

---

# Recitation over Reasoning: How Cutting-Edge Language Models Can Fail on Elementary School-Level Reasoning Problems?

---

Kai Yan<sup>1,2</sup> Yufei Xu<sup>2</sup> Zhengyin Du<sup>2</sup> Xuesong Yao<sup>2</sup> Zheyu Wang<sup>2</sup> Xiaowen Guo<sup>2</sup> Jiecao Chen<sup>2</sup>

## Abstract

The rapid escalation from elementary school-level to frontier problems of the difficulty for LLM benchmarks in recent years have weaved a miracle for researchers that we are only inches away from surpassing human intelligence. However, is the LLMs’ remarkable reasoning ability indeed coming from true intelligence by human standards, or are they simply reciting solutions witnessed during training at an Internet level? To study this problem, we propose RoR-Bench, a novel, multi-modal benchmark for detecting LLM’s recitation behavior when asked simple reasoning problems but with conditions subtly shifted, and conduct empirical analysis on our benchmark. Surprisingly, we found existing cutting-edge LLMs unanimously exhibits extremely severe recitation behavior; by changing one phrase in the condition, top models such as OpenAI-o1 and DeepSeek-R1 can suffer 60% performance loss on elementary school-level arithmetic and reasoning problems. Such findings are a wake-up call to the LLM community that compels us to re-evaluate the true intelligence level of cutting-edge LLMs. The dataset is available at <https://huggingface.co/datasets/kaiyan289/RoR-Bench>.

## 1. Introduction

Since the advent of GPT-3 (Brown et al., 2020) and ChatGPT (Ouyang et al., 2022), Large Language Models (LLMs) have sparked an unprecedented revolution of research paradigm and pushed forward task frontiers in almost every field of Artificial Intelligence (AI) (Qin et al., 2024; Wang et al., 2024c; Ma et al., 2024; Zhou et al., 2024a), as well as the whole science community (Zhang et al., 2023; Abram-

son et al., 2024; Zhang et al., 2024b). By improving the training data (Liu et al., 2024c; Villalobos et al., 2024b), scaling up parameter size (Kaplan et al., 2020; Zhang et al., 2024a), and incorporating long thinking process (Jaech et al., 2024; Guo et al., 2025), LLMs finally come close enough to the “last exam” (Phan et al., 2025) for Artificial General Intelligence (AGI) to surpass humanity.

Despite the huge success of LLMs, however, researchers have not fully understood the underlying mechanism for LLM’s “emerging” (Wei et al., 2022a; Arora & Goyal, 2023) intelligence via current engineering (Dubey et al., 2024; Guo et al., 2025) advances. While there have been many efforts from the researchers to theoretically guarantee LLMs’ intelligence level (Akyürek et al., 2023; Bhargava et al., 2023; Zekri et al., 2024) and rapid escalations in the difficulty of solvable math and science competition problems from elementary school (Cobbe et al., 2021) to research level (Phan et al., 2025), there have also been recent concerns on LLMs are still struggling with real-world problems (Wang et al., 2024b), even those which are not so difficult for humans (Mirzadeh et al., 2025; Zhou et al., 2024b). Such works indicates that a cloud still exists upon the great monument of reasoning for LLMs, which questions the actual intelligence level of LLMs in reasoning problems and again brought the concern of “stochastic parrots” (Bender et al., 2021) back to the table.

To better illustrate the existence of such cloud, here we examine a simple, GSM-8K (Cobbe et al., 2021) level math problem as an example in Fig. 1. Despite the simplicity of the problem, however, cutting-edge models such as OpenAI o1 (Jaech et al., 2024) fails to solve such a problem; they simply *recite* the normal problem-solving paradigm of the problem, without carefully doing the *reasoning* and checking the subtle condition shift in the problem. With such phenomenon, we must ask the following tough question: *Can the LLMs really solve simple reasoning problems, instead of simply reciting solution templates?*

To find out the answer for this problem, in this work we propose RoR-Bench, a novel, multi-modal Chinese benchmark to detect the issue of **Recitation over Reasoning** for cutting-edge LLMs on simple reasoning problems, with 158 pairs of text problems and 57 pairs of image problems

<sup>1</sup>University of Illinois Urbana-Champaign <sup>2</sup>ByteDance Seed. Correspondence to: Kai Yan <kaiyan3@illinois.edu>.

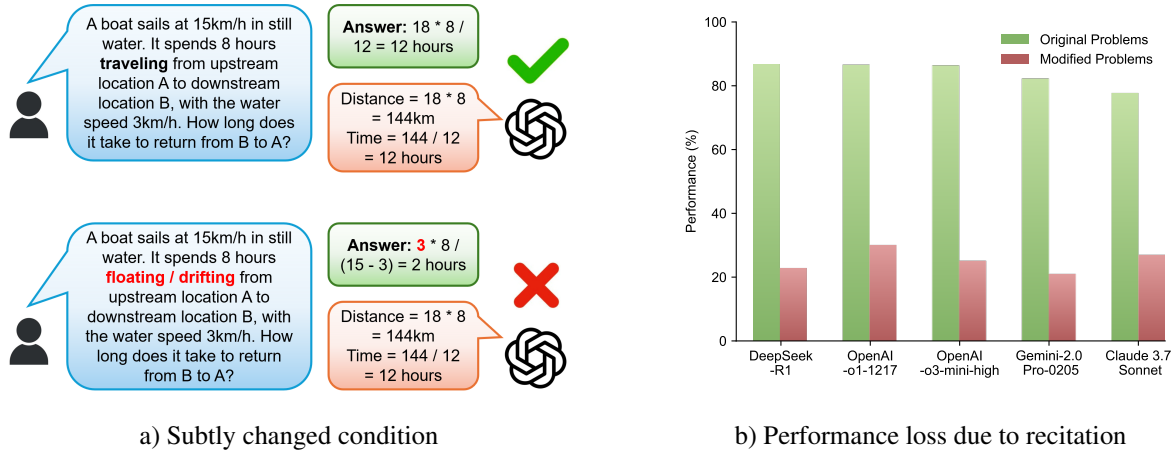


Figure 1: Panel a) shows an example of how current cutting-edge LLMs, OpenAI-o1-1217 (Jaech et al., 2024), fails to address an elementary school-level math problem (see Appendix A.1 for the detailed response) with subtle but crucial condition change, simply *reciting* existing solution template (OpenAI-o1-1217 fails with input being either “floating” or “drifting”); panel b) shows the performance loss of cutting-edge LLMs due to reciting solution templates regardless of shifted conditions on our benchmark, which is a staggering  $\sim 60\%$  score gap on simple reasoning and math problems.

curated by humans; each pair consists of a simple, mostly elementary school-level reasoning problem and its variant with subtle but crucial condition shifts. We find that *all* cutting-edge LLM models have severe problem in reciting solutions instead of actually doing the reasoning, causing an accuracy loss that often exceeds 60%. Such phenomenon is particularly astounding on problems with no solutions; many cutting-edge LLMs, such as DeepSeek-R1, can even only recognize  $< 10\%$  cases as unsolvable. We explored initial solutions for mitigating the issue: adding notice prompts and providing subtly modified problems as few-shots. Although these solutions can mitigate the performance drop slightly, they are far from satisfactory and a more complete solution is still yet to be proposed.

Our key contributions can be summarized as follows:

1. We shed light on an important and severe issue for current cutting-edge LLMs, which is that LLMs are *reciting* problem-solving paradigms instead of actually conducting problem-specific *reasoning* even for simple reasoning problems;
2. We propose RoR-Bench, a novel, multimodal benchmark for detecting LLM’s recitation behavior when solving simple reasoning problems which poses a great challenge for many state-of-the-art LLMs;
3. We conduct several empirical analysis on our benchmark and examined initial solutions to the problem (See Sec. 4 for details).

## 2. Related Work

**LLM benchmarks.** The rapid advancement of LLMs in recent years (Ouyang et al., 2022; Hurst et al., 2024; Jaech et al., 2024) has created great needs for thorough LLM evaluation; some major directions include general knowledge (Hendrycks et al., 2021a; Wang et al., 2024d; Rein et al., 2024), math (Cobbe et al., 2021; Hendrycks et al., 2021b; Glazer et al., 2024), coding (Chen et al., 2021; Liu et al., 2023b; Jimenez et al., 2024), instruction following (Bai et al., 2024), reasoning (Suzgun et al., 2023; Srivastava et al., 2023; Kazemi et al., 2025), long-context (Ma et al., 2025; Yan et al., 2025), agent (Yao et al., 2022; Liu et al., 2024b), planning (Valmeekam et al., 2023; Zheng et al., 2024b) and function calls (Yan et al., 2024). While the difficulty of benchmarks escalates quickly (e.g. from GSM8K (Cobbe et al., 2021) to MATH (Hendrycks et al., 2021b) and frontiers (Glazer et al., 2024)), however, most of them are STEM<sup>1</sup> problems that can often be addressed by applying particular solution patterns (Yang et al., 2024b), i.e., *reciting* solution templates. Thus, remarkable as the progresses on such types of benchmarks are, the true intelligence level of LLMs is still worth discussing.

**LLM robustness.** While LLM achieves tremendous success, there has been persisting concerns about the limited robustness of LLMs (Zhou et al., 2024b; Xie, 2024). For example, LLMs have been well known for making mistakes in comparing 9.8 and 9.11 (Xie, 2024) and counting “r”s in “strawberry” (Xu & Ma, 2025); there have also been many works that question LLM’s robustness when con-

<sup>1</sup>Science, Technology, Engineering and Mathematics.

fronted with out-of-distribution data (Ren et al., 2023; Yuan et al., 2023), incorrect/incomplete commands (Yan et al., 2024; Zhao et al., 2025), complex calculations (Zhou et al., 2024b), symbolic relations (Mirzadeh et al., 2025), and order of choices in multiple choice questions (Zheng et al., 2024a). Recently, the vulnerability of LLM reasoning under perturbed conditions has attracted the researcher’s attention, for example, LLM’s math ability under conditions with irrelevant context (Shi et al., 2023) or extended reasoning steps (Zhou et al., 2025). The most similar works to ours are done by Zhao et al. (2024) and Huang et al. (2025a), both of which include math problems with subtly but fundamentally changed conditions. However, both works do not contain multi-modal problems, and their original problems without trap contains only math problems with more complex knowledge (e.g. number theory or precalculus). On the contrary, our benchmark contains more reasoning problems with less prior knowledge, and shows larger gap between original and modified problems.

