 176 by either prompting methods or fine-tuning methods. For prompting methods, they can correct their 177 responses when interacting with external tools (Gou et al., 2023a; Chen et al., 2023b), such as search 178 engines and calculators, or designing complex prompts (Chen et al., 2024b). Notably, prompting 179 methods can be orthogonally combined with finetuning methods. For finetuning methods, previous 180 works (Yu et al., 2024; An et al., 2023) only model single-round correction during the training stage, and perform verification in a straightforward manner during the inference stage. 181

2.3 DIFFERENCE BETWEEN THE PROPOSED METHOD AND EXISTING SELF-CORRECTION TECHNIQUES

185 The existing self-correction research can be broadly classified into two categories: prompt-based methods (Gou et al., 2023a; Chen et al., 2023b; 2024b) and SFT-based methods (Yu et al., 2024; 187 An et al., 2023). Our approach falls under the SFT-based category and fundamentally differs from 188 prompt-based methods. Specifically, prompt-based methods leverage the intrinsic capabilities of 189 large language models (LLMs) for mathematical reasoning. In contrast, our method embeds the chain 190 of self-correction, a strong reasoning mechanism, as an inherent capability in the LLMs through 191 parameter fine-tuning. Furthermore, compared to existing SFT-based methods, our approach models 192 mathematical reasoning as a multi-round procedure during the training stage. In addition, these SFT-193 based works (Yu et al., 2024; An et al., 2023) perform the verification process in a straightforward manner, whereas our work conducts verification in a step-by-step manner by generating intermediate 194 verification steps. Similar to the essence of the CoT approach, our model first verifies whether the 195 generated code aligns with the problem description and then verifies whether the obtained result 196 is consistent with the problem description. This step-by-step verification strategy significantly 197 improves the verification process in self-correction. The excellent performance of our method on the challenging mathematical problem tasks, such as MATH and GSM8k, clearly demonstrates its 199 effectiveness compared to the existing SFT-based methods (Yu et al., 2024; An et al., 2023). 200

201

3 Method

202 203 204

3.1 OVERVIEW

205 We propose a Chain of Self-Correction (CoSC) to address mathematical questions in a self-correction 206 manner. Given a mathematical question denoted as q, the CoSC mechanism initiates the process 207 by generating a program p that comprises a Python function. The program p is then executed 208 using a Python interpreter to produce an output o. However, the program p may contain bugs or 209 incorrect intermediate code due to potential misinterpretations of the question q, leading to erroneous 210 outputs. Unlike existing methods that either base their reasoning on incorrect intermediate content 211 or halt reasoning upon encountering an error, the CoSC mechanism introduces a self-correction 212 rationale during the reasoning process. After the generation of program p and its output o, the 213 CoSC model analyzes them and generates a verification v. This verification includes suggestions for improvements to ensure the consistency of the program p and its output o with the mathematical 214 question q, respectively. Following this, the CoSC mechanism draws a conclusion c based on these 215 suggestions, which is used to either refine the program p or generate the final answer. This process is



Figure 2: The training of Chain of Self-Correction (CoSC) consists of two phases. The first phase, (a) 231 **CoSC Foundational Learning**, trains LLMs with seeding data generated from proprietary models, 232 equipping them with a baseline proficiency in the CoSC methodology. In particular, we prompt 233 GPT-4 with training questions from MATH (Hendrycks et al., 2021) and GSM8K (Cobbe et al., 2021) 234 datasets to generate mathematical reasoning trajectories that adhere to the CoSC protocol. The second 235 phase, (b) CoSC Self Enhancement, further adapts the seed model obtained from the previous phase 236 with self-generated trajectories. These trajectories are produced by the seed model trained in the 237 foundational phase, thereby enabling the generation of a substantial volume of data without the need 238 for additional GPT-4 intervention. In both phases, we only retain trajectories whose answers match 239 the ground-truth label. 240

repeated until a conclusive answer or a maximum limit of self-correction stages is reached, which can be summarized as $answer = povcpovc \cdots povc$. By employing this reasoning method, we can annotate public mathematical datasets and use the annotated multi-round, self-correction data to fine-tune LLMs.

3.2 TRAINING

To enable the CoSC mechanism, we propose a two-phase finetuning method. (1) CoSC foundational learning with seeding data generated from proprietary models. (2) CoSC self enhancement with self-generated data using the seed model obtained in the first training phase. A summary for training our CoSC is shown in Fig. 2.

251

253

245

246

3.2.1 COSC FOUNDATIONAL LEARNING

Existing mathematical reasoning datasets, such as chain-of-thoughts (CoT) (Wei et al., 2022) and 254 program-of-thoughts (PoT) (Chen et al., 2023a), primarily contain single-round annotations without 255 multi-round, self-correction solutions for the reasoning process. This makes it challenging to fine-tune 256 models to inherently possess self-correction capabilities. To address this issue, we utilize GPT-4 257 (\mathcal{G}) and a few-shot approach to annotate questions from publicly available mathematical datasets, 258 including MATH (Hendrycks et al., 2021) and GSM8K (Cobbe et al., 2021), to generate seeding data 259 with our Chain of Self-Correction (CoSC) mechanism. Table 1 compares our seeding data of CoSC 260 with recently proposed mathematical reasoning datasets. The algorithm of the generation with our 261 CoSC for each question from the training set is shown in Algorithm 1. Detailed instructions and 262 example cases of our CoSC can be found in Appendix A and Appendix B.

Given a mathematical question q, the few-shot prompt used to call GPT-4 is defined as \wp . The corresponding trajectory, denoted as τ , is generated through the following process. Firstly, we feed both the few-shot prompt \wp and the question q into GPT-4, which then generates a Python function code p to solve the question q. It can be represented as:

- 267 $\wp \oplus q \oplus \tau_{i-1} \to p_i, \tag{1}$
- where the symbol \oplus represents concatenation and τ_i indicates the trajectory in the *i* interaction round. After obtaining the code p_i , we execute it using a Python interpreter to acquire the runtime result

270 Table 1: The comparison to the previous mathematical reasoning datasets. Note that our CoSC 271 only employs questions from the original training set of MATH and GSM8K datasets. Particularly, 272 we provide the number of annotation data as #Annotation in each mathematical reasoning method. 273 Additionally, Correction signifies whether the self-correction is to be utilized.

Methods	#Annotation	Tool	Correction	LLM Used	Source
RFT (Yuan et al., 2023)	>100k	X	×	LLaMA-2	GSM8K
Open-Platypus (Lee et al., 20)	23) 25k	×	×	GPT-4	11 datasets with MATH
WizardMath (Luo et al., 2023) >96k	×	×	ChatGPT	MATH & GSM8K
Lila (Mishra et al., 2022)	134k	1	×	-	20 datasets with MATH & GSM8K
MathInstruct (Yue et al., 2023	3) 260k	1	×	GPT-4	14 datasets with MATH & GSM8K
ToRA (Gou et al., 2023b)	16k	\checkmark	×	GPT-4	MATH & GSM8K
CoSC (Ours)	37k	1	1	GPT-4	MATH & GSM8K

 o_i . Then, we enclose the result o within "```output" and feed it, along with the previous few-shot prompt \wp , question q, the previous trajectory τ_i , and generated code p_i to the GPT-4. This produces a verification v_i that analyzes the consistency between the question q and the program code p_i , as well as between the question q and the program output o_i . Additionally, it produces a conclusion c_i on whether to proceed with the next round of reasoning to perform adjustment, denoted as:

$$\oplus q \oplus \tau_{i-1} \oplus p_i \oplus o_i \to v_i \oplus c_i.$$
⁽²⁾

Depending on the conclusion c_i , it decides whether to terminate the reasoning process or proceed to 289 the next round. If there is no inconsistency found during the verification step, it generates a natural language conclusion c_i and encloses the answer to the problem with the special symbol "boxed{}". Otherwise, it should further revise the generated code p_i in a new reasoning round by updating the trajectory τ_i , denoted as: 293

$$\tau_i = \tau_{i-1} \oplus p_i \oplus o_i \oplus v_i \oplus c_i. \tag{3}$$

295 Based on the above steps, CoSC iteratively generates PoT code, program output, verification, and conclusion. At final, the reasoning process is stopped by checking whether the generated result 296 contains an answer enclosed in "boxed{}". 297

298 To ensure that the self-correction can concisely and clearly analyze the code and its runtime results, 299 we design the self-correction in a step-by-step format as follows: 300

- The first step involves verifying whether the generated program code p is consistent with the question q, such as checking the variables and their relationships.
- The second step involves verifying whether the code runtime results o meet the requirements of the question q, such as checking the reasonableness of numerical values.

