

RCoT: DETECTING AND RECTIFYING FACTUAL INCONSISTENCY IN REASONING BY REVERSING CHAIN-OF-THOUGHT

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language Models (LLMs) have achieved promising performance on arithmetic reasoning tasks by incorporating step-by-step chain-of-thought (CoT) prompting. However, LLMs face challenges in maintaining factual consistency during reasoning, exhibiting tendencies to condition overlooking, question misinterpretation, and condition hallucination over given problems. Existing methods use coarse-grained feedback (e.g., whether the answer is correct) to improve factual consistency. In this work, we propose RCoT (**R**eversing **C**hain-of-**T**hought), a novel method to improve LLMs’ reasoning abilities by automatically detecting and rectifying factual inconsistency in LLMs’ generated solutions. To detect factual inconsistency, RCoT first asks LLMs to reconstruct the problem based on generated solutions. Then fine-grained comparisons between the original problem and the reconstructed problem expose the factual inconsistency in the original solutions. To rectify the solution, RCoT formulates detected factual inconsistency into fine-grained feedback to guide LLMs in revising solutions. Experimental results demonstrate improvements of RCoT over standard CoT, Self-Consistency and Self-Refine across seven arithmetic datasets. Moreover, we find that manually written fine-grained feedback can dramatically improve LLMs’ reasoning abilities (e.g., ChatGPT reaches 94.6% accuracy on GSM8K), encouraging the community to further explore the fine-grained feedback generation methods.

1 INTRODUCTION

Large language models (LLMs) (Brown et al., 2020b; Zhang et al., 2022; Chowdhery et al., 2022; Touvron et al.) have showcased strong reasoning capabilities using chain-of-thought (CoT) (Wei et al., 2022; Chowdhery et al., 2022; Fung et al., 2022), where LLMs are prompted to generate intermediate steps before the final answer. Despite the impressive performance of CoT prompting across various reasoning tasks (Dua et al., 2019; Miao et al., 2021; Cobbe et al., 2021a; Yu et al., 2020; Bhagavatula et al., 2019; Talmor et al., 2019), LLMs still struggle to maintain factual consistency in reasoning. Specifically, each reasoning problem usually consists of several conditions and a question, and LLMs exhibit tendencies to hallucinate, overlook conditions and misinterpret questions (Golovneva et al., 2022).

While previous research has proposed various methods to enhance Chain-of-Thought performance (Shum et al., 2023; Fu et al., 2022; Diao et al., 2023; Shum et al., 2023; Zhou et al., 2022; Wang et al., 2022; Gao et al., 2023; Chen et al., 2022; Weng et al., 2022; Paul et al., 2023; Shinn et al., 2023), there remains a noticeable absence of explicit studies addressing the issue of factual inconsistency. The most relevant work is probably Self-Verification (Weng et al., 2022), which verifies answers by swapping conditions and answers. However, it can only tell whether answers are correct and fail to give fine-grained feedback on factual inconsistency to guide LLMs in revising solutions. Figure 1 shows an instance of factual inconsistent solutions generated by ChatGPT, where the critical condition of "tomorrow, 10/16/1924" is mistakenly overlooked. Despite the incorrect solution, ChatGPT could be rectified if we explicitly point out that it overlooked the "tomorrow, 10/16/1924." condition. Motivated by this observation, we introduce Reverse Chain-of-Thought (RCoT), a novel method for detecting and rectifying condition hallucination, overlooking and question misinterpretation (Figure 1) in CoT. More examples could be found in Appendix B.1. The schematics

of our proposed method are illustrated in Figure 2. To detect factual inconsistency, RCoT begins with reconstructing new problems from the original solutions. Next, RCoT detects inconsistencies by conducting a fine-grained comparison between the reconstructed problems and the original problems. To rectify inconsistencies, the detected inconsistencies will be formulated as fine-grained feedback to guide the LLMs in revising their solutions. Fine-grained feedback can not only be used to revise solutions but can also offer more interpretability of reasoning errors, allowing humans to inspect solutions easier.

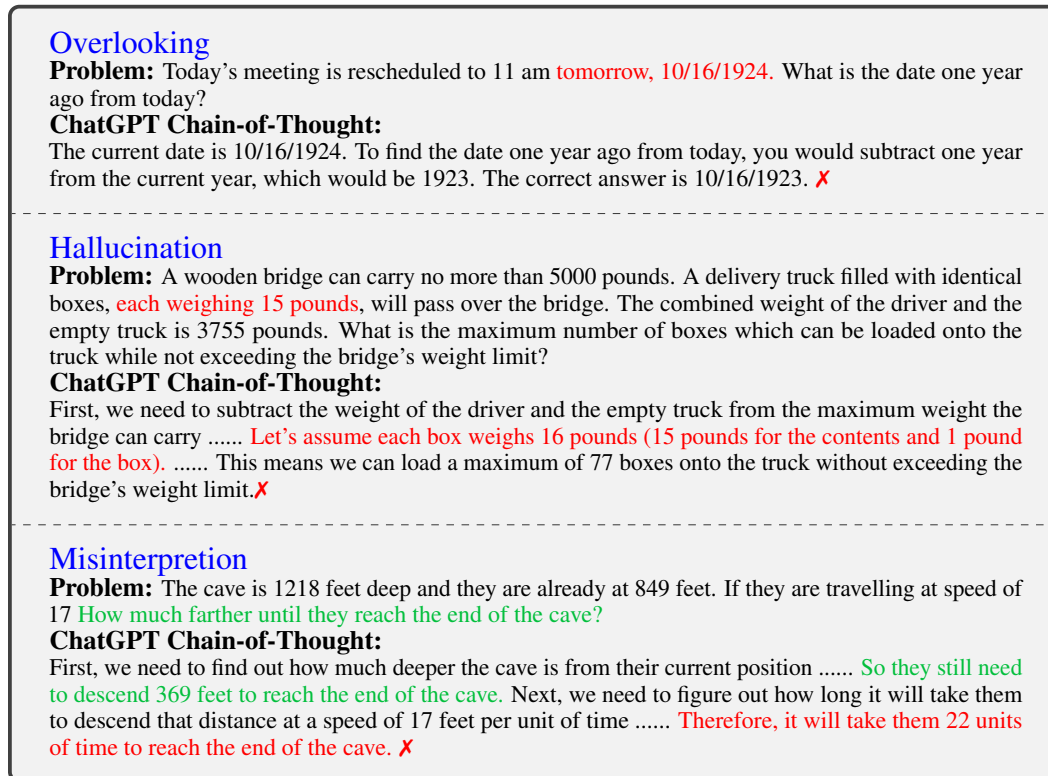


Figure 1: The examples of factual inconsistency. **Overlooking:** The problem mentions that tomorrow is 10/16/1924, whereas LLMs overlook the condition of "tomorrow", thinking that today is 10/16/1924. **Hallucination:** The problem mentions that each box weighs 15 pounds. However, LLMs assume each box weighs 16 pounds, which contradicts real conditions. **Misinterpretation:** LLMs have correctly answered the question in the problem, but it continues answering "how long", resulting in a wrong answer. This is a typical misinterpretation of the problem.

We evaluate RCoT on seven arithmetic reasoning datasets, including GSM8k (Cobbe et al., 2021b), AQuA (Ling et al., 2017), SVAMP (Patel et al., 2021), AddSub (Hosseini et al., 2014), ASDiv (Miao et al., 2021), Date (Srivastava et al., 2022) and SingelEq (Koncel-Kedziorski et al., 2016). Experimental results demonstrate the effectiveness of RCoT, outperforming competitive baselines in both zero-shot and few-shot settings. In-depth analysis and human evaluation suggest that fine-grained feedback on factual inconsistency is crucial for LLMs to revise solutions for arithmetic problems. For example, ChatGPT could achieve 94.6% accuracy on GSM8k with manually written fine-grained feedback. Moreover, we conduct comprehensive ablation studies examining the impact of individual modules. Our findings encourage the community to further explore detecting and rectifying factual inconsistency to enhance LLMs’ reasoning ability.

Our contributions are summarized as follows:

- We propose a novel prompting method, Reversing Chain-of-Thought (RCoT) to effectively detect and rectify the factual inconsistency of LLMs in arithmetic reasoning, focusing on overlooked, hallucinated conditions and misinterpreted questions. RCoT outperforms competitive baseline models consistently across seven arithmetic reasoning tasks.

- Prompting with fine-grained feedback on factual inconsistency shows encouraging results on improving LLMs’ reasoning abilities. Though automatically generated feedback by RCoT shows consistent improvement compared with standard CoT, we find that human-written ground-truth feedback can further improve the LLMs’ reasoning ability (e.g., ChatGPT reaches 94.6% accuracy on GSM8k). The gap between RCoT’s feedback and human-written feedback encourages the community to further explore the automatic generation of fine-grained feedback.
- RCoT offers more interpretability to the reasoning errors with fine-grained feedback on factual inconsistency, allowing humans to inspect solutions easier.

2 RELATED WORK

Language Model for Reasoning Reasoning ability is a critical skill to solve complex problems, such as arithmetic reasoning (Koncel-Kedziorski et al., 2016; Roy & Roth, 2016; Miao et al., 2021; Cobbe et al., 2021a; Dua et al., 2019), logical reasoning (Yu et al., 2020), commonsense reasoning (Bhagavatula et al., 2019; Talmor et al., 2019; Zellers et al., 2018), and tabular reasoning (Zhu et al., 2021). Recently, Large Language Models (e.g., GPT3 (Brown et al., 2020a), ChatGPT, PaLM (Chowdhery et al., 2022) and LLaMA (Touvron et al.)) have demonstrated promising reasoning capability with Chain-of-Thought methods. However, large language models exhibit tendencies to generate intermediate steps that are factually inconsistent, rendering them incapable of solving complex problems requiring multi-step reasoning. In this work, we focus on the detection and rectification of factually inconsistent errors in the intermediate reasoning steps, including question misinterpretation, condition hallucination and condition overlooking.

Prompt Engineering Some prompting methods can elicit useful knowledge in large language models to better solve complex tasks, two representative examples of which are In-context Learning (Brown et al., 2020a) and Chain-of-Thought (Wei et al., 2022). In-Context Learning encourages the language models to learn from a few input-output examples as prompts (Liu et al., 2022; Rubin et al., 2021; Min et al., 2022). Chain-of-Thought prompting improves reasoning performance by prompting LLMs to think of intermediate steps. Inspired by the promising performance of CoT, many methods have explored how to further improve standard CoT. Least-to-most (Zhou et al., 2022) prompting proposes to decompose a complex problem into a series of subproblems and solve them sequentially. Self-Consistency prompting (Wang et al., 2022) improves performance through majority voting on multiple solutions. Similarly, Complex CoT (Fu et al., 2022) emphasizes the importance of prompt complexity and selects the most complex examples as prompts. Auto-CoT (Shum et al., 2023) is proposed to reduce the workload of manual labeling. Active Prompting (Diao et al., 2023) selects the most uncertain questions as demonstration examples to further improve performance. However, these methods fail to address the factual inconsistency problem. Probably the most relevant work are Self-Verification (Weng et al., 2022), REFINER (Paul et al., 2023), and Reflexion (Shinn et al., 2023). These approaches focus on correcting LLMs outputs. However, Self-Verification can only generate binary feedback and fail to get fine-grained feedback, REFINER needs externally trained models, and Reflexion requires environmental feedback, which cannot be easily obtained in arithmetic reasoning. Compared to these methods, RCoT entirely relies on the LLM itself to generate fine-grained feedback on factual consistency.

