
Not All LLM Reasoners Are Created Equal

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We study the depth of problem-solving capabilities of LLMs, and to what extent
2 they perform mathematical reasoning in a compositional manner. To this end,
3 we create a new benchmark by composing pairs of existing math word problems
4 together so that the answer to the second problem depends on correctly answering
5 the first problem. We measure the difference between the performance of solving
6 each question independently and solving the compositional pairs as the reasoning
7 gap of a model. Our findings reveal a significant reasoning gap in most frontier
8 LLMs. This gap is more pronounced in smaller and more cost-efficient models.
9 The objective of this study is not to introduce yet another benchmark, but rather to
10 provide a case study aimed at gaining deeper insights into current models' reasoning
11 abilities, and to reassess existing established training methods and benchmarks.

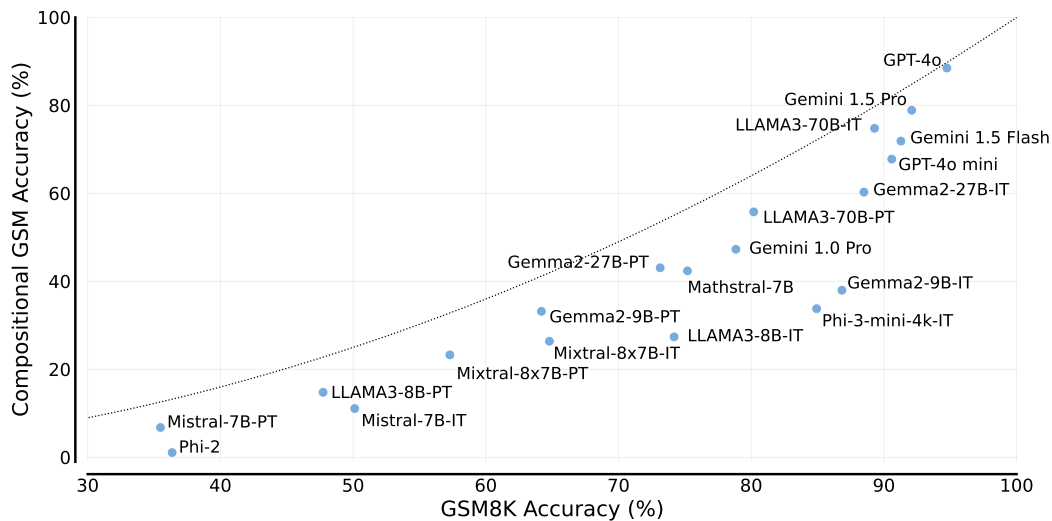


Figure 1: **Reasoning Gap:** Pairs of GSM8K test questions are chained together so that the answer of the first question (Q_1) is a variable in the second one (Q_2). The model is required to correctly answer both questions to solve the problem. If a model has an accuracy of S_1 on the Q_1 set, and S_2 on Q_2 set, then the expected Compositional GSM accuracy is $S_1 \times S_2$. The x-axis corresponds to the geometric mean $\sqrt{S_1 \times S_2}$, labeled GSM8K accuracy for simplicity. The trend-line $y = x^2$ is the expected Compositional GSM accuracy.

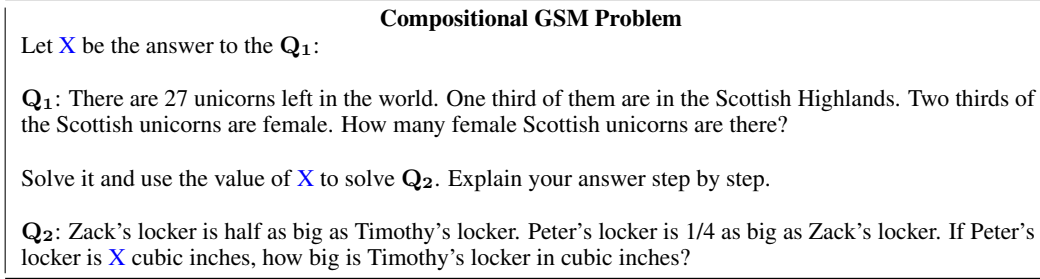


Figure 2: **Example Problem from the Compositional GSM benchmark.** The answer of Question-1 (Q_1) is a variable X in Question-2 (Q_2). Therefore, the model has to be able to solve the first question correctly in order to solve the second question. The new final answer of Question-2 is calculated by modifying its code-form solution and executing it. Question-1 and the number to modify in Question-2 are chosen to have a new final answer which is a positive integer not too far from the old answer of Question-2.

1 Introduction

The strong performance of large language models (LLMs) on high-school and college-level math reasoning benchmarks (Dubey et al., 2024; Google, 2024; OpenAI, 2023b), has led to the common belief that LLMs have “mastered” grade-school math, particularly as measured by the GSM8K benchmark (Cobbe et al., 2021a). This apparent mastery of grade-school math problems raises a deeper question: do LLMs truly grasp the underlying concepts or do they mostly rely on dataset contamination or memorization (Srivastava et al., 2024b)? For example, a recent examination on private “held-out” grade-school problems (Zhang et al., 2024a) reveals that while frontier closed-source LLMs show minimal signs of overfitting, some open-weights models show systematic overfitting, possibly due to test data contamination.

In this work, we perform a case study to evaluate how well LLMs can combine learned concepts to solve unseen problems, to probe the brittleness of their reasoning abilities. To do so, we introduce *Compositional GSM*, a two-hop version of GSM8K with higher difficulty, where each problem chains two test questions together such that the answer to the first question is used as a variable in the second question (Figure 2). As LLMs can easily solve grade-school math problems, they should also be capable of solving combinations of those problems. As such, we measure the gap between their performance on solving the questions individually and on Compositional GSM. Specifically, we benchmark frontier open-weights and closed LLMs, including Gemini (Google, 2023, 2024), Gemma2 (Gemma Team et al., 2024), Llama-3 (AI@Meta, 2024), GPT (OpenAI, 2023a), Phi (Abdin et al., 2024) and Mistral families (Jiang et al., 2024).

Here are our key findings:

- Most models exhibit a gap between their performance on GSM8K test set and Compositional GSM (Figure 1).
- This reasoning gap is particularly larger in small and more cost-efficient models (Figure 3 and Figure 5).
- Instruction-following tuning of LLMs heavily favours the original GSM8K split (Figure 6).
- Finetuning with human data and synthetic data results in a similar reasoning gap trend (Figure 8).
- Math-specialized LLMs exhibit reasoning gap and task-specific overfitting (Figure 7).
- Smaller models benefit more from generating code rather than natural language Chain-of-Thought (CoT) to solve Compositional GSM problems (Figure 11).

2 Compositional Grade-School Math (GSM)

Each question in compositional GSM consists of two questions, Question-1 and Question-2, from a subset of 1200 examples of the original GSM8K test set. The final answer of Question-1 is referred to as X which is a variable in Question-2 (Figure 2). The final answer of Question-2 is obtained by substituting X and solving it. The choice of Question-1 and the number to modify and replace with X in Question-2 was made in a way such that the new final answer of Question-2 is different from its old final answer, and is a positive integer not too far from the old final answer. To obtain

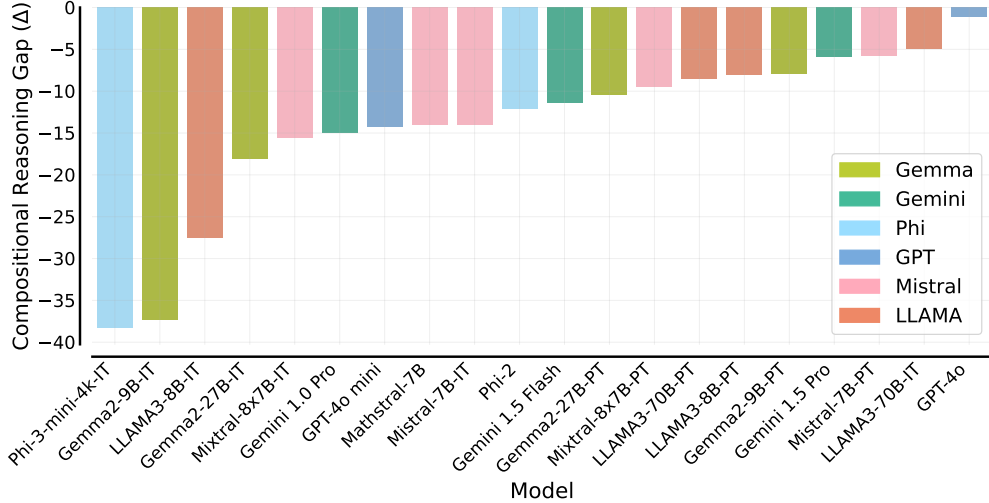


Figure 3: **Compositional Reasoning Gap** of notable open-weights and closed-source LLMs. Smaller, more cost-efficient and math specific models have a bigger gap.

