

LogicPro: Improving Complex Logical Reasoning via Program-Guided Learning

Jin Jiang^{1,2*}, Yuchen Yan^{2,3}, Yang Liu², Yonggang Jin², Shuai Peng¹
Mengdi Zhang², Xunliang Cai², Yixin Cao⁴, Liangcai Gao¹, Zhi Tang¹

¹Peking University ²Meituan Group ³Zhejiang University ⁴Fudan University
jiangjin@stu.pku.edu.cn, yanyuchen@zju.edu.cn,
{liuyang509, zhangmengdi02}@meituan.com, yxcao@fudan.edu.cn, gaoliangcai@pku.edu.cn

Abstract

In this paper, we present a novel approach, called **LogicPro**, to enhance Large Language Models (LLMs) complex Logical reasoning through **Program Examples**. We do this effectively by simply utilizing widely available algorithmic problems and their code solutions. First, we constructed diverse test samples input based on algorithmic questions and code solutions. Then, we designed different complex reasoning questions based on algorithmic problems and test samples. Finally, combining the intermediate variable outputs of the code solutions and the complex reasoning questions, we derived the reasoning process and the final answer. With this approach, we can construct a dataset that is sufficiently **difficult** (all models are ineffective), **diverse** (synthesized from 2,360 different algorithmic questions), and **scalable** (building different test samples and collecting more algorithmic questions). In addition, we obtain a **high-quality reasoning process** guided by the values of intermediate variables. As a result, our approach achieves significant improvements in multiple models for the *BBH*²⁷, *GSM8K*, *HellSwag*, *Logicqa*, *Reclor*, and *RTE* datasets, outperforming a wide range of existing reasoning datasets.

Introduction

Large language models (LLMs) have shown remarkable capabilities across various domains (Achiam et al. 2023; Schaeffer, Miranda, and Koyejo 2024). Utilizing chain-of-thought (CoT) prompting (Wei et al. 2022), LLMs demonstrate strong reasoning abilities (Bubeck et al. 2023). In fields like mathematics (Shao et al. 2024) and coding (Guo et al. 2024), LLMs have achieved significant progress. Nonetheless, they still encounter challenges in complex logical reasoning tasks (Suzgun et al. 2023), underscoring the need for more high-quality reasoning data in this area.

In recent years, many studies have focused on constructing high-quality reasoning data. Some data come from real-world sources, such as the LOGIQA (Liu et al. 2023a) from civil service exams and AR-LSAT (Zhong et al. 2021), RECLOR (Yu et al. 2020) from law exams. These questions form **diverse** and **difficult** datasets, but collecting real data is time-consuming and labor-intensive for further training. In contrast, synthetic data offers an efficient method to generate **large-scale** datasets. Llama3.1 (Dubey et al. 2024) has

shown strong capabilities when trained on massive synthetic data. The synthetic mathematical data (Lu et al. 2024; Huang et al. 2024) significantly improved the model’s mathematical abilities. Additionally, some research has synthesized reasoning data in the logical reasoning (Nie et al. 2020; Saeed et al. 2021) field while providing the reasoning process (Dalvi et al. 2021) during synthesis. However, existing complex logical synthetic data are often constructed based on propositional logic and its combinations $((p \vee q) \vdash (q \vee p))$, limiting their difficulty and diversity. Consequently, large models show limited improvement in complex reasoning when learning from such data, highlighting the need for more challenging and diverse datasets. Additionally, a **high-quality reasoning process** is a core component of superior reasoning datasets. Luo et al. (2024) optimized the reasoning process in mathematical reasoning tasks using the Monte Carlo Tree Search (MCTS) method, followed by reinforcement learning. Some studies (Zhang et al. 2023; Yao et al. 2024; Besta et al. 2024) have also employed the “X-of-Thought” (XoT) methods to achieve better reasoning outcomes like code, math, etc. However, these methods primarily rely on the model’s own improvements, which have inherent limitations. Another approach is to leverage external models to enhance the reasoning process. Liu et al. (2023c) utilized stronger models (such as GPT-4) in reasoning tasks to obtain better results, but such approaches often struggle to ensure that the reasoning process and outcomes are entirely correct. Therefore, obtaining high-quality reasoning process data remains a crucial and challenging task.

In this paper, we propose a data synthesis method that provides high-quality reasoning questions and reasoning processes utilizing widely available algorithm problems and program solutions. First, we construct multiple test case inputs based on the algorithm problem statements and their corresponding Python code. Then, we use model rewriting to generate logical reasoning problems by combining different test case inputs with the corresponding problem statements. Next, we consider the test case inputs and the code solution to construct a solution for the current example, obtaining intermediate variable values and final results from the code. Finally, by integrating the problem statements with the intermediate variable outputs, we derive the final reasoning path. The reasoning questions constructed by this method meet the following three criteria: Sufficient **Difficulty**: As

*Contribution during internship at Meituan Group.

shown in Table 4, all model results are below 50%; **Diversity**: Our data is synthesized from 2360 different LeetCode algorithm problems; **Scalability**: Our data can be further expanded by creating more test samples and collecting more algorithmic question. Moreover, the method is able to generate **high-quality reasoning processes**: guided by intermediate variable values, we construct reasoning processes that are accurate at the step level.

Our main contributions are summarized as follows:

1. We propose a novel complex reasoning enhancement method that utilizes widely available algorithmic questions and their codes to generate reasoning data.
2. This approach allows us to construct a dataset that is sufficiently difficult, diverse and scalable, while also obtaining a high-quality reasoning process guided by the values of intermediate variables.
3. The experimental results show that our approach achieves significant improvements on multiple models for the *BBH*²⁷, *GSM8K*, *HellSwag*, *Logicqa*, *Reclor*, and *RTE* datasets, outperforming a wide range of existing reasoning datasets.

Approach

Our method, as illustrated in Figure 2, is divided into four steps: (1) Building test samples by using algorithm problems and corresponding code, (2) Constructing complex reasoning problems based on these test samples and algorithm problems, (3) Deriving intermediate variable results from the code and test samples, and (4) Forming final question-answer pairs by integrating complex reasoning problems with the intermediate variable results. Through these steps, we generate complex logical reasoning question-answer pairs, incorporating reasoning processes from algorithmic problem statements and Python solution code.

Step 1: Build Test Sample Inputs

In the first step, our input consists of LeetCode algorithmic questions and their corresponding Python code solutions. We provide these questions and code to the open-source LLM (Llama3-70B-Instruct), prompting the model to generate appropriate test sample inputs for each question (see Appendix: Complete Prompts, Figure 8). We require the model to generate 30 test sample inputs in a single response. Specifically, we set the temperature to 0.7, perform multiple inferences, extract the test sample inputs from these results, and de-duplicate the generated samples to obtain the final set of test samples.

As in the example in Figure 2, we enter the description of LeetCode Question 70: Climbing Stairs¹ and the corresponding Python solution. The model gives inputs of test samples that match the requirements of the question, e.g., $n = 17$, as inputs for subsequent steps.

