

Arithmetic Reasoning with LLM: Prolog Generation & Permutation

Xiaocheng Yang and Bingsen Chen and Yik-Cheung Tam

Shanghai Frontiers Science Center of Artificial Intelligence and Deep Learning

New York University Shanghai

{xy2128,bc3088,yt2267}@nyu.edu

Abstract

Instructing large language models (LLMs) to solve elementary school math problems has shown great success using Chain of Thought (CoT). However, the CoT approach relies on an LLM to generate a sequence of arithmetic calculations which can be prone to cascaded calculation errors. We hypothesize that an LLM should focus on extracting predicates and generating symbolic formulas from the math problem description so that the underlying calculation can be done via an external code interpreter. We investigate using LLM to generate Prolog programs to solve mathematical questions. Experimental results show that our Prolog-based arithmetic problem-solving outperforms CoT generation in the GSM8K benchmark across three distinct LLMs. In addition, given the insensitive ordering of predicates and symbolic formulas in Prolog, we propose to permute the ground truth predicates for more robust LLM training via data augmentation.

1 Introduction

Large language models (LLMs), with their scaling of model size and data size, have demonstrated impressive performance across various understanding and generation tasks (Brown et al., 2020; Chowdhery et al., 2022; Rae et al., 2021; Thoppilan et al., 2022; Touvron et al., 2023; Almazrouei et al., 2023; Jiang et al., 2023). Nevertheless, such LLMs fall short in addressing mathematical problems that involves arithmetic, commonsense, and symbolic reasoning – topics that may appear deceptively simple to humans (Rae et al., 2021). Existing works leveraged Chain-of-Thought (CoT) reasoning that asks language models to generate both the answer and the step-by-step reasoning chain, which helps break down a complex reasoning task into a sequential thought process (Wei et al., 2022b). Particularly, arithmetic reasoning with CoT is shown to be an emergent ability that language models acquired during the scaling process (Wei et al., 2022a).

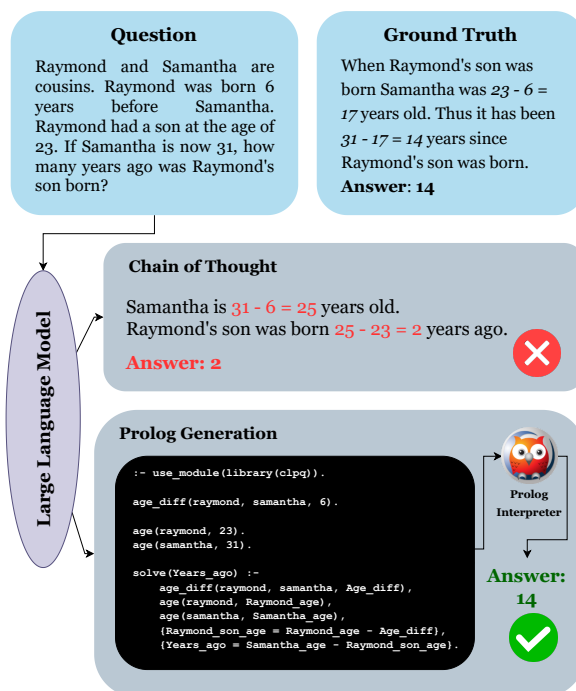


Figure 1: Overview of Prolog generation for arithmetic reasoning with large language models.

Yet, natural language reasoning is not native to mathematical operations and symbolic manipulations. A line of work has focused on augmenting language models with deterministic computation resources like a calculator (Schick et al., 2023) or program-based tools (Gao et al., 2023; Gou et al., 2023). However, all such methods require a sequential reasoning trajectory, where models need to translate the natural language questions into sequential mathematical or logical operations. Our research probes into the application of Prolog, a logic programming language, in solving the arithmetic reasoning task. Prolog solves arithmetic reasoning tasks by defining an unordered set of predicates and running queries over them. We further explain the unique properties of Prolog in Section 2. In Prolog code generation for arithmetic reasoning, LLMs extract facts and rules in mathematical ques-

tions and formulate them into Prolog code. If the facts and rules are accurately captured, a Prolog interpreter can precisely solve for a correct answer in a deterministic way.

Our research has the following contributions: 1) We curate and open-source the GSM8K-Prolog dataset with a semi-automatic approach, which contains arithmetic reasoning problems and their corresponding Prolog code solutions. 2) Our experiments show that Prolog code generation is consistently better than CoT on the arithmetic reasoning task, indicating that LLM can focus on predicate extractions and rely on an external tool to calculate and perform the logical induction to address mathematical problems. 3) Given the non-sequential nature of predicates in Prolog code, we propose predicate permutation as a data augmentation method and demonstrate its efficacy in robust LLM training.