Reverse Engineering. RCoT is inspired by the concept of Reverse Engineering, which has various applications in machine learning research. (Fredrikson et al., 2014) proposes a reverse method for linear models to evaluate models’ privacy safety. (Fredrikson et al., 2015) introduces a model inversion method for shallow neural networks, which can reconstruct the face information. (Geva et al., 2022) unveils the internal prediction construction process of Transformer-based language models by reverse engineering the operations of the feed-forward network (FFN) layers. Inverting model hyperparameters is another application of reverse engineering techniques. (Bhagavatula et al., 2019) reverses network parameters by repeatedly requesting the predicted label from the target model. (Tramèr et al., 2016) develops an avatar method to estimate training data and model architectures, while (Oh et al., 2019) trains a set of white-box models to estimate model hyperparameters. (Hua et al., 2018) estimates both the structure and the weights of a CNN model on a hardware accelerator from information leaks of memory access patterns. Different from their goal of opening up the

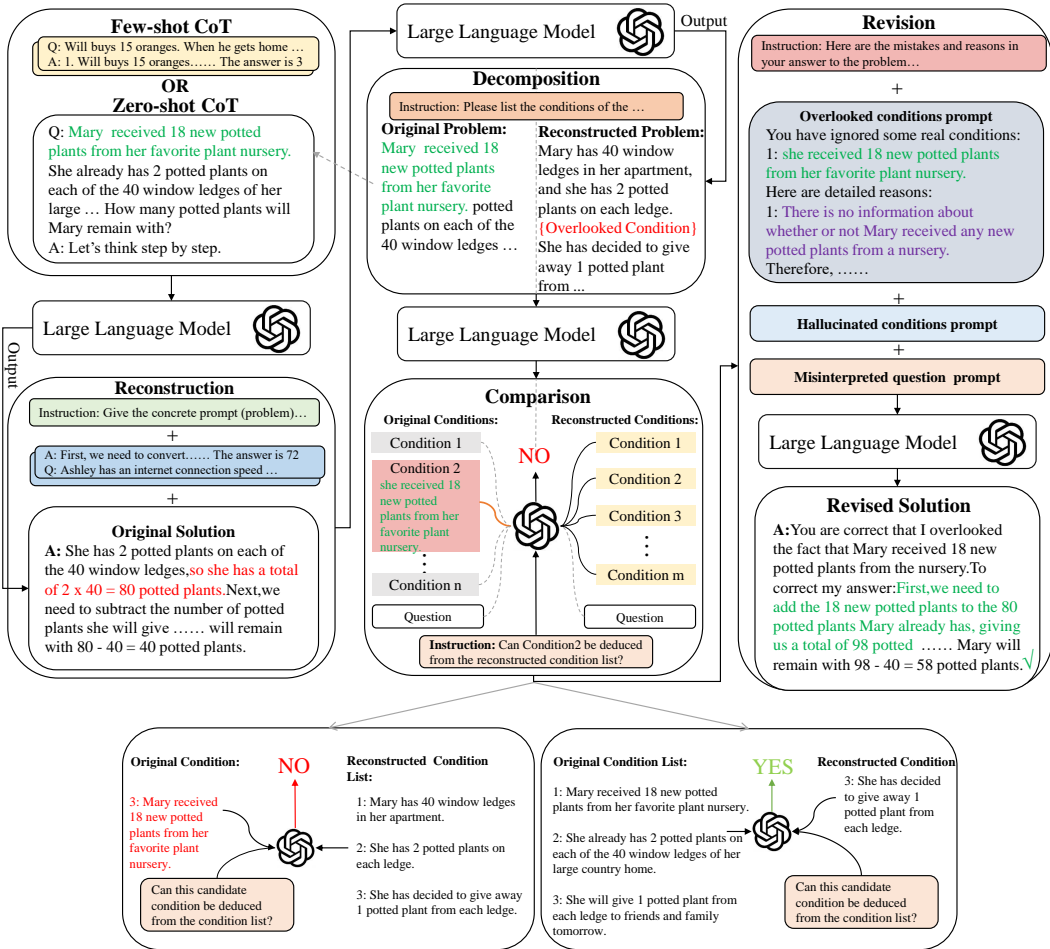


Figure 2: The framework of RCoT. (1) **Reconstruction**: Ask LLMs to reconstruct the problem according to the original solution with instruction and demonstration examples. (2) **Decomposition**: Decomposing the original problem and reconstructed problem into fine-grained condition lists. (3) **Comparison**: Compare both lists of sub-conditions and questions to verify whether there are hallucinations, overlookings and misinterpretations. (4) **Revision**: Gathering all factual inconsistencies into fine-grained feedback to instruct LLMs to revise solutions.

black-box of deep learning models, our work focuses on automatically detecting and rectifying factual inconsistencies that appeared in the solutions generated by LLMs.

3 REVERSING CHAIN-OF-THOUGHT (RCoT)

We propose RCoT for detecting and rectifying factual inconsistency (i.e., condition hallucinations, overlookings, and question misinterpretation) in CoT to enhance LLMs’ reasoning ability. Specifically, given a complex reasoning problem Q and original solution c generated by the LLM, we first ask LLMs to detect factual inconsistency: (i) **Problem Reconstruction**: Reconstruct the problem \hat{Q} based on the generated solution c . (ii) **Fine-grained Comparison**: Conduct a fine-grained comparison between the original problem Q and the reconstructed problem \hat{Q} to detect condition hallucinations, overlookings, and question misinterpretation. Then we rectify LLMs using detected factual inconsistency: (iii) **Fine-grained Feedback and Revision**: The fine-grained comparison reveals the factual inconsistency in original solutions. The detected factual inconsistencies are formulated into fine-grained feedback to guide LLMs in revising their solution accordingly. The overall schematic illustrations of our proposed approach are illustrated in Figure 2, and an example of RCoT is shown in Appendix B.3.

3.1 PROBLEM RECONSTRUCTION

Intuitively, if the generated step-by-step solution of an arithmetic problem is logically and factually correct and complete, it is more likely for a human to infer what is the original problem. Similarly, we ask the LLM to reconstruct the problem to get \hat{Q} based on its own solution c , in order to verify whether it truly understands the problem. We manually write instructions and in-context examples as the reconstruction prompt. We find that the factual inconsistencies such as *condition hallucinations* (e.g., the LLM uses conditions that are not mentioned in the problem Q), *condition overlookings* (e.g., the LLM overlooks some important conditions in the problem Q), and *question misinterpretations* (e.g., the LLM misunderstand the question of Q) can be effectively exposed by comparing the reconstructed problem \hat{Q} with the original problem Q (§ 3.2), as shown in Figures 9, 6, and 15 in Appendix B.1, respectively. The prompt template can be found in Figure 21.

3.2 FINE-GRAINED COMPARISON

To detect condition hallucinations and overlookings, as well as question misinterpretations in the solution c from the reconstructed problem \hat{Q} , a naive approach is to ask the LLM to directly compare Q with \hat{Q} . However, such comparisons usually fail to produce high-quality detection results (Figure 3), which is unsurprising because Q and \hat{Q} contain rich information, and the coarse-grained comparison will inevitably ignore some vital information, causing a sub-optimal result. Therefore, we use fine-grained step-by-step comparisons to improve the detection quality. All prompt templates are shown in Figure 21. The process is as follows:

Problem Decomposition. Q and \hat{Q} are unstructured texts, which are hard to be compared in an organized manner. To overcome this issue, we ask the LLM to decompose the problem into a list of conditions $L_Q = [L_Q^1, \dots, L_Q^m]$, $L_{\hat{Q}} = [L_{\hat{Q}}^1, \dots, L_{\hat{Q}}^n]$. The structured condition list will then be used in fine-grained comparison.

Condition Comparison To find the differences between Q and \hat{Q} , we first check whether their condition lists L_Q and $L_{\hat{Q}}$ are the same. Specifically, the LLM is required to answer whether each L_Q^i can be inferred from $L_{\hat{Q}}$. If L_Q^i cannot be inferred from $L_{\hat{Q}}$, then L_Q^i is either (1) overlooked in the solution or (2) hallucinated by the LLM as a different condition. Similarly, we ask the LLM to tell whether $L_{\hat{Q}}^j$ can be inferred from L_Q for every j . If $L_{\hat{Q}}^j$ cannot be inferred from L_Q , then $L_{\hat{Q}}^j$ is hallucinated. Apparently, we need to conduct comparisons for nm times in total.

Question Comparison The LLM sometimes will also misinterpret the question (Figure 1). Therefore, we also ask LLM to compare the questions being asked in Q and \hat{Q} . If LLMs find the two questions are different, then LLMs misinterpret the question in their solutions. This comparison only needs to be done once since there is one question per problem in most cases.

After these comparisons, we detect hallucinated conditions, overlooked conditions, and misinterpreted of questions. We then use them to formulate our fine-grained feedback to guide the LLM in revising its solution.

3.3 FINE-GRAINED FEEDBACK AND REVISION

We assume the original solution is correct if we do not detect any factual inconsistency through fine-grained comparison. On the contrary, we formulate fine-grained feedback to guide the LLM in revising its solution if any factual inconsistency is detected. Specifically, the fine-grained feedback will first state that the solution is incorrect, then list the detected factual inconsistency, and finally ask the LLM to revise its solution. Figure 21 shows the template we use to formulate the feedback. We take the answer of the revised solution as the final output for evaluation.

4 EXPERIMENT

Our extensive experiments aim to show that (1) RCoT benefits arithmetic reasoning by automatically detecting and rectifying condition hallucination and overlooking, and question misinterpretation;

(2) fine-grained feedback of factual consistency is critical for LLMs to self-revise the solution. (3) fine-grained comparison is essential for constructing high-quality fine-grained feedback.

4.1 EXPERIMENT SETTING

We used closed-source ChatGPT and open-source LLaMA-13B-Chat (Touvron et al.) as the backbone LLMs for solution generation and set the temperature to 0 to improve reproducibility. We evaluate RCoT on seven arithmetic datasets with different difficulties, including GSM8k (Cobbe et al., 2021b), AQuA (Ling et al., 2017), SVAMP (Patel et al., 2021), AddSub (Hosseini et al., 2014), ASDiv (Miao et al., 2021), Date (Srivastava et al., 2022) and SingelEq (Koncel-Kedziorski et al., 2016). Due to the high time cost of API calls, we do not use the whole test set but randomly sample test sub-sets. To reduce the randomness caused by test set sampling and make our results more convincing, we sample three test sub-sets that each contains 256 inputs. We report the average accuracy with deviation on the three test sub-sets. For the dataset that has less than 256 test inputs, we still evaluate three times since ChatGPT’s outputs may change and report the average accuracy with deviation. A detailed description of each dataset is shown in Appendix B.5.

We consider both zero-shot and few-shot settings. For the zero-shot setting, we add the prompt "Let’s think step by step" to encourage LLMs to think intermediate steps but without any demonstration examples (Kojima et al., 2022). For the few-shot setting, we use four-shot CoT prompts that consist of problems, solutions, and final answers.