Step 2: Construct Complex Reason Problem

For the second step, our inputs are the LeetCode algorithmic question and the constructed test sample input. As the

¹<https://leetcode.com/problems/climbing-stairs/description/>

Step 2: Construct Complex Reasoning Problem
<p>I have an <u>algorithmic question</u> and a corresponding <u>test input</u>: please rewrite the algorithmic question as a text-only <u>logical reasoning question</u> based on the test input.</p> <p>Instructions:</p> <ol style="list-style-type: none"> 1. Please incorporate the test input into the description ... 2. Please first give the name of this logical reasoning ... <p>Reference case:</p> <ul style="list-style-type: none"> - <u>Algorithmic Question</u>: Given a sequence containing only (,), {, }, [,], <, >, complete the rest of the sequence, making sure that all the parentheses are properly closed and in the right order. - <u>Test Input</u>: "<> ([[(({ }))] [<>]]" - <u>Text-Only logical reasoning question</u>: Title: Correctly close a Dyck-n word. <p>Q: Complete the rest of the sequence, making sure that the parentheses are closed properly. Input: <> ([[(({ }))] [<>]]</p> <p>Refer to the above example of rewriting an algorithmic question into a text-only logical reasoning question based on test input:</p> <ul style="list-style-type: none"> - <u>Algorithmic Question</u>: {LeetCode 70: Climbing Stairs} - <u>Test Input</u>: n=17 - <u>Text-Only logical reasoning question</u>: ...

Figure 1: The second step of constructing a complex logical reasoning problem prompt example, the full prompt is in Appendix: Complete Prompts, Figure 9.

prompt in Figure 1, we input the test sample input and the description of the LeetCode algorithmic question into the model and ask the model to generate a concrete complex reasoning question. To help the model better understand our instructions, we provide an example *close to a Dyck-n word* (from the *BBH*²⁷ subset) as a reference.

Combined with the example in Figure 2, the model reformulates the "Climbing Stairs" algorithm question into a specific complex reasoning query based on the test sample input ($n = 17$). The question becomes: "How many paths are there to reach the top of the staircase when climbing 17 flights of stairs?"

Step 3: Obtain Intermediate Variable From Program

In this step, we start with the Python code solution and the input of the constructed test sample. We prompt the model to rewrite the original code to both fit the test sample input and print the values of key intermediate variables (Prompt in Appendix: Complete Prompts 11). This approach ensures that the rewritten code can run directly and represents the specific logic of the algorithm through the printed intermediate variables.

For example, for the *Climbing Stairs* problem in Figure 2, the model rewrites the code to output not only the number of ways to climb 17 stairs (Result: *stairs* = 2584) but also the values of important intermediate variables. Specifically, the problem can be solved using the Fibonacci series, where the variables 'a' and 'b' represent the number of ways to reach stairs $n - 2$ and $n - 1$, respectively. The initial values $a = 1$ and $b = 2$ mean that there is 1 way to reach level 1 and 2 ways to reach level 2. Printing these variables helps illustrate the logical progression of the algorithm.

Step 4: Program-Guided Reasoning

For step 4, our input consists of the complex reasoning question constructed in step 2 and the intermediate variable output of the rewritten code in step 3. As shown in Figure 3, the prompt in Step 4 is primarily divided into three parts:

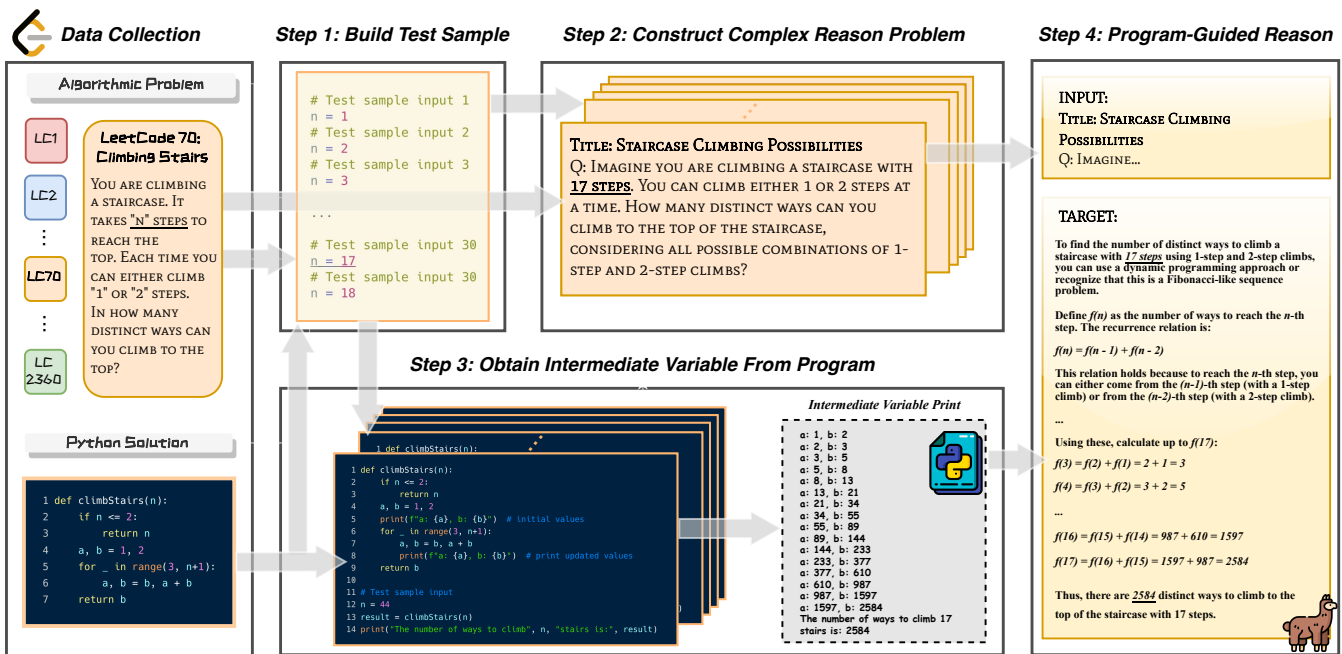


Figure 2: An overview of LogicPro (Example: LeetCode 70): 1. Build Test Sample (Section Step 1), using algorithm problems and corresponding code to construct test samples. 2. Construct Complex Reason Problem (Section Step 2), based on algorithm problems and test samples to create complex reasoning problems. 3. Obtain Intermediate Variable From Program (Section Step 3), deriving intermediate variable results from code and test samples. 4. Program-Guided Reason (Section Step 4), forming final question-answer pairs from complex reasoning problems and intermediate variable results (complete question-answer pairs are in Figure 7).

Step 4: Program-Guided Reasoning

There is a logical reasoning question and the intermediate variable output of its code solution. Please answer this logical reasoning question based on the intermediate variable output of the code.

