
R2P: Reformulate–Retrieve–Program for Robust Mathematical Reasoning in LLMs

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Large language models (LLMs) remain brittle on mathematical word problems:
2 small surface-form changes can shift answer distributions and degrade solve rates,
3 while multi-step computation is error-prone. We present R2P, a three-stage infer-
4 ence framework that (1) Reformulates each problem into diverse paraphrases
5 to reduce surface-form bias, (2) Retrieves domain-aligned few-shot exemplars
6 from a curated bank via lightweight embeddings, and (3) Programs explicit inter-
7 mediate Python code (Program-of-Thoughts) to execute symbolic computations.
8 Under a fixed path budget, R2P samples multiple reasoning trajectories across
9 reformulations and aggregates outcomes by voting, improving both accuracy and
10 consistency. Evaluations on GSM8K, AQuA, and SVAMP with an off-the-shelf
11 9B-parameter LLM (zero/few-shot, no fine-tuning) show consistent gains over
12 Chain-of-Thought, self-consistency, and vanilla Program-of-Thoughts. Ablations
13 varying the number of reformulations and comparing naïve vs. in-context reformu-
14 lation demonstrate that (i) exposing the model to multiple surface forms reliably
15 improves solve rates, and (ii) domain-aware retrieval further boosts robustness.
16 We analyze typical failure modes—misinterpretation of quantities and arithmetic
17 slips—and show how reformulation plus code execution mitigates both. R2P offers
18 a simple, model-agnostic recipe for more reliable mathematical reasoning without
19 additional training.

20 1 Introduction

21 Mathematical reasoning is a cornerstone of problem-solving, with applications spanning diverse
22 fields such as physics, engineering, economics, and computer science. However, despite their success
23 in general natural language processing tasks, existing Large Language Models (LLMs) such as GPT-4
24 struggle with mathematical problems that demand precision, logical reasoning, and step-by-step
25 computation [1]. This gap arises because LLMs rely heavily on statistical patterns in natural language,
26 which often fail to capture the formal structure and symbolic complexity of mathematical problems.
27 [2].

28 Current advancements in mathematical reasoning with LLMs have primarily focused on methods
29 like Chain-of-Thought (CoT) prompting [3, 4, 5], which encourages models to break problems into
30 reasoning steps. Self-Consistency (SC) methods further refine CoT by introducing multiple reasoning
31 paths to identify consistent answers. While effective, these approaches are still limited when faced
32 with problems requiring complex computation or diverse logical forms. As shown in Fig 1(a), the
33 LLM fails to provide the correct answer directly when posed with a complex calculation problem.

34 This limitation is handled by another notable approach, Program of Thoughts (PoT), which introduces
35 intermediate code generation to represent reasoning steps explicitly, improving the interpretability of
36 solutions and disentangle computation from the reasoning process. However, PoT alone cannot resolve