49 the new final answer of Question-2 automatically, we replace a number in the code-form solution of
 50 Question-2. We used a slightly modified version of the code-form solutions from Gao et al. (2023).
 51 The new final answer is the result of executing the code with the new number. We put significant
 52 efforts into ensuring that the modification to Question-2 is sensible. Figure 10 shows the distribution
 53 of final answers (magnitude) of the original test set of GSM8K and compositional GSM. Both test
 54 sets have a similar distribution of final answers.

55 **Quality Checks** To make sure that the modified questions are sensible and logical, we generated
 56 16 candidate solutions per modified question from GPT-4o and Gemini 1.5 Pro. We filtered those
 57 questions for which less than 4 (out of 16) agree with the expected final answer from code execution.
 58 We checked these questions manually and modified them if needed so that they are logical (about
 59 25% of questions).

60 **Reasoning Gap** Question-1 and Question-2 in our compositional queries are from the original
 61 test split $\mathcal{D}_{\text{original}}$, and the modified test split $\mathcal{D}_{\text{modified}}$ respectively. Assuming that a model has an
 62 accuracy of S_1 on $\mathcal{D}_{\text{original}}$ and S_2 on $\mathcal{D}_{\text{original}}$, it is expected for it to have an accuracy of $S_1 \times S_2$
 63 on the compositional split $\mathcal{D}_{\text{comp}}$. We report the following as the compositional reasoning gap score,

$$\text{Compositional reasoning gap} : \Delta = S_c - S_1 \times S_2 \quad (1)$$

64 where S_c is the performance of the model on $\mathcal{D}_{\text{comp}}$.

65 3 Experiments & Results

66 **Setup** We evaluate each model on three test sets: **1) the original GSM8K test split**, **2) the modified**
 67 **GSM8K test split** which are the questions with X being substituted, and **3) the compositional GSM**
 68 test set. Each test set has 1200 examples. Following Zhang et al. (2024a), we evaluate all models with
 69 an 8-shot prompt (Appendix A) for the original and modified GSM8K test splits. We also created
 70 a similar 8-shot prompt (Appendix D) for the compositional GSM questions. We evaluate GPT-4o,
 71 GPT-4o mini, LLaMA3-70B and 8B (PT and IT), Phi 2, Phi-3-mini-instruct, Gemini 1.0, 1.5 Flash
 72 and 1.5 Pro, Gemma2 9B and 27B (PT and IT), Mistral-7B (PT and IT), Mixtral-8x7B (PT and IT)
 73 and Mathstral-7B. All models are sampled with temperature 0, and pass@1 (Chen et al., 2021) is
 74 used to measure the performance on each test split. Some of the models require a preamble prefixed
 75 to the 8-shot prompt in order to output in a consistent format (Appendix B). We test both cases and
 76 report the best performance for each model.

77 The distance to the trend-line in Figure 1 shows the reasoning gap of models. The x-axis corresponds
 78 to $\sqrt{S_1 \times S_2}$, which is the geometric mean of the accuracies on the set of Q_1 and Q_2 independently.

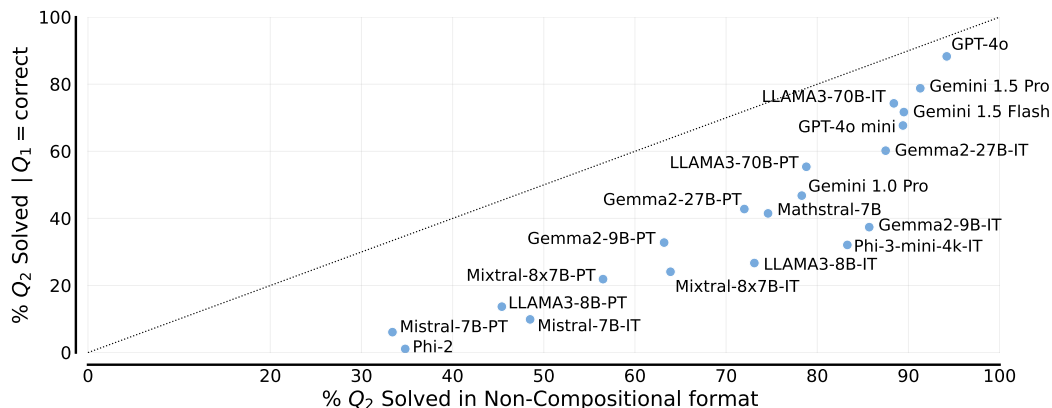


Figure 4: **Can models answer the second question if they have correctly answered the first one?** Here, we compare how often models are able to solve a question independently to how often they are able to solve them in the compositional format given that the first question is solved correctly. This is an alternate measurement of the compositional reasoning gap. If a model can solve a question independently, it should be able to solve it in a compositional setting given that the prerequisites are met. The gap from the diagonal line suggests that some models have overfit to the format of GSM8K type questions. While models may correctly answer the first question, they frequently makes subtle errors and miss key details when solving the second question.

79 We find that most models fall below expectation on Compositional GSM. Specifically, it is evident
 80 that cost-efficient models have a larger gap than more expensive models. Our analysis shows that
 81 tuning models for instruction following introduces more gap in their reasoning capabilities. In the
 82 following sections, we will further examine the models and discuss their shortcomings.

83 3.1 Does Solving Question-1 Guarantee Solving Question-2?

84 Correctly solving Question-1 is a prerequisite to solve Question-2 in the compositional format.
 85 In Figure 4, we look at how often models are able to solve a question independently versus how
 86 often can they solve it given they have correctly solved the previous question in the compositional
 87 format. What remains for the model to do here is to substitute X and solve Q_2 . The deviation from
 88 the diagonal line indicates that certain models may have become too specialized in handling GSM8K-
 89 style questions, and are unable to answer a second question having generated the solution to the first
 90 question. Our qualitative analysis shows that when given two questions, the model might answer the
 91 first one correctly, but often makes subtle errors and overlooks details, leading to inaccurate reasoning
 92 and solution for the second question.

93 3.2 Cost-Efficient LLMs Reason Differently

94 The reasoning abilities of cost-efficient LMs has been rapidly improving over time, as evaluated using
 95 standard benchmarks (Bansal et al., 2024). For example, GPT-4o mini and Gemini 1.5 Flash both
 96 achieve above 90% accuracy on GSM, while costing 25 – 35× cheaper than GPT-4o and Gemini
 97 1.5 Pro respectively. This progress could be attributed to several factors, such as better pretraining
 98 data (AI@Meta, 2024), and knowledge distillation (Agarwal et al., 2024; Team et al., 2024). To this
 99 end, we investigate whether these reasoning gains on GSM8K still persist on Compositional GSM.

100 We study four family of models, each comprising both a high-cost and low-cost option, where cost
 101 is measured via parameter count or API pricing. Figure 5 shows the original GSM8K test split
 102 performance and Compositional GSM performance for all models. The numbers above the bars
 103 represents the reasoning gap defined in Eq 1. While cheaper models perform comparably on the
 104 original GSM8K test, they exhibit a notable drop in performance on the Compositional GSM test set.
 105 Furthermore, our results demonstrates that the additional cost does not translate into a proportional
 106 gain in performance on the Compositional GSM test. Specifically, Gemini 1.5 Pro is 35x more
 107 expensive than 1.5 Flash, but the gap narrows from -11.3% to -5.8%. Similarly, GPT-4o costs 25x
 108 more than GPT-4o mini, while the gap is reduced from -14.2% to -1.1%.