We compare our method with five baselines: (1) **Chain-of-thought (CoT)** (Wei et al., 2022) (2) **Active-Prompting** (Diao et al., 2023), a method that selects the most uncertainty problems as demonstration examples. (3) **Double-Check** asks LLMs to check their answers but does not point out whether the answer is correct. In our experiment, we use the prompt "You should double-check your answers". (4) **Self-Consistency (SC)** (Wang et al., 2022) through majority voting on multiple solutions to improve the performance. (5) **Self-Refine** (Madaan et al., 2023) uses iterative feedback and refinement to revise the answer. We use Tiktoken from Openai to calculate the cost of average tokens.¹

4.2 RCoT BENEFITS ARITHMETIC REASONING

Table 1 shows the results of RCoT on seven arithmetic datasets. Our method consistently outperforms the standard CoT and the double-check methods in the zero-shot setting. Moreover, LLMs benefit more from our method on more challenging tasks that require complex reasoning. For example, the AQuA dataset contains diverse problems, and the Date dataset requires multi-hop reasoning and common sense date knowledge. Both ChatGPT and LLaMA achieve lower accuracy scores on AQuA and Date (51.3% and 66.7% for ChatGPT and 27.2% and 52.4% for LLaMA) among all seven tasks. Meanwhile, we observe that our method helps LLMs improve by apparent margins on AQuA and Date (4.1%, 5.0% and 4.7%, 2.9% for ChatGPT and LLaMA), the highest gains in all seven tasks. Our method also remains effective for easier tasks. For example, RCoT enhances the performance of the SVAMP dataset, which contains problems that usually only need one-step calculation, by 2.8% and 2.5%. Moreover, we also observe greater improvements from our method on ChatGPT than LLaMA, potentially due to the stronger abilities of ChatGPT to detect and correct errors.

We can observe similar results in the few-shot setting to the zero-shot setting. Although selecting the most uncertain problems for LLMs as demonstration examples is helpful for reasoning (Diao et al., 2023), RCoT still improves the accuracy. It is worth noting that the performance of Double-Check method in the few-shot CoT setting decreases immensely. On the AQuA and GSM8K datasets, its performance drops by 27.0% and 4.0%, suggesting that few-shot examples may increase the risk of revising correct solutions to the incorrect ones. LLaMA exhibits a lower degree of susceptibility compared to ChatGPT.

We also compared RCoT with other stronger baselines (i.e., Self-Consistency, SC for short, and Self-Refine). Specifically, We conducted 30 trials per problem for SC and 3 trials per problem for RCoT in the zero-shot setting (set temperature to 0.7 Wang et al. (2022)), which uses similar costs. Due to the extremely high cost, we do not experiment with the few-shot setting and leave it as our future work. We set max attempt to 5 for Self-Refine. Table 2 has shown the results. RCoT could achieve

¹<https://github.com/openai/tiktoken>

Table 1: Average accuracy and standard deviation on seven arithmetic reasoning datasets. **Bold** denotes the best result. **Green**: The performance improvement compared with Standard CoT and Active-Prompting in Zero-shot and Few-shot settings, respectively. * denotes the LLM that uses Manual-CoT. - denotes that Active-Prompting (Diao et al., 2023) does not support the dataset in their source codes.

Model	Method	Arithmetic						
		GSM8K	AQuA	AddSub	Date	SingleEq	ASDiv	SVAMP
UL2-20B*	Standard	4.4	23.6	18.2	14.4	20.2	16.9	12.5
LaMDA-137B*	Standard	14.3	20.6	51.9	26.8	58.7	46.6	37.5
Text-davinci-002*	Standard	46.9	24.8	81.3	52.1	86.6	71.3	68.9
Zero-shot CoT								
ChatGPT	Standard	79.0 \pm 0.95	51.3 \pm 0.6	85.2 \pm 1.2	66.7 \pm 1.4	90.3 \pm 0.6	84.3 \pm 0.4	76.7 \pm 4.1
	+Double-Check	79.3 \pm 2.1	42.7 \pm 0.6	85.6 \pm 1.2	60.5 \pm 6.5	88.8 \pm 0.8	82.8 \pm 1.4	77.6 \pm 2.0
	+RCoT	82.0\pm0.3	55.5\pm0.8	87.1\pm1.1	71.7\pm1.3	91.4\pm0.8	86.0\pm0.3	79.6\pm4.1
Δ RCoT-Standard (Acc)		(+3.1 \pm 0.6)	(+4.1 \pm 0.2)	(+1.8 \pm 0.1)	(+5.0 \pm 0.4)	(+1.1 \pm 0.2)	(+1.7 \pm 0.3)	(+2.8 \pm 0.2)
LLaMA-13B-Chat	Standard	36.9 \pm 0.8	27.2 \pm 0.0	66.7 \pm 0.5	52.4 \pm 1.5	62.6 \pm 2.6	52.2 \pm 3.7	38.6 \pm 1.1
	+Double-Check	35.6 \pm 1.1	24.8 \pm 0.0	62.0 \pm 0.7	27.0 \pm 0.9	62.1 \pm 3.2	53.2 \pm 3.6	41.1 \pm 0.2
	+RCoT	39.8\pm0.8	31.9\pm0.0	67.4\pm0.5	55.3\pm2.0	63.5\pm2.1	53.0\pm3.7	41.1\pm0.8
Δ RCoT-Standard (Acc)		(+2.9 \pm 0.4)	(+4.7 \pm 0.0)	(+0.7 \pm 0.3)	(+2.9 \pm 1.0)	(+0.9 \pm 0.5)	(+0.8 \pm 0.0)	(+2.5 \pm 0.4)
Few-shot CoT								
ChatGPT	Active-Prompting	81.8 \pm 0.6	53.3 \pm 0.6	87.2 \pm 1.2	-	91.7 \pm 0.4	87.9 \pm 0.8	82.5 \pm 0.6
	+Double-Check	77.8 \pm 0.7	26.3 \pm 0.5	86.0 \pm 1.6	-	91.5 \pm 0.2	85.7 \pm 2.4	82.2 \pm 0.8
	+RCoT	84.6\pm0.6	57.1\pm0.3	88.2\pm1.5	-	93.0\pm0.8	89.3\pm0.5	84.9\pm1.3
Δ RCoT-Active (Acc)		(+2.7 \pm 0.1)	(+3.7 \pm 0.9)	(+1.0 \pm 0.4)	-	(+1.2 \pm 0.4)	(+1.4 \pm 0.5)	(+2.3 \pm 1.0)
LLaMA-13B-Chat	Active-Prompting	37.9 \pm 0.6	29.1 \pm 0.0	68.4 \pm 0.7	-	67.9 \pm 2.2	53.3 \pm 0.6	49.4 \pm 0.4
	+Double-Check	36.2 \pm 0.1	23.2 \pm 0.0	61.9 \pm 2.1	-	64.9 \pm 1.3	50.3 \pm 3.5	47.4 \pm 0.8
	+RCoT	40.1\pm0.4	30.7\pm0.0	68.8\pm0.9	-	68.1\pm2.3	53.6\pm0.4	51.2\pm0.3
Δ RCoT-Active (Acc)		(+2.1 \pm 0.3)	(+1.6 \pm 0.0)	(+0.4 \pm 0.3)	-	(+0.2 \pm 0.1)	(+0.3 \pm 0.2)	(+1.8 \pm 0.2)

comparable performance to SC at nearly one-third of the cost (e.g., AddSub, SingleEq, SVAMP) and even outperforms SC on the GSM8K dataset. However, the performance significantly drops on AQuA and Date datasets. That is because there are multiple-choice tasks, making it exceedingly simple for the model to approximate the answer by employing multiple guesses with incorrect logical steps. Combining RCoT with SC, our method can further improve the performance and surpass all baselines, reaching a high accuracy of 84.5% across seven arithmetic datasets. Our experiments demonstrate the same conclusion as Wang et al. (2022) that Self-Refine is not good at arithmetic reasoning. It’s worth noting that Self-Refine achieves the highest accuracy on SingleEq and AddSub datasets. Nevertheless, the improvement does not come from refinement but the usage of code in the Self-Refine implementation, reducing a large number of calculation errors. The real improvements brought by refinement are actually 0.8 and 0.4 in the AddSub and SingleEq datasets, respectively. Another phenomenon is that self-refine does not bring more token cost even if we give it more refinement budget. This is because self-refine tends to state that the solution is correct after the second refinement.

4.3 FINE-GRAINED FEEDBACK IS CRITICAL FOR SOLUTION REVISION

The success of our method comes from fine-grained feedback that points out detailed factual inconsistency (condition hallucination and overlooking, and question misinterpretation). In this section, we show that coarse-grained feedback will lead to worse performance to prove the necessity of fine-grained feedback. We replace our fine-grained feedback with two kinds of coarse-grained feedback: (1) w/o reasons: we do not tell LLMs the detected factual inconsistency by RCoT and only give a high-level judgment. Therefore, if RCoT detects no factual inconsistency, we take the original solution as the final output for evaluation. Otherwise, we use the prompt "Your answer is wrong. You should double-check your answer" to guide LLMs in revising solutions. (2) w/o judgment+reasons (i.e., Double-Check): We further remove the high-level judgment from the prompts. Therefore, we

Table 2: Average accuracy on seven arithmetic reasoning datasets among Self-Consistency (Wang et al., 2022), RCoT and Self-Refine (Madaan et al., 2023). **Bold** denotes the best result.

Method	GSM8K	AQuA	AddSub	Date	SingleEq	ASDiv	SVAMP	Avg Acc	Avg Tokens
SC (30 trials per problem)	81.6	70.8	88.6	80.0	92.9	90.2	80.4	83.5	5615.0
RCoT (1 trial per problem)	82.0	56.3	87.2	71.9	92.4	86.3	79.7	79.4	1831.0
RCoT (3 trials per problem)	83.2	72.8	89.8	78.9	93.8	91.8	81.2	84.5	5453.3
attempt 0	79.1	45.2	90.6	51.3	97.6	83.5	75.2	74.7	190.2
attempt 1	80.7	49.2	91.4	52.7	98.0	84.3	76.8	76.1	3108.4
attempt 2	80.7	49.2	91.4	52.7	98.0	84.3	76.8	76.1	3324.9
Self-Refine									
attempt 3	80.7	49.2	91.4	52.7	98.0	84.3	76.8	76.1	3359.6
attempt 4	80.7	49.2	91.4	52.7	98.0	84.3	76.8	76.1	3367.7
attempt 5	80.7	49.2	91.4	52.7	98.0	84.3	76.8	76.1	3367.7

Table 3: The performance of RCoT using fine-grained feedback and coarse-grained feedback. **w/o reasons** : remove explanations of specific mistakes from the original fine-grained feedback. The prompt becomes "Your answer is wrong. You should double-check your answer.". **w/o judgment+reasons**: further remove the high-level judgment. The prompt becomes "You should double-check your answer." **Red**: The performance drops compared with RCoT method.

Method	GSM8K	AQUA	SVAMP
Standard CoT	79.0	51.3	76.7
RCoT(ours)	82.0	55.5	79.6
- w/o reasons	80.0 (-2.0)	52.3 (-3.2)	78.9 (-0.7)
- w/o judgment+reasons	79.3 (-2.7)	42.7 (-12.8)	77.6 (-2.0)

Table 4: The performance without question comparison and condition comparison, as well as the performance with coarse-grained comparison. **coarse-grained**: We directly ask LLMs to compare the original problem with the reconstructed problem. Results show that (1) fine-grained comparison is important to get fine-grained feedback, and (2) both question comparison and condition comparison are important in the fine-grained comparison. **Red**: The performance drops compared with the RCoT method.