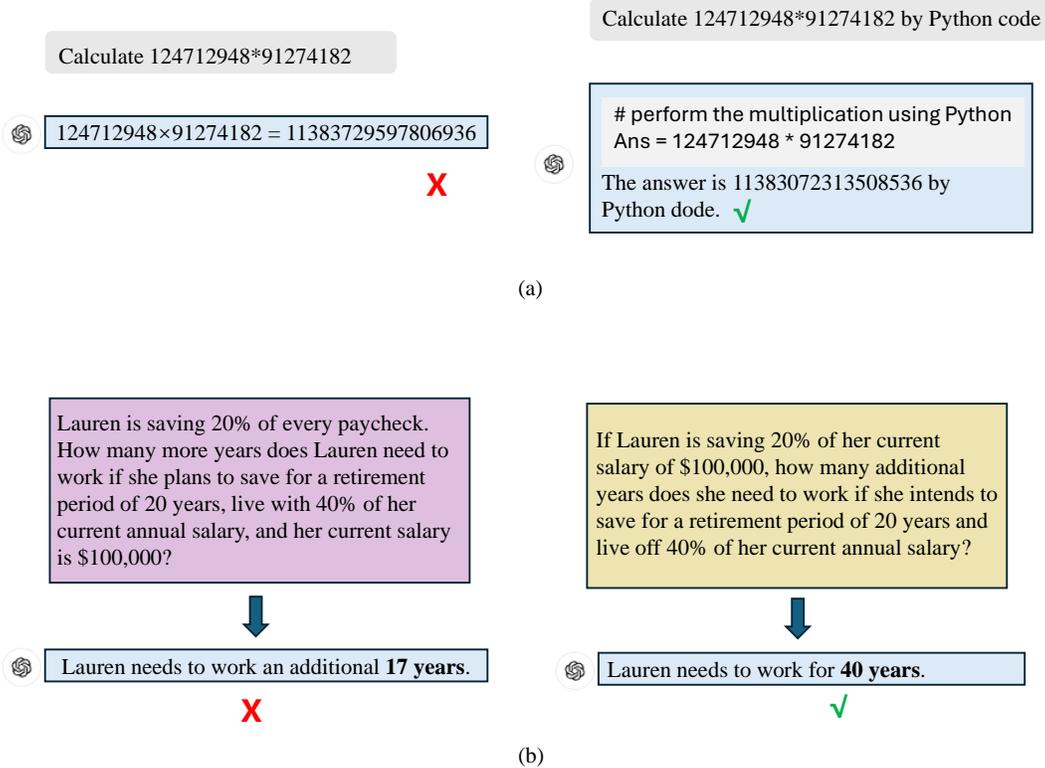


Figure 1: Limitations for GPT-4 solving MWPs. (a) When GPT-4 is asked to directly provide the result of a complex arithmetic problem, the answer is incorrect.(b) When the original problem is reformulated, answers are inconsistent.

37 inconsistencies arising from ambiguous structured problem forms. As shown in Fig 1(b), we can
 38 observe that small changes in the surface form of the mathematical problem can lead to significantly
 39 different outcomes.

40 To address these challenges, we propose R2P, a novel two-stage framework that integrates Refor-
 41 mulation of Mathematical Problems (RM) and Program of Thoughts (PoT). The RM stage prompts
 42 the LLM to generate multiple paraphrased versions of the same problem, thereby mitigating the
 43 adverse effects of poorly structured problem formulations. By generating multiple reformulations of
 44 a given problem and exposing the model to various formulations of the same task, LLMs can uncover
 45 the core mathematical structure underlying the problem [6]. This diversity in problem presentation
 46 improves the LLM’s accuracy and consistency in solving the problem. In the PoT stage, the model
 47 generates Python code as intermediate reasoning steps, decoupling symbolic computations from
 48 natural language reasoning. A voting mechanism ensures consistency across reformulated solutions,
 49 enhancing accuracy and robustness. Domain-specific context activation can retrieve semantically
 50 coherent ⟨question-solution⟩ pairs from a pre-built question bank, enabling LLM to reason based on
 51 specific domains.

52 During the deduction phase, we introduce a novel integration of Few-Shot Learning module to further
 53 improve the accuracy and robustness of LLMs in mathematical problem-solving which operates as
 54 follows : Given a problem p , we (1) classify its mathematical domain D (e.g., algebra, arithmetic)
 55 using a lightweight embedding model, (2) retrieve the top-5 ⟨question,solution⟩ pairs from D ’s bank
 56 via cosine similarity over sentence-BERT embeddings, and (3) inject these examples into the prompt
 57 to prime the LLM with domain-specific reasoning patterns. This process ensures semantic alignment
 58 between the input problem and the few-shot exemplars, enhancing both reformulation diversity and
 59 solution accuracy. Additionally, we analyze the performance gains from the inclusion of few-shot
 60 examples, showing that incorporating these examples allows for more accurate model predictions,
 61 even when faced with unfamiliar or highly variable problem structures.

62 The main contributions of this paper are as follows:

- 63 • We introduce R2P, a novel framework that combines surface-level problem reformulation
64 and explicit intermediate code generation to improve mathematical reasoning in LLMs.
- 65 • We conduct comprehensive experiments on widely-used benchmarks, including GSM8K,
66 AQuA, and MATH, demonstrating that R2P outperforms baseline approaches across different
67 datasets.
- 68 • We provide analysis of the effectiveness of RM and PoT stages, highlighting their contribu-
69 tions to solving mathematical problems.

70 2 Related Work

71 2.1 Mathematical Reasoning in LLMs

72 In recent years, numerous studies have explored using the System-2 reasoning approach to solve
73 mathematical problems with LLMs [3, 7, 8]. As a prominent framework, chain-of-thought(CoT)
74 is proposed by [3]. Rather than generate the answer directly, it prompts the LLMs to produce a
75 sequence of intermediate reasoning steps. [7] further extended CoT by self-consistency. They
76 generate multiple reasoning steps from different angles and potentially lead to the same answer. [9]
77 proposes program-of-thoughts, using intermediate codes to represent the reasoning process.