Instructions

1. Refer to the code's intermediate variable outputs. Use information ...
2. First, outline ... Then, provide the exact reasoning process step by step.
3. Do not use code to solve this logical reasoning ...
4. Avoid phrases like "From the intermediate variables" ...
5. Please ensure that the logical reasoning questions are answered in a way that reflects each step of the reasoning process in the code.

Reference case:

- Logical reasoning question: {logical_question_ref}
- Code intermediate variables: {code_print_ref}
- Logical Reasoning Question Answer: {answer_with_print_ref}

Refer to the above case to give a solution to a reasoning question:

- Logical reasoning question: {logical_question}
- Code intermediate variables: {code_print}
- Logical Reasoning Question Answer: ...

Figure 3: The 4 step in constructing the final answer is illustrated in the example, with the full prompt shown in Appendix: Complete Prompts, Figure 12.

instruction, reference, and the current question. The first part, the instruction, outlines the requirements of the current task, specifically instructing the model to generate an improved answer based on the values of intermediate variables. The second part provides a reference case. We observed that open-source models often struggle to follow instructions and generate a complete reasoning process directly from intermediate variables. To address this, we initially used GPT-4 (gpt-4-turbo-2024-04-09) to construct a reference for each algorithm problem (Prompts are detailed in the Appendix:

Complete Prompts, Figure 13 and Figure 14). The third part focuses on the current question. In the reference case, the reasoning process, intermediate variables, and the final answer for a specific test case of the same algorithm problem are presented as a one-shot example. The model then uses this example to reason through and provide answers to other test cases. Through this approach, we can achieve high-quality code-driven reasoning processes at a low cost.

For the question *climb stairs* in Fig. 2, we first let GPT-4 construct a good answer based on a certain test sample (e.g., in $n = 10$), and then use the answer in $n = 10$ as a reference sample to let the open source model rewrite it on other test samples of the same topic (e.g., at $n = 11, 12, 13...17$) to get the final answer for all final answers for the data.

Dataset-LogicPro

Based on the above process, we constructed the LogicPro dataset. First, we collected 2,360 LeetCode algorithm problems and their corresponding solutions. Next, following step 1, we created test cases for all the algorithm problems, resulting in a dataset of 117,793 entries. In step 3, we filtered out data where code execution failed, with 104,824 entries running successfully and 12,969 failing. Finally, considering the varying number of test samples for different algorithm problems in previous steps, we set a threshold of 30, extracting up to 30 samples per problem from past data. After these steps, our final training dataset size is 69,866. In Table 5 of Appendix, we compare LogicPro with other datasets.

Base Model	SFT Data	BBH ²⁷	GSM8K	HellaSwag	Logiqa	Reclor	RTE	Average
GPT-4	-	74.9	94.2	83.0	69.0	88.5	77.5	76.1
ChatGPT	-	50.3	65.3	73.0	44.5	61.0	90.0	52.8
Qwen1.5-7B	General-100k	45.0	65.5	69.5.0	50.0	62.5	82.5	48.3
	Gen-Logic	44.8	65.2	71.0	54.0	63.5	84.0	48.3
	+ Code	45.6	65.6	74.5	55.0	61.5	84.0	49.1
	+ LogicPro_{ours}	47.9	66.2	71.0	55.5	64.0	84.0	51.1
Llama-2-7B	General-100k	34.9	23.4	29.5	40.0	40.0	52.5	35.2
	Gen-Logic	35.0	24.6	40.0	37.5	45.0	52.5	35.8
	+ Code	34.0	25.2	35.0	38.5	44.5	52.5	34.8
	+ LogicPro_{ours}	35.8	25.7	35.0	40.0	49.5	53.0	36.6
Llama-3-8B	General-100k	43.3	63.8	38.8	57.5	60.0	76.0	45.9
	Gen-Logic	44.7	63.3	41.5	59.0	64.0	77.5	47.2
	+ Code	45.2	63.7	41.3	57.0	67.0	77.5	47.7
	+ LogicPro_{ours}	49.3	65.7	44.3	59.0	67.5	78.0	51.4
Llama-2-13B	General-100k	40.9	34.2	35.5	44.0	56.0	75.5	42.2
	Gen-Logic	38.8	35.4	38.3	43.0	59.5	77.0	40.6
	+ Code	38.8	35.9	35.8	45.0	57.5	76.5	40.6
	+ LogicPro_{ours}	40.2	44.0	37.3	46.5	60.0	81.0	42.4
Qwen1.5-14B	General-100k	50.8	69.8	47.0	62.0	68.0	84.0	53.2
	Gen-Logic	51.2	71.7	63.5	65.0	69.0	85.0	54.2
	+ Code	51.4	70.3	64.0	64.0	68.0	86.0	54.4
	+ LogicPro_{ours}	51.8	71.8	66.0	65.0	66.5	86.0	54.8

Table 1: The results of LogicPro data on different models. General-100k and Gen-Logic serve as baseline results (Section Train Datasets). + Code refers to training on Gen-Logic combined with 2,360 original algorithm problems. + **LogicPro_{ours}** is trained on Gen-Logic combined with our LogicPro data. The complete results can be found in Table 6 of the Appendix.

Experiments

Train Datasets

We collected open-source general and reasoning data to train our baseline model. The general data are mainly from OpenHermes-2.5 (Teknum 2023), where we extracted all the *alpaca* data and randomly sampled from other data to bring the total number of data to 100,000 (**General-100k**). The logical reasoning data were mainly from LogiCoT (Liu et al. 2023c), totaling 68,983 items. We mixed the above generalized and logic data to construct the SFT dataset (General-mix-Logic, abbreviated as **Gen-Logic**) to train the baseline model. Next, the constructed LogicPro data was mixed with other data and trained to verify its effectiveness.

To further validate the effectiveness of LogicPro data, we categorized the open source reasoning data into three dimensions for further comparison. The first dimension is real-world reasoning data, including the larger LOGIQA (Mill 2013), RECLOR (Yu et al. 2020) and AR-LSAT (Wang et al. 2022). The second dimension is synthetic data, including Entailmentbank (Dalvi et al. 2021), ProntoQA (Saparov and He 2023) and RuleTaker (Clark, Tafjord, and Richardson 2021). The third dimension is mathematical reasoning data, including GSM8K (Cobbe et al. 2021), ASDIV (Miao, Liang, and Su 2020), and SVAMP (Patel, Bhattamishra, and Goyal 2021). Based on these data, we first conduct experiments on each dataset individually. Then, we experiment

with the data mixed within the three specified dimensions. Finally, we experiment with all the data combined. Specifically, each dataset is combined with Gen-Logic to train a separate model for parallel comparison with our LogicPro.

Eval Datasets

We evaluated the model in a Zero-Shot CoTs format on the BBH (Suzgun et al. 2023), GSM8K (Cobbe et al. 2021), Hellswag (Zellers et al. 2019), Logiqa (Liu et al. 2023a), Reclor (Yu et al. 2020) and RTE (Wang et al. 2020) datasets. BBH is the core benchmark for evaluating the model’s complex reasoning ability, containing 27 challenging reasoning tasks. Therefore, we computed averages by weighting subsets of the data. Specifically, for LogicPro data, these evaluation datasets are **out-of-domain**.