Original Ground Truth

```

:- use_module(library(clpqr)).

earn(weng, 12).
work(weng, 50).

solve(Total_salary) :-
    earn(weng, Salary_per_hour),
    work(weng, Working_minutes),
    (Total_salary = Salary_per_hour * Working_minutes / 60).

```

Permuted Ground Truth

```

:- use_module(library(clpqr)).

solve(Total_salary) :-
    work(weng, Working_minutes),
    (Total_salary = Salary_per_hour * Working_minutes / 60),
    earn(weng, Salary_per_hour).

work(weng, 50).
earn(weng, 12).

```

Figure 2: Prolog and permuted Prolog code samples.

2 Preliminaries: Prolog Language

Prolog is a logic programming language, which was initially designed for artificial intelligence and computational linguistics (Clocksin and Mellish, 2003; Bratko, 2012; Covington, 2002). As shown in the upper graph of Figure 2, a Prolog program defines a set of predicates that contains facts and goals. In the example, facts include `earn(weng, 12)` that declares the hourly salary of Weng, and `work(weng, 50)` that defines the working minutes of Weng; the goals constitute a rule in the form of `solve<answer>:-<goal_1>, <goal_2>, ... A`

rule is true when all the goals are satisfied. Having all the facts and goals defined in the program, users can make a query to obtain the solutions that make the rule true given all the facts. Moreover, Prolog codes are not sequential like Python, meaning that the order of facts and rules does not alter the result of the program. The lower graph in Figure 2, shows an equivalent sample that permutes the order of the predicates, which produces the same result as the original program.

3 Method

3.1 GSM8K-Prolog Dataset

To our knowledge, there has not been a dataset for solving mathematical questions with Prolog. We hence curated a dataset based on GSM8K (Cobbe et al., 2021), a popular benchmark of diverse grade school math word problems, in a semi-automatic manner with OpenAI’s Text Completion API¹. In particular, we used the same dataset splits and questions in GSM8K and prompted GPT-4 to generate the Prolog programs to solve the questions. We then manually corrected some malfunctioning samples. In this manner, we obtained a high-quality corpus with 100% accuracy in terms of the code results. Algorithm 1 describes the detailed pseudocode for creating this dataset. We open-sourced this dataset to the research community with the MIT license.²

3.2 PROPER: Prolog Permutation

Since Prolog predicates are permutable, inspired by XLNet (Yang et al., 2020) that performs a token-wise permutation via attention masking, we decided to also use the permutation technique. The XLNet, via the permutation, can attend to tokens on both sides during training and thus can partially obtain the property of autoencoding while maintaining the property of autoregressive modeling. Similarly, PROPER takes advantage of the permutative property of facts and goals in the Prolog programs as indicated in Figure 2. For each original program, we sample n of its permutations and mix them into the dataset. In this way, models can learn to extract predicates in the mathematical questions based on any other predicates regardless of the ordering, which more precisely reflects the nature

¹<https://platform.openai.com/docs/guides/text-generation/chat-completions-api>

²<https://huggingface.co/datasets/Thomas-X-Yang/gsm8k-prolog>

Method	Llama-2		CodeLlama		Mistral	
	GSM8K	GSM-HARD	GSM8K	GSM-HARD	GSM8K	GSM-HARD
CoT	33.8%	12.0%	37.5%	13.9%	58.9%	30.8%
Prolog	41.5%	32.4%	55.0%	41.6%	66.3%	50.6%
PROPER	51.0%	37.4%	59.0%	45.9%	70.2%	54.4%

Table 1: Accuracy results on the GSM8K and GSM-HARD datasets. We compare regular Prolog generation (Prolog) and PROPER Prolog generation with the CoT baseline (supervised finetuning with LoRA using CoT ground truth labels in the original GSM8K dataset).

of the Prolog language. We describe the practical details of permutation in Appendix A.2.

4 Experiments

4.1 Setup

Dataset We used the GSM8K-Prolog described in Section 3.1. We denote the corpus as D . The training set is D_{train} and the test set is D_{test} . The total corpus size is 8792, where 7473 samples belong to the training set and 1319 belong to the test set. During training, 100 samples were selected from the training set to constitute the validation set. The input format follows the instruction prompt used in Stanford Alpaca (Taori et al., 2023) (See sample prompts in Appendix A.3). We discarded samples that exceeded 512 tokens. Notably, when we used PROPER to augment the dataset, we used slightly altered input prompts for permuted samples because we found that using the same instruction for both the original ground truth codes and the permuted ones degraded the performance of the model. A likely reason is that having multiple correct output tokens for the same input instruction confuses the model. In addition, besides the GSM8K’s test set, GSM-HARD (Gao et al., 2023), which replaces the numbers in the GSM8K test set with large numbers and thus makes questions hard for language models, was also used for evaluation.