Finally, for summary, CoSC will generate a conclusion c to determine whether we should start the 305 next round of reasoning or provide the final answer. 306

307 With CoSC, we use GPT-4 to annotate the MATH and GSM8K training datasets. To explore diverse data, we apply nucleus sampling (Holtzman et al., 2019) during GPT-4 annotation. Each question is 308 sampled in 3 times. However, for some complex questions, we are unable to obtain even a single 309 correct solution. For these problems, we apply 10 more samplings and retained up to 4 correct data. 310 Finally, we filtered out incorrect answers and constructed 37k pieces of data using GPT-4. 311

312 Based on the 37k data constructed by GPT-4 as the seeding data, we apply CoSC foundational 313 learning to train models. For a given question q, the response trajectory τ generated by GPT-4, 314 denoted as $\tau = povcpovc \cdots povc$, the minimized negative log-likelihood loss used for training the model can be represented as: 315

$$\arg\min_{\theta} \sum_{q,\tau} \sum_{i=1}^{n_q-1} -\log \mathbb{P}_{\theta}(p_{i+1}o_{i+1}v_{i+1}c_{i+1}|q, p_i \cdots o_i v_i c_i),$$
(4)

318 319 320

316 317

281 282

283

284

285

286

287 288

290

291

292

301

302

303

where n_q is the iteration rounds of the question q in our CoSC for the mathematical reasoning.

321 3.2.2 COSC SELF ENHANCEMENT 322

After the completion of the CoSC foundational learning, the seed model gains the ability to self-323 correct during inference and perform multi-round reasoning. Subsequently, we employ the seed

Require: question q , model \mathcal{G} , prompt \wp , external too	ls \mathcal{E} , stop condition $Stop(\cdot)$, iteration rounds n_q
1: $\tau_0 \leftarrow$ ""	Trajectory Initialization
2: for $i \leftarrow 1$ to n_q do	
3: $p_i \sim \mathbb{P}_{\mathcal{G}}(\cdot \wp \oplus q \oplus \tau_{i-1})$	▷ Program Generation (Eq. 1)
4: $o_i \leftarrow \mathcal{E}(p_i)$	▷ Tool Execution
5: $v_i \oplus c_i \sim \mathbb{P}_{\mathcal{G}}(\cdot \wp \oplus q \oplus \tau_{i-1} \oplus p_i \oplus o_i)$	▷ Verification and Conclusion Generation (Eq. 2)
$6: \tau_i \leftarrow \tau_{i-1} \oplus p_i \oplus o_i \oplus v_i \oplus c_i$	▷ Trajectory Adjustment (Eq. 3)
7: if $Stop(c_i)$ then	▷ Stopping Criteria
8: return τ_i	
9: end if	
10: end for	
11: return τ_n	

model, after CoSC foundational learning, to apply dense solution sampling and dense question sampling, which enables the generation of more self-generated data with self-correction mechanisms.

Dense solution sampling. For mathematical questions, there are usually multiple solutions, but using
 GPT-4 for annotating them is expensive. Therefore, we use the seed model after CoSC foundational
 learning to resample the questions in the datasets multiple times, which further improves the model
 generalization ability in answering questions. Specifically, we use the CodeLLaMA-34B model
 which has the best performance after CoSC foundational learning to perform dense nucleus-sampling
 on 16k training data questions. Each question in the MATH and GSM8K training datasets is sampled
 times. We filter out the correct answers based on whether they match the ground-truth.

Dense question sampling. Similarly, there are multiple ways to ask a mathematical question. To
 improve the model generalization ability of question understanding, we use the data from Meta Math (Yu et al., 2023) to generalize the questions. Each question in MetaMATH is rewritten by simply
 rephrasing it, self-verifying the conditions in the question, adding if-then questions to reverse the
 question conditions, *etc.* Specifically, we use the CodeLLaMA-34B model with CoSC foundational
 learning to perform nucleus-sampling on the questions, sampling each question once, and filtering
 out data points based on whether they match the ground-truth.

In the end, we obtain a total of 339k data points, including 37k seeding data generated from GPT-4 and 302k generated from the CodeLLaMA-34B model with CoSC foundational learning. Then, we adopt them to train models from scratch to obtain our final CoSC model.

359 360

22/

338

3.3 IMPLEMENTATION DETAILS

361 By using 339k data points, including 37k seeding data and 302k self-generated data, we fine-tune 362 the base models of LLaMA-2 (Touvron et al., 2023) and CodeLLaMA (Roziere et al., 2023) to 363 obtain our CoSC and CoSC-Code, respectively. They have different parameter sizes, such as 7B, 364 13B, and 34B. All models use full-scale fine-tuning. We use the AdamW optimizer with a learning rate of 2e-5 for all models, with a batch size set to 128, training for 1 epoch. To enable training, 366 we use DeepSpeed ZeRO stage 3 (Rajbhandari et al., 2021) and Flash-Attention 2 (Dao, 2023) to optimize the model's memory usage. During inference, we set a maximum of 3 calls to the Python 367 interpreter and a maximum token length of 2048. The GPT-4 version for CoSC data generation is 368 gpt-4-0613. Our experiments train the models in 7B size with 8 NVIDIA A800 80GB GPUs and 369 train the models in 13B and 34B with 16 NVIDIA A800 80GB GPUs. 370

371 372

373

4 EXPERIMENTS

374 4.1 EVALUATION SETUP

375
376 Datasets. We evaluated models on the most widely used mathematical problem tasks, MATH (Hendrycks et al., 2021) and GSM8K (Cobbe et al., 2021). The MATH dataset encompasses a total of 12,500 problems, partitioned into 7,500 for training and 5,000 for testing. The GSM8K dataset

Table 2: Accuracy results (%) on MATH and GSM8K datasets. Vanilla models are tested with
 CoT. ZS indicates the zero-shot inference without demonstrations. PAL refers to the Program-Aided
 Language model prompting (Gao et al., 2023). The best results in each section are in **bold** and
 the second-best results are <u>underlined</u>. Our CoSC can outperform open-sourced LLMs, proprietary
 non-multi-modal LLMs and some of multi-modal LLMs, on the more challenging MATH dataset.

Model	Base	Size	ZS	MATH	GSM8K	AVG
	Prop	rietary	Mod	els		
GPT-40	-	-	X	76.6 (0-shot CoT)	96.1 (8-shot CoT)	86.4
GPT-4V	-	-	X	52.9 (4-shot)	92.0 (5-shot CoT)	72.5
GPT-4 (PAL)	-	-	X	51.8 (PAL)	94.2 (PAL)	73.0
GPT-4	-	-	X	42.5 (Cot)	92.0 (5-shot CoT)	67.3
ChatGPT (PAL)	-	-	X	38.7 (PAL)	78.6 (pal)	58.7
ChatGPT	-	-	X	35.5 (Сот)	80.8 (5-shot CoT)	58.2
Claude-3.5 Sonnet	-	-	X	71.1 (0-shot CoT)	96.4 (0-shot CoT)	83.8
Claude-3 Opus	-	-	X	60.1 (0-shot CoT)	95.0 (0-shot CoT)	77.6
Claude-2	-	-	X	32.5 (0-shot CoT)	85.2 (0-shot CoT)	58.9
Gemini-1.5 Pro	-	-	X	67.7 (4-shot Minerva)	90.8 (11-shot)	79.3
Gemini-1.5 Flash	-	-	X	54.9 (4-shot Minerva)	86.2 (11-shot)	70.6
Gemini-1.0 Ultra	-	-	X	53.2 (4-shot Minerva)	88.9 (11-shot)	71.1
Gemini-1.0 Pro	-	- 5 40D	X	32.6 (4-shot Minerva)	77.9 (11-shot)	55.3
PaLM-2	-	540B	×	34.3 (4-shot CoT)	80. / (8-shot CoT)	57.5
	Open-	Source	e Mo	dels		
LLaMA-2	LLaMA-2	7B	X	4.1 (CoT)	13.3 (Сот)	8.7
LLaMA-2 SFT	LLaMA-2	7B	1	7.2	41.3	24.3
LLaMA-2 RFT	LLaMA-2	7B	1	-	51.2	-
CodeLLaMA (PAL)	CodeLLaMA	7B	X	16.6 (PAL)	34.0 (PAL)	25.3
Platypus-2 (Lee et al., 2023)	LLaMA-2	7B	X	5.4 (CoT)	14.4 (CoT)	9.9
WizardMath (Luo et al., 2023)	LLaMA-2	7B	1	10.7	54.9	32.8
MetaMath (Yu et al., 2023)	LLaMA-2	7B	1	19.8	66.5	43.2
ToRA (Gou et al., 2023b)	LLaMA-2	7B	1	40.1	68.8	54.5
CoSC (Ours)	LLaMA-2	7B	1	42.7	70.5	56.6
ToRA-Code (Gou et al., 2023b)	CodeLLaMA	7B	1	<u>44.6</u>	<u>72.6</u>	<u>58.6</u>
CoSC-Code (Ours)	CodeLLaMA	7B	✓	47.6	74.7	61.2 (+2.6)
LLaMA-2	LLaMA-2	13B	X	6.3 (Cot)	24.3 (Сот)	15.3
LLaMA-2 SFT	LLaMA-2	13B	1	9.2	51.1	30.2
LLaMA-2 RFT	LLaMA-2	13B	I	-	55.3	-
CodeLLaMA (PAL)	CodeLLaMA	13B	×	19.9 (PAL)	39.9 (PAL)	29.9
Platypus-2 (Lee et al., 2023)	LLaMA-2	13B	X	7.1 (CoT)	23.7 (Cot)	15.4
WizardMath (Luo et al., 2023)	LLaMA-2	13B		14.0	63.9	39.0
MetaMath (Yu et al., 2023)	LLaMA-2	13B	 Image: A start of the start of	22.4	72.3	47.4
ToRA (Gou et al., 2023b)	LLaMA-2	13B	 Image: A start of the start of	43.0	72.7	57.9
CoSC (Ours)	LLaMA-2	13B	 Image: A start of the start of	45.3	73.9	59.6
ToRA-Code (Gou et al., 2023b)	CodeLLaMA	13B	 Image: A start of the start of	$\frac{48.1}{59.2}$	75.8	$\frac{62.0}{2}$
CoSC-Code (Ours)	CodeLLaMA	13B	✓	50.3	77.3	63.8 (+1.8)
LLaMA-1 RFT	LLaMA-1	34B	 I 	-	57.9	-
CodeLLaMA (PAL)	CodeLLaMA	34B	X	23.9 (PAL)	53.3 (PAL)	38.6
TokA-Code (Gou et al., 2023b)	CodeLLaMA	54B	 Image: A start of the start of	<u>50.8</u>	80.7	65.8
CoSC-Code (Ours)	CodeLLaMA	34B	✓	53.5	82.3	67.9 (+2.1)