Metrics

We report accuracy metrics for all evaluation tasks. For GSM8K, results are obtained through rule extraction, with exact matching used to compute the corresponding metrics. For other datasets, such as BBH and HellSwag, we employ an internal scoring model for evaluation. This scoring model is trained on the open-source LLaMA model, using historical data from GPT-4 and human ratings. Our prompt model assesses the accuracy of the current answer by comparing the model’s response to the reference answer, outputting a

score in JSON format (either 0 or 1).

Baselines

For proprietary models, we present results from state-of-the-art large language models (LLMs) such as OpenAI’s GPT-4 (gpt-4-0613) and ChatGPT (gpt-3.5-turbo). For open source models, we performed full parameter fine-tuning experiments on Llama-2 (7B-13B), Llama-3 (8B), Qwen1.5 (7B-14B), Yi1.5-9B and Deepseek-7B (llm-math encoder). All of our experiments were performed on base models without supervised fine-tuning (SFT).

We employed a learning rate of $3e-06$ with a cosine learning rate strategy. The training process included a global batch size of 128, a warmup period of 3%, and was carried out over 3 epochs. All models were trained using the Megatron framework. During the inference phase, we utilized the vllm tool with a greedy decoding strategy, and the maximum output length was set to 2048 tokens.

Results

LogicPro on Different Base Models

Table 1 shows the results of different base models on General-100k, Gen-Logic, Gen-Logic+Code and Gen-Logic+LogicPro. Overall, the LogicPro method consistently enhances model performance, especially in GSM8K, Logiqa, Reclor, and RTE tasks. Categorizing models by size reveals key insights: For smaller models: Qwen1.5-7B and Deepseek-7B show significant gains with LogicPro, improving from 48.3 to 51.1 and 36.3 to 38.7, respectively, particularly in BBH and GSM8K tasks. This suggests that LogicPro’s high-quality reasoning data effectively compensates for the limited capacity of smaller models. For Llama-2-7B and Llama-3-8B benefit from LogicPro, with average scores rising from 35.8 to 36.6 and 47.2 to 51.4, particularly in Logiqa and Reclor. This indicates that LogicPro enhances logical reasoning even in models with strong generalization. For larger models (9B-14B): Larger models like Yi-1.5-9B and Qwen1.5-14B see more modest gains, from 54.5 to 56.1 and 54.2 to 54.8, suggesting that while general and code data may saturate performance, LogicPro still offers an edge in complex tasks. The slight decrease in Llama-2-7B performance with code data (from 35.8 to 34.8) may result from interference or overfitting, highlighting that code data can sometimes introduce noise. However, LogicPro overcomes this, providing a net gain. General logic and code data can enhance performance, especially in smaller models, but their impact diminishes with larger models. LogicPro’s curated data continues to drive improvements, underscoring its value in complex reasoning tasks.

Comparison of LogicPro with Other Datasets

Table 2 presents the comparison results of LogicPro against other reasoning datasets. Considering its affordable size with best performance, and significant community impact, we chose Qwen1.5-7B as the model for data comparison. Overall, LogicPro achieves the highest overall performance, particularly excelling in BBH²⁷ and GSM8K, where it matches or surpasses Realistic Logical^{mix3} and SVAMP.

It also leads in Reclor and performs strongly in HellaSwag. Mathematical logical datasets perform well in GSM8K but generally lag in other areas. The **R-S-M-Mix-Logic** dataset excels in RTE, showing that mixed logical data can produce strong results in specific tasks. While specialized datasets show strengths in certain areas, LogicPro’s diverse and challenging nature enables it to enhance overall reasoning abilities more effectively. This makes LogicPro a superior choice for becoming a reasoning expert.

Ablation Study

As shown in Table 3, we conducted an ablation study on the Qwen1.5-7B model to evaluate the impact of different SFT data on model performance. Using **General-100k** and **Gen-Logic**, the model achieved an average accuracy of 48.3% on both the BBH²⁷ and GSM8K tasks. When the original 2360 algorithm code data was introduced, the model’s performance improved slightly to an average accuracy of 49.1%, but it remained significantly lower than that of LogicPro. Furthermore, when the code data was repeated 30 times (+ **Code*30**) to match the scale of LogicPro data, the average accuracy actually decreased to 46.1%. This suggests that the performance gains achieved by LogicPro over code are not simply due to an increase in data volume. Regarding the ablation study on intermediate variables, the model trained with LogicPro data without intermediate variables (**LogicPro_{w/o} Inter-Var**) achieved an average accuracy of 49.7%, already surpassing the performance of code-based training. However, when intermediate variables were incorporated into the LogicPro data construction (**LogicPro_w Inter-Var**), the model achieved the highest performance, with accuracies of 47.9% and 66.2% on the BBH²⁷ and GSM8K tasks, respectively, and an overall average accuracy of 51.1%. This demonstrates that the inclusion of intermediate variables is crucial for enhancing the model’s reasoning capabilities, leading to a significant improvement in complex task performance.

Case Study

<p>Question: Vina tells the truth. Alejandro says Vina lies. Maybelle says Alejandro tells the truth. Millicent says Maybelle lies. Elanor says Millicent tells the truth. Does Elanor tell the truth?</p>	
<p>Answer (Qwen1.5-7B: Gen_Logic+ LogicPro):</p> <ol style="list-style-type: none"> Vina tells the truth. Alejandro says Vina lies. Maybelle says Alejandro tells the truth. Millicent says Maybelle lies. Elanor says Millicent tells the truth. <p>Now, let's analyze the statements in relation to each other:</p> <ul style="list-style-type: none"> - Vina tells the truth, so Alejandro is lying. - Maybelle says Alejandro tells the truth, but Alejandro is lying, so Maybelle is lying. - Millicent says Maybelle lies, but Maybelle is lying, so Millicent is telling the truth. - Elanor says Millicent tells the truth, and Millicent is telling the truth, so Elanor is telling the truth. <p>Therefore, <u>Elanor does tell the truth.</u></p>	<p>Answer (Qwen1.5-7B: Gen_Logic):</p> <ol style="list-style-type: none"> Vina tells the truth. Alejandro says Vina lies. Maybelle says Alejandro tells the truth. Millicent says Maybelle lies. Elanor says Millicent tells the truth. <p>Now, let's analyze the statements in relation to each other:</p> <ol style="list-style-type: none"> Vina tells the truth. Alejandro says Vina lies. Maybelle says Alejandro tells the truth. Millicent says Maybelle lies. Elanor says Millicent tells the truth. <p>We can see that Alejandro and Maybelle are contradicting each other. Since we <u>don't have enough information</u> to determine who is telling the truth, we cannot determine if Elanor is telling the truth or not.</p>

Figure 4: Success Case for LogicPro: *bbh web of lies*. The correct answer is: Elanor tell the truth.