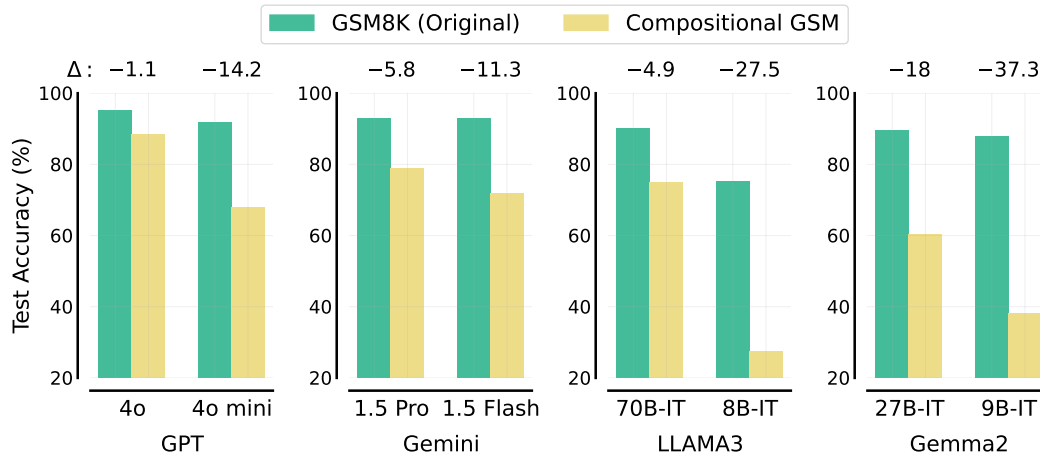


Figure 5: **Cost efficient LLMs reason differently**: showing four family of models, each having a high-cost and low-cost option. Although the cheaper models perform similarly on the original GSM8K test, they show a significant decline in performance on the Compositional GSM test. Furthermore, it is clear that additional cost does not yield a corresponding decrease in the reasoning gap.

109 3.3 Instruction-Tuning Impacts LLM Reasoning Differently

110 We compare pretrained and instruction-following tuned versions of models in three families of Mistral,
 111 LLAMA3 and Gemma2. Figure 6 illustrates this comparison, along with the performance gains from
 112 instruction-tuning, displayed above bars for each test set. On small models (top row), this comparison
 113 shows that current instruction-tuning is heavily optimized for GSM8K questions. Instruction-tuning
 114 leads to a significantly larger improvement on the original GSM8K test set than the Compositional
 115 GSM test across model families. However, this trend does not apply to larger models (bottom row),
 116 where the improvements are inconsistent.

117 3.4 Reasoning Gap in Math-Specialized LLMs

118 Math-specialized models are LLMs tailored to solve problems in the specific domain of mathematical
 119 reasoning. Such LLMs have an extensive coverage of diverse mathematical fields, so can they
 120 generalize to Compositional GSM or do they demonstrate task-specific overfitting? To answer this
 121 question, we evaluated three mathematical LLMs, namely NuminaMath-7B-CoT (Beeching et al.,
 122 2024), Mathstral-7B and Qwen2.5-Math-7B-IT (Yang et al., 2024) on GSM8K and Compositional
 123 GSM (Figure 7). We observe that these math-specialized LLMs exhibit reasoning gaps comparable
 124 to those in other models of similar size within our analysis. For instance, Qwen2.5-Math-7B-IT is
 125 reported to have an accuracy of 83.6% on MATH (Hendrycks et al., 2021), but has a large reasoning
 126 gap of -21.9% .

127 3.5 Supervised Finetuning Can Lead to Task-Specific Overfitting

128 Finetuning models on task specific problems is a common strategy to improve reasoning performance.
 129 In this section, we explore how it impacts the performance on Compositional GSM. We investigate
 130 the performance of Gemma2 27B PT as we finetune it on the original GSM8K training data, and
 131 self-generated rationales (aka synthetic data) to identify any difference in the characteristics of these
 132 two sources. We collect self-generated rationales which results in correct final answers for all GSM8K
 133 training queries.

134 See Appendix C for details of data generation and training for this set of experiments. We evaluated
 135 intermediate checkpoints (at 50, 100 and 400 training steps) from both settings on GSM8K original
 136 test split and Compositional GSM split (Figure 8). We observe a similar pattern for both settings.
 137 The Compositional GSM performance increases with some training (up to 100 steps), but drops with
 138 more training steps while GSM8K test performance keeps increasing, which suggests overfitting.

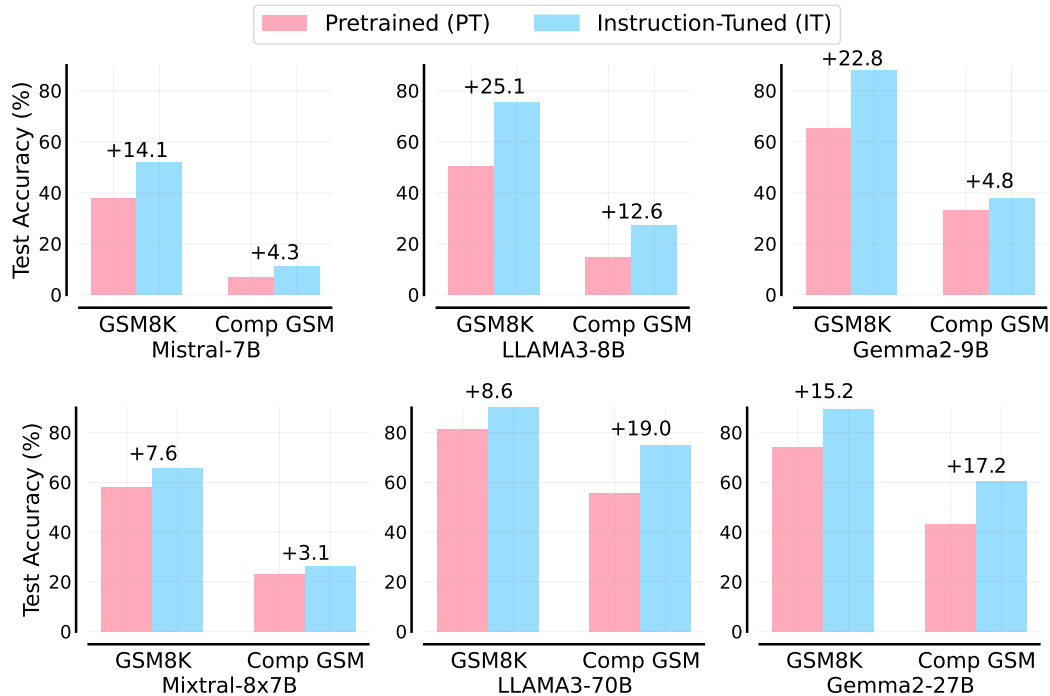


Figure 6: **Impact of Instruction-Tuning on Reasoning Gap:** comparing pretrained and instruction-following tuned variant of models from Mistral, LLAMA3 and Gemma2 families. Numbers above bars represent improvements from instruction-tuning on each set. For smaller models (top), we observe that instruction-tuning is highly optimized for GSM8K questions, which results in a greater improvement on the original GSM8K test set compared to the Compositional GSM test. However, this pattern does not hold for larger models (bottom).

139 Our results show that training on synthetic data generally leads to a higher Compositional GSM
 140 performance. We did not observe further improvements on either test splits after 400 training steps.

141 3.6 Models Get Distracted Easily

142 Assuming an LLM answers a question correctly, it is somewhat expected that it would answer the
 143 same question correctly with additional context. Figure 9 shows how often a model answers a question
 144 (from Q_1 set) correctly on the x-axis, and how often it answers it correctly in our compositional
 145 format, as Q_1 . Ideally, models should be on the $x = y$ line, but we observe that most of the models
 146 fall short of this expectation. Examining the responses from models with greater deviations from the
 147 trendline in Figure 9 reveals that they frequently make subtle errors. They often overlook important
 148 details, such as missing a reasoning step related to *each* in the question or omitting a multiplication
 149 step when the question specifies *in a month*. The models generally adhere well to the output format
 150 provided in the 8-shot context, resulting in negligible instances of non-extractable answers. This
 151 distraction is caused by the existence of a second question Q_2 in the prompt. Such failures lead to
 152 not being able to correctly answer Q_1 , which subsequently impairs the models' ability to answer Q_2
 153 correctly.