⁴²²

contains 8,500 problems, with 7,500 for training and 1,000 for testing. These datasets collectively
 encompass a broad range of mathematical questions, from basic arithmetic to competition level. More
 details of datasets are provided in Appendix C.1.

Baselines. For proprietary models, we present results from an array of SoTA LLMs, such as OpenAI's ChatGPT (gpt-3.5-turbo), OpenAI's GPT-4, OpenAI's GPT-4V, OpenAI's GPT-4o, Google's PaLM-2, Google's Gemini-1.0 Pro, Google's Gemini-1.0 Ultra, Google's Gemini-1.5 Flash, Google's Gemini-1.5 Pro, Anthropic's Claude-2, Anthropic's Claude-3 Opus, and Anthropic's Claude-3.5
Sonnet. By default, we report CoT prompting results, and include PAL (Program-Aided Language model) prompting (Gao et al., 2023) results for selected models. For open-source models, base models comprise LLaMA-2 (Touvron et al., 2023) and CodeLLaMA (Roziere et al., 2023) with

CoT and PAL prompting (Gao et al., 2023). Supervised Fine-Tuning (SFT) employs CoT rationales
from the original MATH and GSM8K dataset (15k samples) for fine-tuning. Rejection sampling
Fine-Tuning (RFT) (Yuan et al., 2023) leverages multiple models to generate diverse reasoning paths
for fine-tuning. WizardMath (Luo et al., 2023) augments data using ChatGPT, and conducts SFT
and RLHF. Platypus-2 (Lee et al., 2023) is fine-tuned with Open-Platypus reasoning datasets. ToRA
(Gou et al., 2023b) uses GPT-4 to generate tool-used trajectories and finetunes on it. The prompting
methods for each evaluation are shown in Appendix C.2.

439 Metric. We report accuracies of predicted answers. For numerical values, we perform rounding,
 440 while for expressions, we employ the Python library *sympy* for parsing.

441 442

442 4.2 MAIN RESULTS

In our experiments, we observe several findings that demonstrate the effectiveness of our proposed
CoSC model on mathematical datasets. We show the experimental results in Table 2. First and
foremost, our proposed CoSC consistently outperforms previous state-of-the-art open-source LLMs
across all scales. Specifically, our CoSC-Code can achieve an average improvement of 2.6%, 1.8%,
and 2.1%, on 7B, 13B, and 34B size, respectively.

Moreover, to further emphasize the superiority of our CoSC, we conduct a comprehensive comparison
 against multiple proprietary models. The results reveal that our CoSC-Code-34B can outperform
 all the advanced proprietary non-multi-modal LLMs, such as ChatGPT and GPT-4, as well as most
 advanced proprietary multi-modal LLMs, such as GPT-4V, Gemini-1.0 Pro, and Gemini-1.0 Ultra on
 MATH dataset. It is important to note that, unlike these proprietary models, our CoSC performs the
 inference in a zero-shot manner without demonstrations.

455 456

4.3 Ablation Study

In our ablation study, we conduct experiments on MATH dataset with the largest number of 5,000 test
samples with broad spectrum of subjects and difficulty levels. For the consideration of computational
efficiency, we choose the smaller models of CodeLLaMA with size 7B and 13B as the base models.

461 4.3.1 EFFECT OF EACH TRAINING PHASE

463 Our CoSC mechanism comprises two training phases: (1) CoSC foundational learning and (2) CoSC
 464 self-enhancement. In order to assess the individual contributions of each training phase, we evaluate
 465 multiple combinations of the proposed training phases, the results of which are shown in Table 3.

The experiments conducted on CodeLLaMA serve as the baseline for our study. Initially, we employ GPT-4 to generate 37k seeding data for CoSC foundational learning, which equips LLMs with initial self-correction capabilities. As a result, we observe a notable improvement in performance on the MATH dataset, with accuracy increasing from 16.6% to 42.3% for the 7B size and from 19.9% to 470 47.0% for the 13B size.

To further enhance the self-correction performance, we utilize the seed model obtained in the first phase to generate additional data without relying on GPT-4. This approach leads to further improvements in performance on the MATH dataset from 42.3% to 47.6% and 47.0% to 50.3%, for 7B and 13B respectively.

It is exciting to see that the originally weak LLMs, such as the CodeLLaMA 7B base model and the 13B base model, can significantly improve in mathematical reasoning by using the proposed CoSC method, as shown in Table 3. This demonstrates the effectiveness of the CoSC method, which is able to embed self-correction as an inherent capability in LLMs, leading to significant improvements in mathematical reasoning performance.

481

4.3.2 EFFECT OF MULTI-ROUND REASONING IN THE PROPOSED COSC MECHANISM

Our CoSC mechanism integrates a series of self-correction stages to progressively verify and refine output of LLMs. In particular, the conclusion step in our CoSC relies on clues from the verification step to determine whether to proceed to the next round of reasoning or directly provide the final answer. Such iterative multi-round mechanism enables LLMs to self-correct their outputs and

self enhancement. The results show	w the superiority and necessit	y of both training	phases.	
Chain of Self-Corre Foundational Learning	ection (CoSC) Self Enhancement	7B	13B	
×	×	16.6	19.9	
1	X	42.3	47.0	

47.6

50.3

Table 3: Accuracy results (%) on the MATH dataset for the CoSC foundational learning and CoSC
 self enhancement. The results show the superiority and necessity of both training phases.

Table 4: The distribution of reasoning rounds on MATH test dataset across three LLMs in 7B and 13B size. Our CoSC-Code can generate more rounds compared to other models in mathematical reasoning. #Round indicates the number of reasoning rounds during inference.

		7B			13B	
Models	#Round=1	#Round=2	#Round=3	#Round=1	#Round=2	#Round=3
CodeLLaMA	100%	0	0	100%	0	0
ToRA-Code	99.9%	0.1%	0	100%	0	0
CoSC-Code (Ours)	78.3%	12.7%	9.0%	79.3%	13.1%	7.6%

Table 5: Accuracy results (%) for the test samples on the MATH dataset of our CoSC-Code with single-round reasoning and multi-round reasoning during the inference stage.

Models	7B	13B
CoSC-Code with single-round reasoning	40.2	42.4
CoSC-Code with multi-round reasoning (Ours)	47.6	50.3

improve accuracy. This experiment quantitatively investigates effect of this multi-round mechanism.
 As shown in Table 4, our CoSC can effectively generate more rounds of mathematical reasoning during inference, confirming efficacy of our CoSC mechanism in enhancing the reasoning process.

In addition, we conduct a comparison between single-round reasoning and multi-round reasoning using our CoSC mechanism on the test samples from the MATH dataset. The single-round results are obtained by extracting the answer solely from the output of the first round, without any self-correction enabled. To illustrate the impact of multi-round reasoning, we present the comparative results in Table 5. The results clearly demonstrate the effectiveness of multi-round reasoning in rectifying errors and improving overall accuracy. With the 7B model, we observe an improvement in accuracy from 40.2% to 47.6%, while with the 13B model, accuracy increases from 42.4% to 50.3%. These findings highlight the significant benefits of employing multi-round reasoning within our CoSC mechanism. More ablation studies are shown in Appendix D.