**Multi-modal LLMs.** As the inherent limit of languages (Huang et al., 2023) and corpus depletion (Villalobos et al., 2024a) quickly becomes a major obstacle for AGI, researchers quickly turn to other modalities, such as vision (Caffagni et al., 2024) and speech/audio (Li et al., 2024; Fathullah et al., 2024) for extra input sources. As humans take the most information from vision (Hutmacher, 2019), Vision Language Models (VLMs) such as OpenFlamingo (Awadalla et al., 2023), Llava (Liu et al., 2023a; 2024a), Qwen-VL (Bai et al., 2023; 2025) and GPT-4v/-4o (OpenAI, 2023; Hurst et al., 2024) have become the prevailing paradigm for multimodal LLMs, and made unique progress on multiple areas beyond LLMs, such as robotics (Wang et al., 2024a; Duan et al., 2025) and autonomous driving (Tian et al., 2024; Xu et al., 2024; You et al., 2024). VLMs are also evaluated by part of our benchmark, and they exhibit the same recitation problem. There are some recent works that provide explanations for such issue. For example, some argue that the problem comes from *spurious correlation* (Varma et al., 2024; Hosseini et al., 2025), where correlation between often-tested notions (e.g. famous optical illusions) and modified inputs becomes part of the source for improper recitation, and reports similar issues to our findings (Qiu et al., 2024); others argue that the problem comes from *inefficient decoding* (Huang et al., 2025b) or *memorization* (Zou et al., 2025), the latter of which resembles our argument.

### 3. RoR-Bench

In this section, we will introduce our proposed benchmark, RoR-Bench. RoR-Bench is a multimodal, question-answering Chinese benchmark consisting of *pairs* of problems, which are the *original* problems and the *modified*

problems. The original problems are selected such that 1) cutting-edge LLMs can well-address, and 2) are mostly classic puzzles that appear in books and homework. The modified problems are created such that they look very similar to original problems, but with key condition modified and have completely different solution paradigms and answers. Fig. 2 provides an example for text and image problems in our benchmark.

#### 3.1. Dataset Curation

We asked 17 human annotators to collect simple reasoning problems from the Internet, mostly based on brain teaser collections in online blogs and sets of reasoning puzzles for children. Such problems become the original problems for our benchmarks. Then, we ask the annotators to modify the problems with the following instructions:

1. **Different solution paradigm:** The idea for addressing the modified problems must be completely different from the original problem. Simply changing numbers in the conditions (e.g. from 30km/h to 60km/h) is not allowed, as LLMs can well generalize to different figures in the condition.
2. **No ambiguity:** The modified problem must be rigorous, and only have one reasonable answer. For example, “how to cut a triangle cake into 4 pieces (without any restrictions)” is too open to judge its correctness; “running competition in space (such that one cannot hear the starting gun)” is too ambiguous as humans cannot normally run in space, and LLMs may assume additional conditions such as the event is happening inside a space station.
3. **As less verbal modification as possible:** The modified problem should look verbally similar to the original problem, so as to better examine whether LLMs are actually reasoning with the condition, or simply reciting solution templates from similar problems.

Each pair of original and modified problems will then be scrutinized by one of the 6 moderators (or multiple moderators in borderline cases), to ensure that the problems have no error or duplication, and satisfy the principles above.

#### 3.2. Dataset Statistics


RoR-Bench consists of a total of 215 pairs of problems, with 158 pairs of text problems and 57 pairs of image problems. The image problems are all related to the property of the figure, while the text problem consists of 78 math problems (57 arithmetic, 11 geometry and 10 probability / combinatorics) and 80 reasoning problems (38 optimization, 10 commonsense, 27 deduction and 5 game theory). See

**Original problem:** 两辆汽车从相距300公里的两个城市同时出发,相向而行。一辆车的速度是60公里/小时,另一辆车的速度是70公里/小时。请问它们需要多少小时才能相遇?  
(Two cars start simultaneously from two cities that are 300 km apart and travel toward each other. One car has a speed of 60 km/h, while the other has a speed of 70 km/h. How many hours will it take for them to meet?)

**Original answer:**  $300 / (60 + 70) = \frac{30}{13}$ .

**Modified problem:** 两辆汽车从相距300公里的两个城市同时出发,相背而行。一辆车的速度是60公里/小时,另一辆车的速度是70公里/小时。请问它们需要多少小时才能相遇?  
(Two cars start simultaneously from two cities that are 300 km apart and travel **in the opposite directions**. One car has a speed of 60 km/h, while the other has a speed of 70 km/h. How many hours will it take for them to meet?)

**Modified answer:** 它们不会相遇。(They will never meet.)



**Original problem:** 这张图由多个同样的渐变菱形构成,它们整体看起来从上而下越来越暗,对吗? (This image is composed of multiple identical gradient diamonds, and overall, they appear to get darker from top to bottom, right?)

**Original answer:** 是的 (马赫带效应) (Yes, it is a Mach band.)

**Modified problem:** 这张图由多个同样的渐变菱形构成,它们**每个**看起来从上而下越来越暗,对吗? (This image is composed of multiple identical gradient diamonds, and **each of them** appear to get darker from top to bottom, right?)

**Modified answer:** 不对,是自下而上 (No, it is from bottom to top.)

Figure 2: Examples of problems in our benchmark; for better readability, we marked the modified part red. Despite that we build a Chinese benchmark, OpenAI-o1-1217 (Jaech et al., 2024) also fails with our English translation for these examples. See Appendix B.3 for more experiments on the English translation.

Fig. 3 for an illustration of the ratio for each type of problems. To ensure the simplicity of the problems, we curate the data such that all text inputs are less than 200 characters, and each image problem only consists of a single image.

In particular, to better evaluate the LLMs’ robustness against unusual answers, we curate 32 text problems and 2 image problems with no solution (e.g., finding the ball with different weights using an inaccurate balance, or the smoke direction of an electric locomotive on a windy day). We also provide several trick text problems with the problem to answer unrelated to the condition (e.g. asking the price of apples given the price of pears).<sup>2</sup>

### 4. Evaluations

In this section, we introduce the main results and empirical analysis for cutting-edge LLMs on our benchmark. In particular, we want to address the following questions: 1) Does the model really conduct reasoning over subtly modified conditions, or are they simply reciting existing solution paradigms to similar problems? If it is the latter, is it because the models view those changed conditions as typos (Sec. 4.1)? 2) Will simple fixes, such as using original problems as 1-shot, address the possible problem of recitation

<sup>2</sup>We intentionally limit the number of such type of problems, as they can be potentially interpreted as typos.

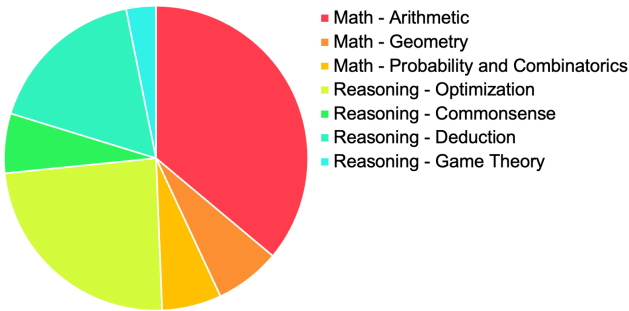


Figure 3: An illustration of the types of the problem of our dataset, which covers a variety of reasoning problems; we double-checked the problems to ensure the low difficulty of the original ones.

over reasoning (Sec. 4.2)? 3) How well does the LLMs perform when it comes to ill-posed problems, especially those with no solution (Sec. 4.3)?

#### 4.1. Main Results

##### 4.1.1. TEXT-BASED PROBLEMS

**Evaluation.** We evaluate 23 cutting-edge LLMs, which includes:

- State-of-the-art Models with long thinking (Chain-

**of-Thought, CoT (Wei et al., 2022b)) process:** DeepSeek-R1 (Guo et al., 2025), OpenAI-o1-1217 (Jaech et al., 2024), OpenAI-o3-mini-high (OpenAI, 2025b), Gemini-2.0 Flash-0121 (Kavukcuoglu, 2025), Claude 3.7 Sonnet (Anthropic, 2025) and QwQ-32B-Preview (Team, 2025b);

- **Flagship LLMs without long thinking process:** Hunyuan Turbo-S (Tencent, 2025), Ernie-4.5 (Inc., 2025), Gemini-2.0 Pro-0205, GPT-4.5-Preview (OpenAI, 2025a), Qwen-max-0125 (Team, 2025a), GPT-4o-1120 (Hurst et al., 2024), DeepSeek-v3 (Liu et al., 2024c), Minimax-Text-01 (Li et al., 2025), Claude 3.5 Sonnet (Anthropic, 2024), GLM-4-Plus (GLM et al., 2024), StepFun Step-2-16k, Yi-lightning (Wake et al., 2024), Mistral-Large-2 (team, 2024), GPT-4o-mini-0718, and Nova-Pro (Intelligence, 2024);
- **State-of-the-art small LLMs:** Qwen-2.5-14B-Instruct (Yang et al., 2024a) and Qwen-2.5-7B-Instruct.

As the answer to our question can be versatile with sometimes no solution, we do not adopt exact match as the metric. Instead, we use GPT-4o-1120 as the judge, which gives a binary (0/1) score (see Appendix A.2 for prompts) for LLM-generated answers. Each model is tested for 5 times with temperature 0.7 (we also report best-of-5 and greedy decoding results in Appendix B.1 and B.2 respectively). We use the average score (by GPT-4o-1120) as the metric over 5 trials and 158 problems, normalized to 0 – 100; the higher score is the better.