154 4 Related Work

155 **Mathematical Reasoning with Language Models.** The performance of large language models on
 156 mathematical reasoning tasks has seen rapid development in recent years (Lightman et al., 2023;
 157 McAleese et al., 2024; Shao et al., 2024). Majority of these improvements are due to vast amount of
 158 pretraining and post-training on math related data (Google, 2024; Lewkowycz et al., 2022; OpenAI,
 159 2023b) and various self-improvement (Gulcehre et al., 2023; Singh et al., 2023; Zelikman et al.,
 160 2022) and specialized prompting techniques. We do not plan to evaluate all of the methods from the
 161 existing literature on Compositional GSM benchmark. These methods could potentially influence

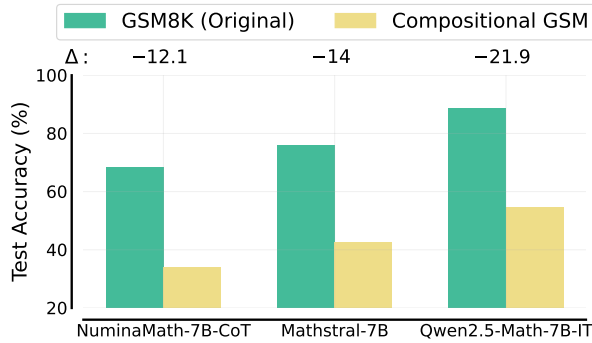


Figure 7: **Mathematical LLMs:** We compare the performance of three models specifically designed for mathematical problem solving on GSM8K and Compositional GSM sets. Aiming to explore whether extensive specialized training in mathematics can bridge the reasoning gap observed among models of similar size or family. Surprisingly, we find that mathematical LLMs exhibit similar reasoning gaps and signs of overfitting to standard benchmarks.

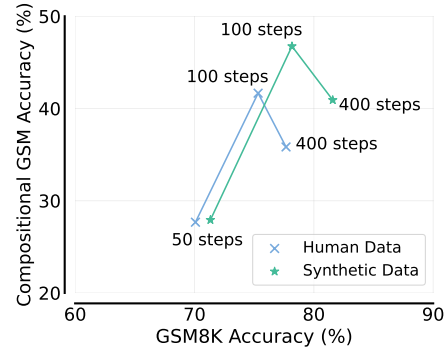


Figure 8: **Human Data v.s. Synthetic Data:** We finetune Gemma2 27B on the original GSM8K training data, and self-generated rationales. In both settings, after 100 training steps, Compositional GSM test performance drops while original GSM8K test performance keeps increasing. No further improvements were observed on either test split after 400 training steps.

162 and reduce the reasoning gap at inference. Future research could investigate how they impact the
 163 reasoning gap and overall robustness.

164 **Prompting Techniques to Improve Reasoning in Language Models.** There is a broad range of
 165 ongoing research aimed at improving reasoning abilities of language models. For instance, specialized
 166 prompting techniques designed to elicit explicit reasoning (Li et al., 2024a; Wei et al., 2022; Yao
 167 et al., 2023a) and methods to provide models with more inference compute (Wang et al., 2023c) or
 168 verification (Cobbe et al., 2021b; Hosseini et al., 2024; Lightman et al., 2023; Wang et al., 2023b;
 169 Zhang et al., 2024b) to solve problems.

170 **Mathematical Reasoning Robustness.** Our work is heavily inspired by the study of mathematical
 171 reasoning gap via rewrites of GSM8K test queries (Zhang et al., 2024a). Srivastava et al. (2024a)
 172 introduce a framework for thoroughly assessing the mathematical reasoning abilities of LLMs by
 173 employing different functional variants of benchmarks. Others have investigated the robustness of
 174 mathematical reasoning abilities of LLMs via adversarial examples (Anantheswaran et al., 2024; Li
 175 et al., 2024b), leakage estimation (Xu et al., 2024), semantic substitutions (Chen et al., 2023; Wang
 176 et al., 2023a) and distractions within the context (Shi et al., 2023).

177 **Compositional Reasoning.** The ability of models to apply learned patterns to novel combinations
 178 of elements and generalize effectively has been studied extensively. Hupkes et al. (2020); Lake and
 179 Baroni (2018) have looked at seq2seq models’ ability to compose known fractions together into novel
 180 combinations in synthetic settings. More recently, the in-context compositional generalization of
 181 LLM reasoners has been further examined (Hosseini et al., 2022; Khot et al., 2023; Press et al., 2023;
 182 Yao et al., 2023b; Yin et al., 2024; Zhou et al., 2023).

183 5 Discussion and Conclusion

184 We designed the Compositional GSM benchmark, which requires solving dependent pairs of math
 185 word problems. These problems are from the original GSM8K test split. We investigate the “System
 186 2” mathematical reasoning capabilities of LLMs by comparing their performance on the original
 187 GSM8K test split and our Compositional GSM test set. Our analysis reveals a notable reasoning gap
 188 in most models. Many leading LLMs exhibit a substantial difference in performance when solving
 189 questions independently versus as part of the compositional pair. Our study indicates that smaller and
 190 more cost efficient models exhibit a larger reasoning gap. Models frequently struggle with pairs of
 191 questions and get distracted likely because they are tuned to handle one question at a time. They often
 192 answer the first question correctly, but lose attention to details and make subtle errors in answering
 193 the second question. We also noticed that learning from human data and self-generated data results in

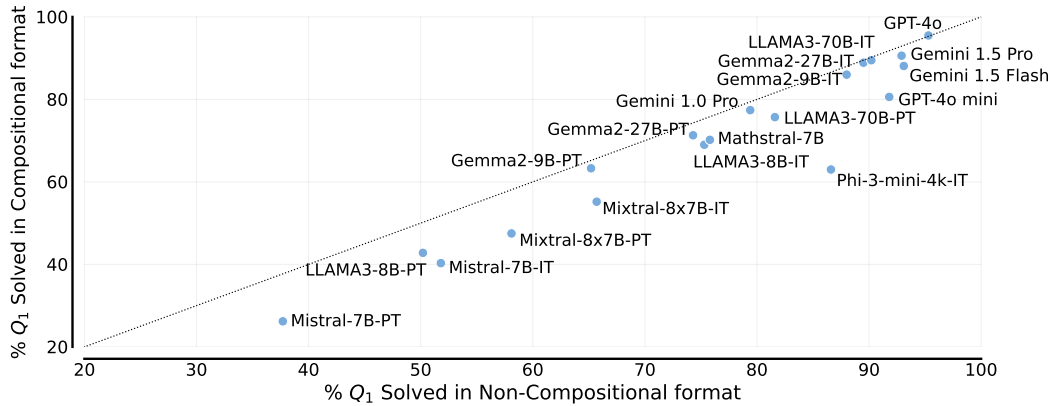


Figure 9: **Some LLMs get distracted easily:** Measuring models’ ability to solve a question in the standard format (non-compositional) versus solving the same question as Q_1 in the compositional format. Models below the trend-line get distracted and cannot answer Q_1 in the compositional format even though solving it does not depend on solving any other question.