78 2.2 Specialized Models for Mathematical Tasks

79 To address the limitations of general-purpose LLMs, specialized approaches have been developed.
80 For instance, MathGPT [10] integrates symbolic computation engines to solve algebraic problems,
81 combining language modeling with symbolic reasoning. However, this hybrid approach still depends
82 heavily on the capabilities of external tools. Other models, such as GeoSolver [11] and MathQA [12],
83 focus on specific domains like geometry or math question-answering, using domain-specific datasets
84 and tailored architectures. While these models perform well in their respective areas, their narrow
85 focus limits generalization to broader mathematical tasks.

86 2.3 Problem Reformulation in Mathematical Problem Solving

87 The role of problem reformulation in enhancing the performance of language learning models (LLMs)
88 in mathematical problem solving has gained increasing attention in recent years. Early work by
89 [13, 14] demonstrated that altering the structure and phrasing of a problem can lead to improved
90 reasoning outcomes in LLMs. This line of research has been extended by [15], who explored how
91 different surface forms of mathematical problems can significantly influence the solve rate of LLMs.
92 Further studies have investigated the impact of reformulation strategies such as paraphrasing [16]
93 and introducing problem variants [17], suggesting that diverse representations help models better
94 understand and process complex tasks. Building on these insights, recent works have also delved into
95 the combination of reformulation with in-context learning techniques to optimize model performance
96 through adaptive problem presentations.

97 3 Method

98 We propose R2P, a framework that leverages the potential of large language models (LLMs) to solve
99 mathematical problems. The overall architecture of R2P is illustrated in Fig. 2, consisting of three
100 key components:

- 101 • **Stage I - Problem Reformulation (RM):** We reformulate the given mathematical problems
102 into diverse surface forms, enabling the LLM to better grasp the underlying structure of the
103 problems and mitigate structural biases that may hinder accurate reasoning.
- 104 • **Stage II - Domain-Aware Few-Shot Learning:** To enhance problem-solving accuracy,
105 we introduce a domain-specific few-shot learning mechanism that retrieves semantically
106 aligned example problems from a pre-constructed question bank. This structured retrieval
107 ensures that the model benefits from contextualized exemplars, improving solution accuracy
108 even for problems with unfamiliar or highly variable structures.

109 • **Stage III - Program of Thoughts (PoT):** We employ the **Program of Thoughts** [9]
 110 approach, which instructs the LLM to generate **intermediate Python code** to execute
 111 symbolic computations and enhance reasoning transparency. By explicitly decoupling the
 112 computational process from the logical reasoning process, this method ensures that numerical
 113 calculations remain precise while maintaining interpretability. To further improve
 114 robustness, we implement a voting mechanism across multiple reformulated problem versions,
 115 aggregating diverse reasoning paths to achieve more consistent and reliable final
 116 answer.

117 The details of these two stages are described in the following subsections.

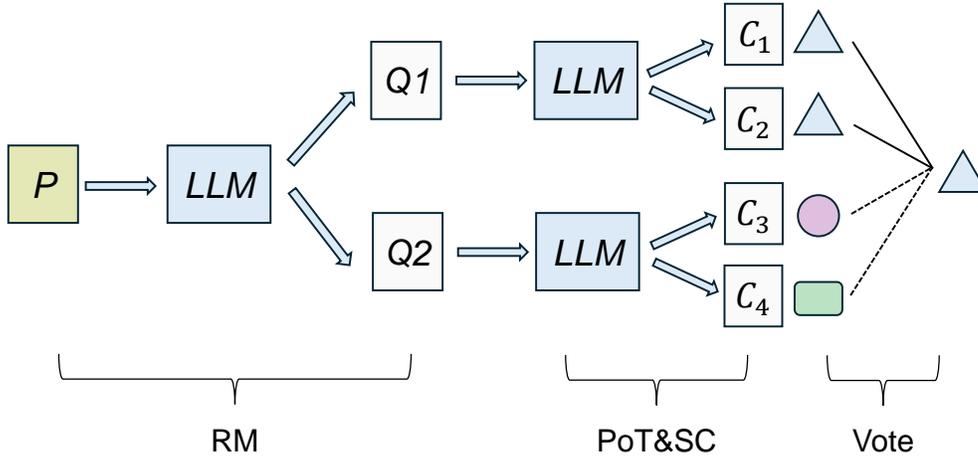


Figure 2: Overview of our proposed R2P framework. It operates in two stages. First, the RM method is employed to reformulate the original problem into various forms using the LLM. Next, the PoT approach is utilized to guide the LLM in generating intermediate code, which not only illustrates the reasoning process but also computes the result. The final answer is then determined through a voting mechanism.