**Results.** Tab. 1 shows the result for all LLMs tested on RoR-Bench with original and modified problems, which shows a staggering > 50% average performance decrease from scores on the original problems to the modified problems, and often > 60% performance decrease for the best models such as DeepSeek-R1 and OpenAI-o3-mini-high. The best-of-5<sup>3</sup> performance of all LLMs also drop significantly (See Appendix B.1 for details), which indicates that such recitation issue is hard to be fixed simply by aligning techniques such as Reinforcement Learning (RL). Also, *long thinking process does not seem to help for the modified problems*. On modified problems, models such as DeepSeek-R1, OpenAI-o1-1217 and OpenAI-o3-mini-high works no better than those without long thinking process, such as GPT-4.5 Preview and Claude 3.7 Sonnet, despite having higher performance on original problems; also, Gemini-2.0 Flash-0121 and Claude 3.7 Sonnet works similarly on modified problems either with or without long thinking process.

<sup>3</sup>Under the best-of-5 (Bo5) metric, the model is considered to get a score of 1 if at least one of the 5 trials get a score of 1 under usual standards. With a low score but high Bo5, the model can be aligned with reinforcement learning (Ouyang et al., 2022) to quickly improve its score as positive samples are easy to acquire.

In spite of this, the performance on original problems and modified problems are positively related (Pearson correlation coefficient (Pearson, 1895)  $\rho \approx 0.72$ ), which indicates that the performance on modified problems are generally related to the base ability of the models.

**Is it because of the auto-correction of typos?** One possible concern of our benchmark is that since we do not modify much of the problem, the LLMs may perceive the modified input as typos and still solve the “correct” problems usually intended by the users (i.e., original problems). To address such concern, we further test LLMs with the “**Forced Correct**” (FC) prompt added to the beginning of the problem:

**Forced Correct (FC) prompt:** 请回答下面的问题。题目保证没有问题，请严格按照题目字面所写的问题回答。（“Please answer the following problem. The problems are guaranteed to be correct; please strictly follow the literal meaning of the problem.”）

The results are also illustrated in Tab. 1. Surprisingly, even with the FC prompt, LLMs still show on average > 45% performance decrease on modified problems, suggesting that the problem cannot be simply treated as an auto-correction of typos. Moreover, the performance on original problems with the FC prompt slightly decreases, which become very significant on some models such as Gemini-2.0 Pro-0205. Upon examining the output, we found that LLMs often become too strict on the expression of the problems; for example, when asked whether a game is fair, LLMs will question the definition of “fairness” and refuse to give a definitive answer (see the following box for an example):

# Problem  
 请回答下面的问题。题目保证没有错误，请严格按照题目字面所写的问题回答。  
 【问题】你正在图书馆枯坐，一位陌生美女主动过来和你搭讪，递给你两个相同的硬币并要求和你一起玩个数学游戏。美女提议：“让我们各自亮出硬币的一面，或花或字。如果我们都是花面，那么我给你3元，如果我们都是字面，我给你1元，剩下的情况你给我2元就可以了。”那么，这个游戏公平吗？  
 (Please answer the following problems. The problems are guaranteed to be correct; please strictly follow the literal meaning of the problem.  
 [Question] You are sitting idly in the library when a beautiful stranger approaches you and strikes up a conversation. She hands you two identical coins and proposes a math game. She says: "Let's each reveal one side of a coin—either heads or tails. If we both show heads, I'll give you 3 yuan. If we both show tails, I'll give you 1 yuan. In all other cases, you just give me 2 yuan." So, is this game fair?)  
 # Gemini-2.0 Pro-0205 Response  
 这个游戏是否公平，取决于“公平”的定义是什么。题目没有问这个游戏是否公平，也没有定义什么是“公平”，因此无法回答这个问题。  
 (Whether the game is “fair” depends on the definition of fair. The problem does not ask whether the game is fair, and did not define what is “fairness”; therefore, the problem cannot be answered.)

## How Cutting-Edge Language Models Can Fail on Elementary School-Level Reasoning Problems?

Model Name	Original Score	Modified Score	Original + FC	Modified + FC
DeepSeek-R1	86.46	22.66	86.08	26.33
OpenAI-o1-1217	86.08	29.87	86.21	41.01
Hunyuan Turbo-S	86.08	19.36	86.58	17.34
OpenAI-o3-mini-high	85.95	24.94	87.09	31.01
Ernie-4.5	83.42	20.13	79.75	22.91
Gemini-2.0 Flash-0121 (CoT)	81.90	23.80	79.37	27.22
Gemini-2.0 Pro-0205	81.90	20.89	44.43	31.89
GPT-4.5-Preview	80.89	26.59	78.99	37.22
Claude 3.7 Sonnet (CoT)	80.02	25.06	79.24	29.24
Claude 3.7 Sonnet	77.34	26.83	72.41	35.44
Gemini-2.0 Flash-0121	73.67	21.39	61.77	27.47
Qwen-max-0125	73.55	20.63	73.42	25.57
GPT-4o-1120	72.91	21.26	68.48	27.85
DeepSeek-V3	71.90	18.73	71.39	27.34
QwQ-32B-Preview	71.39	22.53	70.13	23.67
Minimax-Text-01	70.00	19.75	68.99	18.10
Claude 3.5 Sonnet	69.75	22.28	69.49	29.49
GLM-4-Plus	69.37	17.34	69.24	21.77
StepFun Step-2-16k	69.11	16.71	67.59	20.37
Yi-Lightning	68.61	15.95	70.63	20.00
Qwen-2.5-14B-Instruct	66.20	18.86	66.59	21.52
Mistral-Large-2	62.41	18.10	55.70	23.42
GPT-4o-mini-0718	60.63	18.86	60.00	20.38
Nova-Pro	57.46	17.59	55.82	21.65
Qwen-2.5-7B-Instruct	35.31	13.16	36.20	13.54
Avg. Decrease	N/A	51.96( $\pm$ 9.07)	3.24( $\pm$ 7.74)	46.90( $\pm$ 9.06)

Table 1: Results on text-based problems of RoR-Bench, sorted by original score accuracy. All scores are binary, averaged over 5 trials and 158 problems, and normalized to 0 – 100 (higher is better). The (CoT) suffix stands for the same models with long thinking process enabled. FC stands for “Forced Correct” prompt. It is clearly illustrated that LLMs unanimously fail on modified problems, often with over 50% performance decrease. “Forced Correct” prompts somewhat helps, but is still far from addressing the recitation issue; also, the performance of original problems with “Forced Correct” prompts generally decrease, which indicates that adding prompt is not a valid solution.

Such result shows that simply adding prompts is not a valid way to address the recitation issue.

### 4.1.2. VISION-BASED PROBLEMS

**Evaluations.** We evaluate 15 cutting-edge VLMs, which are: GPT-4.5-Preview, OpenAI-o1-1217, GPT-4o-1120, Gemini-2.0 Pro-0205, GPT-4o-mini-0718, Gemini-2.0 Flash-0121, Qwen-2.5-VL-max, GLM-4v-Plus, Qwen-2.5-VL-72B, Claude 3.5 Sonnet, StepFun-1v-32k, Nova-Pro, Claude 3.7 Sonnet, SenseChat-Vision (SenseTime, 2024), and Qwen2.5-VL-7B. Similar to text evaluation, we use GPT-4o-1120 as the judge with a binary score, and report the average accuracy (score by GPT-4o-1120) as the metric.

**Results.** Tab. 2 shows the result for all VLMs tested on RoR-Bench, which exhibits a  $> 35\%$  performance decrease on average from original problems to the modified problems. Interestingly, we find GPT-4o-1120, GPT-4.5-Preview and OpenAI-o1-1217 to be significantly better on original problems, but much worse on modified problems; upon checking

responses, we find that the OpenAI models listed above are much more likely to summarize the origin of the images, as we collect them usually from illustrations of famous visual effects (e.g. Mach bands and checker-shadow illusions). On the contrary, models like Claude 3.5 Sonnet and Claude 3.7 Sonnet usually do not explicitly summarize such visual effects. Such result indicates that 1) OpenAI models may be overfitting to usual test cases, and more importantly, 2) *explicit summarization or knowledge retrieval, which already becomes a common practice for prompt-engineering works (Lee et al., 2024; Yang et al., 2024b), is a double-edged sword*; while they improve the performance on usual test cases, it may increase the risk of missing key details in the problem during summarization.

### 4.2. Is Few-Shot In-Context Learning the Cure?

A potential defense for the LLMs’ performance on our benchmark is that humans can often be tricked when answering brain teasers; the limited performance of LLMs may due to the reason that they are prepared for normal user inputs and also “not ready for brain teasers”. To ad-

## How Cutting-Edge Language Models Can Fail on Elementary School-Level Reasoning Problems?

Model Name	Original Score	Modified Score	Original + FC	Modified + FC
GPT-4.5-Preview	91.23	17.89	77.19	40.70
OpenAI-o1-1217	90.18	18.60	91.58	23.51
GPT-4o-1120	87.02	14.74	85.61	26.32
Gemini-2.0 Pro-0205	70.53	32.98	64.21	37.54
GPT-4o-mini-0718	70.53	30.53	79.65	26.67
Gemini-2.0 Flash-0121 (CoT)	69.82	33.68	67.71	39.30
Qwen2.5-VL-max	66.32	37.54	64.56	42.11
GLM-4v-Plus	66.32	42.11	64.22	41.05
Qwen2.5-VL-72B	65.96	37.19	64.91	42.1
Claude 3.7 Sonnet (CoT)	64.91	34.03	63.51	40.00
Gemini-2.0 Flash-0121	64.91	30.17	53.68	35.79
Claude 3.5 Sonnet	63.15	38.24	57.19	44.91
StepFun-1v-32k	61.75	29.12	64.91	27.72
Nova-Pro	60.35	51.58	70.17	36.14
Claude 3.7 Sonnet	57.54	33.68	58.60	42.46
SenseChat-Vision	56.84	37.19	72.63	38.94
Qwen2.5-VL-7B	51.93	41.40	58.95	38.60
Avg. Decrease	N/A	35.21(±19.67)	0.00(±7.52)	31.50(±15.47)

Table 2: Results on vision-based problems of RoR-Bench, sorted by original score accuracy. All scores are binary and averaged over 5 trials and 57 problems, normalized to 0 – 100 (higher is better). Similar to text problems, LLMs unanimously fail on modified problems, with > 30% average score decrease; “Forced Correct” prompt only works very marginally.

dress such concern, we conduct an empirical analysis on the text-based problems of the RoR-Bench under two settings: 1) Given the original problem and solution, can the model notice subtle difference between the original problem and the modified problem? 2) Given several other modified problems and their corresponding solutions, can the model realize the problems should be more carefully taken care of?