Method	GSM8K	AQUA	SVAMP
Standard CoT	79.0	51.3	76.7
RCoT	82.0	55.5	79.6
- w/o question comparison	80.9 (-1.1)	54.6 (-0.9)	79.2 (-0.4)
- w/o condition comparison	80.1 (-1.9)	53.5 (-2.0)	78.1 (-1.5)
RCoT (Corase-grained)	74.2 (-7.8)	49.6 (-5.9)	76.1 (-3.5)

always use "You should double-check your answer" to guide LLMs in revising solutions regardless of the detection results of RCoT.

Table 3 shows the results on SVAMP(easy), GSM8K(medium), and AQuA (hard) datasets. We can see consistent performance drops when we remove detected factual inconsistency and only keep a high-level judgment, showing the effectiveness of fine-grained feedback. Moreover, we can observe that further removing judgment will make the performance even worse than standard CoT. This is not surprising because LLMs may mistakenly revise the correct solution to the incorrect one. Appendix B.6 shows an example of RCoT and Double-Check, where we can see that RCoT could help the LLM correct the solution, but Double-Check couldn't.

To further show the power of fine-grained feedback, we conduct a human evaluation. Specifically, we go through the generated solution and write fine-grained feedback on factual inconsistency by ourselves. Amazingly, the LLM could reach 94.6% accuracy on the GSM8K dataset and could only reach 86.3% if we remove explanations of factual inconsistencies from our feedback (i.e., the same setting as "w/o reasons" in Table 3). Appendix B.2 shows examples of manually written and RCoT-generated feedback. This result shows the same observation of Table 3 and reveals the strong power of fine-grained feedback. Since RCoT still has a gap (12.6% accuracy gap) compared with humans, we encourage the community to further explore fine-grained feedback generation.

4.4 FINE-GRAINED COMPARISON LEADS TO FINE-GRAINED FEEDBACK

To get fine-grained feedback, conditions and questions are compared in a fine-grained manner in RCoT. A simpler way is to ask LLMs to generate fine-grained feedback directly by comparing the original problems and reconstructed problems. Appendix A has illustrated that LLMs can

reconstruct questions well when the reasoning maintains factual consistency but not vice versa. Table 4 shows that coarse-grained comparison will cause a significant accuracy drop (even worse than standard CoT), which suggests that it fails to generate high-quality feedback (Figure 3). Therefore, problem decomposition and fine-grained comparison are essential. We also show that both condition comparison and question comparison are important. Removing either of them will lead to worse performance. This is because LLMs may hallucinate/overlook conditions (Figure 1) as well as misinterpret questions (Figure 1).

Original Problem: Rebecca wants to split a collection of eggs into groups of 3. Rebecca has 99 bananas 9 eggs and 27 marbles. How many groups will be created?

Reconstructed Problem: Rebecca has 9 eggs 99 bananas and 27 marbles. If we can group items into sets of 3 how many groups can we make in total?

Coarse-Grained Comparison:
Yes, the reconstructed problem does not add new conditions or lack any conditions that were present in original problem. Original problem and reconstructed problem essentially ask the same question, with the only difference being the way the information is presented. Both problems ask how many groups of 3 items can be made from a collection of 9 eggs, 99 bananas, and 27 marbles. Therefore, reconstructed problem does not introduce any new conditions or leave out any conditions from the original problem. ✗

Fine-Grained Comparison:
No, the candidate condition "Group items into sets of 3" cannot be deduced from the given condition list as it is not applicable to all items in the condition list. The condition only applies to the eggs that Rebecca wants to split into groups of 3. It does not apply to the bananas or the marbles. Additionally, the candidate condition implies that all items must be grouped into sets of 3, but there is no indication in the context that Rebecca wants to group all items into sets of 3. She only wants to split the eggs into groups of 3. ✓

Figure 3: This is an example showing the failure of the coarse-grained comparison. Red: The inconsistent condition between the original problem and the reconstructed problem.

4.5 ADDITIONAL QUANTITATIVE ANALYSIS OF RCoT

To explore the effectiveness of RCoT, we construct further quantitative analysis on 100 problems the ChatGPT answered incorrectly. We manually divided these problems into four categories: condition hallucination, condition overlooking, question misinterpretation and other errors. The statistical results are shown in Table 5. We find that RCoT is better at detecting overlooking and misinterpretation errors than hallucination errors.

Table 5: **Found/Not Found:** RCoT can or cannot find the reasons for errors. **Other errors:** such as computation error, logical error and so on.

Type	Found	Not Found	total
Overlooking	5	1	6
Hallucinating	16	15	31
Misinterpreting	5	3	8
Other errors	0	55	55

5 CONCLUSION

In this paper, we propose RCoT, a method that enables LLMs to detect and rectify factual inconsistency automatically to improve LLMs’ reasoning abilities. RCoT detects factual inconsistency through fine-grained comparison between the reconstructed problems and original questions, and then asks LLMs to rectify inconsistencies through fine-grained feedback. Experimental results on seven arithmetic reasoning datasets demonstrate the effectiveness of RCoT. Our experiments also show encouraging results of LLMs’ reasoning abilities with the help of manually written fine-grained feedback, encouraging the community to further explore fine-grained feedback generation. RCoT could, in principle, be applied to other tasks requiring CoT solutions. We discuss the limitations and future work in Appendix C.

6 REPRODUCIBILITY STATEMENT

The supplementary material includes the code for all experiments and their corresponding running scripts. The dataset (GSM8K, AQuA, AddSub, SingleEq, Date, ASDiv and SVAMP) can be easily accessible on the HuggingFace website or from their official repositories. We explain all the experimental details (temperature, dataset size and so on) in Section 4.1.

REFERENCES

- Chandra Bhagavatula, Ronan Le Bras, Chaitanya Malaviya, Keisuke Sakaguchi, Ari Holtzman, Hannah Rashkin, Doug Downey, Scott Wen-tau Yih, and Yejin Choi. Abductive commonsense reasoning. *arXiv preprint arXiv:1908.05739*, 2019.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020a.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020b.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *arXiv preprint arXiv:2211.12588*, 2022.
- 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. *arXiv preprint arXiv:2204.02311*, 2022.
- 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. *arXiv preprint arXiv:2110.14168*, 2021a.
- 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. *arXiv preprint arXiv:2110.14168*, 2021b.
- Shizhe Diao, Pengcheng Wang, Yong Lin, and Tong Zhang. Active prompting with chain-of-thought for large language models. *arXiv preprint arXiv:2302.12246*, 2023.
- Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. Drop: A reading comprehension benchmark requiring discrete reasoning over paragraphs. *arXiv preprint arXiv:1903.00161*, 2019.
- Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. Model inversion attacks that exploit confidence information and basic countermeasures. In *Proceedings of the 22nd ACM SIGSAC conference on computer and communications security*, pp. 1322–1333, 2015.
- Matthew Fredrikson, Eric Lantz, Somesh Jha, Simon Lin, David Page, and Thomas Ristenpart. Privacy in pharmacogenetics: An end-to-end case study of personalized warfarin dosing. In *23rd {USENIX} Security Symposium ({USENIX} Security 14)*, pp. 17–32, 2014.
- Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and Tushar Khot. Complexity-based prompting for multi-step reasoning. *arXiv preprint arXiv:2210.00720*, 2022.
- Yi R Fung, Tuhin Chakraborty, Hao Guo, Owen Rambow, Smaranda Muresan, and Heng Ji. Norm-sage: Multi-lingual multi-cultural norm discovery from conversations on-the-fly. *arXiv preprint arXiv:2210.08604*, 2022.
- Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. Pal: Program-aided language models. In *International Conference on Machine Learning*, pp. 10764–10799. PMLR, 2023.

- Mor Geva, Avi Caciularu, Kevin Ro Wang, and Yoav Goldberg. Transformer feed-forward layers build predictions by promoting concepts in the vocabulary space. *arXiv preprint arXiv:2203.14680*, 2022.
- Olga Golovneva, Moya Chen, Spencer Poff, Martin Corredor, Luke Zettlemoyer, Maryam Fazel-Zarandi, and Asli Celikyilmaz. Roscoe: A suite of metrics for scoring step-by-step reasoning. *arXiv preprint arXiv:2212.07919*, 2022.
- Mohammad Javad Hosseini, Hannaneh Hajishirzi, Oren Etzioni, and Nate Kushman. Learning to solve arithmetic word problems with verb categorization. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 523–533, Doha, Qatar, October 2014. Association for Computational Linguistics. doi: 10.3115/v1/D14-1058. URL <https://aclanthology.org/D14-1058>.
- Weizhe Hua, Zhiru Zhang, and G Edward Suh. Reverse engineering convolutional neural networks through side-channel information leaks. In *Proceedings of the 55th Annual Design Automation Conference*, pp. 1–6, 2018.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35: 22199–22213, 2022.
- Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate Kushman, and Hannaneh Hajishirzi. MAWPS: A math word problem repository. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 1152–1157, San Diego, California, June 2016. Association for Computational Linguistics. doi: 10.18653/v1/N16-1136. URL <https://aclanthology.org/N16-1136>.
- Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 158–167, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi: 10.18653/v1/P17-1015. URL <https://aclanthology.org/P17-1015>.
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. What makes good in-context examples for GPT-3? In *Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures*, pp. 100–114, Dublin, Ireland and Online, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.deelio-1.10. URL <https://aclanthology.org/2022.deelio-1.10>.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. *arXiv preprint arXiv:2303.17651*, 2023.
- Shen-Yun Miao, Chao-Chun Liang, and Keh-Yih Su. A diverse corpus for evaluating and developing english math word problem solvers. *arXiv preprint arXiv:2106.15772*, 2021.
- Sewon Min, Mike Lewis, Luke Zettlemoyer, and Hannaneh Hajishirzi. MetaICL: Learning to learn in context. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 2791–2809, Seattle, United States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.201. URL <https://aclanthology.org/2022.naacl-main.201>.
- Seong Joon Oh, Bernt Schiele, and Mario Fritz. Towards reverse-engineering black-box neural networks. *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, pp. 121–144, 2019.
- Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are NLP models really able to solve simple math word problems? In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 2080–2094, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.168. URL <https://aclanthology.org/2021.naacl-main.168>.

Debjit Paul, Mete Ismayilzada, Maxime Peyrard, Beatriz Borges, Antoine Bosselut, Robert West, and Boi Faltings. Refiner: Reasoning feedback on intermediate representations. *arXiv preprint arXiv:2304.01904*, 2023.

Subhro Roy and Dan Roth. Solving general arithmetic word problems. *arXiv preprint arXiv:1608.01413*, 2016.

Ohad Rubin, Jonathan Herzig, and Jonathan Berant. Learning to retrieve prompts for in-context learning. *arXiv preprint arXiv:2112.08633*, 2021.

Noah Shinn, Beck Labash, and Ashwin Gopinath. Reflexion: an autonomous agent with dynamic memory and self-reflection. *arXiv preprint arXiv:2303.11366*, 2023.

KaShun Shum, Shizhe Diao, and Tong Zhang. Automatic prompt augmentation and selection with chain-of-thought from labeled data. *arXiv preprint arXiv:2302.12822*, 2023.