118 3.1 Reformulate Mathematical Problems(RM)

119 As demonstrated in Section 1, the surface form of a problem significantly influences the performance
 120 of LLMs [18]. Therefore, we prompt the LLM to generate K different surface forms of the original
 121 problem and use a voting mechanism to determine the answer. The intuition behind this approach
 122 is that if a problem exhibits a low solve rate and ineffective reasoning paths due to its original
 123 surface form, introducing diversity in its surface forms can enhance the likelihood of finding a correct
 124 solution.

125 It is important to note that the original problem is reformulated into different forms using the same
 126 LLM that is used to solve the problem. This ensures that the improvement in model performance is
 127 due to the diversity of the reformulated problem forms, rather than the sharing of knowledge with
 128 other LLMs.

129 In this paper, we explore two ways of prompting LLM to generate reformulated problems. The naive
 130 prompt template is "Reformulate the following math problem, try to change the sentence structure of
 131 the problem: {input problem}" . The In-Context method begins by identifying effective examples: We
 132 conduct experiments to determine which pairs of (Original Problem, Reformulated Problem) exhibit
 133 the largest margin in solve rates. Then, we use these examples to enable in-context learning for the
 134 LLM, as depicted in Fig. 3.

135 3.2 Program of Thoughts(PoT)

136 Program of Thoughts (PoT) is a method aimed at enhancing the reasoning capabilities of large
 137 language models (LLMs) [9]. It works by decomposing complex tasks into a series of intermediate
 138 steps, or "thoughts", which guide the model through a structured reasoning process. Each thought

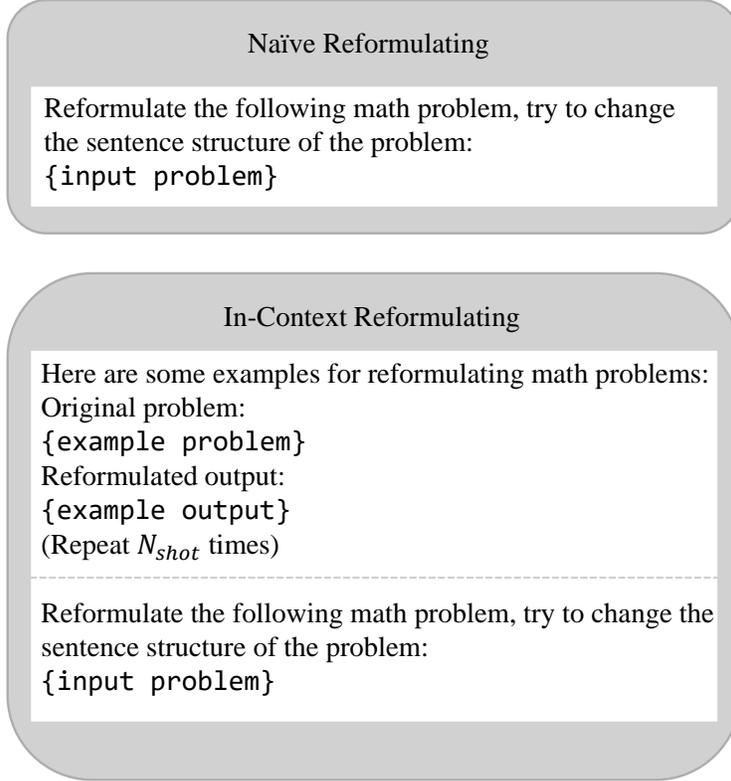


Figure 3: Naive and In-Context Reformulating prompts for LLM.

139 represents a logical step that incrementally leads to the final answer, enabling the model to tackle
 140 intricate problems using a step-by-step approach.

141 This method improves the model’s ability to solve tasks that require multi-step reasoning, such as
 142 mathematical or logical problems, by fostering transparency in the reasoning process and increasing
 143 accuracy in the final result. PoT has proven effective in scenarios where direct answers are difficult,
 144 allowing LLMs to perform more reliably in problem-solving tasks.

145 In our proposed R2P, we aim to instruct the LLM to generate intermediate Python code to solve
 146 mathematical problems. This approach can clearly show the reasoning process of LLM and improve
 147 accuracy.

148 3.3 Self-Consistency(SC) and Voting

149 In our proposed R2P framework, we reformulate the original problems into K different surface forms.
 150 Moreover, we utilize the Self-Consistency(SC) [7] method. Specifically, for each formulated problem,
 151 we let LLM generate $\frac{K}{N}$ reasoning paths, and thus the total number of generated answer is N . We
 152 then vote for the final answer.