To further understand the impact of LogicPro on enhancing complex reasoning capabilities, we provide case study examples in this section. As shown in Figure 4, BBH Web of Lies is a subset of BBH used to evaluate the truth value of

Logic Data	BBH ²⁷	GSM8K	HellaSwag	Logiqa	Reclor	RTE	Average
-	44.8	65.2	71.0	54.0	<u>63.5</u>	84.0	48.3
Realistic Logical ^{mix3}	47.9	65.4	<u>75.0</u>	54.0	60.5	84.5	<u>51.0</u>
AR-LSAT (Zhong et al. 2021)	44.9	64.2	73.0	55.0	62.5	84.0	48.5
Logiqa (Liu et al. 2023a)	44.6	64.2	74.5	58.0	63.0	84.0	48.4
Reclor (Yu et al. 2020)	47.3	64.7	72.0	53.0	<u>63.5</u>	84.0	50.5
Synthetic Logical ^{mix3}	47.3	64.4	76.5	53.0	61.0	83.5	50.5
Entailmentbank (Dalvi et al. 2021)	44.0	65.4	74.5	54.0	60.5	84.5	47.7
ProntoQA (Saparov and He 2023)	46.9	64.7	73.0	53.5	63.0	84.0	50.2
RuleTaker (Clark et al. 2021)	46.4	67.1	74.5	54.5	<u>63.5</u>	84.5	49.9
Mathematical Logical ^{mix3}	46.7	65.7	74.5	56.0	<u>63.5</u>	84.5	50.1
ASDiv (Miao, Liang, and Su 2020)	45.7	66.7	73.5	<u>57.0</u>	63.0	83.5	49.3
GSM8K (Cobbe et al. 2021)	46.3	65.3	72.0	55.0	62.0	<u>85.0</u>	49.7
SVAMP (Patel et al. 2021)	47.1	<u>66.2</u>	74.0	54.0	61.5	84.5	50.4
R-S-M-Mix-Logic	<u>47.6</u>	65.0	69.5	57.0	62.5	86.0	50.7
LogicPro	47.9	<u>66.2</u>	71.0	55.5	64.0	84.0	51.1

Table 2: Results on LogicPro-Train vs. Different Logic Data. Base Model: Qwen1.5-7B. Baseline Data: Gen-Logic. ^{mix3} represents the result of mixing the data within each dimension, and **R-S-M-Mix-Logic** represents the result of mixing the data from all three dimensions.

SFT Data	BBH ²⁷	GSM8K	Average
General-100k	45.0	65.5	48.3
Gen-Logic	44.8	65.2	48.3
+ Code	45.6	65.6	49.1
+ Code*30	43.5	62.4	46.1
+ LogicPro _{w/o Inter-Var}	<u>46.2</u>	<u>65.9</u>	<u>49.7</u>
+ LogicPro _{w. Inter-Var}	47.9	66.2	51.1

Table 3: Ablation Study on different *SFT* data. + denotes the mixing of different data based on Gen-Logic. "code" refers to the source 2360 LeetCode algorithm question. "Code*30" indicates repeating the "code" 30 times to match the calculation of LogicPro. LogicPro_{w/o Inter-Var} and LogicPro_{w. Inter-Var} respectively whether the data construction process uses intermediate variables. Base Model: Qwen1.5-7B.

a random Boolean function expressed as a natural language word problem. For this problem, the baseline model merely organized and listed the conditions in the problem with repetition, failing to provide the correct answer. In contrast, the model incorporating LogicPro data offered a clear and accurate reasoning path, systematically analyzing the lying situations of Alejandro, Maybelle, and Millicent, and ultimately correctly judged the truthfulness of Elanor’s statement. Appendix provides more case studies on the results of BBH subsets. Furthermore, we provide case analyses for the construction of test samples and LogicPro data in Appendix and Appendix , respectively.

Analysis

LogicPro under different types base models As shown in Figure 5, we observe significant differences in the per-

formance of various model types in certain subsets. Firstly, the math model shows a noticeable decline in performance on the *object counting* task compared to other models, likely due to its original strength in numerical capabilities. Secondly, the coder model exhibits unusual trends on several subsets. The coder model’s average improvement (+12.5%) is significantly higher than that of the LLM model (+2.7%) and the math model (-1.9%). On the *ruin names* and *dyck_language* subsets, despite only slight changes in other models, the coder model shows significant improvement. This anomaly might be due to the logicpro data originating from synthesized code data. Additionally, Figure 4 illustrates the changing trends of more general models in the BBH subset, providing a basis for comparative analysis.

Different Proportions of LogicPro As shown in Figure 6, as the data proportion increases, the performance of most datasets improves. In the BBH²⁷ and GSM8K datasets, the performance improvement is particularly notable, especially when the data volume increases from 10% to 30%. Some datasets, like Reclor and RTE, show similar performance across different data proportions, lacking a clear upward trend, which may indicate a lower sensitivity to data volume. Overall, the "average" performance gradually improves with increased data volume, indicating that when integrating multiple datasets, data volume positively impacts performance enhancement.

Discussion of Difficulty and Diversity of LogicPro In the previous section, we mentioned that LogicPro is a sufficiently challenging, diverse, and scalable dataset. Regarding difficulty, we extracted five different examples from each algorithmic problem in LogicPro to create LogicPro-eval (11800) and conducted evaluations on several models (the evaluation method is the same as in Section Met-