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew Dai, Andrew La, Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes, Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut Erdem, Ayla Karakaş, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph, Bartłomiej Bojanowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno Stein, Berk Ekmekci, Bill Yuchen Lin, Blake Howald, Cameron Diao, Cameron Dour, Catherine Stinson, Cedrick Argueta, César Ferri Ramírez, Chandan Singh, Charles Rathkopf, Chenlin Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Chris Waites, Christian Voigt, Christopher D. Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette, Dan Hendrycks, Dan Kilman, Dan Roth, Daniel Freeman, Daniel Khashabi, Daniel Levy, Daniel Moseguí González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra, Dilyar Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Ekaterina Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, Ellie Pavlick, Emanuele Rodola, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang, Ethan A. Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Engfu Manyasi, Evgenii Zheltonozhskii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francesca Happé, Francois Chollet, Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germán Kruszewski, Giambattista Parascandolo, Giorgio Mariani, Gloria Wang, Gonzalo Jaimovitch-López, Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh Hajishirzi, Harsh Mehta, Hayden Bogar, Henry Shevlin, Hinrich Schütze, Hiromu Yakura, Hongming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack Geissinger, Jackson Kernion, Jacob Hilton, Jaehoon Lee, Jaime Fernández Fisac, James B. Simon, James Koppel, James Zheng, James Zou, Jan Kocoń, Jana Thompson, Jared Kaplan, Jarema Radom, Jascha Sohl-Dickstein, Jason Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova, Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen Taal, Jesse Engel, Jesujoba Alabi, Jiacheng Xu, Jiaming Song, Jillian Tang, Joan Waweru, John Burden, John Miller, John U. Balis, Jonathan Berant, Jörg Froberg, Jos Rozen, Jose Hernandez-Orallo, Joseph Boudeman, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chua, Kamil Kanclerz, Karen Livescu, Karl Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Katja Markert, Kaustubh D. Dhole, Kevin Gimpel, Kevin Omondi, Kory Mathewson, Kristen Chiafullo, Ksenia Shkaruta, Kumar Shridhar, Kyle McDonell, Kyle Richardson, Laria Reynolds, Leo Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-Ochando, Louis-Philippe Morency, Luca Moschella, Lucas Lam, Lucy Noble, Ludwig Schmidt, Luheng He, Luis Oliveros Colón, Luke Metz, Lütfi Kerem Şenel, Maarten Bosma, Maarten Sap, Maartje ter Hoeve, Maheen Farooqi, Manaal Faruqi, Mantas Mazeika, Marco Baturan, Marco Marelli, Marco Maru, Maria Jose Ramírez Quintana, Marie Tolkiehn, Mario Giulianelli, Martha Lewis, Martin Potthast, Matthew L. Leavitt, Matthias Hagen, Mátyás Schubert, Medina Orduna

- Baitemirova, Melody Arnaud, Melvin McElrath, Michael A. Yee, Michael Cohen, Michael Gu, Michael Ivanitskiy, Michael Starritt, Michael Strube, Michał Swędrowski, Michele Bevilacqua, Michihiro Yasunaga, Mihir Kale, Mike Cain, Mímee Xu, Mirac Suzgun, Mo Tiwari, Mohit Bansal, Moin Aminnaseri, Mor Geva, Mozhdah Gheini, Mukund Varma T, Nanyun Peng, Nathan Chi, Nayeon Lee, Neta Gur-Ari Krakover, Nicholas Cameron, Nicholas Roberts, Nick Doiron, Nikita Nangia, Niklas Deckers, Niklas Muennighoff, Nitish Shirish Keskar, Niveditha S. Iyer, Noah Constant, Noah Fiedel, Nuan Wen, Oliver Zhang, Omar Agha, Omar Elbaghdadi, Omer Levy, Owain Evans, Pablo Antonio Moreno Casares, Parth Doshi, Pascale Fung, Paul Pu Liang, Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao, Percy Liang, Peter Chang, Peter Eckersley, Phu Mon Htut, Pinyu Hwang, Piotr Miłkowski, Piyush Patil, Pouya Pezeshkpour, Priti Oli, Qiaozhu Mei, Qing Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta Rudolph, Raefer Gabriel, Rahel Habacker, Ramón Risco Delgado, Raphaël Millière, Rhythm Garg, Richard Barnes, Rif A. Saurous, Riku Arakawa, Robbe Raymaekers, Robert Frank, Rohan Sikand, Roman Novak, Roman Sitelew, Ronan LeBras, Rosanne Liu, Rowan Jacobs, Rui Zhang, Ruslan Salakhutdinov, Ryan Chi, Ryan Lee, Ryan Stovall, Ryan Teehan, Rylan Yang, Sahib Singh, Saif M. Mohammad, Sajant Anand, Sam Dillavou, Sam Shleifer, Sam Wiseman, Samuel Gruetter, Samuel R. Bowman, Samuel S. Schoenholz, Sanghyun Han, Sanjeev Kwatra, Sarah A. Rous, Sarik Ghazarian, Sayan Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schuster, Sepideh Sadeghi, Shadi Hamdan, Sharon Zhou, Shashank Srivastava, Sherry Shi, Shikhar Singh, Shima Asaadi, Shixiang Shane Gu, Shubh Pachchigar, Shubham Toshniwal, Shyam Upadhyay, Shyamolima, Debnath, Siamak Shakeri, Simon Thormeyer, Simone Melzi, Siva Reddy, Sneha Priscilla Makini, Soo-Hwan Lee, Spencer Torene, Sriharsha Hatwar, Stanislas Dehaene, Stefan Divic, Stefano Ermon, Stella Biderman, Stephanie Lin, Stephen Prasad, Steven T. Piantadosi, Stuart M. Shieber, Summer Misherggi, Svetlana Kiritchenko, Swaroop Mishra, Tal Linzen, Tal Schuster, Tao Li, Tao Yu, Tariq Ali, Tatsu Hashimoto, Te-Lin Wu, Théo Desbordes, Theodore Rothschild, Thomas Phan, Tianle Wang, Tiberius Nkinyili, Timo Schick, Timofei Kornev, Timothy Telleen-Lawton, Titus Tunduny, Tobias Gerstenberg, Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay Ramasesh, Vinay Uday Prabhu, Vishakh Padmakumar, Vivek Srikumar, William Fedus, William Saunders, William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi Wu, Xudong Shen, Yadollah Yaghoobzadeh, Yair Lakretz, Yangqiu Song, Yasaman Bahri, Yejin Choi, Yichi Yang, Yiding Hao, Yifu Chen, Yonatan Belinkov, Yu Hou, Yufang Hou, Yuntao Bai, Zachary Seid, Zhuoye Zhao, Zijian Wang, Zijie J. Wang, Zirui Wang, and Ziyi Wu. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models, 2022.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4149–4158, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1421. URL <https://aclanthology.org/N19-1421>.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. corr, abs/2302.13971, 2023. doi: 10.48550. *arXiv preprint arXiv:2302.13971*.
- Florian Tramèr, Fan Zhang, Ari Juels, Michael K Reiter, and Thomas Ristenpart. Stealing machine learning models via prediction {APIs}. In *25th USENIX Security Symposium (USENIX Security 16)*, pp. 601–618, 2016.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022.
- Yixuan Weng, Minjun Zhu, Shizhu He, Kang Liu, and Jun Zhao. Large language models are reasoners with self-verification. *arXiv preprint arXiv:2212.09561*, 2022.

Weihao Yu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. Reclor: A reading comprehension dataset requiring logical reasoning. *arXiv preprint arXiv:2002.04326*, 2020.

Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. SWAG: A large-scale adversarial dataset for grounded commonsense inference. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 93–104, Brussels, Belgium, October–November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1009. URL <https://aclanthology.org/D18-1009>.

Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. *ArXiv preprint*, abs/2205.01068, 2022. URL <https://arxiv.org/abs/2205.01068>.

Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, et al. Least-to-most prompting enables complex reasoning in large language models. *arXiv preprint arXiv:2205.10625*, 2022.

Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng, and Tat-Seng Chua. Tat-qa: A question answering benchmark on a hybrid of tabular and textual content in finance. *arXiv preprint arXiv:2105.07624*, 2021.

A THE QUALITY OF RECONSTRUCTED PROBLEM

We measured the Rouge1, Rouge2, RougeL, RougeSum, and sentence embedding similarity (using `sentence-transformers/all-mpnet-base-v2`) between original problems and reconstructed problems. We can observe from table 6 that higher CoT accuracies correspond to higher similarities between original problems and reconstructed problems. This is expected since correctly solved problems' reconstruction should ideally be itself, and incorrectly solved problems' reconstruction should differ from itself.

Table 6: The metrics between original problem and reconstructed problem

Dateset	Rouge1	Rouge2	RougeL	RougeSum	Similarity	Standard CoT Acc
GSM8K	71.4622	49.3915	58.8991	58.8974	93.57	79.0
AQuA	54.2383	33.5828	43.5771	43.7721	84.44	51.3
AddSub	78.0939	57.5594	66.3973	66.4099	94.05	85.2
Date	46.9414	28.7005	39.2268	39.3934	79.12	66.7
SingleEq	72.8212	53.5232	64.8562	64.7956	94.01	90.3
ASDiv	68.1849	46.5488	59.7907	59.7911	92.16	84.3
SVAMP	75.0074	55.3162	65.1539	65.1973	93.71	76.7

B EXAMPLES

B.1 EXAMPLES FOR DIFFERENT KINDS OF FACTUAL INCONSISTENCIES FOUND BY RCoT

Original Solution:

First, we need to find out how many potted plants Mary currently has. She has 2 potted plants on each of the 40 window ledges, **so she has a total of $2 \times 40 = 80$ potted plants.**

Next, we need to subtract the number of potted plants she will give away. She has decided to give away 1 potted plant from each of the 40 window ledges, so she will give away a total of $1 \times 40 = 40$ potted plants.

Therefore, Mary will remain with $80 - 40 = 40$ potted plants after giving away 1 potted plant from each of the 40 window ledges. **X**

Original Problem: Mary is an avid gardener. **Yesterday, she received 18 new potted plants from her favorite plant nursery.** She already has 2 potted plants on each of the 40 window ledges of her large country home. Feeling generous, she has decided that she will give 1 potted plant from each ledge to friends and family tomorrow. How many potted plants will Mary remain with?

Reconstructed Problem: Mary has 40 window ledges in her apartment, and she has 2 potted plants on each ledge. She has decided to give away 1 potted plant from each ledge. How many potted plants will Mary have left after giving away the plants?

Figure 4: A example for overlooking

Original Solution:
First, we need to find the date that is one month before the deadline of June 1, 2021. To do this, we can subtract one month from the deadline date.
June 1, 2021 - 1 month = May 1, 2021
So, the date that is one month before the deadline is May 1, 2021.
Next, we need to write this date in the format MM/DD/YYYY.
May 1, 2021 = 05/01/2021
Therefore, the answer is (C) 05/01/2021.✗

Original Problem: The deadline is Jun 1, 2021, which is 2 days away from now. What is the date a month ago in MM/DD/YYYY? Answer Choices: (A) 04/30/2021 (B) 02/19/2021 (C) 05/01/2021 (D) 05/09/2021 (E) 07/06/2021 (F) 08/03/2021

Reconstructed Problem: What is the date that is one month before the deadline of June 1, 2021, written in the format MM/DD/YYYY? Answer Choices: (A) 04/30/2021 (B) 02/19/2021 (C) 05/01/2021 (D) 05/09/2021 (E) 07/06/2021 (F) 08/03/2021

Figure 5: A example for overlooking

Original Solution:
Let's first calculate how much money Tobias saved up in three months: He gets a \$5 allowance per month, so in three months he has saved up:
 $\$5/\text{month} * 3 \text{ months} = \15
In addition to his allowance, Tobias also earns money by mowing lawns and shoveling driveways. If he mows 4 lawns, he earns:
 $4 \text{ lawns} * \$15/\text{lawn} = \60
Let's say he shovels x driveways. Then, he earns:
 $x \text{ driveways} * \$7/\text{driveway} = \$7x$
In total, he has saved up:
 $\$15 (\text{allowance}) + \$60 (\text{lawn mowing}) + \$7x (\text{shoveling}) = \$75 + \$7x$
He spent \$95 on shoes, so he has:
 $\$75 + \$7x - \$95 = \0
Simplifying this equation:
 $\$7x - \$20 = \$0$
 $\$7x = \20
 $x = 2.85$
Therefore, Tobias shoveled about 2.85 driveways. However since he can't shovel a fraction of a driveway, we can round this up to 3 driveways.✗

Original Problem: Tobias is buying a new pair of shoes that costs \$95. He has been saving up his money each month for the past three months. He gets a \$5 allowance a month. He also mows lawns and shovels driveways. He charges \$15 to mow a lawn and \$7 to shovel. After buying the shoes, he has \$15 in change. If he mows 4 lawns, how many driveways did he shovel?