5 CONCLUSION

In conclusion, our Chain of Self-Correction (CoSC) mechanism equips Large Language Models (LLMs) with the ability to autonomously validate and refine their outputs. This mechanism facilitates a sequence of self-correction stages that progressively refine the reasoning process, leading to enhanced accuracy in mathematical reasoning. Our two-phase finetuning approach, involving initial training with proprietary model-generated seeding data and subsequent training with self-generated data, effectively enables the implementation of CoSC at a low cost. Through extensive experiments, we have demonstrated the remarkable performance improvement that CoSC brings to traditional mathematical datasets when compared to existing open-source LLMs, even advanced close-sourced non-multi-modal LLMs, such as ChatGPT, GPT-4 and multi-modal LLMs like GPT-4V, Gemini-1.0 Pro, and Gemini-1.0 Ultra. We believe that our CoSC mechanism can provide valuable insights for future research and contribute to the ongoing advancement of LLMs across various domains.

Ethics statement. We experiment on two mathematical datasets, including GSM8K and MATH, both
 of which use MIT License code. The prompts used in these experiments are listed in Appendix A,
 and we want to emphasize that none of the prompts contain any words that discriminate against any
 individual or group. Furthermore, prompts would not negatively impact anyone's safety in this work.

Reproducibility statement. We provide part of the codes and some seeding data. We will provide
 the remaining codes and data of CoSC for reproducibility upon the acceptance of the paper.

540 REFERENCES

- Shengnan An, Zexiong Ma, Zeqi Lin, Nanning Zheng, Jian-Guang Lou, and Weizhu Chen. Learning
 from mistakes makes llm better reasoner. arXiv preprint arXiv:2310.20689, 2023.
- Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos,
 Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. <u>arXiv</u>
 preprint arXiv:2305.10403, 2023.
- 547 548 Anthropic. https://www.anthropic.com/news/claude-3-family/. 2023.
- Zhangir Azerbayev, Hailey Schoelkopf, Keiran Paster, Marco Dos Santos, Stephen McAleer, Albert Q
 Jiang, Jia Deng, Stella Biderman, and Sean Welleck. Llemma: An open language model for
 mathematics. arXiv preprint arXiv:2310.10631, 2023.
- Changyu Chen, Xiting Wang, Ting-En Lin, Ang Lv, Yuchuan Wu, Xin Gao, Ji-Rong Wen, Rui Yan, and Yongbin Li. Masked thought: Simply masking partial reasoning steps can improve mathematical reasoning learning of language models. <u>arXiv preprint arXiv:2403.02178</u>, 2024a.
- Sijia Chen, Baochun Li, and Di Niu. Boosting of thoughts: Trial-and-error problem solving with
 large language models. In ICLR, 2024b.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. In <u>TMLR</u>, 2023a.
- Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching large language models to self-debug. arXiv preprint arXiv:2304.05128, 2023b.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. Journal of Machine Learning Research, 24(240):1–113, 2023.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. <u>arXiv preprint arXiv:2110.14168</u>, 2021.
- Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. <u>arXiv</u> preprint arXiv:2307.08691, 2023.
- Jiayi Fu, Lei Lin, Xiaoyang Gao, Pengli Liu, Zhengzong Chen, Zhirui Yang, Shengnan Zhang,
 Xue Zheng, Yan Li, Yuliang Liu, et al. Kwaiyiimath: Technical report. <u>arXiv preprint</u> <u>arXiv:2310.07488</u>, 2023.
- Tommaso Furlanello, Zachary Lipton, Michael Tschannen, Laurent Itti, and Anima Anandkumar.
 Born again neural networks. In <u>ICML</u>, 2018.
- Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. Pal: Program-aided language models. In <u>ICML</u>, 2023.
- Silin Gao, Jane Dwivedi-Yu, Ping Yu, Xiaoqing Ellen Tan, Ramakanth Pasunuru, Olga Golovneva, Koustuv Sinha, Asli Celikyilmaz, Antoine Bosselut, and Tianlu Wang. Efficient tool use with chain-of-abstraction reasoning. <u>arXiv preprint arXiv:2401.17464</u>, 2024.
- Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Nan Duan, and Weizhu Chen.
 Critic: Large language models can self-correct with tool-interactive critiquing. <u>arXiv preprint</u> arXiv:2305.11738, 2023a.
- Zhibin Gou, Zhihong Shao, Yeyun Gong, Yujiu Yang, Minlie Huang, Nan Duan, Weizhu Chen, et al. Tora: A tool-integrated reasoning agent for mathematical problem solving. <u>arXiv preprint</u> arXiv:2309.17452, 2023b.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. <u>arXiv</u> preprint arXiv:2103.03874, 2021.

- 594 Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text 595 degeneration. arXiv preprint arXiv:1904.09751, 2019. 596
- Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alexander Ratner, 597 Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. Distilling step-by-step! outperforming larger 598 language models with less training data and smaller model sizes. In ACL, 2023.
- 600 Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, 601 and Denny Zhou. Large language models cannot self-correct reasoning yet. In ICLR, 2024a. 602
- 603 Yiming Huang, Xiao Liu, Yeyun Gong, Zhibin Gou, Yelong Shen, Nan Duan, and Weizhu Chen. Key-point-driven data synthesis with its enhancement on mathematical reasoning. arXiv preprint 604 arXiv:2403.02333, 2024b. 605
 - Ariel N Lee, Cole J Hunter, and Nataniel Ruiz. Platypus: Quick, cheap, and powerful refinement of llms. arXiv preprint arXiv:2308.07317, 2023.
- 609 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan 610 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. arXiv preprint 611 arXiv:2305.20050, 2023.
- Bingbin Liu, Sebastien Bubeck, Ronen Eldan, Janardhan Kulkarni, Yuanzhi Li, Anh Nguyen, Rachel 613 Ward, and Yi Zhang. Tinygsm: achieving> 80% on gsm8k with small language models. arXiv 614 preprint arXiv:2312.09241, 2023a. 615
- 616 Yixin Liu, Avi Singh, C Daniel Freeman, John D Co-Reyes, and Peter J Liu. Improving large 617 language model fine-tuning for solving math problems. arXiv preprint arXiv:2310.10047, 2023b.
- Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, 619 Qingwei Lin, Shifeng Chen, and Dongmei Zhang. Wizardmath: Empowering mathematical 620 reasoning for large language models via reinforced evol-instruct. arXiv preprint arXiv:2308.09583, 621 2023. 622
- 623 Swaroop Mishra, Matthew Finlayson, Pan Lu, Leonard Tang, Sean Welleck, Chitta Baral, Tanmay 624 Rajpurohit, Oyvind Tafjord, Ashish Sabharwal, Peter Clark, et al. Lila: A unified benchmark for 625 mathematical reasoning. arXiv preprint arXiv:2210.17517, 2022.
- Arindam Mitra, Hamed Khanpour, Corby Rosset, and Ahmed Awadallah. Orca-math: Unlocking the 627 potential of slms in grade school math. arXiv preprint arXiv:2402.14830, 2024. 628
- 629 OpenAI. Gpt-4 technical report. 2023a. 630

606

607

608

612

618

626

641

- 631 OpenAI. https://openai.com/index/gpt-4v-system-card/. 2023b.
- 632 OpenAI. https://openai.com/index/hello-gpt-4o/. 2024. 633

634 Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, 635 Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The refinedweb 636 dataset for falcon llm: outperforming curated corpora with web data, and web data only. arXiv 637 preprint arXiv:2306.01116, 2023. 638