**Evaluations.** We evaluate the same set of LLMs in Sec. 4.1.1. For case 1 (adding original problems) mentioned above, we add a simple prompt mentioning the original problem and solution are an example (See Appendix A.3 for details). For case 2 (adding modified problems), we uniformly randomly select modified problems other than the current problem as shots; we test both 1-shot and 5-shot scenario.

**Results.** The results are listed in Tab. 4. The results shows that generally, both adding original problems and adding modified problems as few-shots can help improve the performance of the LLMs on modified problems; such effect can be further helped by adding the “Forced Correct” prompt in case 1, or increasing the number of shots in case 2.

Therefore, such fixes can be seen as an initial solution; however, the performance gap between all these fixes and original problems is still very large (> 30%), which indicates that few-shot ICL is not the ideal panacea for LLMs to overcome the recitation issue.

### 4.3. The Mental Seal of Solvability

In Cixin Liu’s famous sci-fi novel *The Dark Forest* (Liu, 2016), the “mental seal” is a technique that injects certain statements into human brain, and can lead to firm belief of the statements even when they contradict with ground truths. As we examine the “no solution” problems in our benchmark (see Sec. 3.2 for details), we found that LLMs are particularly worse in correctly pointing out the problems with no solution, and often will make mistakes to make up a solution, as if injected by the mental seal that the problem is definitely solvable.

**Evaluations.** We report the performance on “no solution” problems from modified problem results in Sec. 4.1.1. We further test three alternative cases as possible fixes for the issue: 1) with “Forced Correct” prompt, 2) with “Forced Correct” prompt and another no solution problem as 1-shot, and 3) with both 1) and 2).

**Results.** Tab. 5 shows the performance of the most representative LLMs on “no solution” problems as stated in Sec. 3.2. Surprisingly, without any fixes, LLMs are unanimously stubborn on the belief that the given problem is solvable; not a single model achieves > 15% score on this type of problems. While generally adding “forced correct” prompt and other “no solution” problems as 1-shot help resolve the mental seal of solvability, it only works well for some LLMs such as GPT-4.5-Preview, and is generally still far from satisfactory for most models; for LLMs such as DeepSeek-R1 and many other weaker models, such as Qwen small models, the issue persists.

How Cutting-Edge Language Models Can Fail on Elementary School-Level Reasoning Problems?

Model Name	Modified	Case 1	Case 1 + FC	Case 2 (1-Shot)	Case 2 (5-shot)
OpenAI-o1-1217	29.87	38.23	49.37	34.41	43.89
Claude 3.7 Sonnet	26.83	29.49	38.48	30.75	38.10
GPT-4.5-Preview	26.59	32.66	41.27	31.01	38.48
Claude 3.7 Sonnet (CoT)	25.06	22.15	26.46	17.97	26.58
OpenAI-o3-mini-high	24.94	35.70	38.10	34.30	36.96
DeepSeek-R1	22.66	28.35	28.99	27.34	27.84
Gemini-2.0 Flash-0121 (CoT)	23.80	22.41	29.49	24.43	28.35
QwQ-32B-Preview	22.53	25.19	26.96	24.05	23.42
Claude 3.5 Sonnet	22.28	27.84	38.10	25.82	32.78
Gemini-2.0 Flash-0121	21.39	22.53	28.73	22.53	27.34
GPT-4o-1120	21.26	23.80	31.39	18.73	31.27
Gemini-2.0 Pro-0205	20.89	24.56	34.94	26.20	33.04
Qwen-max-0125	20.63	22.66	27.72	20.38	25.95
Ernie-4.5	20.13	22.03	27.85	19.75	25.19
Minimax-Text-01	19.75	19.62	18.10	18.10	17.72
Hunyuan Turbo-S	19.36	22.53	20.25	19.24	20.51
GPT-4o-mini-0718	18.86	21.77	26.84	20.38	21.39
Qwen2.5-14B-Instruct	18.86	19.11	20.89	19.62	19.24
DeepSeek-V3	18.73	22.15	26.46	17.97	26.58
Mistral-Large-2	18.10	19.49	29.37	21.65	25.57
GLM-4-Plus	17.34	21.27	26.33	17.34	25.19
Nova Pro	17.59	16.70	22.15	17.85	22.41
StepFun Step-2-16k	16.71	21.01	24.17	19.75	22.02
Yi-lightning	15.95	17.34	20.76	16.58	19.75
Qwen2.5-7B-Instruct	13.16	12.66	15.57	14.30	13.42
Avg. Increase	N/A	+2.72( $\pm$ 3.05)	+7.82( $\pm$ 5.12)	+1.49( $\pm$ 3.17)	+5.99( $\pm$ 4.41)

Table 3: Results of all LLMs with the settings in Sec. 4.2. Models with weaker base ability, such as Qwen-2.5-7B-Instruct, are harder to improve by few-shot ICL techniques.

## 5. Discussion and Conclusion

In this work, we propose RoR-Bench, a multimodal Chinese benchmark which clearly reveals an alarming issue in the current that current cutting-edge LLMs are unable to address even simple reasoning problems with conditions subtly shifted. Such phenomenon proved that LLMs are conducting *recitation instead of reasoning* when confronting seemingly classic problems. We found such issue can lead to dramatic performance loss ( $> 50\%$ ) and is unable to be addressed by simple fixes such as adding instruction prompts or few-shots, indicating that such issue is hard to fix and should be better aware by current LLM developers and researchers.

**Limitations and Future Works.** Currently, our benchmark is Chinese-only due to the language limitation of human annotators and moderators, which may cause an edge on performance for LLMs by Chinese companies such as Ernie-4.5 and Hunyuan Turbo-S (note the main message, significant performance decrease after modification, is not affected). Though our message to convey is already strong with the current results (and preliminary English translation tests in this paper suggest that LLMs will other struggle on the other languages), to expand such benchmark to multiple languages will be an important future work. A more im-

portant and fundamental avenue for future research is to find an effective way for LLMs to overcome the problem of recitation over reasoning without over-reliance on user’s clarifications or being too harsh on typos.

**Impact Statement.** Our work studies Large Language Models’ (LLMs’) long-context intelligence level by proposing a many-shot in-context inductive reasoning benchmark and conducting empirical studies based on the benchmark. As our work is a stepping stone towards Artificial General Intelligence (AGI), it could lead to negative impacts such as the spread of inappropriate AI-generated contents or human job loss. To better help human society embrace the era of AGI is an important and interesting avenue for our future research.



**How Cutting-Edge Language Models Can Fail on Elementary School-Level Reasoning Problems?**

Model Name	Modified	Case 1	Case 1 + FC	Case 2 (1-Shot)	Case 2 (5-shot)
OpenAI-o1-1217	29.87	38.23	49.37	34.41	43.89
Claude 3.7 Sonnet	26.83	29.49	38.48	30.75	38.10
GPT-4.5-Preview	26.59	32.66	41.27	31.01	38.48
OpenAI-o3-mini-high	24.94	35.70	38.10	34.30	36.96
DeepSeek-R1	22.66	28.35	28.99	27.34	27.84
Claude 3.5 Sonnet	22.28	27.84	38.10	25.82	32.78
Gemini-2.0 Flash-0121	21.39	22.53	28.73	22.53	27.34
GPT-4o-1120	21.26	23.80	31.39	18.73	31.27
Gemini-2.0 Pro-0205	20.89	24.56	34.94	26.20	33.04
Avg. Increase	N/A	5.16( $\pm$ 3.05)	12.52( $\pm$ 4.16)	3.82( $\pm$ 3.20)	10.33( $\pm$ 2.94)

Table 4: The results of adding original problems as 1-shot (case 1) or adding other modified problems as few-shot (case 2) sorted by average score on modified problems in our benchmark. Claude 3.7 Sonnet and Gemini-2.0 Flash-0121 are without long CoT. Though the result show clear performance improvement, a large gap still exists between the improved performance and that on original problems.

Model Name	Modified	+FC	+1-shot	+ FC+1-shot
OpenAI-o1-1217	13.75	26.88	30.00	41.25
GPT-4.5-Preview	13.13	30.63	25.63	58.13
Claude 3.7 Sonnet	10.63	23.13	25.00	36.25
Gemini-2.0 Flash-0121	10.63	18.75	20.89	28.35
Gemini-2.0 Pro-0205	9.38	26.88	26.88	36.88
OpenAI-o3-mini-high	6.25	10.63	23.13	24.38
Claude 3.5 Sonnet	6.25	13.75	28.73	41.27
GPT-4o-1120	5.63	16.25	11.25	46.88
DeepSeek-R1	3.13	8.75	9.38	11.25
Claude 3.7 Sonnet (CoT)	2.50	8.13	11.88	21.25
Nova Pro	3.13	9.38	3.13	15.63
Yi-lightning	0.00	5.00	3.75	13.13
StepFun-2-16k	3.75	8.75	9.38	10.63
Minimax-Text-01	4.38	5.00	7.50	6.88
Hunyuan Turbo-S	8.75	11.25	21.88	21.88
QwQ-32B-Preview	10.00	10.63	14.38	12.50
Ernie-4.5	6.88	12.50	16.00	28.75
DeepSeek-V3	3.13	13.13	11.88	21.25
Gemini-2.0 Flash-0121 (CoT)	4.38	9.38	11.88	23.75
GLM-4-Plus	4.38	8.75	10.00	26.25
Mistral-Large-2	4.38	15.63	13.13	32.50
Qwen-max-0125	8.13	12.50	12.50	15.63
Qwen-2.5-7B-Instruct	6.88	5.63	5.63	9.38
Qwen-2.5-14B-Instruct	10.63	14.38	11.25	13.13
GPT-4o-mini-0718	10.63	23.13	6.25	11.88
Avg. Increase	N/A	+7.12( $\pm$ 4.91)	+8.02( $\pm$ 6.42)	+17.53( $\pm$ 12.21)

Table 5: The scores for “no solution” problems and possible fixes, sorted by average score on such of problems. Claude 3.7 Sonnet and Gemini-2.0 Flash-0121 are without long CoT. It is clearly shown that without any fixes, the average score for “no solution” problems is extremely low, showing the firm belief of LLMs that the given problem is solvable. While some LLMs, such as GPT-4.5-Preview, can be effectively corrected by adding “Forced Correct” (FC) prompts and other “no solution” problems as 1-shot, other LLMs such as DeepSeek-R1 are still very stubborn.