Reconstructed Problem: How many driveways did Tobias have to shovel to save enough money to buy \$95 shoes, if he gets a \$5 allowance per month for three months, mows 4 lawns for \$15 each, and earns \$7 for every driveway he shovels?

Figure 6: A example for overlooking.

Original Solution:
First, we need to find out how far John traveled during the first 3 hours at a speed of 60 mph. We can use the formula distance = speed x time, so:
distance = 60 mph x 3 hours = 180 miles
Next, we need to find out how far John traveled during the next 2 hours in standstill traffic. Since he didn't move, the distance is 0 miles.
Then, we need to find out how far John traveled during the next half-hour at a speed of 30 mph. Using the same formula, we get:
distance = 30 mph x 0.5 hours = 15 miles
Finally, we need to find out how far John traveled during the last 1.5 hours at a speed of 80 mph. Again, using the same formula, we get:
distance = 80 mph x 1.5 hours = 120 miles
To find out how far John is from home at the end of those 4 hours, we need to add up all the distances he traveled:
180 miles + 0 miles + 15 miles + 120 miles = 315 miles
Therefore, John is 315 miles away from home at the end of those 4 hours.✗

Original Problem: John drives for 3 hours at a speed of 60 mph and **then turns around because he realizes he forgot something very important at home**. He tries to get home in 4 hours but spends the first 2 hours in standstill traffic. He spends the next half-hour driving at a speed of 30mph, before being able to drive the remaining time of the 4 hours going at 80 mph. How far is he from home at the end of those 4 hours?

Reconstructed Problem: John drove for 4 hours. During the first 3 hours, he drove at a speed of 60 mph. During the next 2 hours, he was stuck in standstill traffic. During the next half-hour, he drove at a speed of 30 mph. During the last 1.5 hours, he drove at a speed of 80 mph. How far is John from home at the end of those 4 hours?

Figure 7: A example for overlooking

Original Solution:
First, we need to find out how many books and movies are left unread/watched.
Books left unread: $11 - 7 = 4$
Movies left unwatched: $17 - 21 = -4$
Uh oh, we have a problem. We can't have negative movies left unwatched, so we need to adjust our calculations.
Let's assume that the person has actually watched 17 movies, but some of them were repeats. So, we can say that they have watched 17 unique movies and some of them multiple times.
Now, let's recalculate:
Books left unread: $11 - 7 = 4$
Movies left unwatched: $17 - 17 = 0$
So, the person has read 7 books and watched 17 unique movies. To find out how many more movies than books they have watched, we can subtract:
 $17 - 7 = 10$
Therefore, the person has watched 10 more movies than books they have read.✗

Original Problem: There are 11 different books and 17 different movies in the ' crazy silly school ' series. If you read 7 of the books and **watched 21 of the movies** How many more movies than books have you read?

Reconstructed Problem: A person has 11 books and 21 movies. They have read/watched 7 books **and 17 movies**. How many more movies than books have they watched?

Figure 8: A example for hallucination.

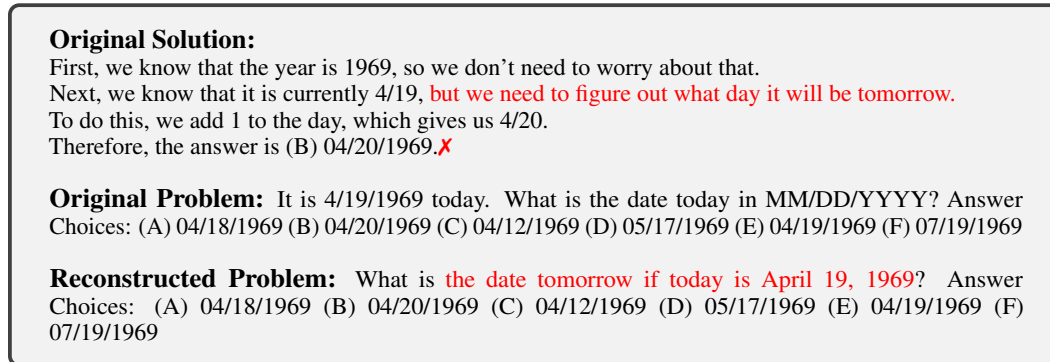


Figure 9: A example for hallucination.

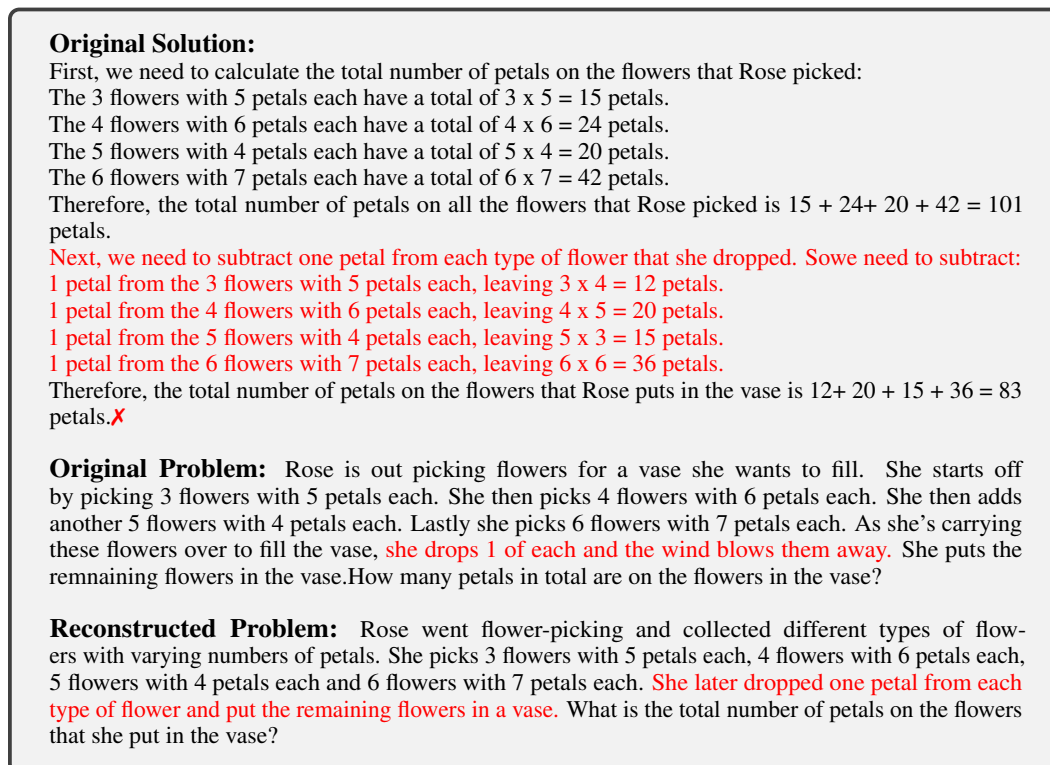


Figure 10: A example for hallucination

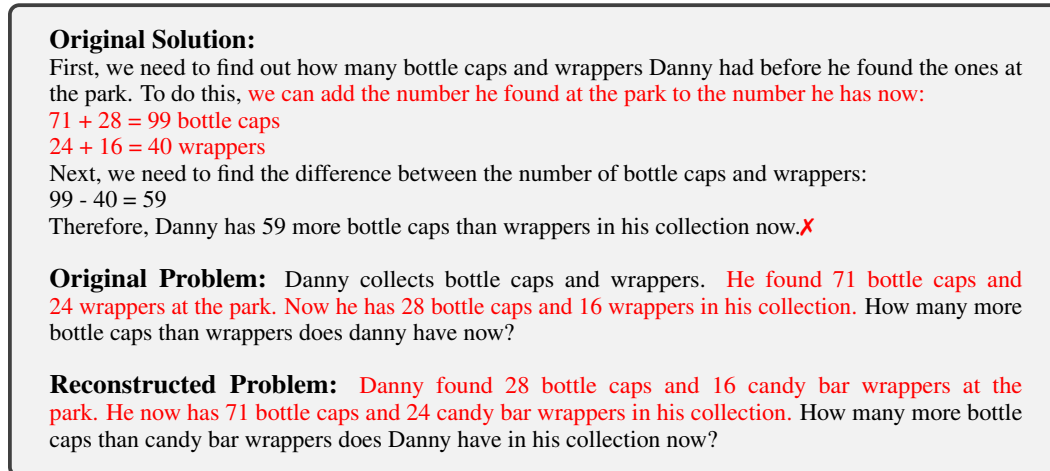


Figure 11: A example for hallucination

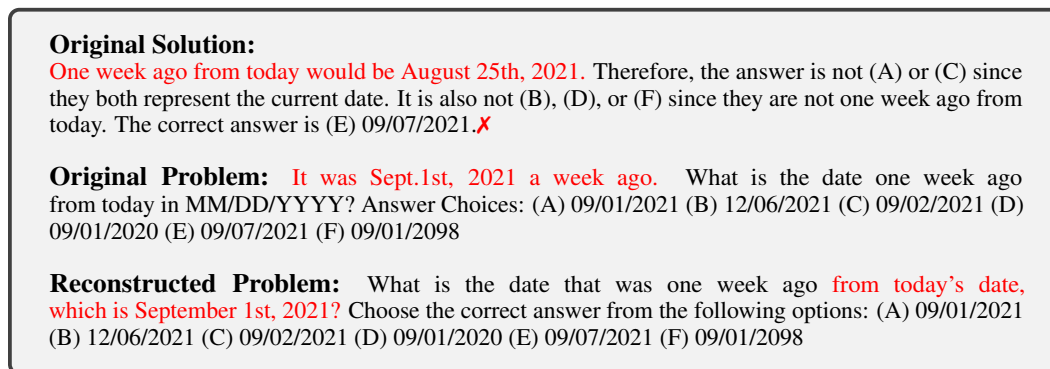


Figure 12: A example for hallucination

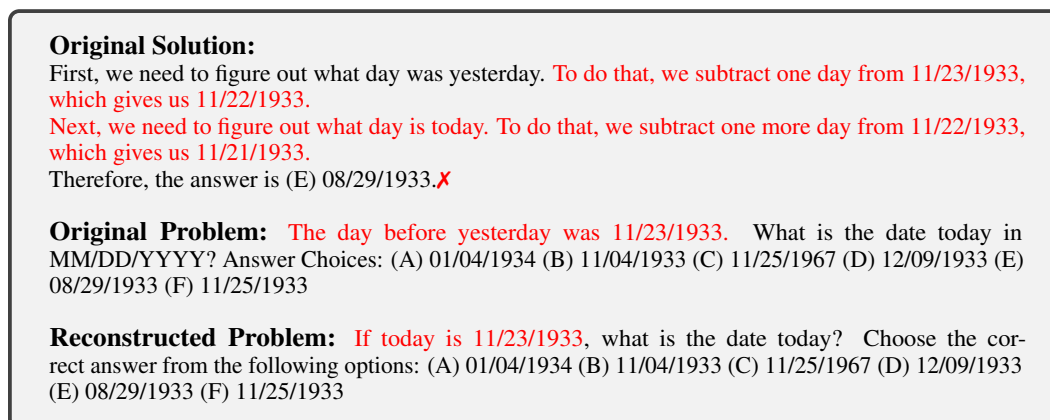


Figure 13: A example for hallucinating

Original Solution:
Martha has been collecting shells since she turned 5 years old, which means she has been collecting shells for 5 years. In one year, Martha collects 12 shells (one per month). So in 5 years, Martha collects $5 \times 12 = 60$ shells. **By her 10th birthday, Martha will have collected $60 + 12 = 72$ shells.** X

Original Problem: Martha has been collecting shells since she turned 5 years old, every month she collects one shell. By her 10th birthday, how many shells will Martha have collected?