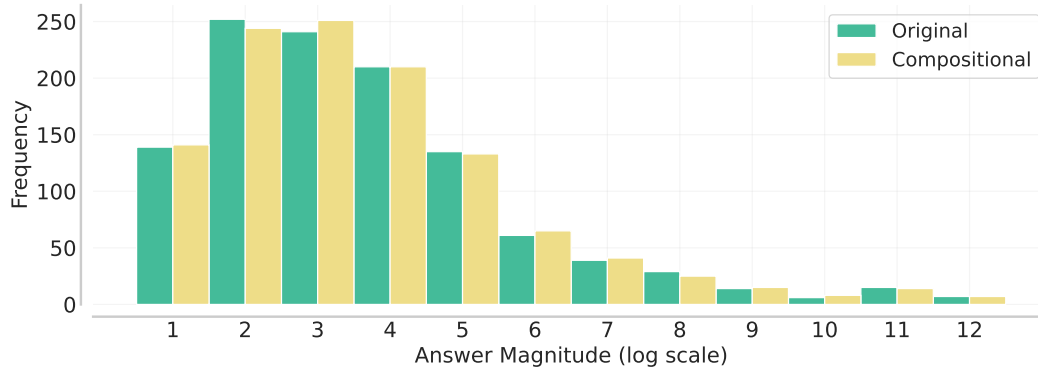


Figure 10: **Distribution of final answers** from the test set of original GSM8K and compositional GSM benchmark. The number modification in the compositional benchmark was done in a way to ensure that the new final answer is a positive integer not too far from the old answer. Our compositional GSM benchmark has a similar distribution of final answers.

194 similar behaviour. In both settings, as training progresses, the model’s performance on the original
 195 test split improves. However, beyond a certain point, performance on the Compositional GSM test
 196 begins to decline.

197 We emphasize that this benchmark should not be viewed as an endpoint or merely as a tool for
 198 generating additional training data, but as a catalyst to gain insights about current models and to re-
 199 evaluate and improve existing benchmarks. Our findings are intended to stimulate further exploration
 200 and provide new perspectives. Future research could build on this setup by incorporating more
 201 challenging questions, such as those from the MATH dataset, or by extending the framework to
 202 multi-modal problems to gain deeper insights into the reasoning capabilities of LLMs.

203 References

204 M. Abdin, S. A. Jacobs, A. A. Awan, J. Aneja, A. Awadallah, H. Awadalla, N. Bach, A. Bahree,
 205 A. Bakhtiari, J. Bao, H. Behl, A. Benhaim, M. Bilenko, J. Bjorck, S. Bubeck, Q. Cai, M. Cai,
 206 C. C. T. Mendes, W. Chen, V. Chaudhary, D. Chen, D. Chen, Y.-C. Chen, Y.-L. Chen, P. Chopra,
 207 X. Dai, A. D. Giorno, G. de Rosa, M. Dixon, R. Eldan, V. Fragoso, D. Iter, M. Gao, M. Gao, J. Gao,
 208 A. Garg, A. Goswami, S. Gunasekar, E. Haider, J. Hao, R. J. Hewett, J. Huynh, M. Javaheripi,
 209 X. Jin, P. Kauffmann, N. Karampatziakis, D. Kim, M. Khademi, L. Kurilenko, J. R. Lee, Y. T.
 210 Lee, Y. Li, Y. Li, C. Liang, L. Liden, C. Liu, M. Liu, W. Liu, E. Lin, Z. Lin, C. Luo, P. Madan,
 211 M. Mazzola, A. Mitra, H. Modi, A. Nguyen, B. Norick, B. Patra, D. Perez-Becker, T. Portet,

- 212 R. Pryzant, H. Qin, M. Radmilac, C. Rosset, S. Roy, O. Ruwase, O. Saarikivi, A. Saied, A. Salim,
213 M. Santacroce, S. Shah, N. Shang, H. Sharma, S. Shukla, X. Song, M. Tanaka, A. Tupini, X. Wang,
214 L. Wang, C. Wang, Y. Wang, R. Ward, G. Wang, P. Witte, H. Wu, M. Wyatt, B. Xiao, C. Xu, J. Xu,
215 W. Xu, S. Yadav, F. Yang, J. Yang, Z. Yang, Y. Yang, D. Yu, L. Yuan, C. Zhang, C. Zhang, J. Zhang,
216 L. L. Zhang, Y. Zhang, Y. Zhang, Y. Zhang, and X. Zhou. Phi-3 technical report: A highly capable
217 language model locally on your phone, 2024. URL <https://arxiv.org/abs/2404.14219>.
- 218 R. Agarwal, N. Vieillard, Y. Zhou, P. Stanczyk, S. R. Garea, M. Geist, and O. Bachem. On-
219 policy distillation of language models: Learning from self-generated mistakes. In *The Twelfth*
220 *International Conference on Learning Representations*, 2024.
- 221 AI@Meta. Llama 3 model card. 2024. URL [https://github.com/meta-llama/llama3/blob/](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md)
222 [main/MODEL_CARD.md](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md).
- 223 U. Anantheswaran, H. Gupta, K. Scaria, S. Verma, C. Baral, and S. Mishra. Investigating the
224 robustness of llms on math word problems. *CoRR*, abs/2406.15444, 2024. doi: 10.48550/ARXIV.
225 2406.15444. URL <https://doi.org/10.48550/arXiv.2406.15444>.
- 226 H. Bansal, A. Hosseini, R. Agarwal, V. Q. Tran, and M. Kazemi. Smaller, weaker, yet better: Training
227 llm reasoners via compute-optimal sampling. *arXiv preprint arXiv:2408.16737*, 2024.
- 228 E. Beeching, S. C. Huang, A. Jiang, J. Li, B. Lipkin, Z. Qina, K. Rasul, Z. Shen, R. Soletskyi, and
229 L. Tunstall. Numinamath 7b cot. <https://huggingface.co/AI-MO/NuminaMath-7B-CoT>,
230 2024.
- 231 M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. de Oliveira Pinto, J. Kaplan, H. Edwards, Y. Burda,
232 N. Joseph, G. Brockman, A. Ray, R. Puri, G. Krueger, M. Petrov, H. Khlaaf, G. Sastry, P. Mishkin,
233 B. Chan, S. Gray, N. Ryder, M. Pavlov, A. Power, L. Kaiser, M. Bavarian, C. Winter, P. Tillet, F. P.
234 Such, D. Cummings, M. Plappert, F. Chantzis, E. Barnes, A. Herbert-Voss, W. H. Guss, A. Nichol,
235 A. Paino, N. Tezak, J. Tang, I. Babuschkin, S. Balaji, S. Jain, W. Saunders, C. Hesse, A. N. Carr,
236 J. Leike, J. Achiam, V. Misra, E. Morikawa, A. Radford, M. Knight, M. Brundage, M. Murati,
237 K. Mayer, P. Welinder, B. McGrew, D. Amodei, S. McCandlish, I. Sutskever, and W. Zaremba.
238 Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 239 W. Chen, X. Ma, X. Wang, and W. W. Cohen. Program of thoughts prompting: Disentangling
240 computation from reasoning for numerical reasoning tasks. *Trans. Mach. Learn. Res.*, 2023, 2023.
241 URL <https://openreview.net/forum?id=YfZ4ZPt8zd>.
- 242 K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton,
243 R. Nakano, C. Hesse, and J. Schulman. Training verifiers to solve math word problems. *arXiv*
244 *preprint arXiv:2110.14168*, 2021a.
- 245 K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton,
246 R. Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*,
247 2021b.
- 248 A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang,
249 A. Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- 250 L. Gao, A. Madaan, S. Zhou, U. Alon, P. Liu, Y. Yang, J. Callan, and G. Neubig. PAL: program-aided
251 language models. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett,
252 editors, *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu,*
253 *Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 10764–10799.
254 PMLR, 2023. URL <https://proceedings.mlr.press/v202/gao23f.html>.
- 255 T. M. Gemma Team, C. Hardin, R. Dadashi, S. Bhupatiraju, L. Sifre, M. Rivière, M. S. Kale,
256 J. Love, P. Tafti, L. Hussenot, and et al. Gemma. 2024. doi: 10.34740/KAGGLE/M/3301. URL
257 <https://www.kaggle.com/m/3301>.
- 258 G. T. Google. Gemini: a family of highly capable multimodal models. *arXiv preprint*
259 *arXiv:2312.11805*, 2023.
- 260 G. T. Google. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context.
261 *arXiv e-prints*, pages arXiv–2403, 2024.