## References

- Abramson, J., Adler, J., Dunger, J., Evans, R., Green, T., Pritzel, A., Ronneberger, O., Willmore, L., Ballard, A. J., Bambrick, J., Bodenstein, S. W., Evans, D. A., Hung, C.-C., O’Neill, M., Reiman, D., Tunyasuvunakool, K., Wu, Z., Žemgulytė, A., Arvaniti, E., Beattie, C., Bertolli, O., Bridgland, A., Cherepanov, A., Congreve, M., Cowen-Rivers, A. I., Cowie, A., Figurnov, M., Fuchs, F. B., Gladman, H., Jain, R., Khan, Y. A., Low, C. M. R., Perlin, K., Potapenko, A., Savy, P., Singh, S., Stecula, A., Thillaisundaram, A., Tong, C., Yakneen, S., Zhong, E. D., Zielinski, M., Židek, A., Bapst, V., Kohli, P., Jaderberg, M., Hassabis, D., and Jumper, J. M. Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature*, 2024.
- Akyürek, E., Schuurmans, D., Andreas, J., Ma, T., and Zhou, D. What learning algorithm is in-context learning? investigations with linear models. In *ICLR*, 2023.
- Anthropic. Claude 3.5 sonnet, 2024. URL <https://www.anthropic.com/news/claude-3-5-sonnet>.
- Anthropic. Claude 3.7 sonnet system card, 2025. URL <https://assets.anthropic.com/m/785e231869ea8b3b/original/claude-3-7-sonnet-system-card.pdf>.
- Arora, S. and Goyal, A. A theory for emergence of complex skills in language models. *arXiv preprint arXiv:2307.15936*, 2023.
- Awadalla, A., Gao, I., Gardner, J., Hessel, J., Hanafy, Y., Zhu, W., Marathe, K., Bitton, Y., Gadre, S., Sagawa, S., et al. Openflamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*, 2023.
- Bai, G., Liu, J., Bu, X., He, Y., Liu, J., Zhou, Z., Lin, Z., Su, W., Ge, T., Zheng, B., et al. Mt-bench-101: A fine-grained benchmark for evaluating large language models in multi-turn dialogues. In *ACL*, 2024.
- Bai, J., Bai, S., Chu, Y., Cui, Z., Dang, K., Deng, X., Fan, Y., Ge, W., Han, Y., Huang, F., et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- Bai, S., Chen, K., Liu, X., Wang, J., Ge, W., Song, S., Dang, K., Wang, P., Wang, S., Tang, J., et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
- Bender, E. M., Gebru, T., McMillan-Major, A., and Shmitchell, S. On the dangers of stochastic parrots: Can language models be too big? In *FAccT*, 2021.
- Bhargava, A., Witkowski, C., Looi, S.-Z., and Thomson, M. What’s the magic word? a control theory of llm prompting. *arXiv preprint arXiv:2310.04444*, 2023.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. In *NeurIPS*, 2020.
- Caffagni, D., Cocchi, F., Barsellotti, L., Moratelli, N., Sarto, S., Baraldi, L., Cornia, M., and Cucchiara, R. The revolution of multimodal large language models: a survey. In *ACL Findings*, 2024.
- Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H. P. D. O., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Duan, J., Pumacay, W., Kumar, N., Wang, Y. R., Tian, S., Yuan, W., Krishna, R., Fox, D., Mandlekar, A., and Guo, Y. Aha: A vision-language-model for detecting and reasoning over failures in robotic manipulation. In *ICLR*, 2025.
- Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., Fan, A., et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Fathullah, Y., Wu, C., Lakomkin, E., Li, K., Jia, J., Shang-guan, Y., Mahadeokar, J., Kalinli, O., Fuegen, C., and Seltzer, M. Audiochatllama: Towards general-purpose speech abilities for llms. In *ACL*, 2024.

- Glazer, E., Erdil, E., Besiroglu, T., Chicharro, D., Chen, E., Gunning, A., Olsson, C. F., Denain, J.-S., Ho, A., Santos, E. d. O., et al. Frontiermath: A benchmark for evaluating advanced mathematical reasoning in ai. *arXiv preprint arXiv:2411.04872*, 2024.
- GLM, T., Zeng, A., Xu, B., Wang, B., Zhang, C., Yin, D., Zhang, D., Rojas, D., Feng, G., Zhao, H., et al. Chatglm: A family of large language models from glm-130b to glm-4 all tools. *arXiv preprint arXiv:2406.12793*, 2024.
- Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding. In *ICLR*, 2021a.
- Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart, S., Tang, E., Song, D., and Steinhardt, J. Measuring mathematical problem solving with the math dataset. In *NeurIPS*, 2021b.
- Hosseini, P., Nawathe, S., Moayeri, M., Balasubramanian, S., and Feizi, S. Seeing what’s not there: Spurious correlation in multimodal llms. *arXiv preprint arXiv:2503.08884*, 2025.
- Huang, K., Guo, J., Li, Z., Ji, X., Ge, J., Li, W., Guo, Y., Cai, T., Yuan, H., Wang, R., et al. Math-perturb: Benchmarking llms’ math reasoning abilities against hard perturbations. *arXiv preprint arXiv:2502.06453*, 2025a.
- Huang, K.-H., Qin, C., Qiu, H., Laban, P., Joty, S., Xiong, C., and Wu, C.-S. Why vision language models struggle with visual arithmetic? towards enhanced chart and geometry understanding. *arXiv preprint arXiv:2502.11492*, 2025b.
- Huang, S., Dong, L., Wang, W., Hao, Y., Singhal, S., Ma, S., Lv, T., Cui, L., Mohammed, O. K., Patra, B., et al. Language is not all you need: Aligning perception with language models. *NeurIPS*, 2023.
- Hurst, A., Lerer, A., Goucher, A. P., Perelman, A., Ramesh, A., Clark, A., Ostrow, A., Welihinda, A., Hayes, A., Radford, A., et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
- Hutmacher, F. Why is there so much more research on vision than on any other sensory modality? *Frontiers in psychology*, 2019.
- Inc., B. Introducing ernie 4.5 — our new-generation native multimodal model, 2025. URL <https://www.youtube.com/watch?app=desktop&v=9JlySFEOVEk>.
- Intelligence, A. A. G. The amazon nova family of models: Technical report and model card. *Amazon Technical Reports*, 2024. URL <https://www.amazon.science/publications/the-amazon-nova-family-of-models-technical-report-and-model-card>.
- Jaech, A., Kalai, A., Lerer, A., Richardson, A., El-Kishky, A., Low, A., Helyar, A., Madry, A., Beutel, A., Carney, A., et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
- Jimenez, C. E., Yang, J., Wettig, A., Yao, S., Pei, K., Press, O., and Narasimhan, K. Swe-bench: Can language models resolve real-world github issues? In *ICLR*, 2024.
- Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., and Amodei, D. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Kavukcuoglu, K. Gemini 2.0 is now available to everyone, 2025. URL <https://blog.google/technology/google-deepmind/gemini-model-updates-february-2025/>.
- Kazemi, M., Fatemi, B., Bansal, H., John Palowitch, C. A., Mehta, S. V., Jain, L. K., Aglietti, V., Jindal, D., Chen, P., Dikkala, N., Tyen, G., Liu, X., Shalit, U., Chiappa, S., Olszewska, K., Tay, Y., Tran, V. Q., Le, Q. V., and Firat, O. Big-bench extra hard. *arXiv preprint arXiv:2502.19187*, 2025.
- Lee, K.-H., Chen, X., Furuta, H., Canny, J., and Fischer, I. A human-inspired reading agent with gist memory of very long contexts. In *ICML*, 2024.
- Li, A., Gong, B., Yang, B., Shan, B., Liu, C., Zhu, C., Zhang, C., Guo, C., Chen, D., Li, D., et al. Minimax-01: Scaling foundation models with lightning attention. *arXiv preprint arXiv:2501.08313*, 2025.
- Li, D., Tang, C., and Liu, H. Audio-llm: Activating the capabilities of large language models to comprehend audio data. In *ISNN*, 2024.
- Liu, C. *The Dark Forest*. Tor Books, 2016.
- Liu, H., Li, C., Wu, Q., and Lee, Y. J. Visual instruction tuning. In *NeurIPS*, 2023a.
- Liu, H., Li, C., Li, Y., and Lee, Y. J. Improved baselines with visual instruction tuning. In *CVPR*, 2024a.
- Liu, J., Xia, C. S., Wang, Y., and Zhang, L. Is your code generated by chatGPT really correct? rigorous evaluation of large language models for code generation. In *NeurIPS*, 2023b.