153 By increasing K , we maintain a fixed total number of reasoning paths N . This approach effectively
 154 isolates the effect of diversifying the reasoning paths from merely increasing their number. It also
 155 ensures a fair comparison with the self-consistency baselines.

156 4 Experiments

157 4.1 Experimental Settings

158 **Datasets** We evaluate our approach on the following public mathematics reasoning benchmarks:

- **GSM8k** [2] contains 8.5K linguistically diverse grade school-level math questions with moderate difficulties.
- **AQuA** [19] consists of 100K algebraic word problems, including the questions, the possible multiple-choice options, and natural language answer rationales from GMAT and GRE.
- **SVAMP** [20] contains 1K arithmetic word problems. It focuses on basic arithmetic operations such as addition, subtraction, multiplication, and division.

Large Language Model We use GLM-4-9B [21] as our base model. All experiments are conducted in zero-shot or few-shot settings, without training or fine-tuning it. For generation configs, We set the temperature $T = 0.7$, Top-p= 0.8 and Top-k= 3. The total number of reasoning paths N we sample for each problem is 16.

Implementation Details For problem solving, the PoT process consists of two steps, as illustrated in Fig. 4. First, we prompt the LLM to generate Python code and store the result in a variable with a fixed name, allowing for convenient extraction. If the problem includes options, we instruct the LLM to identify the closest match among the provided options.

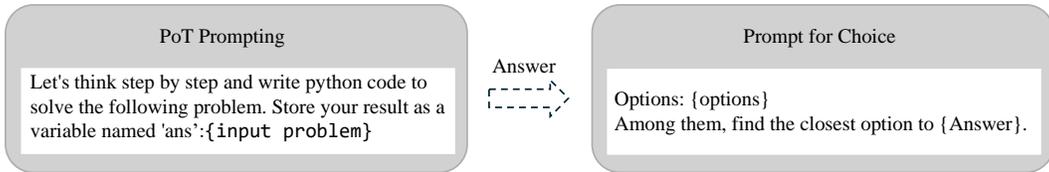


Figure 4: Prompts for solving the problem and selecting the correct choice when options are provided.

172

173 4.2 Effectiveness of Reformulating

174 To verify the effectiveness of RM, we reformulate the selected problems from the AQuA dataset and
 175 calculate the solve rate difference between the original problems and their reformulated versions. An
 176 example is shown in Fig. 5 to provide readers with an intuitive understanding of the RM process.
 177 When the surface form of the original problem is altered, the solve rate increases from 43.8% to
 178 81.3%.

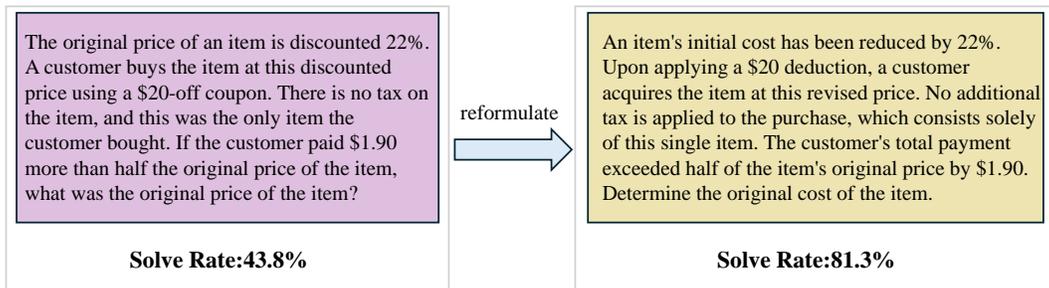


Figure 5: An example from AQuA dataset. When the original problem is reformulated, the solve rate improves significantly.

179 Furthermore, we compute the solve rate difference between the original problems p and the reformu-
 180 lated problems p_r from the AQuA dataset, defined as where $SR(p_r) - SR(p)$, where SR denotes
 181 the solve rate. The results are shown in Fig. 6. we can observe that reformulating the problems leads
 182 to an overall improvement in the solve rate.