Training We experimented with different LLMs’ 7B versions, including Llama2 (Touvron et al., 2023), CodeLlama (Rozière et al., 2023) and Mistral (Jiang et al., 2023). We adopted 8-bit quantization and LoRA (Hu et al., 2021) to finetune models efficiently at a reasonable performance cost. We applied LoRA to finetune query and value weight matrices in the transformer blocks. We experimented with different LoRA rank and alpha settings, including $(r, \alpha) = (8, 16), (16, 32),$ and $(32, 64)$. With more trainable parameters, $r = 32, \alpha = 64$ yielded significantly better results,

which we thereby adopted as the configuration for all the experiments. Note that this setting resulted in training only 0.248% of the 7 billion parameters for Llama2 and CodeLlama, and 0.188% of the 7 billion for Mistral. We document our training details and GPU usage in Appendix A.5.

Evaluation At inference time, we used beam search with a beam size of 4 to generate the Prolog code. We then used the PySwip library³, a foreign interface of Prolog in Python, as the Prolog interpreter to produce the final answer. We used accuracy as the metric for evaluation. It is defined as

$$\text{Acc} = \frac{\sum_{i=1}^{|D_{test}|} \mathbb{1}_{\{\mathcal{P}(a_i^{pred})=P(a_i^{true})\}}}{|D_{test}|} \times 100\%$$

where \mathcal{P} denotes the Prolog interpreter. Notably, since we noticed that the PySwip library cannot handle decimal answers, we only considered the samples with an integer answer.

4.2 Results

Prolog generation performs consistently better than CoT across three models. According to Table 1, generating Prolog to solve mathematical questions yields significantly more accurate results with a 10.9% margin over the CoT baseline on average across all models on GSM8K. This gap further expands to 22.6% on GSM-HARD, indicating exceptional superiority over CoT when large number calculations are involved. Although Llama-2 and Mistral exhibit large performance gaps when applying CoT reasoning, generating Prolog code produces better results than CoT on both models. This observation indicates that Prolog generation works well regardless of the model’s inherent arithmetic reasoning capability. Also, CodeLlama demonstrates a larger performance gain when switching

³Prolog version 9.0.4. PySwip version 0.2.11. <https://github.com/yuce/pyswip>

Ratio	Llama-2	CodeLlama	Mistral
1:0	41.5	55.0	66.3
1:1	50.9 (49.5)	58.7 (56.6)	70.2 (69.1)
1:2	51.0 (49.4)	59.0 (58.3)	68.8 (66.8)

Table 2: Accuracy(%) results on GSM8K with different permutation ratios. We report both the best and average accuracy of 1:1 and 1:2 over three trials with different randomly permuted data in the form of max (avg). Note that the 1:0 case essentially means not applying PROPER.

from CoT to Prolog generation, which is potentially attributed to its pretraining on the code-related corpus. In other words, CodeLlama is specifically trained to generate structured programs better than natural language reasoning.

With a proper permutation ratio, PROPER further enhances LLM’s arithmetic reasoning with Prolog generation. Permutation ratio refers to the ratio between original samples and permuted samples. As shown in Table 2, by adding two permuted samples for each original sample, we observed an increased accuracy of 9.5% and 4.0% of Llama-2 and CodeLlama respectively on the test set. This improvement indicates that learning the non-sequential structure of Prolog predicates is helpful for LLMs to generate correct Prolog programs to solve arithmetic problems. On the other hand, the lowered accuracy of Mistral, compared with its case of one permutation per sample, suggests that PROPER might be limited for models already with high Prolog generation capacity.

Lowered validation loss from PROPER does not lead to higher accuracy. As is shown in Figure 3, increasing the permutation ratio results in significantly lowered validation loss. This is because we first added in permutations and then split a validation set from the training set. Consequently, the permutations of validation samples were included in the training set and the generalization ability of the language models enabled the models to utilize the permutations to improve the performance on the validation set, causing a soft data leakage. Therefore, according to Table 2, the permutation ratio of 1:2 yielded a weakened performance on Mistral although the validation loss was the lowest.