- Samyam Rajbhandari, Olatunji Ruwase, Jeff Rasley, Shaden Smith, and Yuxiong He. Zero-infinity: 639 Breaking the gpu memory wall for extreme scale deep learning. In Proceedings of the international 640 conference for high performance computing, networking, storage and analysis, pp. 1–14, 2021.
- 642 Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi 643 Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. Code llama: Open foundation models for code. 644 arXiv preprint arXiv:2308.12950, 2023. 645
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, YK Li, Y Wu, 646 and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language 647 models. arXiv preprint arXiv:2402.03300, 2024.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu 649 Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable 650 multimodal models. arXiv preprint arXiv:2312.11805, 2023. 651 Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett 652 Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal 653 understanding across millions of tokens of context. arXiv preprint arXiv:2403.05530, 2024. 654 655 Shubham Toshniwal, Ivan Moshkov, Sean Narenthiran, Daria Gitman, Fei Jia, and Igor Git-656 man. Openmathinstruct-1: A 1.8 million math instruction tuning dataset. arXiv preprint 657 arXiv:2402.10176, 2024. 658 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée 659 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and 660 efficient foundation language models. arXiv preprint arXiv:2302.13971, 2023. 661 662 Peiyi Wang, Lei Li, Zhihong Shao, RX Xu, Damai Dai, Yifei Li, Deli Chen, Y Wu, and Zhifang Sui. Math-shepherd: A label-free step-by-step verifier for llms in mathematical reasoning. arXiv 663 preprint arXiv:2312.08935, 2023a. 664 665 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-666 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. 667 In ICLR, 2023b. 668 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny 669 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. In NeurIPS, 670 2022. 671 672 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, 673 Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. arXiv preprint 674 arXiv:2407.10671, 2024. 675 Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo 676 Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for 677 large language models. arXiv preprint arXiv:2309.12284, 2023. 678 679 Xiao Yu, Baolin Peng, Michel Galley, Jianfeng Gao, and Zhou Yu. Teaching language models to self-improve through interactive demonstrations. In NAACL, 2024. 680 681 Zheng Yuan, Hongyi Yuan, Chengpeng Li, Guanting Dong, Chuanqi Tan, and Chang Zhou. Scaling 682 relationship on learning mathematical reasoning with large language models. arXiv preprint 683 arXiv:2308.01825, 2023. 684 685 Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. Mammoth: Building math generalist models through hybrid instruction tuning. arXiv preprint 686 arXiv:2309.05653, 2023. 687 688 Boning Zhang, Chengxi Li, and Kai Fan. Mario eval: Evaluate your math llm with your math llm-a 689 mathematical dataset evaluation toolkit. arXiv preprint arXiv:2404.13925, 2024. 690 691 692 693 694 696 697 699 700

A PROMPT FOR OUR COSC

```
704
       We present specific instructions and example few-shot prompts of our CoSC for querying GPT-4 to
705
       generate the seeding data.
706
707
       A.1 PROMPT FOR OUR COSC ON MATH
708
709
      The prompt for CoSC on MATH is as follows:
710
       Integrate step-by-step reasoning and Python code to solve math problems
711
          using the following guidelines:
712
713
        - Analyze the question and write functions to solve the problem; the
714
            function should not take any arguments.
715
        - Present the final results in LaTeX using a "\boxed{}" without any
716
           units.
717
718
        - Utilize the "pi" symbol and "Rational" from Sympy for $\pi$ and
            fractions, and simplify all fractions and square roots without
719
            converting them to decimal values.
720
721
        - In "verification", verify the consistency between the "Question" and
722
           the "python" code, as well as between the "Question" and "output" in
             "Solution" step-by-step.
723
724
        - According to the "verification", rewrite the code if neccessary.
725
726
       Here are some examples you may refer to:
727
728
       ___
729
       Question: Find the matrix that corresponds to projecting onto the vector
730
           $\begin{pmatrix} 2 \\ -3 \end{pmatrix}.$
731
732
       Solution:
733
       '''python
734
       from sympy import Matrix, simplify
735
736
       def projection_matrix():
737
           v = Matrix([2, -3])
738
739
           # Calculate the projection matrix
           P = (v * v.T) / (v.T * v) [0]
740
741
           # Simplify the elements
742
           P = P.applyfunc(simplify)
743
744
           return P
745
       answer = projection_matrix()
746
       print (answer)
747
748
749
       '''output
      Matrix([[4/13, -6/13], [-6/13, 9/13]])
750
       * * *
751
752
       ```verification
753
 Step 1. Verify the consistency between the question and the code. The "
754
 python" code defines a function binary_sum_diff() that calculates the
755
 projection matrix for the given vector. The vector is defined as 'v
 = Matrix([2, -3])', and then use standard formula '(v \star v.T) / (v.T \star
```

```
756
 v) [0] ' for calculating the projection matrix onto a vector, finally
757
 use 'applyfunc(simplify)' to simplify the elements of the matrix.
758
 Step 2: Verify the consistency between the question and the output. The
759
 answer matrix of 2 dimension vector should be a 2x2 matrix, the "
 output" of "python" code is `[[4/13, -6/13], [-6/13, 9/13]]`, which
760
 is reasonable.
761
 . . .
762
763
       ```conclusion
764
       Therefore, the "python" code and "output" are consistent with "Question".
           The matrix is \boxed{\begin{pmatrix} \frac{4}{13} & -\frac{6}{13} \
765
           -\frac{6}{13} & \frac{9}{13} \end{pmatrix}}.$
766
767
768
       ___
769
770
      Question: Compute $1011_2 + 101_2 - 1100_2 + 1101_2$. Express your
          answer in base 2.
771
772
      Solution:
773
774
       ```python
775
 def binary_sum_diff():
 # Convert binary to decimal
776
 num1 = int("1011", 2)
777
 num2 = int("101", 2)
778
 num3 = int("1100", 2)
779
 num4 = int("1101", 2)
780
 # Perform the operations
781
 result = num1 + num2 - num3 + num4
782
783
 # Convert the result back to binary
784
 result_binary = format(result, "b")
785
 return result_binary
786
787
 answer = binary_sum_diff()
788
 print(answer)
789
 • • •
790
791
 '''output
792
 10001
 ...
793
794
       ```verification
795
       Step 1. Verify the consistency between the question and the code. The "
796
          python" code defines a function binary_sum_diff() that convert the
797
          binary numbers to decimal, performing the operations, and then
          converting the result back to binary. It assigns the decimal form of
798
          values $1011, 101, 1100, 1101$ to $num1, num2, num3, num4$
799
          respectively, and then perform the operation $num1 + num2 - num3 +
800
          num4, finally convert the result back to binary.
801
       Step 2: Verify the consistency between the question and the output. The "
802
          output" of "python" code is $10001$, which is reasonable.
803
804
       ```conclusion
805
 Therefore, the "python" code and "output" are consistent with "Question".
806
 The answer expressed in base 2 is \boxed{10001_2}.
 . . .
807
808
809