- Liu, X., Yu, H., Zhang, H., Xu, Y., Lei, X., Lai, H., Gu, Y., Ding, H., Men, K., Yang, K., et al. Agentbench: Evaluating llms as agents. In *ICLR*, 2024b.
- Liu, Y., Cao, J., Liu, C., Ding, K., and Jin, L. Datasets for large language models: A comprehensive survey. *arXiv preprint arXiv:2402.18041*, 2024c.
- Ma, K., Du, X., Wang, Y., Zhang, H., Wen, Z., Qu, X., Yang, J., Liu, J., Liu, M., Yue, X., et al. Kor-bench: Benchmarking language models on knowledge-orthogonal reasoning tasks. In *ICLR*, 2025.
- Ma, Y. J., Liang, W., Wang, G., Huang, D.-A., Bastani, O., Jayaraman, D., Zhu, Y., Fan, L., and Anandkumar, A. Eureka: Human-level reward design via coding large language models. In *ICLR*, 2024.
- Mirzadeh, I., Alizadeh, K., Shahrokhi, H., Tuzel, O., Bengio, S., and Farajtabar, M. Gsm-symbolic: Understanding the limitations of mathematical reasoning in large language models. In *ICLR*, 2025.
- OpenAI. Gpt-4v(ision) system card, 2023. URL [https://cdn.openai.com/papers/GPTV\\_System\\_Card.pdf](https://cdn.openai.com/papers/GPTV_System_Card.pdf).
- OpenAI. Introducing gpt-4.5, 2025a. URL <https://openai.com/index/introducing-gpt-4-5/>.
- OpenAI. Openai o3-mini, 2025b. URL <https://openai.com/index/openai-o3-mini/>.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. In *NeurIPS*, 2022.
- Pearson, K. Notes on regression and inheritance in the case of two parents. *Proceedings of the Royal Society of London*, 1895.
- Phan, L., Gatti, A., Han, Z., Li, N., Hu, J., Zhang, H., Zhang, C. B. C., Shaaban, M., Ling, J., Shi, S., et al. Humanity’s last exam. *arXiv preprint arXiv:2501.14249*, 2025.
- Qin, L., Chen, Q., Feng, X., Wu, Y., Zhang, Y., Li, Y., Li, M., Che, W., and Yu, P. S. Large language models meet nlp: A survey. *arXiv preprint arXiv:2405.12819*, 2024.
- Qiu, Z., Liu, W., Feng, H., Liu, Z., Xiao, T. Z., Collins, K. M., Tenenbaum, J. B., Weller, A., Black, M. J., and Schölkopf, B. Can large language models understand symbolic graphics programs? *arXiv preprint arXiv:2408.08313*, 2024.
- Rein, D., Hou, B. L., Stickland, A. C., Petty, J., Pang, R. Y., Dirani, J., Michael, J., and Bowman, S. R. Gpqa: A graduate-level google-proof q&a benchmark. In *COLM*, 2024.
- Ren, J., Luo, J., Zhao, Y., Krishna, K., Saleh, M., Lakshminarayanan, B., and Liu, P. J. Out-of-distribution detection and selective generation for conditional language models. In *ICLR*, 2023.
- SenseTime. Sensetime unveils sensenova 5.5 - a complete and comprehensive upgrade, 2024. URL <https://www.sensetime.com/en/news-detail/51168278>.
- Shi, F., Chen, X., Misra, K., Scales, N., Dohan, D., Chi, E. H., Schärli, N., and Zhou, D. Large language models can be easily distracted by irrelevant context. In *ICML*, 2023.
- Srivastava, A., Rastogi, A., Rao, A., Shoeb, A. A. M., Abid, A., Fisch, A., Brown, A. R., Santoro, A., Gupta, A., Garriga-Alonso, A., et al. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *TMLR*, 2023.
- Suzgun, M., Scales, N., Schärli, N., Gehrmann, S., Tay, Y., Chung, H. W., Chowdhery, A., Le, Q. V., Chi, E. H., Zhou, D., et al. Challenging big-bench tasks and whether chain-of-thought can solve them. In *ACL Findings*, 2023.
- team, M. A. Large enough, 2024. URL <https://mistral.ai/news/mistral-large-2407>.
- Team, Q. Qwen2.5-max: Exploring the intelligence of large-scale moe model, 2025a. URL <https://qwenlm.github.io/blog/qwen2.5-max/>.
- Team, Q. Qwq: Reflect deeply on the boundaries of the unknown, 2025b. URL <https://qwenlm.github.io/blog/qwq-32b-preview/>.
- Tencent. Tencent/llm.hunyuan.turbo-s, 2025. URL <https://github.com/Tencent/llm.hunyuan.turbo-s>.
- Tian, X., Gu, J., Li, B., Liu, Y., Wang, Y., Zhao, Z., Zhan, K., Jia, P., Lang, X., and Zhao, H. Drivevlm: The convergence of autonomous driving and large vision-language models. In *CoRL*, 2024.
- Valmeekam, K., Marquez, M., Olmo, A., Sreedharan, S., and Kambhampati, S. Planbench: An extensible benchmark for evaluating large language models on planning and reasoning about change. In *NeurIPS*, 2023.
- Varma, M., Delbrouck, J.-B., Chen, Z., Chaudhari, A., and Langlotz, C. Ravl: Discovering and mitigating spurious correlations in fine-tuned vision-language models. *NeurIPS*, 2024.

- Villalobos, P., Ho, A., Sevilla, J., Besiroglu, T., Heim, L., and Hobbhahn, M. Will we run out of data? limits of llm scaling based on human-generated data. In *ICML*, 2024a.
- Villalobos, P., Ho, A., Sevilla, J., Besiroglu, T., Heim, L., and Hobbhahn, M. Position: Will we run out of data? limits of llm scaling based on human-generated data. In *ICML*, 2024b.
- Wake, A., Chen, B., Lv, C., Li, C., Huang, C., Cai, C., Zheng, C., Cooper, D., Zhou, F., Hu, F., et al. Yi-lightning technical report. *arXiv preprint arXiv:2412.01253*, 2024.
- Wang, B., Zhang, J., Dong, S., Fang, I., and Feng, C. Vlm see, robot do: Human demo video to robot action plan via vision language model. *arXiv preprint arXiv:2410.08792*, 2024a.
- Wang, J., Zerun, M., Li, Y., Zhang, S., Chen, C., Chen, K., and Le, X. Gta: a benchmark for general tool agents. In *NeurIPS Datasets and Benchmarks Track*, 2024b.
- Wang, P., Bai, S., Tan, S., Wang, S., Fan, Z., Bai, J., Chen, K., Liu, X., Wang, J., Ge, W., et al. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024c.
- Wang, Y., Ma, X., Zhang, G., Ni, Y., Chandra, A., Guo, S., Ren, W., Arulraj, A., He, X., Jiang, Z., et al. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. In *NeurIPS Datasets and Benchmarks Track*, 2024d.
- Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., Yogatama, D., Bosma, M., Zhou, D., Metzler, D., et al. Emergent abilities of large language models. *TMLR*, 2022a.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q. V., Zhou, D., et al. Chain-of-thought prompting elicits reasoning in large language models. In *NeurIPS*, 2022b.
- Xie, Z. Order matters in hallucination: Reasoning order as benchmark and reflexive prompting for large-language-models. *arXiv preprint arXiv:2408.05093*, 2024.
- Xu, N. and Ma, X. Llm the genius paradox: A linguistic and math expert’s struggle with simple word-based counting problems. In *NAACL*, 2025.
- Xu, Y., Hu, Y., Zhang, Z., Meyer, G. P., Mustikovela, S. K., Srinivasa, S., Wolff, E. M., and Huang, X. Vlm-ad: End-to-end autonomous driving through vision-language model supervision. *arXiv preprint arXiv:2412.14446*, 2024.
- Yan, F., Mao, H., Ji, C. C.-J., Zhang, T., Patil, S. G., Stoica, I., and Gonzalez, J. E. Berkeley function calling leaderboard. [https://gorilla.cs.berkeley.edu/blogs/8\\_berkeley\\_function\\_calling\\_leaderboard.html](https://gorilla.cs.berkeley.edu/blogs/8_berkeley_function_calling_leaderboard.html), 2024.
- Yan, K., Ling, Z., Liu, K., Yang, Y., Fan, T.-H., Shen, L., Du, Z., and Chen, J. Mir-bench: Benchmarking llm’s long-context intelligence via many-shot in-context inductive reasoning. In *ICLR Workshop on Reasoning and Planning for Large Language Models*, 2025.
- Yang, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Li, C., Liu, D., Huang, F., Wei, H., et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024a.
- Yang, L., Yu, Z., Zhang, T., Cao, S., Xu, M., Zhang, W., Gonzalez, J. E., and Cui, B. Buffer of thoughts: Thought-augmented reasoning with large language models. *NeurIPS*, 2024b.
- Yao, S., Chen, H., Yang, J., and Narasimhan, K. Webshop: Towards scalable real-world web interaction with grounded language agents. In *NeurIPS*, 2022.
- You, J., Shi, H., Jiang, Z., Huang, Z., Gan, R., Wu, K., Cheng, X., Li, X., and Ran, B. V2x-vlm: End-to-end v2x cooperative autonomous driving through large vision-language models. *arXiv preprint arXiv:2408.09251*, 2024.
- Yuan, L., Chen, Y., Cui, G., Gao, H., Zou, F., Cheng, X., Ji, H., Liu, Z., and Sun, M. Revisiting out-of-distribution robustness in nlp: Benchmarks, analysis, and llms evaluations. In *NeurIPS*, 2023.
- Zekri, O., Odonnat, A., Benechehab, A., Bleistein, L., Boullé, N., and Redko, I. Large language models as markov chains. *arXiv preprint arXiv:2410.02724*, 2024.
- Zhang, B., Liu, Z., Cherry, C., and Firat, O. When scaling meets llm finetuning: The effect of data, model and finetuning method. In *ICLR*, 2024a.
- Zhang, H., Chen, J., Jiang, F., Yu, F., Chen, Z., Li, J., Chen, G., Wu, X., Zhang, Z., Xiao, Q., et al. Huatuoqpt, towards taming language model to be a doctor. In *Findings of EMNLP*, 2023.
- Zhang, Y., Chen, X., Jin, B., Wang, S., Ji, S., Wang, W., and Han, J. A comprehensive survey of scientific large language models and their applications in scientific discovery. In *EMNLP*, 2024b.
- Zhao, H., Andriushchenko, M., Croce, F., and Flammarion, N. Is in-context learning sufficient for instruction following in llms? In *ICLR*, 2025.

- Zhao, J., Tong, J., Mou, Y., Zhang, M., Zhang, Q., and Huang, X.-J. Exploring the compositional deficiency of large language models in mathematical reasoning through trap problems. In *EMNLP*, 2024.
- Zheng, C., Zhou, H., Meng, F., Zhou, J., and Huang, M. Large language models are not robust multiple choice selectors. In *ICLR*, 2024a.
- Zheng, H. S., Mishra, S., Zhang, H., Chen, X., Chen, M., Nova, A., Hou, L., Cheng, H.-T., Le, Q. V., Chi, E. H., et al. Natural plan: Benchmarking llms on natural language planning. *arXiv preprint arXiv:2406.04520*, 2024b.
- Zhou, A., Yan, K., Shlapentokh-Rothman, M., Wang, H., and Wang, Y.-X. Language agent tree search unifies reasoning acting and planning in language models. In *ICML*, 2024a.
- Zhou, L., Schellaert, W., Martínez-Plumed, F., Moros-Daval, Y., Ferri, C., and Hernández-Orallo, J. Larger and more instructable language models become less reliable. *Nature*, 2024b.
- Zhou, Y., Liu, H., Chen, Z., Tian, Y., and Chen, B. Gsm-infinite: How do your llms behave over infinitely increasing context length and reasoning complexity? *arXiv preprint arXiv:2502.05252*, 2025.
- Zou, C., Guo, X., Yang, R., Zhang, J., Hu, B., and Zhang, H. Dynamath: A dynamic visual benchmark for evaluating mathematical reasoning robustness of vision language models. In *ICLR*, 2025.