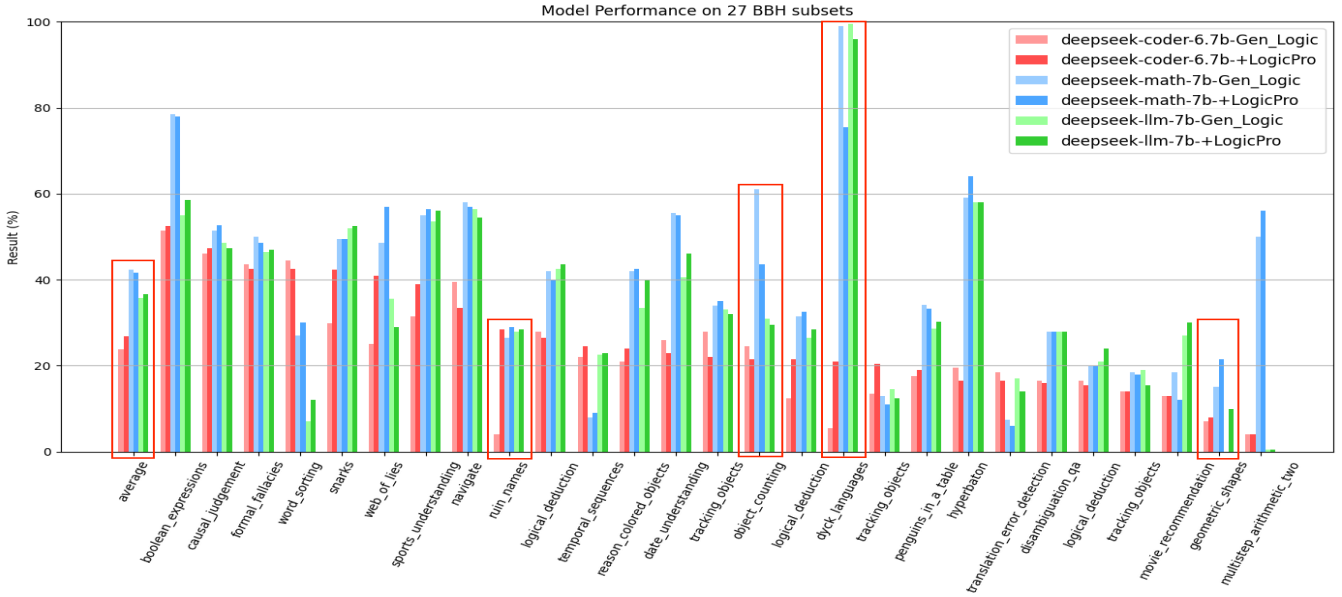


Figure 5: Comparison of results before and after adding LogicPro on different types of base models. (DeepSeek: LLM-Math-Coder). Green, blue, and red represent different models. Light colors indicate baseline results, while dark colors show the results after incorporating LogicPro. Red boxes highlight differences in trends across certain subsets.

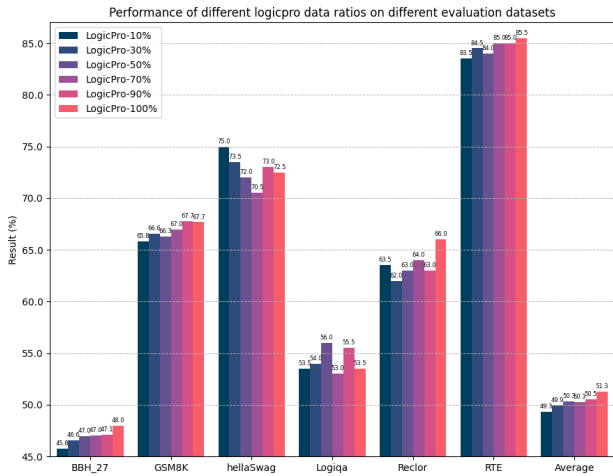


Figure 6: Performance of different logicpro data ratios on different evaluation datasets (Qwen1.5-7B). The performance of BBH 27 subset are shown in Figure 15.

rics). As shown in Table 4, all models performed poorly on this task, with accuracies below 50%. In contrast, according to open source leaderboards (Contributors 2023), these models achieved high scores on common reasoning benchmarks. This indicates that LogicPro is indeed a difficult dataset. When it comes to diversity, finding a unified quantitative measure is challenging. However, unlike datasets constructed from specific rules ($((p \vee q) \wedge \neg p) \vdash q$) or combinations ($((\forall x(p(x) \vee q(x)) \wedge \neg p(a)) \vdash q(a))$) thereof, LogicPro is based on a wide variety of algorithmic problems, suggesting greater diversity in its synthetic origins. Therefore, assessing its diversity warrants further exploration.

Model	LogicPro-Eval
GPT-4 (24-0409)	46.3
Qwen1.5-7B-Chat	28.6
Llama-3-8B-Instruct	27.8
Qwen1.5-14B-Chat	27.6
Llama3-70B-Instruct	37.6
Qwen1.5-72B-Chat	31.6

Table 4: Results on LogicPro-Eval.

Conclusion and Future Work

In this paper, we present LogicPro, which enhances complex logical reasoning through code cases. With this approach, we can construct a dataset that is challenging, diverse, and scalable, while also ensuring a high-quality reasoning process guided by the values of intermediate variables. The experimental results demonstrate that our method significantly enhances performance across multiple models on the datasets: BBH²⁷, GSM8K, HellSwag, Logicqa, Reclor, and RTE.

For future work, considering the vast number of algorithm problems in the real world (far more than the 2630 on LeetCode), such as those from Luogu, ACM competitions, and various online judges, we plan to gather more algorithm problem data on a larger scale and launch a more extensive LogicPro X, which can further enhance the model’s reasoning capabilities. Moreover, our data naturally includes both outcome signals (code execution results) and process signals (intermediate variables during code execution), which might offer new insights for process supervision learning.