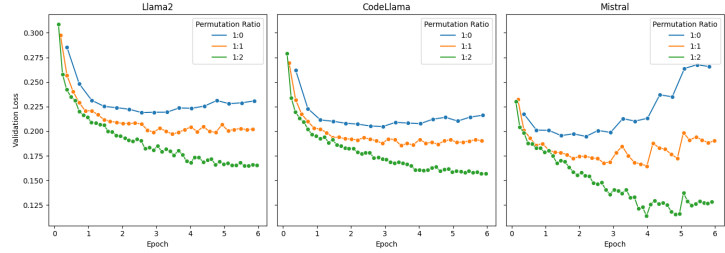


Figure 3: Validation loss curves for training Llama2, CodeLlama, and Mistral with different permutation ratios (We only report the first trial when we use permuted data since the loss curves are very similar across trials).

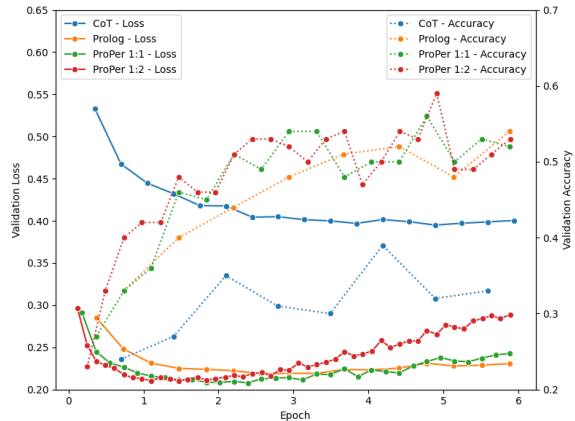


Figure 4: Validation loss curves and validation accuracy curves for training Llama2 with different methods (We only report the first trial when we use permuted data since the loss and accuracy curves are very similar across trials).

Increased validation loss from PROPER does not lead to decreased validation accuracy. Excluding the permutations of validation samples from the training set, we report both the cross-entropy loss and the accuracy on the validation set for Llama2 using different methods in Figure 4. A mismatch between the loss and accuracy is observed. As a loss curve decreases to the minimum and bounces back, the corresponding accuracy curve keeps increasing and then maintains a high level. As is shown in Table 3, by choosing checkpoints based on validation accuracy instead of validation loss, the performance can be improved across all methods. Moreover, the improvement for the Prolog and PROPER method is significantly greater than that of CoT, suggesting a larger divergence between the objective of cross entropy loss and the ultimate accuracy of Prolog generation. Therefore, it is suggested to choose the best checkpoint based on the

Method	Initial	No Leakage (by loss)	No Leakage (by accuracy)
CoT	33.8	33.8	36.5
Prolog	41.5	41.5	47.9
ProPer 1:1	50.9 (49.5)	44.3 (43.4)	50.1 (48.4)
ProPer 1:2	51.0 (49.4)	44.4 (43.6)	51.3 (50.3)

Table 3: Accuracy(%) results of training Llama2 on the GSM8K dataset. We compare the results of avoiding validation sample leakage in the training set and picking the optimal checkpoint based on validation loss and accuracy with the initial results with leakage. The best and average accuracy of 1:1 and 1:2 are in the form of max (avg).

validation accuracy. Nevertheless, the new performance is similar to the initial results where leakage is involved. We notice that late checkpoints yield better performance according to the validation accuracy and the validation loss keeps decreasing in the initial setting. Therefore, both settings happen to pick late checkpoints, resulting in similar performance.

We have also tested Python generation, for which the corpus was generated by the same procedure as Algorithms 1 except that we prepare Python codes instead of Prolog codes. It gives an accuracy of 55.12% on GSM8K using Llama2 as the base model, better than both Prolog and PROPER. One possible reason is that Python now is the prevalent programming language and Llama2 might have been pretrained on a large amount of Python codes. We believe if sufficient Prolog codes are used for training, Prolog generation can at least match up with Python generation due to its essence of symbolic reasoning.

We present some representative error cases of Mistral (1:1) in Appendix A.4.

5 Related Work

Arithmetic Reasoning The Chain-of-Thought (CoT) prompting approach (Wei et al., 2022b) first proposes to prompt the model to generate the reasoning chain step-by-step to reach the final answer. Afterwards, advancements have been made in LLMs’ reasoning capacity via step-by-step methods (Zhou et al., 2023; Zhu et al., 2023; Huang et al., 2022; Liang et al., 2023). However, the natural language generation still performs poorly on complex or multi-step reasoning. Therefore, one trajectory of efforts has been made to leverage reasoning structures like trees (Yao et al., 2023; Long,

2023) and graphs (Besta et al., 2023; Zhang et al., 2023). Another trajectory is to render the reasoning task based on external tools (Cobbe et al., 2021; Mishra et al., 2023; Gou et al., 2023; Gao et al., 2023; Shao et al., 2023; Chen et al., 2023), which is the one that we are following. Besides, Yuan et al.’s (2023) RFT method shares the idea of dataset augmentation, but they compile rejection samples from multiple models to form an augmented training set, which is different from PROPER’s automatic permutation.