183 4.3 Main Results

184 In this subsection, we compare the performance of R2P with Chain of Thoughts (CoT), vanilla
 185 self-consistency (SC), and vanilla Program of Thoughts(PoT). All results are presented in Table 1. In
 186 this evaluation, we apply the naive reformulation of the original problems. the number of reformulated

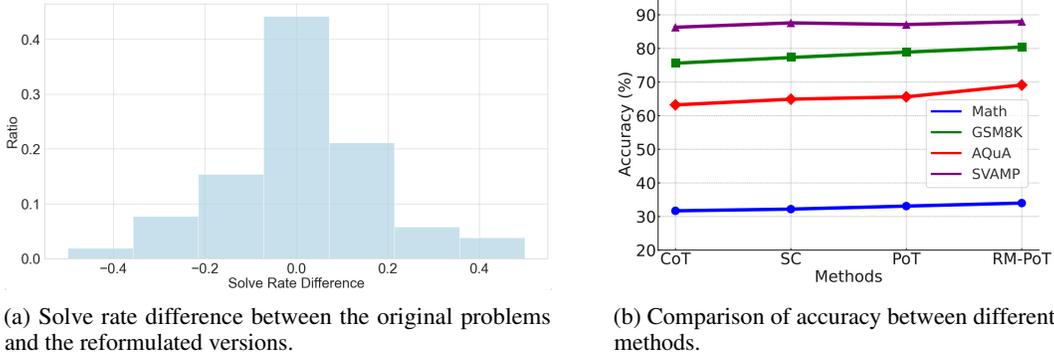


Figure 6: The performance of reformulated versions.

187 problems is set to $K = 4$, and the total number of reasoning paths is $N = 16$. For SC, we also
 188 maintain $N = 16$.

189 The performance metric presented in the table is the accuracy of answers after voting. We can observe
 190 that R2P outperforms the other three baselines across all datasets, despite the considerable likelihood
 191 of generating reformulation with a lower solve rate (see Fig. 6). We assume that this may be because
 192 the model demonstrates more consistent understanding of the problems when exposed to various
 193 reformulations, thereby avoiding misinterpretations of the original problem in certain specific cases.
 194 We will further discuss it in section 5.

Table 1: Comparison of accuracy between different methods.

	Math	GSM8K	AQuA	SVAMP
CoT	31.7	75.6	63.2	86.3
SC	32.2	77.3	64.9	87.6
PoT	33.1	78.9	65.6	87.1
R2P	34.0	80.4	69.1	88.0

195 4.4 Ablation Study

196 In this subsection, we vary the number of reformulated problems K across $\{1,2,4\}$, while keeping
 197 the total reasoning paths fixed at $N = 16$. Additionally, we compare the naive reformulation with the
 In-Context reformulation, as outlined in Fig 3. The results are shown in Table 2.

Table 2: Comparison of accuracy between different settings of K and reformulation.

	K	Math	GSM8K	AQuA	SVAMP
Naive	1	32.2	78.2	66.0	86.9
	2	33.1	79.8	67.7	88.4
	4	34.0	80.4	69.1	88.0
In-Context	1	32.5	78.4	67.6	87.1
	2	33.6	80.1	69.4	89.0
	4	34.3	80.9	72.2	89.6

198
 199 We observe that as K increases, the performance of the LLM improves as well. This further validates
 200 our assumption that exposing the LLM to different surface forms of a problem allows it to better
 201 grasp the underlying structure. Notably, on the SVAMP dataset, the LLM performs better when
 202 $K = 2$ than when $K = 4$. This may be because the problems in this dataset are relatively simple,
 203 and the process of reformulating and solving the problems involves some degree of randomness. By
 204 comparing the naive reformulation with the In-Context reformulation, we can conclude that the LLM
 205 learns the reformulation method more effectively when provided with good examples.