```

```
810
 Question: Convert the point (0, -3 \setminus grt \{3\}, 3) in rectangular
811
 coordinates to spherical coordinates. Enter your answer in the form
812
 813
 $0 \le \phi \le \pi.$
814
 Solution:
815
816
 '''python
817
 from sympy import sqrt, atan2, acos, pi
818
819
 def rectangular_to_spherical():
 x, y, z = 0, -3*sqrt(3), 3
820
 rho = sqrt(x * *2 + y * *2 + z * *2)
821
 theta = atan2(y, x)
822
 phi = acos(z/rho)
823
 return rho, theta, phi
824
 answer = rectangular_to_spherical()
825
 print(answer)
826
 • • •
827
828
       ```output
       (6, -pi/2, pi/3)
829
830
831
       ```verification
832
 Step 1. Verify the consistency between the question and the code. The "
833
 python" code defines a function rectangular_to_spherical() that
834
 convert rectangular coordinates to spherical coordinates. It assigns
 the values 0, -3 \sqrt{3}, 3 to x, y, z respectively, and then
835
 use the following formulas \r = \ x^2 + x^2 + z^2, \t =
836
 \operatorname{\operatorname{V}} = \operatorname{\operatorname{V}} \{x\} 3, \operatorname{D} = \operatorname{\operatorname{V}} \{z\} to
837
 calculates $rho, theta, phi$.
838
 Step 2: Verify the consistency between the question and the output. The
839
 required form of the "output" is (\n, \heat, \bar), where \n, \heat, \bar)
 > 0$, $0 \le \theta < 2\pi$, and $0 \le \phi \le \pi$. The "output"
840
 of "python" code is (6, -pi/2, pi/3), which is not reasonable.
841
 • • •
842
843
       ```conclusion
844
       Therefore, the "python" code and "output" are not consistent with "
845
           Question". Let's rewrite the "python" code based on the "verification
           ":
846
       • • •
847
848
       '''python
849
       from sympy import sqrt, atan2, acos, pi
850
851
       def rectangular_to_spherical():
           x, y, z = 0, -3*sqrt(3), 3
852
           rho = sqrt(x**2 + y**2 + z**2)
853
           theta = atan2(y, x)
854
           phi = acos(z/rho)
855
           theta = (theta + 2 * pi) % (2 * pi)
856
           phi = (phi + pi) % pi
           return rho, theta, phi
857
858
       answer = rectangular_to_spherical()
859
       print(answer)
860
       • • •
861
862
       '''output
863
       (6, 3*pi/2, pi/3)
```

```
864
865
       ```verification
866
 Step 1. Verify the consistency between the question and the code. The "
867
 python" code was rewritten, 'theta = (theta + 2 \star pi) % (2 \star pi) ' and
 'phi = (phi + pi) % pi' were added to meet requirement that $0 \le \
868
 theta < 2 \neq 0, and 0 \leq \phi < \phi.
869
 Step 2: Verify the consistency between the question and the output. The
870
 required form of the "output" is (\no, \heat, \bar), where \no.
871
 > 0$, $0 \le \theta < 2\pi$, and $0 \le \phi \le \pi$. The "output"
872
 of "python" code is (6, 3*pi/2, pi/3), which is reasonable.
873
874
       ```conclusion
875
       Therefore, the "python" code and "output" are consistent with "Question".
876
           Hence, the spherical coordinates for the point (0, -3 \setminus 3)
877
          are $\boxed{(6, \frac{3\pi}{2}, \frac{\pi}{3})}$.
       . . .
878
879
       ____
880
881
       Question: Solve the inequality
       \left[ \frac{x^2 - 25}{x + 5} < 0. \right]
882
883
       Solution:
884
885
       '''python
886
       from sympy import symbols, simplify
887
       from sympy.solvers.inequalities import solve_univariate_inequality
       from sympy.core.relational import LessThan
888
889
       def solve_inequality():
890
           x = symbols('x')
891
           expression = (x * * 2 - 25) / (x + 5)
           inequality = LessThan(expression, 0)
892
           solution = solve_univariate_inequality(inequality, x, relational=
893
               False)
894
           simplified_solution = simplify(solution)
895
896
           return simplified_solution
897
       answer = solve_inequality()
898
       print(answer)
899
900
901
       ''output
902
       Union(Interval.open(-oo, -5), Interval.Lopen(-5, 5))
903
904
       ```verification
905
 Step 1. Verify the consistency between the question and the code. The "
906
 python" code defines a function solve_inequality() that that solves
 the inequality \left[\frac{x^2 - 25}{x + 5} < 0.\right]. It defines a symbolic
907
 variable 'x' and the expression is defined as \left(\frac{x^2 - 25}{x + 2}\right)
908
 5}\]', then use function 'LessThan' to define the inequality and
909
 function 'solve_univariate_inequality' to solve the inequality for 'x
910
 ', finally use 'simplify' to simplify the result.
911
 Step 2: Verify the consistency between the question and the output. The
 required answer 'x' should meet the condition: x + 5 != 0, the "
912
 output" of "python" code is (-\infty, -5) and (-5, 5), which is
913
 reasonable.
914
 • • •
915
916
       ```conclusion
       Therefore, the "python" code and "output" are consistent with "Question".
917
           The answer is \boxed{(-\infty,-5)\cup(-5,5)}
```

```
918
       • • •
919
920
       ___
921
       Question: In the figure, triangles $ABC$ and $BCD$ are equilateral
922
          triangles. What is the value of $AD \div BC$ when expressed in
923
          simplest radical form?
924
925
       \begin{asy}
926
       draw((0,0)--(5,8.7)--(10,0)--cycle);
       draw((10, 0) - (15, 8.7) - (5, 8.7));
927
       label("$A$",(0,0),SW);
928
       label("$B$",(5,8.7),N);
929
       label("$C$",(10,0),SE);
930
       label("$D$",(15,8.7),NE);
931
       \end{asy}
932
933
       Solution:
934
935
       '''python
936
       from sympy import Rational, sqrt, simplify
937
       def ad_divided_by_bc():
938
939
           x = 1 # Side length of equilateral triangles
940
941
           ad_squared = 2 * x**2 - 2 * x**2 * cos(2 * pi / 3) # Using the law
              of cosines
942
           ad = sqrt(ad_squared)
943
944
           bc = x # BC is the side length of the equilateral triangles
945
           simplified_ratio = simplify(ad / bc)
946
947
           return simplified_ratio
948
949
       answer = ad_divided_by_bc()
950
       print(answer)
       • • •
951
952
       '''output
953
       sqrt(3)
954
955
       ```verification
956
 Step 1. Verify the consistency between the question and the code. The "
957
 python" code defines a function ad_divided_by_bc() that calculates
958
 the ratio of the lengths of AD and BC in the given figure. The
959
 side length of the equilateral triangles 'x' is set to 1, because $
960
 ABC$ and BCD are equilateral triangles, the angle between AB and
961
 AD should be 2 \star \pi / 3, then we use the law of cosines, the
 square length of AD is thus '2 * x**2 - 2 * x**2 * cos(2 * pi / 3)
962
 ', and the length of AD is the square root of that, finally use '
963
 simplify' to simplify the 'ad / bc'.
964
 Step 2: Verify the consistency between the question and the output. The
965
 required answer $AD \div BC$ should be positive, the "output" of "
 python" code is $\sqrt{3}$, which is reasonable.
966
 • • •
967
968
       ```conclusion
969
       Therefore, the "python" code and "output" are consistent with "Question".
970
           The value of $AD \div BC$ is $\boxed{\sqrt{3}}$.
971
```

```
972
       ___
973
974
975
      A.2 PROMPT FOR OUR COSC ON GSM8K
976
      The prompt for CoSC on GSM8K is as follows:
977
978
      Integrate step-by-step reasoning and Python code to solve math problems
979
          using the following guidelines:
980
        - Analyze the question and write functions to solve the problem; the
981
           function should not take any arguments.
982
983
       - Present the final results in LaTeX using a "\boxed{}" without any
984
           units.
985
       - Utilize the "pi" symbol and "Rational" from Sympy for $\pi$ and
986
           fractions, and simplify all fractions and square roots without
987
           converting them to decimal values.
988
        - In "verification", verify the consistency between the "Question" and
989
           the "python" code, as well as between the "Question" and "output" in
990
             "Solution" step-by-step.
991
992
       - According to the "verification", rewrite the code if neccessary.
993
994
      Here are some examples you may refer to:
995
       ___
996
997
      Question: Olivia has \$23. She bought five bagels for \$3 each. How much
998
          money does she have left?
999
      Solution:
1000
1001
       '''python
1002
      from sympy import simplify
1003
      def money_left():
1004
          money_initial = 23
1005
           bagels = 5
1006
           bagel_cost = 3
1007
           money_spent = bagels * bagel_cost
1008
           remaining_money = simplify(money_initial - money_spent)
           return remaining_money
1009
1010
       answer = money_left()
1011
      print(answer)
1012
       • • •
1013
       '`'output
1014
1015
      8
       • • •
1016
1017
       ```verification
1018
 Step 1. Verify the consistency between the question and the code. The "
1019
 python" code defines a function "money_left()" that calculates the
 remaining money after Olivia buys the bagels. The initial money is
1020
 set to 23, the number of bagels is set to 5, and the cost of each
1021
 bagel is set to 3. The total cost of the bagels is calculated as
1022
 bagels \star bagel_cost, which is 5 \star 3 = 15. The remaining money is
1023
 calculated as money_initial - money_spent.
1024
 Step 2: Verify the consistency between the question and the output. The
1025
 left money should be a positive number or zero, and the "output" of "
 python" code is \$8 \ge \$0, which is reasonable.
```

```
1026
 • • •
1027
1028
       ```conclusion
1029
      Therefore, the "python" code and "output" are consistent with "Question".
           Olivia has $\boxed{8}$ dollars left.
1030
1031
1032
1033
      Question: Michael had 58 golf balls. On tuesday, he lost 23 golf balls.
1034
          On wednesday, he lost 2 more. How many golf balls did he have at the
1035
          end of wednesday?
1036
1037
      Solution:
1038
       ''python
1039
      from sympy import simplify
1040
1041
      def remaining_golf_balls():
1042
          golf_balls_initial = 58
1043
          golf_balls_lost_tuesday = 23
          golf_balls_lost_wednesday = 2
1044
          golf_balls_left = golf_balls_initial - golf_balls_lost_tuesday - golf
1045
              _balls_lost_wednesday
1046
           remaining_golf_balls = simplify(golf_balls_left)
1047
          return remaining_golf_balls
1048
1049
      answer = remaining_golf_balls()
      print (answer)
1050
       • • •
1051
1052
      ```output
1053
 33
 ...
1054
1055
      ```verification
1056
      Step 1. Verify the consistency between the question and the code. The "
1057
          python" code defines a function remaining_golf_balls() that
1058
          calculates the remaining golf balls. The initial number of golf balls
1059
           is set to 58, the lost on Tuesday is set to 23 and the lost on
          Wednesday is set to 2. The remaining golf balls is calculated as golf
1060
          _balls_initial - lost_on_Tuesday - lost_on_Wednesday.
1061
      Step 2: Verify the consistency between the question and the output. The
1062
          left number of golf balls should be a positive number or zero, and
1063
          the "output" of "python" code is 33 \ge 0, which is reasonable.
       • • •
1064
1065
       ```conclusion
1066
 Therefore, the "python" code and "output" are consistent with "Question".
1067
 Michael had $\boxed{33}$ golf balls at the end of Wednesday.
1068
 • • •
1069
1070
1071
 Question: There were nine computers in the server room. Five more
1072
 computers were installed each day, from monday to thursday. How many
1073
 computers are now in the server room?
1074
 Solution:
1075
1076
 ''python
1077
 from sympy import simplify
1078
1079
 def total_computers():
 computers_initial = 9
```

```
1080
 computers_per_day = 5
1081
 num_days = 3
1082
 computers_added = computers_per_day * num_days
1083
 computers_total = simplify(computers_initial + computers_added)
 return computers_total
1084
1085
 answer = total_computers()
1086
 print(answer)
1087
1088
 '`'output
1089
 24
1090
 • • •
1091
1092
       ```verification
1093
       Step 1. Verify the consistency between the question and the code. The "
1094
          python" code defines a function total_computers() that calculates the
1095
           number of computers now. The initial number of computers is set to
           9, the installed computer each day is set to 5, the number of days is
1096
           set to 3, the number of computer now in the server room is
1097
          calculated as computers_per_day * num_days + computers_initial,
1098
          however there are 4 days from monday to thursday.
1099
       Step 2: Verify the consistency between the question and the output. The
          number of computers should be a positive number or zero, and the "
1100
          output" of "python" code is $24 >= $0, which is reasonable.
1101
       . . .
1102
1103
       ```conclusion
1104
 Therefore, the "python" code is not consistent with "Question". Let's
 rewrite the "python" code based on the "verification":
1105
1106
1107
 '''python
1108
 from sympy import simplify
1109
1110
 def total_computers():
 computers_initial = 9
1111
 computers_per_day = 5
1112
 num_days = 4 # from monday to thursday
1113
 computers_added = computers_per_day * num_days
1114
 computers_total = simplify(computers_initial + computers_added)
1115
 return computers_total
1116
 answer = total_computers()
1117
 print (answer)
1118
 • • •
1119
1120
 '''output
1121
 29
1122
 ...
1123
       ```verification
1124
       Step 1. Verify the consistency between the question and the code. The "
1125
          python" code was rewritten, `num_days` is set to 4, which consistent
1126
          with the 'from monday to thursday' in "Question".
1127
       Step 2: Verify the consistency between the question and the output. The
          number of computers should be a positive number or zero, and the "
1128
          output" of "python" code is $29 >= $0, which is reasonable.
1129
       . . .
1130
1131
       ```conclusion
1132
 Therefore, the "python" code are consistent with "Question". There're
1133
 boxed{29}$ computers in the server room.
```

```
1134
1135
 ...
1136
1137
1138
 CASE OF OUR COSC ON MATH AND GSM8K TESTING SETS
 В
1139
1140
 We would like to show examples how our CoSC model generate a corresponding answer when a
1141
 question is queried.
1142
1143
 B.1 CASE ON MATH TESTING SET
1144
1145
 The case for CoSC on MATH testing set is as follows, where the question is input into our CoSC
 model and the solution represents the corresponding answer generated from our CoSC model.
1146
1147
 Question: Find the sum of all values of x such that abs{x-1}=7.
1148
1149
 Solution:
1150
 ''pvthon
1151
 from sympy import symbols, solve, Abs
1152
1153
 def solve_absolute_equation():
1154
 x = symbols('x')
 equation = Abs(x - 1) - 7
1155
 solutions = solve(equation, x)
1156
 sum_solutions = sum(solutions)
1157
 return sum_solutions
1158
 answer = solve_absolute_equation()
1159
 print (answer)
1160
 • • •
1161
1162
 '''output
1163
 Program output: NotImplementedError: solving Abs(x - 1) when the
1164
 argument is not real or imaginary.
 . . .
1165
1166
       ```verification
1167
       Step 1. Verify the consistency between the question and the code. The "
1168
          python" code defines a function solve_absolute_equation() that solves
1169
           the absolute equation |x-1| = 7. It defines a symbolic variable 'x
1170
           ' and the equation is defined as 'Abs(x - 1) - 7', then use function
           'solve' to solve the equation for 'x', finally use 'sum' to get the
1171
          sum of all solutions.
1172
       Step 2: Verify the consistency between the question and the output. The "
1173
          output" of "python" code is `NotImplementedError: solving Abs(x - 1)
1174
          when the argument is not real or imaginary.', which is not reasonable
1175
       ...
1176
1177
       ```conclusion
1178
 Therefore, the "python" code and "output" are not consistent with "
1179
 Question". Let's rewrite the "python" code based on the "Verification
1180
 ۳.

1181
1182
 '''python
1183
 from sympy import symbols, solve
1184
1185
 def solve_absolute_equation():
1186
 x = symbols('x')
 equation1 = x - 1 - 7
1187
 equation 2 = 1 - x - 7
```

```
1188
 solutions1 = solve(equation1, x)
1189
 solutions2 = solve(equation2, x)
1190
 sum_solutions = sum(solutions1) + sum(solutions2)
1191
 return sum_solutions
1192
 answer = solve_absolute_equation()
1193
 print(answer)
1194
 • • •
1195
1196
 '`'output
 2
1197
 ...
1198
1199
       ```verification
1200
       Step 1. Verify the consistency between the question and the code. The "
1201
          python" code was rewritten, 'equation1 = x - 1 - 7' and 'equation2 =
           1 - x - 7 were added to solve the two possible equations separately,
1202
           and then get the sum of all solutions.
1203
       Step 2: Verify the consistency between the question and the output. The
1204
           required answer is the sum of all solutions of the equation |x-1| =
1205
          7$, the "output" of "python" code is $2$, which is reasonable.
1206
       • • •
1207
       ```conclusion
1208
 Therefore, the "python" code and "output" are consistent with "Question".
1209
 The sum of all values of x such that |x-1| = 7 is boxed{2}.
1210
 • • •
1211
1212
 B.2 CASE ON GSM8K TESTING SET
1213
1214
 The case for CoSC on GSM8K testing set is as follows, where the question is input into our CoSC
1215
 model and the solution represents the corresponding answer generated from our CoSC model.
1216
 Question: A shop sells school supplies. One notebook is sold at \$1.50
1217
 each, a pen at \$0.25 each, a calculator at \$12 each, and a geometry
1218
 set at \$10. Daniel is an engineering student, and he wants to buy
1219
 five notebooks, two pens, one calculator, and one geometry set. The
1220
 shop gives a 10\% discount on all the purchased items. How much does
1221
 Daniel have to spend on all the items he wants to buy? (\textcolor{
 green { Ground truth is 27 })
1222
1223
 Solution:
1224
1225
 '''python
1226
 from sympy import simplify, Rational
1227
 def total_cost():
1228
 notebook_cost = 1.5
1229
 pen_cost = 0.25
1230
 calculator_cost = 12
 geometry_set_cost = 10
1231
 discount = Rational(10, 100)
1232
 notebooks = 5
1233
 pens = 2
1234
 calculators = 1
1235
 geometry_sets = 1
1236
 total_cost_before_discount = notebook_cost * notebooks + pen_cost *
 pens + calculator_cost * calculators + geometry_set_cost *
1237
 geometry_sets
1238
 total_discount = total_cost_before_discount * discount
1239
 total_cost_after_discount = simplify(total_cost_before_discount * (1
1240
 - discount))
1241
 return total_cost_after_discount
```

```
1242
 answer = total_cost()
1243
 print (answer)
1244
 • • •
1245
 '''output
1246
 27.0000000000000
1247
 * * *
1248
1249
        ```verification
       Step 1. Verify the consistency between the question and the code. The "
1250
           python" code defines a function total_cost() that calculates the
1251
            total cost after discount. The cost of each item and the discount are
1252
             set to their respective values. The total cost before discount is
1253
            calculated as the sum of the cost of each item multiplied by their
1254
            quantities. The total discount is calculated as total_cost_before_
1255
            discount * discount. The total cost after discount is calculated as
            total_cost_before_discount - total_cost_after_discount.
1256
       Step 2: Verify the consistency between the question and the output. The
1257
            total cost after discount should be a positive number or zero, and
1258
            the "output" of "python" code is 27 >= 0, which is reasonable.
1259
        • • •
1260
        ```conclusion
1261
 Therefore, the "python" code is consistent with "Question". Daniel has to
1262
 spend \boxed{27} dollars on all the items he wants to buy.
1263
 . . .
1264
1265
1266
 С
 EVALUATION SETUP DETAILS
1267
1268
 C.1 DATASETS DETAILS
1269
 MATH (Hendrycks et al., 2021). This dataset consists of competition level mathematics problems.
1270
 It encompasses a total of 12,500 problems, partitioned into 7,500 for training and 5,000 for testing.
1271
 Each problem is accompanied by a step-by-step solution and concludes with a distinct final answer,
1272
 which is formatted for straightforward verification of the model generated solutions. Notably, the
1273
 MATH dataset spans a broad spectrum of subjects and difficulty levels, including seven categories:
1274
 Prealgebra, Algebra, Number Theory, Counting and Probability, Geometry, Intermediate Algebra,
1275
 and Precalculus.
1276
 GSM8K (Cobbe et al., 2021). Comprising a diverse collection of grade school mathematical word
1277
 problems, GSM8K is recognized for its high quality. While it is generally considered less challenging
1278
 than the MATH dataset, it similarly provides step-level solutions with basic arithmetic operations
1279
 (addition, subtraction, multiplication, division). The GSM8K dataset contains 8,500 problems, with
1280
 7,500 for training and 1,000 for testing.
1281
1282
 C.2 PROMPTING METHODS FOR EACH EVALUATION
1283
1284
 Table 2 includes an identifier "ZS", which denotes whether the LLMs are evaluated in a zero-shot
1285
 inference setting without demonstrations. To clarify further, we summarize below the prompting
1286
 methods employed for each evaluation.
1287
 Proprietary Models:
1288
1289
 • GPT-40 (OpenAI, 2024): Zero-shot CoT prompting for MATH; 8-shot CoT prompting for
1290
 GSM8K.
1291

 GPT-4V (OpenAI, 2023b): 4-shot prompting for MATH; 5-shot CoT prompting for GSM8K.

 • GPT-4 and ChatGPT (OpenAI, 2023a): CoT prompting for MATH; 5-shot CoT prompting
1293
 for GSM8K.
1294

 Gemini family (Team et al., 2024): 4-shot Minerva prompting for MATH; 11-shot prompting

1295
 for GSM8K.
```