Neural Symbolic Reasoning Neural symbolic reasoning (Andreas et al., 2016; Neelakantan et al., 2017; Hudson and Manning, 2019; Gupta et al., 2020; Nye et al., 2021) aims to leverage both neural networks and symbolic reasoning to obtain better reasoning abilities and transparency. Those methods suffer from low scalability of learning and reasoning components. LLMs are hence adopted to generate symbolic representations from natural language (Lyu et al., 2023; Pan et al., 2023; Yang et al., 2023), where deterministic symbolic solvers will process the query and symbolic representations generated by LLMs to conduct reasoning or proofs. Prolog has been a popular candidate for the format of symbolic representations. We are posited on this trajectory and in the specific field of arithmetic reasoning.

6 Conclusion

In conclusion, we aim to enhance the reasoning performance of LLMs. We adopt the pipeline that the model generates Prolog predicates from a mathematical question in natural language and an external Prolog interpreter processes the query for a final result. We contribute an open-sourced corpus named GSM8K-Prolog, which is a high-quality Prolog-annotated version of GSM8K. We show that Prolog generation substantially outperformed CoT generation across all three 7B models for solving arithmetic reasoning problems. We also propose PROPER, a data augmentation method designed specifically for Prolog code generation, which enables the finetuned models to learn the non-sequential nature of Prolog predicates. PROPER further improves the model’s accuracy on GSM8K-Prolog and mitigates early convergence during training. Lastly, due to the gap between cross-entropy loss objective and accuracy, we suggest using validation accuracy instead of validation loss to pick the best checkpoint.

Limitations

Although we have experimentally conducted full-parameter finetuning, the result was not satisfying. We believe it is because of the limited size of the original corpus. Therefore, at the current stage, we cannot have a comparison with other methods like ToRA (Gou et al., 2023) or RFT (Yuan et al., 2023). Future research can look into preparing a larger and more diverse corpus adapted to Prolog code generation. Besides, We did not try scaling the base model to more than 7B parameters. So we do not know the impact of model scaling on the performance of Prolog code generation for arithmetic reasoning. Furthermore, due to the limitation of the PySwip library, solvable questions are restricted to the ones with an integer answer. Future work can expand the domain by using other interpreting tools.

References

- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Hestlow, Julien Launay, Quentin Malartic, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. Falcon-40B: an open large language model with state-of-the-art performance.
- Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. 2016. [Neural module networks](#). In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 39–48.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Michal Podstawski, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoefler. 2023. [Graph of thoughts: Solving elaborate problems with large language models](#).
- Ivan Bratko. 2012. *Prolog programming for Artificial Intelligence*. Addison-Wesley.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2023. [Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks](#).
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pilla, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. [Palm: Scaling language modeling with pathways](#).
- W. F. Clocksin and C. S. Mellish. 2003. *Programming in Prolog*. Springer-Verlag.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#).
- Michael A. Covington. 2002. *Natural language processing for Prolog programmers*. Prentice Hall.
- Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. [Pal: Program-aided language models](#).
- Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujia Yang, Minlie Huang, Nan Duan, and Weizhu Chen. 2023. [Tora: A tool-integrated reasoning agent for mathematical problem solving](#).
- Nitish Gupta, Kevin Lin, Dan Roth, Sameer Singh, and Matt Gardner. 2020. [Neural module networks for reasoning over text](#). In *International Conference on Learning Representations*.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. [Lora: Low-rank adaptation of large language models](#).
- Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2022. [Large language models can self-improve](#).
- Drew Hudson and Christopher D Manning. 2019. [Learning by abstraction: The neural state machine](#). In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao,

- Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. [Mistral 7b](#).
- Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Zhaopeng Tu, and Shuming Shi. 2023. [Encouraging divergent thinking in large language models through multi-agent debate](#).
- Jieyi Long. 2023. [Large language model guided tree-of-thought](#).
- Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. 2023. [Faithful chain-of-thought reasoning](#).
- Swaroop Mishra, Matthew Finlayson, Pan Lu, Leonard Tang, Sean Welleck, Chitta Baral, Tanmay Rajpurohit, Oyvind Tafford, Ashish Sabharwal, Peter Clark, and Ashwin Kalyan. 2023. [Lila: A unified benchmark for mathematical reasoning](#).
- Arvind Neelakantan, Quoc V. Le, Martin Abadi, Andrew McCallum, and Dario Amodei. 2017. [Learning a natural language interface with neural programmer](#). In *International Conference on Learning Representations*.
- Maxwell Nye, Michael Henry Tessler, Joshua B. Tenenbaum, and Brenden M. Lake. 2021. [Improving coherence and consistency in neural sequence models with dual-system, neuro-symbolic reasoning](#).
- Liangming Pan, Alon Albalak, Xinyi Wang, and William Yang Wang. 2023. [Logic-lm: Empowering large language models with symbolic solvers for faithful logical reasoning](#).
- Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. 2021. [Scaling language models: Methods, analysis & insights from training gopher](#). *arXiv preprint arXiv:2112.11446*.
- Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2023. [Code llama: Open foundation models for code](#).
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. [Toolformer: Language models can teach themselves to use tools](#). *arXiv preprint arXiv:2302.04761*.
- Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie Huang, Nan Duan, and Weizhu Chen. 2023. [Synthetic prompting: Generating chain-of-thought demonstrations for large language models](#). In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 30706–30775. PMLR.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. [Stanford alpaca: An instruction-following llama model](#). https://github.com/tatsu-lab/stanford_alpaca.
- Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. 2022. [Lamda: Language models for dialog applications](#). *arXiv preprint arXiv:2201.08239*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. [Llama 2: Open foundation and fine-tuned chat models](#).
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022a. [Emergent abilities of large language models](#).
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022b. [Chain-of-thought prompting elicits reasoning in large language models](#). *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Sen Yang, Xin Li, Leyang Cui, Lidong Bing, and Wai Lam. 2023. [Neuro-symbolic integration brings causal and reliable reasoning proofs](#).
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V. Le. 2020. [Xlnet: Generalized autoregressive pretraining for language understanding](#).
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. [Tree of thoughts: Deliberate problem solving with large language models](#).

Zheng Yuan, Hongyi Yuan, Chengpeng Li, Guanting Dong, Chuanqi Tan, and Chang Zhou. 2023. Scaling relationship on learning mathematical reasoning with large language models. *arXiv preprint arXiv:2308.01825*.

Yifan Zhang, Jingqin Yang, Yang Yuan, and Andrew Chi-Chih Yao. 2023. [Cumulative reasoning with large language models](#).

Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, and Ed Chi. 2023. [Least-to-most prompting enables complex reasoning in large language models](#).

Xinyu Zhu, Junjie Wang, Lin Zhang, Yuxiang Zhang, Yongfeng Huang, Ruyi Gan, Jiaying Zhang, and Yujie Yang. 2023. [Solving math word problems via cooperative reasoning induced language models](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics.

A Appendix

A.1 Generation Procedure of GSM8K-Prolog

Below is the detailed pseudo-code for the GSM8K-Prolog dataset generation.

Algorithm 1 Procedure of GSM8K-Prolog Generation

Input: The original GSM8K dataset, denoted as set $\mathcal{X} = \{(q_i, a_i^{\text{CoT}})\}_{i=1}^N$, where each sample consists of one question q_i and one Chain-of-Thought answer a_i^{CoT} ; A Prolog interpreter \mathcal{P} that returns the output of a Prolog program; A Chain-of-Thought answer retriever \mathcal{C} that parses out the final answer of a natural language reasoning chain.

Output: GSM8K-Prolog dataset $\mathcal{D} = \{(q_i, a_i^{\text{Prolog}})\}_{i=1}^N$

Initialize a set of indices $\mathcal{I} \leftarrow \{1, \dots, N\}$, a static instruction prompt in the new dataset p_{ins} , and an initial question for querying OpenAI API q_{gen} .

Manually craft 10 correct Prolog codes $\{a_i^{\text{Prolog}}\}_{i=1}^{10}$ that correctly solve $\{q_i\}_{i=1}^{10}$ in \mathcal{X} to initialize \mathcal{D}

for $i \in \mathcal{I}$ **do**

 Retrieve a sample $(q_i, a_i^{\text{CoT}}) \in \mathcal{X}$

 Prompt GPT-4 with $\{q^{\text{gen}}\} \cup \{(q_k, a_k^{\text{CoT}}, a_k^{\text{Prolog}})_{k=1}^{10}\} \cup \{q_i, a_i^{\text{CoT}}\}$ to obtain a_i^{Prolog}

if $\mathcal{P}(a_i^{\text{Prolog}}) = \mathcal{C}(a_i^{\text{CoT}})$ **then**

$\mathcal{D} \leftarrow \mathcal{D} \cup \{(p_{\text{ins}}, q_i, a_i^{\text{Prolog}})\}$