- 262 Z. Gou, Z. Shao, Y. Gong, Y. Yang, M. Huang, N. Duan, W. Chen, et al. Tora: A tool-integrated
263 reasoning agent for mathematical problem solving. *arXiv preprint arXiv:2309.17452*, 2023.
- 264 C. Gulcehre, T. L. Paine, S. Srinivasan, K. Konyushkova, L. Weerts, A. Sharma, A. Siddhant,
265 A. Ahern, M. Wang, C. Gu, et al. Reinforced self-training (rest) for language modeling. *arXiv*
266 *preprint arXiv:2308.08998*, 2023.
- 267 D. Hendrycks, C. Burns, S. Kadavath, A. Arora, S. Basart, E. Tang, D. Song, and J. Steinhardt.
268 Measuring mathematical problem solving with the math dataset. *NeurIPS*, 2021.
- 269 A. Hosseini, A. Vani, D. Bahdanau, A. Sordoni, and A. C. Courville. On the compositional
270 generalization gap of in-context learning. In J. Bastings, Y. Belinkov, Y. Elazar, D. Hup-
271 kes, N. Saphra, and S. Wiegrefe, editors, *Proceedings of the Fifth BlackboxNLP Workshop*
272 *on Analyzing and Interpreting Neural Networks for NLP, BlackboxNLP@EMNLP 2022, Abu*
273 *Dhabi, United Arab Emirates (Hybrid), December 8, 2022*, pages 272–280. Association for
274 Computational Linguistics, 2022. doi: 10.18653/V1/2022.BLACKBOXNLP-1.22. URL
275 <https://doi.org/10.18653/v1/2022.blackboxnlp-1.22>.
- 276 A. Hosseini, X. Yuan, N. Malkin, A. Courville, A. Sordoni, and R. Agarwal. V-star: Training verifiers
277 for self-taught reasoners. *arXiv preprint arXiv:2402.06457*, 2024.
- 278 D. Hupkes, V. Dankers, M. Mul, and E. Bruni. Compositionality decomposed: How do neural
279 networks generalise? *J. Artif. Intell. Res.*, 67:757–795, 2020. doi: 10.1613/JAIR.1.11674. URL
280 <https://doi.org/10.1613/jair.1.11674>.
- 281 A. Q. Jiang, A. Sablayrolles, A. Roux, A. Mensch, B. Savary, C. Bamford, D. S. Chaplot, D. de las
282 Casas, E. B. Hanna, F. Bressand, G. Lengyel, G. Bour, G. Lample, L. R. Lavaud, L. Saulnier, M.-A.
283 Lachaux, P. Stock, S. Subramanian, S. Yang, S. Antoniak, T. L. Scao, T. Gervet, T. Lavril, T. Wang,
284 T. Lacroix, and W. E. Sayed. Mixtral of experts, 2024. URL <https://arxiv.org/abs/2401.04088>.
- 286 T. Khot, H. Trivedi, M. Finlayson, Y. Fu, K. Richardson, P. Clark, and A. Sabharwal. Decomposed
287 prompting: A modular approach for solving complex tasks. In *The Eleventh International Confer-*
288 *ence on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net,
289 2023. URL https://openreview.net/forum?id=_nGgzQjzaRy.
- 290 B. M. Lake and M. Baroni. Generalization without systematicity: On the compositional skills of
291 sequence-to-sequence recurrent networks. In J. G. Dy and A. Krause, editors, *Proceedings of the*
292 *35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018*,
293 volume 80 of *Proceedings of Machine Learning Research*, pages
294 2879–2888. PMLR, 2018. URL <http://proceedings.mlr.press/v80/lake18a.html>.
- 295 A. Lewkowycz, A. Andreassen, D. Dohan, E. Dyer, H. Michalewski, V. V. Ramasesh, A. Slone,
296 C. Anil, I. Schlag, T. Gutman-Solo, Y. Wu, B. Neyshabur, G. Gur-Ari, and V. Misra. Solving quan-
297 titative reasoning problems with language models. In S. Koyejo, S. Mohamed, A. Agarwal, D. Bel-
298 grave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems 35: Annual*
299 *Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA,*
300 *USA, November 28 - December 9, 2022*, 2022. URL [http://papers.nips.cc/paper_files/](http://papers.nips.cc/paper_files/paper/2022/hash/18abbeef8cfe9203fdf9053c9c4fe191-Abstract-Conference.html)
301 [paper/2022/hash/18abbeef8cfe9203fdf9053c9c4fe191-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/18abbeef8cfe9203fdf9053c9c4fe191-Abstract-Conference.html).
- 302 C. Li, J. Liang, A. Zeng, X. Chen, K. Hausman, D. Sadigh, S. Levine, L. Fei-Fei, F. Xia, and
303 B. Ichter. Chain of code: Reasoning with a language model-augmented code emulator. In *Forty-*
304 *first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27,*
305 *2024*. OpenReview.net, 2024a. URL <https://openreview.net/forum?id=vKtomqlSxm>.
- 306 Q. Li, L. Cui, X. Zhao, L. Kong, and W. Bi. Gsm-plus: A comprehensive benchmark for evaluating
307 the robustness of llms as mathematical problem solvers. In L. Ku, A. Martins, and V. Srikumar,
308 editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*
309 *(Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 2961–2984.
310 Association for Computational Linguistics, 2024b. URL [https://aclanthology.org/2024.](https://aclanthology.org/2024.acl-long.163)
311 [acl-long.163](https://aclanthology.org/2024.acl-long.163).

- 312 H. Lightman, V. Kosaraju, Y. Burda, H. Edwards, B. Baker, T. Lee, J. Leike, J. Schulman, I. Sutskever,
313 and K. Cobbe. Let’s verify step by step. *arXiv preprint arXiv:2305.20050*, 2023.
- 314 N. McAleese, R. M. Pokorny, J. F. C. Uribe, E. Nitishinskaya, M. Trebacz, and J. Leike. Llm critics
315 help catch llm bugs, 2024. URL <https://arxiv.org/abs/2407.00215>.
- 316 OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023a. doi: 10.48550/ARXIV.2303.08774.
317 URL <https://doi.org/10.48550/arXiv.2303.08774>.
- 318 OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023b.
- 319 O. Press, M. Zhang, S. Min, L. Schmidt, N. A. Smith, and M. Lewis. Measuring and narrow-
320 ing the compositionality gap in language models. In H. Bouamor, J. Pino, and K. Bali, edi-
321 tors, *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, De-*
322 *cember 6-10, 2023*, pages 5687–5711. Association for Computational Linguistics, 2023. doi:
323 10.18653/V1/2023.FINDINGS-EMNLP.378. URL [https://doi.org/10.18653/v1/2023.](https://doi.org/10.18653/v1/2023.findings-emnlp.378)
324 [findings-emnlp.378](https://doi.org/10.18653/v1/2023.findings-emnlp.378).
- 325 Z. Shao, P. Wang, Q. Zhu, R. Xu, J. Song, M. Zhang, Y. K. Li, Y. Wu, and D. Guo. Deepseekmath:
326 Pushing the limits of mathematical reasoning in open language models. *CoRR*, abs/2402.03300,
327 2024. doi: 10.48550/ARXIV.2402.03300. URL [https://doi.org/10.48550/arXiv.2402.](https://doi.org/10.48550/arXiv.2402.03300)
328 [03300](https://doi.org/10.48550/arXiv.2402.03300).
- 329 F. Shi, X. Chen, K. Misra, N. Scales, D. Dohan, E. H. Chi, N. Schärli, and D. Zhou. Large
330 language models can be easily distracted by irrelevant context. In A. Krause, E. Brunskill,
331 K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett, editors, *International Conference on Machine*
332 *Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of*
333 *Machine Learning Research*, pages 31210–31227. PMLR, 2023. URL [https://proceedings.](https://proceedings.mlr.press/v202/shi23a.html)
334 [mlr.press/v202/shi23a.html](https://proceedings.mlr.press/v202/shi23a.html).
- 335 A. Singh, J. D. Co-Reyes, R. Agarwal, A. Anand, P. Patil, P. J. Liu, J. Harrison, J. Lee, K. Xu,
336 A. Parisi, et al. Beyond human data: Scaling self-training for problem-solving with language
337 models. *arXiv preprint arXiv:2312.06585*, 2023.
- 338 S. Srivastava, A. M. B. A. P. V, S. Menon, A. Sukumar, A. S. T, A. Philipose, S. Prince, and
339 S. Thomas. Functional benchmarks for robust evaluation of reasoning performance, and the
340 reasoning gap. *CoRR*, abs/2402.19450, 2024a. doi: 10.48550/ARXIV.2402.19450. URL <https://doi.org/10.48550/arXiv.2402.19450>.
- 341 [//doi.org/10.48550/arXiv.2402.19450](https://doi.org/10.48550/arXiv.2402.19450).
- 342 S. Srivastava, A. PV, S. Menon, A. Sukumar, A. Philipose, S. Prince, S. Thomas, et al. Functional
343 benchmarks for robust evaluation of reasoning performance, and the reasoning gap. *arXiv preprint*
344 *arXiv:2402.19450*, 2024b.
- 345 G. Team M. Riviere, S. Pathak, P. G. Sessa, C. Hardin, S. Bhupatiraju, L. Hussenot, T. Mesnard,
346 B. Shahriari, A. Ramé, et al. Gemma 2: Improving open language models at a practical size. *arXiv*
347 *preprint arXiv:2408.00118*, 2024.
- 348 H. Wang, G. Ma, C. Yu, N. Gui, L. Zhang, Z. Huang, S. Ma, Y. Chang, S. Zhang, L. Shen, X. Wang,
349 P. Zhao, and D. Tao. Are large language models really robust to word-level perturbations?
350 *CoRR*, abs/2309.11166, 2023a. doi: 10.48550/ARXIV.2309.11166. URL [https://doi.org/10.](https://doi.org/10.48550/arXiv.2309.11166)
351 [48550/arXiv.2309.11166](https://doi.org/10.48550/arXiv.2309.11166).
- 352 P. Wang, L. Li, Z. Shao, R. Xu, D. Dai, Y. Li, D. Chen, Y. Wu, and Z. Sui. Math-shepherd: A label-
353 free step-by-step verifier for llms in mathematical reasoning. *arXiv preprint arXiv:2312.08935*,
354 2023b.
- 355 X. Wang, J. Wei, D. Schuurmans, Q. V. Le, E. H. Chi, S. Narang, A. Chowdhery, and D. Zhou. Self-
356 consistency improves chain of thought reasoning in language models. *International Conference on*
357 *Learning Representations (ICLR)*, 2023c.
- 358 J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. H. Chi, Q. V. Le, and
359 D. Zhou. Chain-of-thought prompting elicits reasoning in large language models. In