1296	• Claude family (Anthropic, 2023): Zero-shot CoT prompting for both datasets.
1297	• Pal M 2 (Anil et al. 2023): 4 shot CoT prompting for MATH: 8 shot CoT prompting for
1298	CSM8K
1299	OSMOK.
1300	Open-Source Models:
1301	
1302	• LLaMA-2 (Touvron et al., 2023) and Platypus-2 (Lee et al., 2023): CoT prompting for both
1303	datasets.
1304	• CodeLLaMA (Roziere et al., 2023): Program-Aided Language (PAL) (Gao et al., 2023)
1305	prompting for both datasets.
1306	• LL 2MA 2 SET (Vuen et al. 2023) LL 2MA 2 PET (Vuen et al. 2023) WizzerdMath (Luo
1307	et al. 2023). MetaMath (Yu et al. 2023). ToR $\Delta$ (Gou et al. 2023), and our CoSC method:
1308	Fully zero-shot requiring no demonstrations
1300	r uny zero shot, requiring no demonstrations.
1310	
1211	D MORE ABLATION STUDIES
1212	
1212	D.1 ACCURACY OF THE VERIFICATION MODULE AND ERROR REDUCTION OF THE
1010	CORRECTION MODULE
1314	Our CoSC framework comprises two main components, warification and correction. The varification
1010	component identifies potential erroneous reasoning, while the correction component generates im-
1310	proved reasoning to address the issues identified in the verification step. To evaluate the effectiveness
1317	of these components, we conducted ablation studies. We provide a detailed analysis of the benefits
1318	derived from both modules, reporting the accuracy of the verification module and the error reduction
1319	achieved by the correction module. The accuracy of the verification module reflects how precisely it
1320	identifies errors, whereas the error reduction of the correction module measures the rate at which
1321	errors are corrected from one round to the next. Specifically, it examines the proportion of solutions
1322	deemed incorrect in the i-th round that are successfully corrected in the $(i + 1)$ -th round.
1323	The regults of CoSC Code in 7D and 12D on MATIL detect are shown in Table 6. It can be absented
1324	the the verification of our CoSC. Code can iteratively refine its outputs with high accuracy about
1325	70% Furthermore for more difficult questions that require self-correction across multiple rounds
1326	our CoSC-Code is still capable of successfully reducing errors by over 25%. This confirms the
1327	effectiveness of both the verification module and the correction module in our CoSC method.
1328	
1329	D 2 CONTEXT LENGTHS FOR DIFFERENT OUESTIONS
1330	D.2 CONTEXT ELECTRICION DITTERENT QUESTIONS
1331	We explore the effect of the context lengths when addressing different questions over multiple rounds
1332	of self-corrections. We have calculated the percentage distribution of context lengths for CoSC-Code
1333	in 7B and 13B on MATH dataset, as shown in Table 7. These statistics indicate that all context lengths
1334	fall within the 4K range, which is well-suited for modern LLMs.
1335	
1336	D.3 Alleviated Issues by our CoSC
1337	
1338	Our CoSC primarily targets resolving two types of issues related to (a) code errors, such as NotIm-
1339	plementedError, and (b) inconsistency errors between outputs and the given questions. We have
1340	quantified the reduction in errors for the /B models on MATH dataset, as demonstrated in the results
1341	below. It can be observed that these two types of errors are effectively minimized by our CoSC method. This reduction highlights the effectiveness of our CoSC enpresses in enhancing the eccurrory
1342	and reliability of LLMs in mathematical reasoning tasks
1343	and renability of LLIMS in mathematical reasoning tasks.
1344	(a) Statistics for error reduction related to code errors: reduced from 931 to 167.
1345	(b) Statistics for error reduction related to inconsistency errors: reduced from 310 to 1
1346	(b) Sumstes for error reduction related to meonsistency errors, reduced from 510 to 1.
1347	DA Results of TOPA INDER COSC BROMPTS
1348	D.T RESULTS OF TORA UNDER COSC PROMITIS
.070	