Figure 14: A example for hallucinating

Original Solution:
First, we need to find the total number of houses in all 81 blocks. To do this, we can multiply the number of houses in each block (12) by the number of blocks (81):
 $12 \times 81 = 972$
So there are 972 houses in total.
Next, we need to find out how many pieces of junk mail the mailman needs to give to each house. **We know that he has to give 4 pieces of junk mail to each house, so we can multiply this by the total number of houses:**
 $4 \times 972 = 3,888$
Therefore, the mailman needs to give 3,888 pieces of junk mail in each block. X

Original Problem: A mailman has to give 4 pieces of junk mail to each house in each of the 81 blocks. If there are 12 houses in each block **How many pieces of junk mail should he give in each block?**

Reconstructed Problem: There are 81 blocks in a neighborhood, and each block has 12 houses. The mailman needs to deliver 4 pieces of junk mail to each house. **How many pieces of junk mail does the mailman need to deliver in total?**

Figure 15: A example for misinterpretation

B.2 EXAMPLES THAT COMPARE RCoT FEEDBACK WITH HUMAN FEEDBACK

Original Problem: Kelly has 22 nintendo games. How many does she need to buy so that she will have 140 games left?

Original Solution: To calculate how many Nintendo games Kelly needs to buy so that she will have 140 games left, we need to subtract the number of games she currently has from the desired total number of games. The desired total number of games is $140 + 22 = 162$ games. Therefore, Kelly needs to buy $162 - 22 = 140$ Nintendo games. Therefore, the final answer is 140.

Reconstructed Problem: Kelly currently has 22 Nintendo games. She wants to have 140 more games than she currently has. How many Nintendo games does Kelly need to buy to reach her goal?

RCoT

Fine-grained Feedback:

Here are the mistakes and reasons in your answer to the problem.

You use some wrong candidate conditions:

- 1: She wants to have 140 more games than she currently has. They all can not be deduced from the true condition list.

Here are detailed reasons:

- 1: The true problem only provides information about how many games Kelly currently has and she needs to reach a certain number of games. It does not provide any information that Kelly wants to have 140 more games than she currently has. Therefore, the candidate condition of "She wants to have 140 more games than she currently has" cannot be deduced from your thinking.

The real problem has the conditions:

- 1: Kelly has 22 nintendo games.
- 2: She wants to have 140 games left.

You should consider all real conditions in the problem.

Revised Answer:

Thank you for pointing out the mistakes in my previous answer. Here's a corrected response:

To calculate how many Nintendo games Kelly needs to buy so that she will have 140 games left, we need to subtract 140 from the total number of games she wants to have. The total number of games she wants to have is the current number of games she has plus the number of games she wants to buy. So, we can set up an equation:

$$22 + x = 140$$

where x is the number of games Kelly needs to buy.

Solving for x , we get:

$$x = 140 - 22 = 118$$

Therefore, Kelly needs to buy 118 Nintendo games so that she will have 140 games left. ✓

Reasons By Humans

Fine-grained Feedback:

Your answer is wrong. The problem mentions that she will have 140 games left, but it doesn't mean that she desires 140 more games. You misunderstand it.

Revised Answer:

I apologize for the mistake in my previous response. To calculate how many Nintendo games Kelly needs to buy so that she will have 140 games left, we need to subtract 140 from the number of games she currently has.

So, if Kelly currently has 22 Nintendo games and wants to have 140 games left, we can calculate the number of games she needs to buy as follows:

$$140 \text{ (desired number of games left)} - 22 \text{ (current number of games)} = 118$$

Therefore, Kelly needs to buy 118 Nintendo games to have 140 games left. ✓

Figure 16: This example shows that feedback generated by RCoT sometimes has the same quality as the human-written feedback.

Original Problem: There are 11 different books and 17 different movies in the 'crazy silly school' series. If you read 7 of the books and watched 21 of the movies How many more movies than books have you read?

Original Solution: First, we need to find out how many books and movies are left unread/watched **Let's assume that the person has actually watched 17 movies, but some of them were repeats. So, we can say that they have watched 17 unique movies and some of them multiple times.** Therefore, the person has watched 10 more movies than books they have read.

Reconstructed Problem: A person has 11 books and 21 movies. They have read/watched 7 books and 17 movies. How many more movies than books have they watched?

RCoT

Fine-grained Feedback:
Here are the mistakes and reasons in your answer to the problem.

You use some wrong candidate conditions:
1: The person has 11 books and 21 movies.
They all can not be deduced from the true condition list.
Here are detailed reasons:
1: The candidate condition states that the person has 11 books and 21 movies, but the condition list only tells us that there are 11 books and 17 movies in the series, and that the person has read 7 of the books and watched 21 of the movies. We cannot deduce that the person has all 11 books just because there are 11 books in the series. Therefore, the candidate condition cannot be deduced from the condition list.

The real problem has the conditions:
1: There are 11 different books and 17 different movies in the 'crazy silly school' series.
2: You read 7 of the books and watched 21 of the movies.
You should consider all real conditions in the problem.

Revised Solution:
I apologize for the mistakes in my previous answer. Here is the corrected answer: The true problem has the condition:
There are 11 different books and 17 different movies in the 'crazy silly school' series.
To find out how many more movies than books you have read **Let's assume that you have actually watched 17 movies,** but some of them were repeats. So, we can say that you have watched 17 unique movies and some of them multiple times.
Therefore, you have watched 10 more movies than books you have read. ✗

Reasons By Humans

Fine-grained Feedback:
The books and movies in crazy silly school are not relevant to calculate how many more movies than books you have read. You only need to consider the number of books and movies you have read/watched.

Revised Solution:
You are correct. I apologize for my mistake. To find out how many more movies than books you have read, we simply need to subtract the number of books you have read from the number of movies you have watched:
Number of movies watched - Number of books read = $21 - 7 = 14$
Therefore, you have watched 14 more movies than books you have read. ✓

Figure 17: This example shows RCoT sometimes cannot perform as well as humans.

B.3 A CONCRETE EXAMPLE THAT SHOWS THE WHOLE PROCESS OF RCoT

Original Problem: Lana had 72 extra nickels. If she put them into stacks with 8 in each stack, how many stacks could she make?

Original Solution: To calculate how many stacks Lana could make with 72 extra nickels, we need to divide the total number of nickels by the number of nickels in each stack. **Since there are 8 nickels in each stack, 72 extra nickels is equal to $72/8 = 9$ stacks.** However, Lana cannot have a fraction of a stack, so we need to round down to the nearest whole number. Therefore, Lana could make 9 stacks with the 72 extra nickels. Therefore, the final answer is 9. ✗

Step1: Problem Reconstruction

Input = [instruction, original solution]

Output = [reconstructed problem]

Instruction: Give the concrete prompt (problem) that can generate this answer. The problem should contain all basic and necessary information and correspond to the answer. The problem can only ask for one result.

Output:

Reconstructed Problem: Lana has 72 extra nickels and **wants to stack them in piles of 8**. How many stacks of nickels can she make?

Step2: Fine-grained Comparison

(i) Problem Decomposition:

Input = [instruction, original problem, reconstructed problem]

Output = [original condition list, reconstructed condition list]

Instruction: Please list the conditions of the problem. There may be multiple conditions.....

Output:

Original condition list:

1: Lana had 72 extra nickels.

2: **She put them into stacks with 8 in each stack.**

Reconstructed condition list:

1: Lana has 72 extra nickels.

2: **She wants to stack them in piles of 8.**

(ii) Condition Comparison:

Input = [instruction, sub-condition, condition list]

Output = [yes or no, reason]

Overlooked conditions:

Instruction: Given a candidate condition: "**She put them into stacks with 8 in each stack.**"

Here is a condition list:

"1: Lana has 72 extra nickels. 2: She wants to stack them in piles of 8."

From a mathematical point of view, can this candidate condition be deduced from the condition list?

Please illustrate your reason and answer "yes" or "no".

Output: No, You have ignored a real condition: **She put them into stacks with 8 in each stack.** The question has mentioned it.

hallucinated conditions:

Instruction: Given a candidate condition: "**She wants to stack them in piles of 8.**"

Here is a condition list:

"1: Lana had 72 extra nickels. 2: She put them into stacks with 8 in each stack."

From a mathematical point of view, can this candidate condition be deduced from the condition list?

Please illustrate your reason and answer "yes" or "no".

Output: No, The context only provides information about the number of extra nickels Lana has and the number of nickels she wants to put in each stack. There is no mention of her preference or intention to stack them in piles of 8.

(iii) Question Comparison:

Input = [instruction, original problem, reconstructed problem]

Output = [yes or no, reason]

Q1: Lana had 72 extra nickels. If she put them into stacks with 8 in each stack, how many stacks could she make?
Q2: Lana has 72 extra nickels and wants to stack them in piles of 5. How many stacks of nickels can she make?

Instruction: From a mathematical point of view, are these problems asking the same thing in the end? Please illustrate your reason and answer "yes" or "no".

Output: Yes, the problems are asking the same thing from a mathematical point of view. Both problems are essentially asking how many stacks Lana can make with a given number of extra nickels when arranging them in a certain quantity per stack.

Step3: Fine-grained Feedback and Revision

Input = [instruction, original solution]
Output = [revised solution]

Instruction:
Here are the mistakes and reasons in your answer to the problem.

You have ignored some real conditions:
1: She put them into stacks with 8 in each stack.
You should consider all real conditions in the problem.
Here are detailed reasons:
1: There is no information about whether or not She put them into stacks with 8 in each stack. Therefore, the real condition of "She put them into stacks with 8 in each stack" cannot be logically derived or deduced from your thinking.

You use some wrong candidate conditions:
1: She wants to stack them in piles of 5.
They all can not be deduced from the true condition list.
You should consider all real conditions in the problem.
Here are detailed reasons:
1: The condition list only provides information about the number of extra nickels Lana has and she wants to put 8 nickels in each stack. There is no mention of her preference or intention to stack them in piles of 5.

Revised Solution:
I apologize for my mistake in the previous answer. Here is the corrected answer:
To calculate how many stacks Lana could make with 72 extra nickels, we need to divide the total number of nickels by the number of nickels in each stack. **Since there are 8 nickels in each stack, Lana could make $72/8 = 9$ stacks with the 72 extra nickels.** Therefore, the final answer is 9. ✓

Figure 18: This is an example showing the whole process of RCoT. **Red:** The overlooked condition in the original problem. **Blue:** The hallucinated condition in the original problem. **Green:** The correct solution.