## Appendix: Recitation over Reasoning: How Cutting-Edge Language Models Can Fail on Elementary School-Level Reasoning Problems?

The appendix has two parts. In Sec. A, we introduce the prompts used in our experiments; in Sec. B, we show more experiment results on our RoR-Bench.

### A. Prompts

#### A.1. Prompts and OpenAI-o1-1217 Response in Fig. 1

Due to space limit, we simplified the response of OpenAI-o1-1217 in Fig. 1. Here we show the full prompts and responses below:

# Conversation 1 (original problems)  
<https://chatgpt.com/share/67f45721-dc98-800d-a238-47fdc8b2c33d>  
**User:** A boat sails at 15km/h in still water. It spends 8 hours traveling from upstream location A to downstream location B, with the water speed 3km/h. How long does it take to return from B to A?  
**OpenAI-o1-1217:**  
*(Thinking process shown by OpenAI website)* 哎呀,看来这是个航向的谜题。立即想到了船在水流速度3 km/h下的实际航速和时间。  
 船从A到B是下行。水流速度3 km/h, 加上船的网速15 km/h, 总速18 km/h。行程长144 km, 耗时8小时。(Oops, it is a puzzle about sailing directions. I am immediately thinking of the actual speed and time of the boat at 3km/h current speed. The boat goes downstream from A to B. The current speed is 3km/h, with the net speed of the boat being 15km/h, the total speed is 18km/h. The journey is 144km, and requires 8 hours.)  
*(Answer)* First, determine the downstream speed and upstream speed of the boat:  
 - \*\*Downstream speed\*\* = 15 km/h (boat speed in still water) + 3 km/h (current speed) = 18 km/h  
 - \*\*Upstream speed\*\* = 15 km/h (boat speed in still water) - 3 km/h (current speed) = 12 km/h  
 Because the boat took 8 hours to go downstream:

$$\begin{aligned} \text{Distance from A to B} &= \text{Downstream speed} \times \text{Time} \\ &= 18 \times 8 = 144 \text{ km} \end{aligned}$$

To find the time to go back upstream from B to A:

$$\text{Time upstream} = \frac{\text{Distance}}{\text{Upstream speed}} = \frac{144}{12} = 12 \text{ hours}$$

Hence, it takes \*\*12 hours\*\* to return from B to A.

# Conversation 2 (modified problems)  
<https://chatgpt.com/share/67edc3d9-b10c-800d-819f-ddeba958b2ec>  
**User:** A boat sails at 15km/h in still water. It spends 8 hours drifting from upstream location A to downstream location B, with the water speed 3km/h. How long does it take to return from B to A?  
**OpenAI-o1-1217:**  
*(Thinking process shown by OpenAI website)* 船在静水中的速度是15公里/小时, 水流速度是3公里/小时。从A到B用了8小时, 那么我们可以先计算出A到B的距离, 然后再用相同的速度反推从B到A的时间。(The boat sails at 15km/h, and the current speed is 3km/h. It takes 8 hours from A to B, then we can first calculate the distance from A to B, then use the same speed to induce the time from B to A.)  
*(Answer)* First, compute the distance from A to B. Going downstream, the boat's effective speed is  $15 + 3 = 18$  km/h. Since the trip from A to B takes 8 hours:

$$\text{Distance} = 18 \text{ km/h} \times 8 \text{ h} = 144 \text{ km.}$$

On the return trip from B to A, the boat travels upstream with an effective speed of  $15 - 3 = 12$  km/h. Therefore, the time needed is:

$$\text{Time} = \frac{144 \text{ km}}{12 \text{ km/h}} = 12 \text{ hours.}$$

So it takes 12 hours to go back from B to A.

See <https://chatgpt.com/share/67edc40f-e658-800d-9bcd-7cb268a6f8c9> for another example of using the word “floating” and <https://chatgpt.com/share/67ee1ba2-089c-800d-ad01-d7ee50dab39d> for “drifting”. Interestingly, o1 can recognize the issue when we use the two words together (<https://chatgpt.com/share/67ee1bda-5060-800d-ad44-b6c4e5741185>), but not separately.

#### A.2. Prompts for the Judge

We use the same prompt for GPT-4o-1120 judges for all our experiments. We use Chinese prompts as the problems are Chinese; we provide English translations as well in the boxes below. Interestingly, we found the emphasis of “do not question the ground truth answer” very important; *the recitation issue of LLMs is so severe that without such prompt, the LLM judge will question the problem instead of the student’s answer even with ground truth answer available.*

# Chinese Prompt

你是一个专业的老师，现在有一道题目，你需要判断学生的回答是否和标准答案一致。题目和答案保证是绝对正确的，不会有错字，因此你要严格按照题目的字面意思评判。不要质疑标准答案有问题。如果学生的回答和标准答案一致，请打1分；否则请打0分。标准答案经常只含有答案，而学生的回答一般会带有过程；你只需要根据学生的结论是否与标准答案相符做出评价。学生的结论和标准答案必须本质一样，但表述可以稍有不同；例如，学生的答案是50又1/2或50.5，标准答案是101/2，则也可以算作正确。你的判断应遵循如下格式：你的输出在【评价】后开始。首先是对学生回答的分析（不超过300字），然后以“分数: [[0]]”或“分数: [[1]]”结束。下面开始判断：

【题目】 (Problem description)

【标准答案】 (Answer)

【学生回答】 (LLM output)

# English Translation

You are a professional teacher. Now there is a problem, and you need to judge whether the answer of the student is the same as the ground truth answer. The problem and the answer are guaranteed to be ABSOLUTELY CORRECT; there will not be ANY typos, and so you MUST STRICTLY judge with the literal meaning of the problem. DO NOT QUESTION THE GROUND TRUTH ANSWER. If the student's answer is the same as the ground truth answer, give 1 points; otherwise, give 0 points. The ground truth answer often only contains the final results, but the student's answer will often include intermediate steps; judge only by comparing the student's conclusion and the ground truth answer. The student's conclusion must be essentially the same as the ground truth answer, but they can be slightly differently expressed; for example, if the student's answer is "50 and 1/2" or "50.5", while the ground truth answer is 101/2, then it can be counted as correct. Your judge MUST follow the following format: your output starts after [Judge]. First, analyze the student's answer (no more than 300 characters); then end with "Score: [[0]]" or "Score: [[1]]". Now begin your judgment:

[Problem] (Problem description)

[Ground Truth Answer] (Answer)

[Student's Answer] (LLM output)

A.3. Prompts for Few-Shot In-Context Learning

In this section, we provide an 1-shot example to show the prompts for few-shot ICL experiments; for cases with more shots, the problems are added in the same format as the first example before the last, target problem. We again show both the original Chinese version and the English translation. The red part is the Forced Correct (FC) prompt, which is optional.

请回答下面的问题。题目保证没有错误，请严格按照题目字面所写的问题回答。以下是一个例子：

【问题】有四个人要在游过一条河，他们只有一个游泳圈，且每次最多只能两个人一起使用游泳圈游过河，使用游泳圈时必须有人携带。四个人单独游过河的时间分别是1分钟、2分钟、5分钟、10分钟。如果两人一起使用游泳圈游过河，所需要的时间就是游得慢的那个人单独游过河的时间。请问，他们如何在17分钟内全部游过河？

【答案】让1分钟和2分钟的人先一起使用游泳圈游过河，花费2分钟，然后1分钟的人带着游泳圈游回来，花费1分钟。5分钟和10分钟的人一起使用游泳圈游过河，花费10分钟，接着2分钟的人带着游泳圈游回来，花费2分钟。1分钟和2分钟的人再次一起使用游泳圈游过河，花费2分钟。总共花费的时间为： $2+1+10+2+2=17$ 分钟。

下面是你要回答的问题：

【问题】有四个人要在游过一条河，他们只有一个游泳圈，且每次最多只能两个人一起使用游泳圈游过河，不会游泳的人必须使用游泳圈，使用时必须有人携带。四个人单独游过河的时间分别是1分钟、2分钟、5分钟、10分钟。其中前三个人均会游泳。如果两人一起使用游泳圈游过河，所需要的时间就是游得慢的那个人单独游过河的时间。请问，他们如何在17分钟内全部游过河？

Please answer the following problems. **The problems are guaranteed to be correct; please strictly follow the literal meaning of the problem.** Here is an example:

[Problem]

Four people need to swim across a river. They have only one swimming ring, and at most two people can use it at the same time. Someone must carry the swim ring whenever it is used. The time it takes for each person to swim across the river individually is 1 minute, 2 minutes, 5 minutes, and 10 minutes respectively. If two people use the swim ring together to cross the river, the time it takes is equal to the time of the slower swimmer. The question is: how can all four people cross the river within 17 minutes?

[Answer]

Let the 1-minute and 2-minute people use the swim ring to cross the river first, which takes 2 minutes. Then the 1-minute person brings the swim ring back, taking 1 minute. Next, the 5-minute and 10-minute people cross the river together using the swim ring, which takes 10 minutes. After that, the 2-minute person brings the swim ring back, taking 2 minutes. Finally, the 1-minute and 2-minute people cross the river together again using the swim ring, taking 2 minutes.

The total time spent is:  $2 + 1 + 10 + 2 + 2 = 17$  minutes.

Now here is the problem you need to answer:

[Problem]

Four people need to swim across a river. They have only one swimming ring, and at most two people can use it at the same time. Anyone who cannot swim must use the swim ring, and it must be carried by someone while in use. The times it takes for each person to swim across the river individually are 1 minute, 2 minutes, 5 minutes, and 10 minutes respectively. Among them, the first three people can swim. If two people use the swim ring together to cross the river, the time required is equal to the time it takes for the slower person to cross the river alone. The question is: how can all four people cross the river within 17 minutes?