- 360 S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Ad-*
361 *vances in Neural Information Processing Systems 35: Annual Conference on Neural Infor-*
362 *mation Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - De-*
363 *cember 9, 2022*, 2022. URL [http://papers.nips.cc/paper_files/paper/2022/hash/](http://papers.nips.cc/paper_files/paper/2022/hash/9d5609613524ecf4f15af0f7b31abca4-Abstract-Conference.html)
364 [9d5609613524ecf4f15af0f7b31abca4-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/9d5609613524ecf4f15af0f7b31abca4-Abstract-Conference.html).
- 365 R. Xu, Z. Wang, R. Fan, and P. Liu. Benchmarking benchmark leakage in large language models.
366 *CoRR*, abs/2404.18824, 2024. doi: 10.48550/ARXIV.2404.18824. URL [https://doi.org/10.](https://doi.org/10.48550/arXiv.2404.18824)
367 [48550/arXiv.2404.18824](https://doi.org/10.48550/arXiv.2404.18824).
- 368 A. Yang, B. Zhang, B. Hui, B. Gao, B. Yu, C. Li, D. Liu, J. Tu, J. Zhou, J. Lin, K. Lu, M. Xue, R. Lin,
369 T. Liu, X. Ren, and Z. Zhang. Qwen2.5-math technical report: Toward mathematical expert model
370 via self-improvement. *arXiv preprint arXiv:2409.12122*, 2024.
- 371 S. Yao, D. Yu, J. Zhao, I. Shafran, T. Griffiths, Y. Cao, and K. Narasimhan. Tree of thoughts: Deliber-
372 ate problem solving with large language models. In A. Oh, T. Naumann, A. Globerson, K. Saenko,
373 M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems 36: Annual*
374 *Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA,*
375 *USA, December 10 - 16, 2023*, 2023a. URL [http://papers.nips.cc/paper_files/paper/](http://papers.nips.cc/paper_files/paper/2023/hash/271db9922b8d1f4dd7aaef84ed5ac703-Abstract-Conference.html)
376 [2023/hash/271db9922b8d1f4dd7aaef84ed5ac703-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/271db9922b8d1f4dd7aaef84ed5ac703-Abstract-Conference.html).
- 377 S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. R. Narasimhan, and Y. Cao. React: Synergizing
378 reasoning and acting in language models. In *The Eleventh International Conference on Learning*
379 *Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023b. URL
380 https://openreview.net/forum?id=wE_vluYUL-X.
- 381 Y. Yin, L. Fu, Y. Li, and Y. Zhang. On compositional generalization of transformer-based neural
382 machine translation. *Inf. Fusion*, 111:102491, 2024. doi: 10.1016/J.INFFUS.2024.102491. URL
383 <https://doi.org/10.1016/j.inffus.2024.102491>.
- 384 E. Zelikman, Y. Wu, J. Mu, and N. D. Goodman. Star: Bootstrapping reasoning with reasoning.
385 *Neural Information Processing Systems (NeurIPS)*, 2022.
- 386 H. Zhang, J. Da, D. Lee, V. Robinson, C. Wu, W. Song, T. Zhao, P. Raja, D. Slack, Q. Lyu, S. Hendryx,
387 R. Kaplan, M. Lunati, and S. Yue. A careful examination of large language model performance on
388 grade school arithmetic. *CoRR*, abs/2405.00332, 2024a. doi: 10.48550/ARXIV.2405.00332. URL
389 <https://doi.org/10.48550/arXiv.2405.00332>.
- 390 L. Zhang, A. Hosseini, H. Bansal, M. Kazemi, A. Kumar, and R. Agarwal. Generative verifiers: Re-
391 ward modeling as next-token prediction, 2024b. URL <https://arxiv.org/abs/2408.15240>.
- 392 D. Zhou, N. Schärli, L. Hou, J. Wei, N. Scales, X. Wang, D. Schuurmans, C. Cui, O. Bousquet,
393 Q. V. Le, and E. H. Chi. Least-to-most prompting enables complex reasoning in large language
394 models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali,*
395 *Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=WZH7099tgfM)
396 [WZH7099tgfM](https://openreview.net/forum?id=WZH7099tgfM).

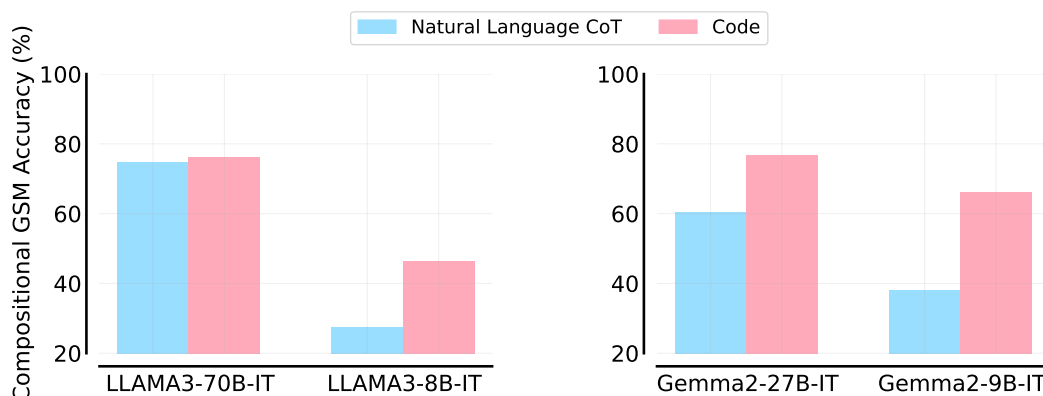


Figure 11: **Natural Language CoT v.s. Code:** Generating code to solve the problems helps in both settings of original test split and Compositional GSM split. Smaller models benefit more from generating code rather than natural language Chain-of-Thought (CoT) to solve Compositional GSM questions.