5 Discussion

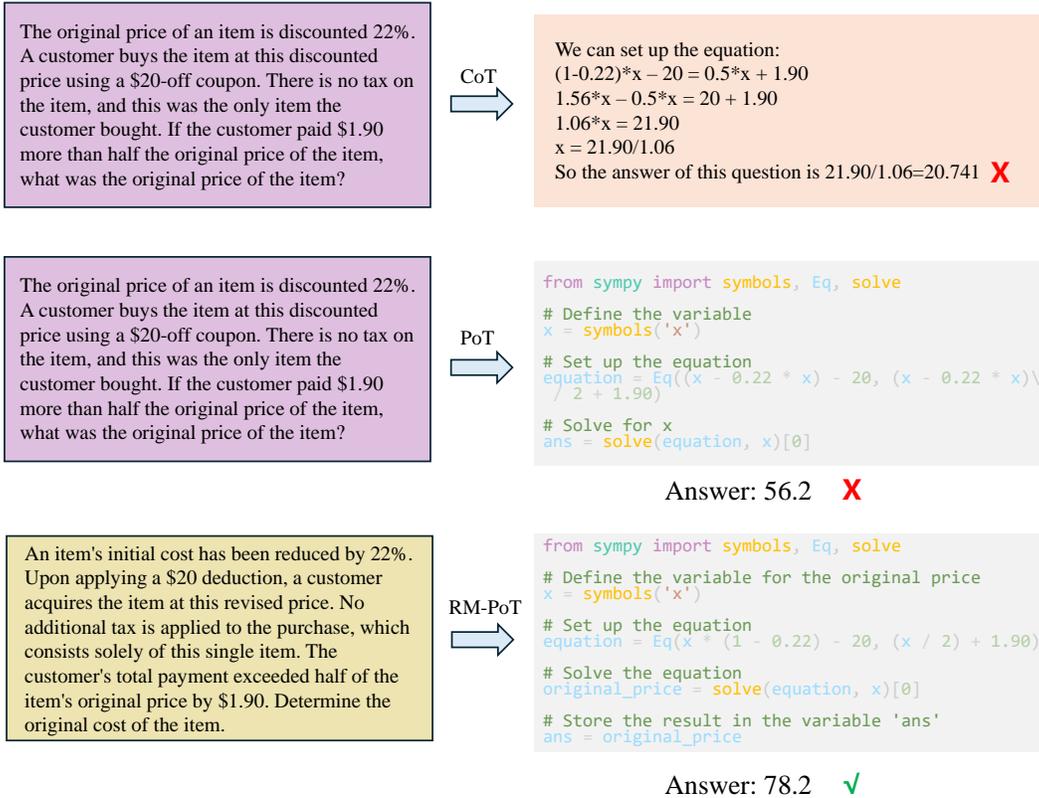


Figure 7: An example showing how CoT, PoT, and R2P solve the same problem. Both CoT and PoT provide incorrect answers, while R2P gives the correct answer.

207 In this section, we discuss how R2P helps the LLM solve math problems through the example shown
 208 in Fig. 7. When applying CoT, the LLM makes two mistakes. First, there is an error in rearranging
 209 the terms of the equation in the second step. Second, the result of $21.90/1.60$ is actually 20.660, not
 210 20.741. This demonstrates that the LLMs has certain shortcomings in computation and reasoning.

211 To disentangle computation from reasoning process, we then apply the PoT method. We find that
 212 the LLM's answer is generally correct, but it misinterprets the original price as the discounted price
 213 during the understanding of the problem. However, after reformulating the original problem, the
 214 LLM correctly understands the question and provides the correct answer. Although both versions
 215 of the problem contain the term 'original price,' the LLM does not fully grasp this in the original
 216 formulation. The deeper reasons behind this remain unclear and will be explored in future work.

217 6 Conclusion

218 In this paper, we present R2P++, a framework that advances mathematical reasoning in LLMs through
 219 domain-aware few-shot learning, problem reformulation, and Program of Thoughts. By retrieving
 220 semantically aligned examples from a pre-constructed question bank, our method primes LLMs to
 221 interpret problems consistently, while reformulation and code generation mitigate structural ambi-
 222 guities and computational errors. Experiments across datasets demonstrate statistically significant
 223 accuracy gains and robustness to linguistic variations.

224 However, the underlying reasons why naive reformulation can enhance performance are still unclear.
 225 In some cases, the reformulated problem even exhibits a lower solve rate than the original. Future
 226 work could explore more effective ways to reformulate problems and integrate fine-tuning methods
 227 to further improve the performance of LLMs. At the same time, this method can be improved in the
 228 following ways:

- 229 • Generalization: Extend to geometry and calculus by curating domain-specific banks with
230 diagram-to-code mappings.
- 231 • Adaptive Retrieval: Dynamically adjust the number of few-shot examples (K) based on
232 problem complexity.
- 233 • Human-in-Loop: Integrate user feedback to refine the question bank and domain classifier.