B.4 AN EXAMPLE COMPARING RCOT AND DOUBLE-CHECK

Original Problem: At the trip to the county-level scavenger hunt competition 90 people were required to split into groups for the competition to begin. To break people up into smaller groups with different leaders 9-person groups were formed. If $3/5$ of the number of groups each had members bring back 2 seashells each how many seashells did they bring? **GSM8K**

Original Solution: Since there were 90 people, and 9 people were in each group, there were $90/9 = 10$ groups. $3/5$ of those groups brought back seashells, which is $3/5 * 10 = 6$ groups. **Each group brought back 2 seashells**, so the total number of seashells brought back is 6 groups * 2 seashells per group = 12 seashells. The final answer is 12. **X**

Reconstructed Problem: There were 90 people on a beach trip, and they split up into groups of 9. $3/5$ of the groups brought back seashells, and each group brought back 2 seashells. How many seashells were brought back in total?

Double-Check

Coarse Feedback:
You should double-check your answer.

Revised Solution:
I apologize for the mistake in my previous answer. Here's the correct solution:
To form groups of 9 people, divide the total number of people (90) by 9:
 $90/9 = 10$ groups
If $3/5$ of the groups brought back 2 seashells each, that means:
 $(3/5) * 10$ groups = 6 groups
So, 6 groups brought back 2 seashells each:
 6 groups \times 2 seashells/group = 12 seashells
Therefore, the correct answer is 12 seashells. **X**

RCoT

Fine-grained feedback:
You used some wrong candidate conditions:
1: $3/5$ of the groups brought back seashells, each group brought back 2 seashells.
They all can not be deduced from the true conditions.
You should consider all real conditions in the problem.

Revised Solution:
Thank you for pointing out my mistake. Here is the corrected answer:
Given that 90 people were required to split into groups for the competition to begin and 9-person groups were formed, we can find the total number of groups as:
 $90/9 = 10$. So, there were 10 groups formed.
Now, we know that $3/5$ of the groups had members bring back 2 seashells each. So, the number of groups that brought back seashells is:
 $3/5 * 10 = 6$
Therefore, 6 groups brought back seashells.
Since each member brought back 2 seashells, the total number of seashells brought back is:
 6 groups \times 9 members per group \times 2 seashells per member = 108 seashells
Therefore, they brought back a total of 108 seashells. **✓**

Figure 19: This is an example that Double-Check fails to correct the solution, whereas RCoT succeeded in correcting the solution. **Green:** The correct solution of the problem. **Red:** The wrong intermediate step of the original solution. **Brown:** The factual inconsistencies found by RCoT.

B.5 DATASETS

Table 7: Examples of each reasoning task and detailed description of each dataset.

Dataset	Answer Format	Train	Test	Example
GSM8K	Number	7473	1319	Joseph had 3 times as many notebooks as Martha. Martha decided she needed more notebooks and then bought 5 more for a total of 7 notebooks. How many more than Joseph does she now have?
AQuA	Multiple choice	97467	254	A man spends 70% of his income. If his income increases by 20%, then what will be his new expenditure? Answer Choices: (A) 58.3% (B) 62.5% (C) 63.5% (D) 64.5% (E) 65.5%
AddSub	Number	-	395	Mary is baking a cake . The recipe wants 8 cups of flour . She already put in 2 cups . How many cups does she need to add ?
SVAMP	Number	-	1000	Bobby ate 28 pieces of candy. Then he ate 42 more. He also ate 63 pieces of chocolate. How many pieces of candy did Bobby eat?
SingleEq	Number	-	508	There were 28 bales of hay in the barn. Tim stacked more bales in the barn today. There are now 54 bales of hay in the barn. How many bales did he store in the barn ?
ASDiv	Number	-	2096	The following week, they decided to go to Lake Huron and Lake Michigan. During their stay there, they caught a total of 30 pikes, 40 sturgeons and 75 herrings. How many fishes did they catch from the two lakes?
Date	Multiple Choices	-	370	Today is 9/7. Jane is watching NFL 2003. What is the date tomorrow in MM/DD/YYYY? Answer Choices: (A) 08/18/2003 (B) 09/08/1916 (C) 09/13/2003 (D) 09/15/2003 (E) 09/01/2003 (F) 09/08/2003

B.6 AN EXAMPLE COMPARING RCOT AND DOUBLE-CHECK

Original Problem: At the trip to the county-level scavenger hunt competition 90 people were required to split into groups for the competition to begin. To break people up into smaller groups with different leaders 9-person groups were formed. If $3/5$ of the number of groups each had members bring back 2 seashells each how many seashells did they bring? **GSM8K**

Original Solution: Since there were 90 people, and 9 people were in each group, there were $90/9 = 10$ groups. $3/5$ of those groups brought back seashells, which is $3/5 * 10 = 6$ groups. Each group brought back 2 seashells, so the total number of seashells brought back is $6 \text{ groups} * 2 \text{ seashells per group} = 12$ seashells. The final answer is 12. **X**

Reconstructed Problem: There were 90 people on a beach trip, and they split up into groups of 9. $3/5$ of the groups brought back seashells, and each group brought back 2 seashells. How many seashells were brought back in total?

Double-Check

Coarse Feedback:
You should double-check your answer.

Revised Solution:
I apologize for the mistake in my previous answer. Here's the correct solution:
To form groups of 9 people, divide the total number of people (90) by 9:
 $90/9 = 10$ groups
If $3/5$ of the groups brought back 2 seashells each, that means:
 $(3/5) * 10 \text{ groups} = 6 \text{ groups}$
So, 6 groups brought back 2 seashells each:
 $6 \text{ groups} * 2 \text{ seashells/group} = 12 \text{ seashells}$
Therefore, the correct answer is 12 seashells. **X**

RCoT

Fine-grained feedback:
You used some wrong candidate conditions:
1: $3/5$ of the groups brought back seashells, each group brought back 2 seashells.
They all can not be deduced from the true conditions.
You should consider all real conditions in the problem.

Revised Solution:
Thank you for pointing out my mistake. Here is the corrected answer:
Given that 90 people were required to split into groups for the competition to begin and 9-person groups were formed, we can find the total number of groups as:
 $90/9 = 10$. So, there were 10 groups formed.
Now, we know that $3/5$ of the groups had members bring back 2 seashells each. So, the number of groups that brought back seashells is:
 $3/5 * 10 = 6$
Therefore, 6 groups brought back seashells.
Since each member brought back 2 seashells, the total number of seashells brought back is:
 $6 \text{ groups} * 9 \text{ members per group} * 2 \text{ seashells per member} = 108 \text{ seashells}$
Therefore, they brought back a total of 108 seashells. **✓**

Figure 20: This is an example that Double-Check fails to correct the solution, whereas RCoT succeeded in correcting the solution. **Green:** The correct solution of the problem. **Red:** The wrong intermediate step of the original solution. **Brown:** The factual inconsistencies found by RCoT.

B.7 TEMPLATE

Figure 21 shows the template prompts of RCoT.

C LIMITATIONS AND FUTURE WORK

RCoT can not detect all possible reasoning errors. For example, it is hard for RCoT to detect computational errors. However, RCoT could be combined with other prompting techniques such as

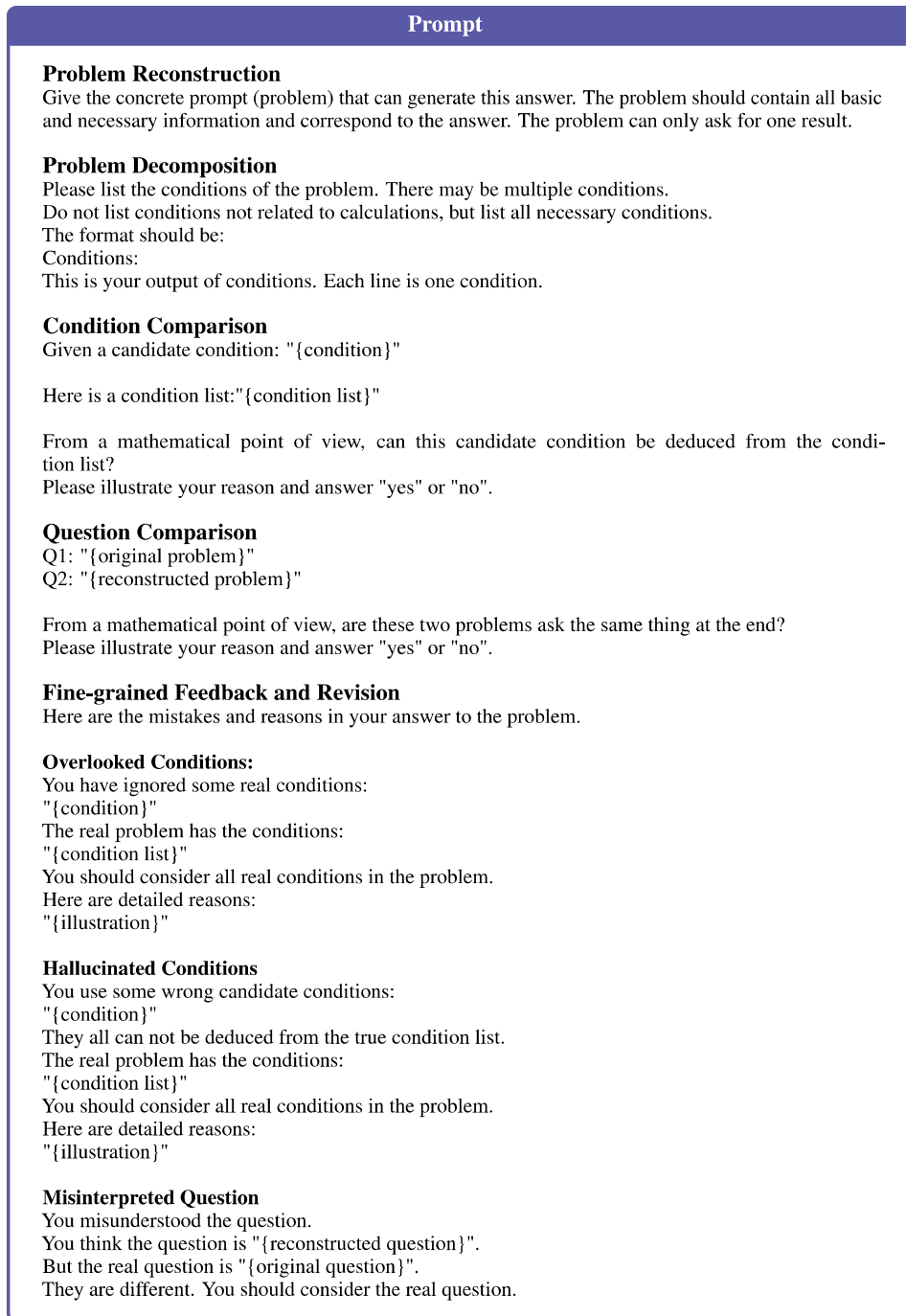


Figure 21: All prompts used in RCoT

Program-of-Thought (Chen et al., 2022), a method to reduce computational errors through disentangling reasoning and computations. Besides, there is still a significant gap between revising the solutions with RCoT-generated feedback and human feedback, which encourages further exploration in the generation of fine-grained feedback with higher quality. RCoT requires multiple conversations with LLMs (e.g., ChatGPT in our paper) and may thus slow down the inference speed due to the low bandwidth of API calls. Nevertheless, a locally deployed model may alleviate such a problem.

In the future, we will explore other applications of RCoT, such as logical reasoning and symbolic reasoning.