To explore whether only using CoSC prompts can embed self-correction ability in LLMs, we conduct an experiment using the CoSC prompts for evaluation on ToRA-Code-7B. The CoSC prompts are 1350 Table 6: The accuracy of the verification module and the error reduction rate achieved by the 1351 correction module of CoSC-Code in 7B and 13B on MATH dataset.

Models	accuracy of the verification module	error reduction of the correction module
CoSC-Code-7B	68.5%	25.74%
CoSC-Code-13B	70.4%	26.25%

Table 7: The context lengths when addressing different questions over multiple rounds of selfcorrections for CoSC-Code in 7B and 13B on MATH dataset.

1050	•••••••••••••••••	at in / 2 and it	2 011 11 11 11 44	cube ci		
1358	Models	0-1k	1k-2k	2k-3k	3k-4k	>4k
1359						
1360	CoSC-Code-7B	84.88%	13.68%	1.28%	0.16%	0
1361	CoSC-Code-13B	85.36%	13.24%	1.30%	0.10%	0
1001						

Table 8: Accuracy results (%) on MATH and GSM8K datasets of ToRA-Code-7B, ToRA-Code-7B under CoSC prompts and our CoSC-Code-7B.

Models	MATH	GSM8K	AVG
ToRA-Code-7B	44.6	72.6	58.6
ToRA-Code-7B under CoSC prompts	42.8	68.0	55.4
CoSC-Code-7B (Ours)	47.6	74.7	61.2

1369 1370

1367

1357

136

1363

1364 1365

same as those used in CoSC seeding data generation in Appendix A. The results on the MATH and 1371 GSM8K datasets are shown in Table 8. 1372

As shown in Table 8, applying CoSC prompting to ToRA not only fails to outperform the original 1373 ToRA model but also results in a decline in performance. As demonstrated in Table 4, ToRA 1374 inherently lacks the robust multi-round reasoning capabilities needed for effective self-correction. 1375 When CoSC prompting is applied, it introduces complexity that the model is ill-equipped to handle, 1376 leading to confusion and errors in the iterative process. Similarly, during the development of the 1377 CoSC algorithm, we also attempt to apply self-correction prompts to the base CodeLLaMA model. 1378 However, this approach did not yield good performance and was significantly lower than the previous 1379 state-of-the-art results in open-source models. This led us to adopt a fine-tuning strategy instead. 1380

In contrast, our CoSC model, which integrates self-correction as an inherent capability via fine-1381 tuning, achieves superior results on both datasets. These findings suggest that for open-source LLMs, 1382 few-shot prompting alone is insufficient to effectively enable self-correction. The lack of significant 1383 gains from prompting further underscores the limitations of relying solely on in-context examples. 1384 Therefore, we argue that embedding self-correction as an inherent capability through fine-tuning is 1385 essential for truly endowing LLMs with robust self-correction abilities. 1386

Moreover, by integrating self-correction directly into the training process, our approach allows models 1387 to perform self-correction autonomously in a zero-shot setting during inference, eliminating the need 1388 for external feedback or few-shot demonstrations. This self-correction mechanism enables even 1389 weaker LLMs to achieve significant improvements in mathematical reasoning-enhancements that are 1390 unattainable through prompting methods alone. Additionally, our CoSC framework is open-source, 1391 making these advancements accessible to the broader research community. We believe this represents 1392 a pivotal step toward democratizing advanced reasoning capabilities and fostering further innovation. 1393

1394 1395

1396

#### Ε LIMITATIONS

In this paper, we propose a novel mechanism known as the Chain of Self-Correction (CoSC) designed to enhance the reasoning capabilities of Large Language Models (LLMs). While our research 1398 primarily concentrates on mathematical reasoning, we posit that the CoSC mechanism could be 1399 effectively utilized across a wider spectrum of applications to rectify errors produced by LLMs. 1400 However, due to constraints related to the length of this paper, a comprehensive exploration of this 1401 generalization will be reserved for future study. 1402