234 **References**

- 235 [1] Dan Hendrycks, Collin Burns, Steven Basart, Spencer Zou, and Mantas Mazeika. Measur-
236 ing mathematical problem solving with the math dataset. *Advances in Neural Information*
237 *Processing Systems (NeurIPS)*, 2021.
- 238 [2] Karl Cobbe, Vineet Kosaraju, and Mohammad Bavarian. Training verifiers to solve math word
239 problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 240 [3] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le,
241 Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models.
242 *Advances in neural information processing systems*, 35:24824–24837, 2022.
- 243 [4] Shizhe Diao, Pengcheng Wang, Yong Lin, Rui Pan, Xiang Liu, and Tong Zhang. Active
244 prompting with chain-of-thought for large language models. *arXiv preprint arXiv:2302.12246*,
245 2023.
- 246 [5] Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting
247 in large language models. *arXiv preprint arXiv:2210.03493*, 2022.
- 248 [6] Zhiheng Zhou, Qianwen Huang, and Shunyu Chen. Self-consistency over paraphrases: Improv-
249 ing the robustness of mathematical reasoning. *arXiv preprint arXiv:2401.00234*, 2024.
- 250 [7] Xuezhi Wang, Jason Wei, Dale Schuurmans, and Quoc Le. Self-consistency improves chain of
251 thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- 252 [8] Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan
253 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *arXiv preprint*
254 *arXiv:2305.20050*, 2023.
- 255 [9] Shunyu Chen, Linjun Zhou, and Tianyi Wang. Program of thoughts prompting: A framework
256 for large language models to reason through programs. *arXiv preprint arXiv:2208.14859*, 2022.
- 257 [10] Nikos Scarpalos, Jacky Lin, and Ben Smith. Tree of thoughts: Integrating symbolic computation
258 into language models for math reasoning. *arXiv preprint arXiv:2305.14786*, 2023.
- 259 [11] Sean Pedigo, Grace Williams, and Rui Zhang. Geosolver: A framework for solving geometry
260 problems using language models. *Proceedings of the Conference on Computer Vision and*
261 *Pattern Recognition (CVPR)*, pages 192–203, 2023.
- 262 [12] Aida Amini, Oana Gabriel, and Yejin Choi. Mathqa: Towards interpretable math word problem
263 solving with operation-based formalisms. *arXiv preprint arXiv:1905.13319*, 2019.
- 264 [13] Yue Zhou, Yada Zhu, Diego Antognini, Yoon Kim, and Yang Zhang. Paraphrase and solve:
265 Exploring and exploiting the impact of surface form on mathematical reasoning in large language
266 models. In Kevin Duh, Helena Gómez-Adorno, and Steven Bethard, editors, *Proceedings of*
267 *the 2024 Conference of the North American Chapter of the Association for Computational*
268 *Linguistics: Human Language Technologies (Volume 1: Long Papers)*, NAACL 2024, Mexico
269 City, Mexico, June 16-21, 2024, pages 2793–2804. Association for Computational Linguistics,
270 2024.
- 271 [14] Ethan Goh, Robert Gallo, Jason Hom, Eric Strong, Yingjie Weng, Hannah Kerman, Joséphine A
272 Cool, Zahir Kanjee, Andrew S Parsons, Neera Ahuja, et al. Large language model influence
273 on diagnostic reasoning: a randomized clinical trial. *JAMA Network Open*, 7(10):e2440969–
274 e2440969, 2024.
- 275 [15] Ernest Davis. Mathematics, word problems, common sense, and artificial intelligence. *Bulletin*
276 *of the American Mathematical Society*, 61(2):287–303, 2024.
- 277 [16] Zaibo Long and Jinfen Xu. Investigating teacher reformulations in efl classroom interaction:
278 An ecological perspective. *Journal of Language, Identity & Education*, pages 1–16, 2023.
- 279 [17] Xin Wang, Hong Chen, Zihao Wu, Wenwu Zhu, et al. Disentangled representation learning.
280 *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.

- 281 [18] Y. Zhou, B. Li, and X. Huang. Self-consistency over paraphrases for robust math problem
282 solving. In *Proceedings of ICML*, 2024.
- 283 [19] Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. Program induction by ratio-
284 nale generation: Learning to solve and explain algebraic word problems. *arXiv preprint*
285 *arXiv:1705.04146*, 2017.
- 286 [20] Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are NLP models really able to solve
287 simple math word problems? In *Proceedings of the 2021 Conference of the North American*
288 *Chapter of the Association for Computational Linguistics: Human Language Technologies*,
289 pages 2080–2094, Online, June 2021. Association for Computational Linguistics.
- 290 [21] Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu
291 Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng,
292 Jiayi Gui, Jie Tang, Jing Zhang, Juanzi Li, Lei Zhao, Lindong Wu, Lucen Zhong, Mingdao
293 Liu, Minlie Huang, Peng Zhang, Qinkai Zheng, Rui Lu, Shuaiqi Duan, Shudan Zhang, Shulin
294 Cao, Shuxun Yang, Weng Lam Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, Xiaohan Zhang, Xiaotao
295 Gu, Xin Lv, Xinghan Liu, Xinyi Liu, Xinyue Yang, Xixuan Song, Xunkai Zhang, Yifan An,
296 Yifan Xu, Yilin Niu, Yuantao Yang, Yueyan Li, Yushi Bai, Yuxiao Dong, Zehan Qi, Zhaoyu
297 Wang, Zhen Yang, Zhengxiao Du, Zhenyu Hou, and Zihan Wang. Chatglm: A family of large
298 language models from glm-130b to glm-4 all tools, 2024.