$\mathcal{I} \leftarrow \mathcal{I} \setminus \{i\}$

end if

end for

Manually select the top 10 clean and logical Prolog code from the current \mathcal{D} to form a new few-shot sample set $Q^{\text{fixed}} = \{(q_k, a_k^{\text{CoT}}, a_k^{\text{Prolog}})_{k \notin \mathcal{I}}\}$, $|Q^{\text{fixed}}| = 10$.

for $j = 1, \dots, M$ **do**

// M trial attempts

for $i \in \mathcal{I}$ **do**

 Retrieve a sample $(q_i, a_i^{\text{CoT}}) \in \mathcal{X}$

 Sample $Q^{\text{random}} \leftarrow \{(q_k, a_k^{\text{CoT}}, a_k^{\text{Prolog}})_{k \notin \mathcal{I}}\}$, $|Q^{\text{random}}| = 10$ from \mathcal{D}

 // Adding 10 dynamic samples and 10 fixed samples into the 20-shot prompt.

 Prompt GPT-4 with $\{q^{\text{gen}}\} \cup Q^{\text{fixed}} \cup Q^{\text{random}} \cup \{q_i, a_i^{\text{CoT}}\}$ to obtain a_i^{Prolog}

if $\mathcal{P}(a_i^{\text{Prolog}}) = \mathcal{C}(a_i^{\text{CoT}})$ **then**

$\mathcal{D} \leftarrow \mathcal{D} \cup \{(p_{\text{ins}}, q_i, a_i^{\text{Prolog}})\}$

$\mathcal{I} \leftarrow \mathcal{I} \setminus \{i\}$

end if

end for

end for

if $\mathcal{I} \neq \emptyset$ **then**

 Manually correct Prolog codes $\{a_i^{\text{Prolog}}\}_{i \in \mathcal{I}}$ that solve $\{q_i\}_{i \in \mathcal{I}}$

$\mathcal{D} \leftarrow \mathcal{D} \cup \{(p_{\text{ins}}, q_i, a_i^{\text{Prolog}})_{i \in \mathcal{I}}\}$

end if

A.2 Permutation procedures

Permutations can be performed both on the level of facts or rules and on the level of goals in a rule. In practice, for each piece of code, we first permute the goals in the solve<answer>: -<goal_1>, <goal_2>, ... predicate. Since the total number of permutations is sensitive to the number of goals and can easily grow to a large magnitude, thus running out of memory, we used the permutation method in the itertools library to yield an iterator over the permutations. Then, we took up to 10 goal permutations from the iterator. If there were less than 10 goal permutations in total because the code was concise and there were not many goals, we took as many goal permutations as possible. Then, in the same manner, we took up

to 10 fact and rule permutations. In principle, there would be at most 100 permuted samples generated for one original sample. Then, for each sample, while conducting an experiment that required a certain number of permutations, we randomly sampled permutations from the set of permutations of size up to 100. For some sample, if the target number of permutations exceeded the total permutations it had, we took all its permutations instead.

A.3 Instruction Prompt Samples

Below are the instruction prompts we used for different training settings (CoT, Prolog, and Permuted Prolog).

Setting	Prompt Template
CoT	<p>Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.</p> <p>### Instruction: Please generate an explanatory answer to solve the given math problem.</p> <p>### Input: <Question></p> <p>### Output: <CoT Reasoning></p>
Prolog Generation	<p>Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.</p> <p>### Instruction: Please generate a piece of Prolog code to solve the given math problem.</p> <p>### Input: <Question></p> <p>### Output: <Prolog Code></p>
Permuted Prolog Generation	<p>Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.</p> <p>### Instruction: Please generate a piece of Prolog code in non-sequential order to solve the given math problem.</p> <p>### Input: <Question></p> <p>### Output: <Prolog Code></p>

A.4 Error Analysis

In this section, we present some typical errors made by our best-performing model to understand the behavior and bottleneck of LLM generation of prolog programs to solve arithmetic reasoning tasks. We sampled 70 problematic Prolog codes generated by finetuned Mistral using the 1:1 permutation setting, which had yielded the highest accuracy, that had errors and manually checked the type of errors using the Swi-Prolog interpreter. 100% of them had semantic errors, meaning their logic of solution was wrong. 7% of them had syntax errors, meaning the Swi-Prolog yields syntax errors on those codes. Some examples of typical syntax errors are listed below. The lines causing errors are in bold. The details of the errors are explained in the comments.