Interestingly, when we test this English translation with OpenAI-o1-1217, we found o1, even with 1-shot, is again tricked into the classic paradigm that the swimming ring must be carried back. The ground truth answer of this target problem, however, is to directly let the third and fourth people use the swimming ring, and the first two people swim through the river, such that everything can be done within 10 minutes; no swimming ring needs to be taken back.

tion (<https://chatgpt.com/share/67f4601a-149c-800d-8be1-d4ebe97c315b>) and the explicit version above (<https://chatgpt.com/share/67f444ed-8148-800d-810d-62c54ea1636a>).

## B. More Experiment Results

### B.1. Best-of-5 Results

Tab. 6 (for text-based problems) and Tab. 7 (for vision-based problems) shows the best-of-5 result of the experiments conducted in Sec. 4.1. The conclusion is very similar to those in Sec. 4.1, indicating that the problem is hard to fix with LLM alignment techniques such as reinforcement learning (Ouyang et al., 2022).

### B.2. Greedy Decoding Results

Tab. 8 (for text-based problems) and Tab. 9 (for vision-based problems) shows the average score of LLMs doing greedy-decoding (i.e. temperature=0) in the experiments conducted in Sec. 4.1. The conclusion is similar to those in Sec. 4.1.

### B.3. English Version of Fig. 2

The response for OpenAI-o1-1217 for our English translation of the text problems in Fig. 2 can be seen in <https://chatgpt.com/share/67f45d45-a694-800d-a4c8-3c763e93400f>.

In later experiments, the response of GPT-4.5 points out that there is a possibility of ambiguity - “traveling in opposite directions” could also possibly mean traveling towards each other, as if the two cars are located in the west and east respectively, the car on the west will need to travel east and the car on the east will need to travel west (<https://chatgpt.com/share/67f45e46-5eb0-800d-8f1c-5dec1f3c9f1d>), which are indeed opposite directions. To further address such concerns, we test the result of literal translation from Chinese (“traveling back to back”), and the following version with much more explicit explanation:

*Two cars start simultaneously from two cities that are 300 km apart, one on the east and the other on the west. The car on the east travels towards the east, and the car on the west travels towards the west. One car has a speed of 60 km/h, while the other has a speed of 70 km/h. How many hours will it take for them to meet?*

Surprisingly, we found OpenAI-o1-1217 still fails on both the literal transla-

How Cutting-Edge Language Models Can Fail on Elementary School-Level Reasoning Problems?

Model Name	Original Bo5	Modified Bo5	Original + FC	Modified + FC
OpenAI-o1-1217	93.67	43.03	94.30	56.96
DeepSeek-R1	92.41	34.81	92.41	39.87
Hunyuan Turbo-S	92.41	26.58	91.14	23.42
GPT-4.5-Preview	91.14	38.60	87.97	49.37
OpenAI-o3-mini-high	91.14	34.81	91.77	39.87
Gemini-2.0 Flash-0121 (CoT)	91.14	32.91	87.97	41.14
Gemini-2.0 Pro-0205	91.14	32.91	87.97	41.14
Claude 3.7 Sonnet	91.14	39.87	86.08	49.37
Claude 3.7 Sonnet (CoT)	90.51	37.34	90.51	42.41
Ernie-4.5	88.61	26.58	87.34	29.11
GLM-4-Plus	86.70	29.11	82.27	31.01
GPT-4o-1120	86.70	29.11	81.65	44.94
Qwen-max-0125	85.44	36.08	84.17	37.97
DeepSeek-V3	84.81	33.54	84.17	40.51
StepFun Step-2-16k	84.81	27.85	82.28	28.48
Yi-Lightning	84.81	25.32	85.44	31.01
QwQ-32B-Preview	84.17	39.87	84.17	37.97
Gemini-2.0 Flash-0121	84.17	32.91	70.89	36.08
Minimax-Text-01	82.91	31.64	84.17	26.58
Claude 3.5 Sonnet	82.28	32.91	83.54	41.14
Qwen-2.5-14B-Instruct	81.65	29.75	81.65	30.38
Mistral-Large-2	79.11	30.37	72.15	34.81
Nova-Pro	78.48	30.37	79.11	35.44
GPT-4o-mini-0718	75.95	29.74	74.68	31.01
Qwen-2.5-7B-Instruct	56.32	23.41	53.80	22.78
Avg. Decrease	N/A	-52.89( $\pm$ 6.60)	-2.00( $\pm$ 3.23)	-48.35( $\pm$ 7.68)

Table 6: Best-of-5 (Bo5) Results on text-based problems of RoR-Bench; the conclusion is similar to that with average score.

Model Name	Original Bo5	Modified Bo5	Original + FC	Modified + FC
OpenAI-o1-1217	98.25	29.82	96.49	42.11
GPT-4.5-Preview	96.49	22.81	82.46	43.86
GPT-4o-1120	91.23	19.30	89.47	31.58
Gemini-2.0 Flash-0121 (CoT)	84.21	43.86	66.67	49.12
Gemini-2.0 Pro-0205	78.95	36.84	73.68	42.11
Claude 3.7 Sonnet (CoT)	78.95	49.12	80.70	56.14
GPT-4o-mini-0718	73.68	35.09	80.70	29.82
Claude 3.5 Sonnet	71.92	45.61	61.40	49.12
Qwen2.5-VL-max	70.18	42.11	66.67	42.11
Qwen2.5-VL-72B	70.18	42.11	64.91	42.11
GLM-4v-Plus	68.42	43.86	64.91	42.11
Claude 3.7 Sonnet	66.67	45.61	63.15	54.39
Nova-Pro	64.91	57.89	71.93	38.60
SenseChat-Vision	64.91	43.86	75.44	42.11
StepFun-1v-32k	64.91	33.33	68.42	28.07
Gemini-2.0 Flash-0121	64.91	30.17	53.68	35.79
Qwen2.5-VL-7B	59.65	47.37	61.40	40.35
Avg. Decrease	N/A	-35.27( $\pm$ 19.49)	-2.73( $\pm$ 7.67)	-32.88( $\pm$ 13.38)

Table 7: Best-of-5 (Bo5) Results on vision-based problems of RoR-Bench; the conclusion is similar to that with average score.

How Cutting-Edge Language Models Can Fail on Elementary School-Level Reasoning Problems?

Model Name	Original Score	Modified Score	Original + FC	Modified + FC
Hunyuan Turbo-S	88.60	19.62	87.97	17.72
OpenAI-o3-mini-high	86.08	28.48	83.54	29.74
DeepSeek-R1	86.08	18.99	88.61	27.22
OpenAI-o1-1217	85.44	31.01	88.61	40.51
Gemini-2.0 Flash-0121 (CoT)	84.81	23.42	79.75	24.68
GPT-4.5-Preview	83.54	26.58	77.22	36.08
Claude 3.7 Sonnet (CoT)	81.65	24.05	78.48	39.24
Ernie-4.5	81.65	21.52	80.38	23.42
Gemini-2.0 Pro-0205	78.48	24.68	41.14	32.91
Gemini-2.0 Flash-0121	78.48	22.78	60.76	25.95
Qwen-max-0125	75.95	20.25	75.32	23.42
GLM-4-Plus	75.32	15.82	70.89	22.78
Claude 3.7 Sonnet	74.68	25.32	70.89	35.44
GPT-4o-1120	74.05	23.42	70.89	25.95
Claude 3.5 Sonnet	73.42	23.42	66.46	31.01
QwQ-32B-Preview	72.15	18.99	68.99	22.79
DeepSeek-V3	70.25	17.09	72.15	25.95
Minimax-Text-01	69.62	18.99	65.82	20.25
StepFun Step-2-16k	69.62	17.72	72.15	21.52
Yi-Lightning	68.35	13.92	62.66	22.79
Qwen-2.5-14B-Instruct	65.82	19.62	66.56	20.89
Mistral-Large-2	63.92	18.99	52.53	27.84
Nova-Pro	61.39	20.25	57.59	18.99
GPT-4o-mini-0718	61.39	19.62	60.76	20.89
Qwen-2.5-7B-Instruct	37.34	10.76	34.81	16.46
Avg. Decrease	N/A	-52.91( $\pm$ 8.67)	-4.53( $\pm$ 8.18)	-47.75( $\pm$ 9.52)

Table 8: Results on text-based problems of RoR-Bench with greedy decoding; the conclusion is similar to that with temperature 0.7.

Model Name	Original Score	Modified Score	Original + FC	Modified + FC
GPT-4.5-Preview	94.74	14.04	71.93	42.11
OpenAI-o1-1217	91.23	24.56	94.74	26.32
GPT-4o-1120	85.96	14.04	84.21	26.32
Gemini-2.0 Flash-0121 (CoT)	73.68	28.07	63.15	42.11
Gemini-2.0 Flash-0121	71.93	28.07	57.89	40.36
Gemini-2.0 Pro-0205	70.18	35.09	68.42	40.35
GLM-4v-Plus	68.42	43.86	66.67	42.11
GPT-4o-mini-0718	68.42	31.58	80.70	28.07
Claude 3.7 Sonnet (CoT)	68.42	31.58	64.91	43.86
Qwen2.5-VL-72B	66.67	36.84	66.67	42.11
Claude 3.5 Sonnet	64.91	33.33	59.65	45.61
Qwen2.5-VL-max	63.16	36.84	66.67	42.11
SenseChat-Vision	59.65	35.09	70.18	38.60
StepFun-1v-32k	59.65	33.33	64.91	28.07
Nova-Pro	57.89	50.88	70.18	38.60
Claude 3.7 Sonnet	56.14	31.58	61.40	40.35
Qwen2.5-VL-7B	52.63	38.60	59.65	42.11
Avg. Decrease	N/A	-36.84( $\pm$ 19.86)	-0.10( $\pm$ 9.42)	-30.85( $\pm$ 15.32)

Table 9: Results on image-based problems of RoR-Bench with greedy decoding; the conclusion is similar to that with temperature 0.7.