References

- Achiam, J.; Adler, S.; Agarwal, S.; Ahmad, L.; Akkaya, I.; Aleman, F. L.; Almeida, D.; Altenschmidt, J.; Altman, S.; Anadkat, S.; et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Besta, M.; Blach, N.; Kubicek, A.; Gerstenberger, R.; Podstawski, M.; Gianinazzi, L.; Gajda, J.; Lehmann, T.; Niewiadomski, H.; Nyczyk, P.; et al. 2024. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 17682–17690.
- Bubeck, S.; Chandrasekaran, V.; Eldan, R.; Gehrke, J.; Horvitz, E.; Kamar, E.; Lee, P.; Lee, Y. T.; Li, Y.; Lundberg, S.; et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.
- Chen, W.; Ma, X.; Wang, X.; and Cohen, W. W. 2022. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *arXiv preprint arXiv:2211.12588*.
- Clark, P.; Tafjord, O.; and Richardson, K. 2021. Transformers as soft reasoners over language. In *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, 3882–3890.
- Cobbe, K.; Kosaraju, V.; Bavarian, M.; Chen, M.; Jun, H.; Kaiser, L.; Plappert, M.; Tworek, J.; Hilton, J.; Nakano, R.; et al. 2021. Training Verifiers to Solve Math Word Problems. *arXiv preprint arXiv:2110.14168*.
- Contributors, O. 2023. OpenCompass: A Universal Evaluation Platform for Foundation Models. <https://github.com/open-compass/opencompass>.
- Dalvi, B.; Jansen, P.; Tafjord, O.; Xie, Z.; Smith, H.; Pipatanangkura, L.; and Clark, P. 2021. Explaining Answers with Entailment Trees. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 7358–7370.
- Dubey, A.; Jauhri, A.; Pandey, A.; Kadian, A.; Al-Dahle, A.; Letman, A.; Mathur, A.; Schelten, A.; Yang, A.; Fan, A.; et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Gou, Z.; Shao, Z.; Gong, Y.; Yang, Y.; Huang, M.; Duan, N.; Chen, W.; et al. 2023. Tora: A tool-integrated reasoning agent for mathematical problem solving. *arXiv preprint arXiv:2309.17452*.
- Guo, D.; Zhu, Q.; Yang, D.; Xie, Z.; Dong, K.; Zhang, W.; Chen, G.; Bi, X.; Wu, Y.; Li, Y.; et al. 2024. DeepSeek-Coder: When the Large Language Model Meets Programming—The Rise of Code Intelligence. *arXiv preprint arXiv:2401.14196*.
- Huang, Y.; Liu, X.; Gong, Y.; Gou, Z.; Shen, Y.; Duan, N.; and Chen, W. 2024. Key-point-driven data synthesis with its enhancement on mathematical reasoning. *arXiv preprint arXiv:2403.02333*.
- Le, H.; Wang, Y.; Gotmare, A. D.; Savarese, S.; and Hoi, S. C. H. 2022. CodeRL: Mastering Code Generation through Pretrained Models and Deep Reinforcement Learning. *arXiv:2207.01780*.
- Lightman, H.; Kosaraju, V.; Burda, Y.; Edwards, H.; Baker, B.; Lee, T.; Leike, J.; Schulman, J.; Sutskever, I.; and Cobbe, K. 2023. Let’s Verify Step by Step. *arXiv preprint arXiv:2305.20050*.
- Liu, H.; Liu, J.; Cui, L.; Teng, Z.; Duan, N.; Zhou, M.; and Zhang, Y. 2023a. Logiqa 2.0—an improved dataset for logical reasoning in natural language understanding. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*.
- Liu, H.; Ning, R.; Teng, Z.; Liu, J.; Zhou, Q.; and Zhang, Y. 2023b. Evaluating the logical reasoning ability of chatgpt and gpt-4. *arXiv preprint arXiv:2304.03439*.
- Liu, H.; Teng, Z.; Cui, L.; Zhang, C.; Zhou, Q.; and Zhang, Y. 2023c. LogiCoT: Logical Chain-of-Thought Instruction Tuning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, 2908–2921.
- liu, H.; Teng, Z.; Ning, R.; Liu, J.; Zhou, Q.; and Zhang, Y. 2023. GLoRE: Evaluating Logical Reasoning of Large Language Models. *arXiv:2310.09107*.
- Lu, Z.; Zhou, A.; Ren, H.; Wang, K.; Shi, W.; Pan, J.; Zhan, M.; and Li, H. 2024. Mathgenie: Generating synthetic data with question back-translation for enhancing mathematical reasoning of llms. *arXiv preprint arXiv:2402.16352*.
- Luo, L.; Liu, Y.; Liu, R.; Phatale, S.; Lara, H.; Li, Y.; Shu, L.; Zhu, Y.; Meng, L.; Sun, J.; et al. 2024. Improve Mathematical Reasoning in Language Models by Automated Process Supervision. *arXiv e-prints*, arXiv–2406.
- Miao, S.-Y.; Liang, C.-C.; and Su, K.-Y. 2020. A Diverse Corpus for Evaluating and Developing English Math Word Problem Solvers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 975–984.
- Mill, J. S. 2013. *A system of Logic, Ratiocinative and Inductive: Being a Connected View of the Principles of Evidence, and the Methods of Scientific Investigation*. Harper and Brothers, Publishers.
- Nie, Y.; Williams, A.; Dinan, E.; Bansal, M.; Weston, J.; and Kiela, D. 2020. Adversarial NLI: A New Benchmark for Natural Language Understanding. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 4885–4901.
- Paster, K.; Santos, M. D.; Azerbayev, Z.; and Ba, J. 2023. OpenWebMath: An Open Dataset of High-Quality Mathematical Web Text. *arXiv:2310.06786*.
- Patel, A.; Bhattamishra, S.; and Goyal, N. 2021. Are NLP Models really able to Solve Simple Math Word Problems? In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2080–2094.
- Saeed, M.; Ahmadi, N.; Nakov, P.; and Papotti, P. 2021. RuleBERT: Teaching Soft Rules to Pre-Trained Language Models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 1460–1476.
- Saparov, A.; and He, H. 2023. Language models can (kind of) reason: A systematic formal analysis of chain-of-thought. In *International Conference on Learning Representations*.

- Schaeffer, R.; Miranda, B.; and Koyejo, S. 2024. Are emergent abilities of large language models a mirage? *Advances in Neural Information Processing Systems*, 36.
- Shao, Z.; Wang, P.; Zhu, Q.; Xu, R.; Song, J.; Zhang, M.; Li, Y.; Wu, Y.; and Guo, D. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.
- Srivastava, A.; Rastogi, A.; Rao, A.; Shoeb, A. A. M.; Abid, A.; Fisch, A.; Brown, A. R.; Santoro, A.; Gupta, A.; Garriga-Alonso, A.; et al. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint arXiv:2206.04615*.
- Suzgun, M.; Scales, N.; Schärli, N.; Gehrmann, S.; Tay, Y.; Chung, H. W.; Chowdhery, A.; Le, Q.; Chi, E.; Zhou, D.; et al. 2023. Challenging BIG-Bench Tasks and Whether Chain-of-Thought Can Solve Them. In *Findings of the Association for Computational Linguistics: ACL 2023*, 13003–13051.
- Taylor, R.; Kardas, M.; Cucurull, G.; Scialom, T.; Hartshorn, A.; Saravia, E.; Poulton, A.; Kerkez, V.; and Stojnic, R. 2022. Galactica: A Large Language Model for Science. *arXiv:2211.09085*.
- Teknum. 2023. OpenHermes 2.5: An Open Dataset of Synthetic Data for Generalist LLM Assistants.
- Wang, A.; Pruksachatkun, Y.; Nangia, N.; Singh, A.; Michael, J.; Hill, F.; Levy, O.; and Bowman, S. R. 2020. SuperGLUE: A Stickier Benchmark for General-Purpose Language Understanding Systems. *arXiv:1905.00537*.
- Wang, P.; Li, L.; Shao, Z.; Xu, R.; Dai, D.; Li, Y.; Chen, D.; Wu, Y.; and Sui, Z. 2023. Math-shepherd: A label-free step-by-step verifier for llms in mathematical reasoning. *arXiv preprint arXiv:2312.08935*.
- Wang, S.; Liu, Z.; Zhong, W.; Zhou, M.; Wei, Z.; Chen, Z.; and Duan, N. 2022. From lsat: The progress and challenges of complex reasoning. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*.
- Wang, S.; Wei, Z.; Choi, Y.; and Ren, X. 2024. Can LLMs Reason with Rules? Logic Scaffolding for Stress-Testing and Improving LLMs. *arXiv preprint arXiv:2402.11442*.
- Wei, J.; Wang, X.; Schuurmans, D.; Bosma, M.; Xia, F.; Chi, E.; Le, Q. V.; Zhou, D.; et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35: 24824–24837.
- Xin, H.; Guo, D.; Shao, Z.; Ren, Z.; Zhu, Q.; Liu, B.; Ruan, C.; Li, W.; and Liang, X. 2024. DeepSeek-Prover: Advancing Theorem Proving in LLMs through Large-Scale Synthetic Data. *arXiv:2405.14333*.
- Yao, S.; Yu, D.; Zhao, J.; Shafran, I.; Griffiths, T.; Cao, Y.; and Narasimhan, K. 2024. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36.
- Yao, S.; Zhao, J.; Yu, D.; Du, N.; Shafran, I.; Narasimhan, K.; and Cao, Y. 2022. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*.
- Yu, W.; Jiang, Z.; Dong, Y.; and Feng, J. 2020. Reclor: A reading comprehension dataset requiring logical reasoning. *arXiv preprint arXiv:2002.04326*.
- Yuan, L.; Cui, G.; Wang, H.; Ding, N.; Wang, X.; Deng, J.; Shan, B.; Chen, H.; Xie, R.; Lin, Y.; Liu, Z.; Zhou, B.; Peng, H.; Liu, Z.; and Sun, M. 2024. Advancing LLM Reasoning Generalists with Preference Trees. *arXiv:2404.02078*.
- Yue, X.; Qu, X.; Zhang, G.; Fu, Y.; Huang, W.; Sun, H.; Su, Y.; and Chen, W. 2023. Mammoth: Building math generalist models through hybrid instruction tuning. *arXiv preprint arXiv:2309.05653*.
- Zellers, R.; Holtzman, A.; Bisk, Y.; Farhadi, A.; and Choi, Y. 2019. HellaSwag: Can a Machine Really Finish Your Sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 4791–4800.
- Zhang, Y.; Yang, J.; Yuan, Y.; and Yao, A. C.-C. 2023. Cumulative reasoning with large language models. *arXiv preprint arXiv:2308.04371*.
- Zhong, W.; Wang, S.; Tang, D.; Xu, Z.; Guo, D.; Wang, J.; Yin, J.; Zhou, M.; and Duan, N. 2021. AR-LSAT: Investigating Analytical Reasoning of Text. *arXiv:2104.06598*.
- Zhou, K.; Zhang, B.; Wang, J.; Chen, Z.; Zhao, W. X.; Sha, J.; Sheng, Z.; Wang, S.; and Wen, J.-R. 2024. JiuZhang3.0: Efficiently Improving Mathematical Reasoning by Training Small Data Synthesis Models. *arXiv preprint arXiv:2405.14365*.