Question	Prolog Code
<p>Julia’s boat sprang a leak while she was out on the lake. The boat was taking on two liters of water for every ten feet she rowed back towards shore. It took her sixteen seconds to row twenty feet. The shore was 64 seconds away. How much water had the boat taken on by the time Julia reached shore?</p>	<pre>:- use_module(library(clpqq)). liters_per_10_feet(julia, 2). time_per_20_feet(julia, 16). distance_to_shore(julia, 64). solve(Water_taken_on) :- liters_per_10_feet(julia, Liters_per_10_feet), time_per_20_feet(julia, Time_per_20_feet), distance_to_shore(julia, Distance_to_shore), {Seconds_per_10_feet = Time_per_20_feet / 2}, {Total_seconds = Distance_to_shore / 10 * Seconds_per_10_feet}, {Total_feet_rowed = Total_seconds / Seconds_per_20_feet * 20}, % The variable has never been defined. {Total_liters_taken_on = Total_feet_rowed * Liters_per_10_feet}, {Water_taken_on = Total_liters_taken_on}.</pre>
<p>Bill bakes 300 rolls, 120 chocolate croissants, and 60 baguettes every day. Each roll is 4 inches long, each croissant is 6 inches long, and each baguette is two feet long. If Bill puts all the baked goods end to end, how long will they be in feet?</p>	<pre>:- use_module(library(clpqq)). bake(bill, rolls, 300). bake(bill, chocolate_croissants, 120). bake(bill, baguettes, 60). length(roll, 4). # The predicate has already been reserved by Prolog. length(chocolate_croissant, 6). length(baguette, 2). solve(Total_length) :- bake(bill, rolls, Num_rolls), bake(bill, chocolate_croissants, Num_chocolate_croissants), bake(bill, baguettes, Num_baguettes), length(roll, Roll_length), length(chocolate_croissant, Chocolate_croissant_length), length(baguette, Baguette_length), {Roll_total_length = Num_rolls * Roll_length}, {Chocolate_croissant_total_length = Num_chocolate_croissants * Chocolate_croissant_length}, {Baguette_total_length = Num_baguettes * Baguette_length}, {Total_length = Roll_total_length + Chocolate_croissant_total_length + Baguette_total_length}. % This code also contains a semantic error because it fails to convert the unit of inch to feet.</pre>
<p>In a race with 240 Asians, 80 were Japanese, and the rest were Chinese. If the number of boys on the Chinese team was 60, how many girls were on the Chinese team?</p>	<pre>:- use_module(library(clpqq)). total_participants(240). japanese_participants(80). chinese_boys(60). solve(Girls_on_chinese_team) :- total_participants(Total_participants), Japanese_participants, % The correct syntax: japanese_participants(Japanese_participants), {Chinese_participants = Total_participants - Japanese_participants}, chinese_boys(Chinese_boys), {Girls_on_chinese_team = Chinese_participants - Chinese_boys}.</pre>
<p>Brinley is in Mr. Bert’s math class. Mr. Bert gives six tests each semester. He allows the students to remove the lowest score from the average each semester. If Brinley has scores of 89, 71, 92, 100 and 86 on the first five tests, what score does she need on the sixth test to get an average of 93?</p>	<pre>:- use_module(library(clpqq)). num_tests(mr_bert, 6). lowest_score_removal(mr_bert, 1). target_average(mr_bert, 93). scores(brinley, [89, 71, 92, 100, 86]). solve(Test_score) :- num_tests(mr_bert, Num_tests), lowest_score_removal(mr_bert, Lowest_score_removal), target_average(mr_bert, Target_average), scores(brinley, Scores), Length is Num_tests - Lowest_score_removal, {Total_score = sum(Scores)}, % The built-in predicate is misused. {Average_score = Total_score / Length}, {Test_score = (Target_average * Length) - Total_score}.</pre>

A.5 Training Details and Computational Budget

During finetuning, we controlled the number of epochs to be 6, batch size to be 128, and learning rate to be 3×10^{-4} . For a single training run, we used 2 NVIDIA RTX 4090 GPUs to finetune Llama2 and CodeLlama and 2 NVIDIA RTX 8000 GPUs to finetune Mistral. We adopted Distributed Data Parallelism to speed up training. Training on the original CoT data in GSM8K or the non-permuted Prolog code data took around 2 hours on 2 NVIDIA RTX 4090 GPUs and around 10 hours on 2 NVIDIA RTX 8000 GPUs. When we added in permuted samples, the training time grew proportionally with the dataset size since we controlled the number of epochs and batch size. During inference on the test set, we used a batch size of 2

on an RTX 4090 GPU, which took around 6 hours to finish a full inference round, and a batch size of 3 on one RTX 8000 GPU, which took around 7 hours to finish a full inference round.