398 Breaking down the natural language problem into executable code steps has been shown to improve
 399 models' reasoning and generalization abilities (Gao et al., 2023; Gou et al., 2023). To this end, we
 400 evaluate whether the compositional problem-solving ability of LLMs improves when generating
 401 natural language CoT rationales compared to generating executable Python code. For code generation,
 402 we utilize a compositional 8-shot prompt(Appendix E), where the answers are written as two functions,
 403 one which solves the first question *solve_q1()*, and *solution()* which solves the second question with a
 404 $X = solve_q1()$ line at the beginning.

405 Our results are shown in Figure 11 for two families of open-weight instruction-tuned models:
 406 LLAMA3-8B and 70B, and Gemma2-9B and 27B. Notably, generating code generally improves
 407 performance on Compositional GSM problems, albeit not uniformly. Specifically, the smaller models,
 408 LLAMA3-8B and Gemma2-9B, benefit significantly more from generating code solutions compared
 409 to generating natural language CoTs. This variability in the performance gains from generating code
 410 solutions indicates that such the in-context learning behaviors of LLMs are inconsistent.

Q: There are 15 trees in the grove. Grove workers will plant trees in the grove today. After they are done, there will be 21 trees. How many trees did the grove workers plant today?

A: There are 15 trees originally. Then there were 21 trees after some more were planted. So there must have been $21 - 15 = 6$. The final answer is 6.

Q: If there are 3 cars in the parking lot and 2 more cars arrive, how many cars are in the parking lot?

A: There are originally 3 cars. 2 more cars arrive. $3 + 2 = 5$. The final answer is 5.

Q: Leah had 32 chocolates and her sister had 42. If they ate 35, how many pieces do they have left in total?

A: Originally, Leah had 32 chocolates. Her sister had 42. So in total they had $32 + 42 = 74$. After eating 35, they had $74 - 35 = 39$. The final answer is 39.

Q: Jason had 20 lollipops. He gave Denny some lollipops. Now Jason has 12 lollipops. How many lollipops did Jason give to Denny?

A: Jason started with 20 lollipops. Then he had 12 after giving some to Denny. So he gave Denny $20 - 12 = 8$. The final answer is 8.

Q: Shawn has five toys. For Christmas, he got two toys each from his mom and dad. How many toys does he have now?

A: Shawn started with 5 toys. If he got 2 toys each from his mom and dad, then that is 4 more toys. $5 + 4 = 9$. The final answer is 9.

Q: There were nine computers in the server room. Five more computers were installed each day, from monday to thursday. How many computers are now in the server room?

A: There were originally 9 computers. For each of 4 days, 5 more computers were added. So $5 * 4 = 20$ computers were added. $9 + 20$ is 29. The final answer is 29.

Q: Michael had 58 golf balls. On tuesday, he lost 23 golf balls. On wednesday, he lost 2 more. How many golf balls did he have at the end of wednesday?

A: Michael started with 58 golf balls. After losing 23 on tuesday, he had $58 - 23 = 35$. After losing 2 more, he had $35 - 2 = 33$ golf balls. The final answer is 33.

Q: Olivia has 23. *She bought five bagels for 3 each.* How much money does she have left?

A: Olivia had 23 dollars. 5 bagels for 3 dollars each will be $5 * 3 = 15$ dollars. So she has $23 - 15$ dollars left. $23 - 15$ is 8. The final answer is 8.

Q: {question}

A:

412 **B Prmopt Preambles**

GSM8K Preamble

I am going to give you a series of demonstrations of math Problems and Solutions. When you respond, respond only with the Solution of the final Problem, thinking step by step. At the end of the Solution, when you give your final answer, write it in the form "The final answer is ANSWER."

Compositional GSM Preamble

I am going to give you a series of demonstrations of compositional math questions and solutions. Respond by thinking step by step. Solve the first question and write the intermediate answer as "The Q1 answer is ANSWER1." Then solve Q2. At the end of the solution, when you give your final answer, write it in the form "The final answer is ANSWER2."

413 **C Rejection Finetuning Details**

414 Synthetic data was generated by prompting Gemma2 27B PT model with the 8-shot prompt in Ap-
415 pendix A to solve GSM8K training questions. We generated 10 solutions for each question in the
416 original GSM8K training data, and only kept those solutions with a correct final answer. These model
417 generation solutions were used to train the model.

418 **D Compositional 8-shot Prompt**

Let X be the answer to Q1:

Q1: There are 15 trees in the grove. Grove workers will plant trees in the grove today. After they are done, there will be 21 trees. How many trees did the grove workers plant today?

solve it and use the value of X to solve Q2. Explain your answer step by step.

Q2: There are X students in Marissa's class. Each student started the year with 10 pencils. After two months, $\frac{1}{5}$ of the total pencils in class were used. At the end of the year, only $\frac{1}{3}$ of the remaining pencils were left. How many pencils were left?

A: There are 15 trees originally. Then there were 21 trees after some more were planted. So there must have been $21 - 15 = 6$. The Q1 answer is 6. Therefore $X=6$. So there were $6 \times 10 = 60$ pencils in the class at the start of the year. After two months, $60 \times \frac{1}{5} = 12$ pencils were used. Thus, $60 - 12 = 48$ pencils were left unused after two months. Therefore, $48 \times \frac{1}{3} = 16$ pencils were left at the end of the year. The final answer is 16.

Let X be the answer to Q1:

Q1: If there are 3 cars in the parking lot and 2 more cars arrive, how many cars are in the parking lot?

solve it and use the value of X to solve Q2. Explain your answer step by step.

Q2: Ingrid drinks X cups of water every day. If there are 16 cups in a gallon, how many gallons of water does she drink in 32 days?

A: There are originally 3 cars. 2 more cars arrive. $3 + 2 = 5$. The Q1 answer is 5. Therefore $X=5$. So Ingrid drinks 5 cups of water a day so after 32 days she drinks $5 * 32 = 160$ cups of water. There are 16 cups in 1 gallon so she drinks $160 / 16 = 10$ gallons of water in 30 days. The final answer is 10.

⋮

Let X be the answer to Q1:

Q1: {QUESTION_1}

solve it and use the value of X to solve Q2. Explain your answer step by step.

Q2: {QUESTION_2}

A:

419 Some examples in the prompt are omitted due to space constraints. The remaining question-and-
 420 answer pairs follow the same format.

421 **E Compositional Code Prompt**

422 Some examples in the prompt are omitted due to space constraints. The remaining question and
423 code solutions are written with the same format. The answers are provided through two functions:
424 *solve_q1()*, which addresses the first question, and *solution()*, which answers the second question.
425 The *solution()* function begins with a line $X = solve_q1()$ to incorporate the result from the first
426 function.

```
Write two functions 'solve_q1' and 'solution' to solve Q1 and Q2 problems.

Let X be the answer to Q1:

Q1: There are 15 trees in the grove. Grove workers will plant trees in the
grove today. After they are done, there will be 21 trees. How many trees did
the grove workers plant today?

Q2: There are X students in Marissa's class. Each student started the year
with 10 pencils. After two months, 1/5 of the total pencils in class were used.
At the end of the year, only 1/3 of the remaining pencils were left. How many
pencils were left?

A: The answer is
...
def solve_q1():
    """There are 15 trees in the grove. Grove workers will plant trees
in the grove today. After they are done, there will be 21 trees. How many trees
did the grove workers plant today?"""
    trees_initial = 15
    trees_after = 21
    trees_added = trees_after - trees_initial
    result = trees_added
    return result

def solution():
    """There are X students in Marissa's class. Each student started the
year with 10 pencils. After two months, 1/5 of the total pencils in class were
used. At the end of the year, only 1/3 of the remaining pencils were left. How
many pencils were left?"""
    X = solve_q1()
    num_students = X
    pencils_per_student = 10
    total_pencils = num_students * pencils_per_student
    pencils_left_after_two_months = total_pencils * (4/5)
    remaining_pencils = pencils_left_after_two_months * (1/3)
    result = remaining_pencils
    return result
...
:

Let X be the answer to the following question:

Q1: {QUESTION_1}

Q: {QUESTION_2}

A: The answer is
```