Appendix

Related work

Reasoning of LLMs Reasoning, which serves as a fundamental ability of LLMs, determines the strength to solve complex real-world problems. Enhancing the reasoning ability of LLM can be divided mainly into two ways.

Improve Reasoning of LLMs by Prompting. The reasoning ability of LLMs can be significantly stimulated by giving them different prompts, such as Chain-of-Thought (Wei et al. 2022), Plan-and-Solve (Wang et al. 2023), etc. It is also possible to assist the model in reasoning by providing it with some external tools (Yao et al. 2022; Chen et al. 2022; Gou et al. 2023). These methods do not require parameter modification on LLMs, but do some control during LLM’s reasoning to get a more reliable reasoning process and a better final result.

Improve Reasoning of LLMs by Training. Continuing pre-training provides a means to enhance the internal reasoning ability of LLMs from a knowledge perspective (Taylor et al. 2022; Paster et al. 2023). The ability of reasoning could be further enhanced by fine-tuning with instruction pairs related reasoning (Yue et al. 2023; Yuan et al. 2024). Reinforcement learning with two types of reward model: Outcome Reward Model (ORM) (Le et al. 2022; Shao et al. 2024) and Process Reward Model (PRM) (Lightman et al. 2023; Wang et al. 2023), have also been used to improve the model’s reasoning accuracy at various granularity. In addition, synthesised data from LLMs (Xin et al. 2024; Zhou et al. 2024) demonstrates the possibility of improving reasoning of LLMs themselves.

Logic Reasoning of LLMs Logical reasoning epitomizes the art of deducing new insights from existing knowledge by adhering to specific principles and laws. This process does not require a robust knowledge base. Instead, it emphasizes the precision and meticulousness with which conclusions are inferred from one piece of information to another.

Training Data of Logic Reasoning. There are various open-source available datasets for different types of logical reasoning tasks. LogiQA2.0 (Liu et al. 2023a) is a complex logical reasoning dataset built from Chinese Civil Service Exam questions. ReClor (Yu et al. 2020), a data set built on standardized graduate admission examinations, contains reading comprehension tasks that require logical reasoning. ULogic (Wang et al. 2024) is logical reasoning dataset constructed from diverse inferential rules, which could improve various commonsense reasoning tasks.

Evaluation of Logic Reasoning. LogiEval (Liu et al. 2023b) and GLoRE (Liu et al. 2023) combines several logical reasoning datasets, evaluating the logical reasoning of LLMs from multiple dimensions. Big-Bench Hard (Suzgun et al. 2023; Srivastava et al. 2022) is a diverse evaluation set that incorporates logical reasoning tasks such as logical deduction and logical fallacy detection.

Approach

Complete Prompts

Step 1 Prompt The complete prompt for the first step to construct test sample inputs is shown in Figure 8.

Step 2 Prompt The complete prompt for the second step in building a logical reasoning problem is shown in Figure 9.

Step 3 Prompt The full prompt for the 3-1 step in constructing the test sample code is illustrated in Figure 9. The complete prompt for the 3-2 step in rewriting the code to print intermediate variables is shown in Figure 11.

Step 4 Prompt The complete prompt for the fourth step in constructing the final answer is shown in Figure 12. The prompt for GPT-4 is shown in Figures 13 and 14.

Data Comparison

As shown in Table. 5, we compare four types of data with LogicPro in terms of three dimensions: data size, data source and reasoning rule complexity. The results show that our method performs well in terms of data size, reasoning rule complexity and scalability.

Result

The complete data results of LogicPro on different models are in Table 6.

Additional Analysis

LogicPro under different types base models

Figure 15: Comparison of LogicPro results on different types of base models (deepseek-Yi-Qwen).

Different Proportions of LogicPro Figure 16: Results of Different Proportions of LogicPro-Train Data on BBH 27 subsets.

Examples

Cases of BBH Evaluation Figure 17, 18, 4, 19, 20 and 21: cases of BBH Evaluation.

Cases of Step1: Construct Test Sample Inputs Figure 22, 23, 24 and 25: cases of constructing Test Sample Inputs.

Cases of LogicPro We open source LogicPro data and show a small number of samples here (Figure 26 and 27).

Dataset	Size	Synthetic	Text Source
Realistic Logical Reasoning			
LOGIQA	8,678	not	China Civil Service Exam
RECLOR	6,138	not	GMAT and LSAT
FOLIO	1,435	not	First-order logic
E-KAR	1155	-	Analogical Reasoning
Synthetic Logical Reasoning			
ProofWriter	20,192	yes	Entailment Tree
DEER	1200	not	Inductive reasoning
RuleTaker	27363	yes	soft rule
ProntoQA	200	yes	Deductive
Mathematical Reasoning			
GSM8K	8,792	not	Multi-step math reasoning
ASDiv	2,305	not	Multi-step math reasoning
SVAMP	1,000	not	Multi-step math reasoning
Commonsense Reasoning			
CommonsenseQA	12,247	-	ConceptNet
OpenBookQA	5,957	-	Open-book knowledges
LogicPro (our)	81,026	yes	complex (Logic from Code)

Table 5: Comparison of four types of datasets and LogicPro. For synthetic data, where the reasoning logic is singular and requires no domain knowledge, it is simple reasoning. For knowledge-based reasoning data, which requires domain knowledge and some ability to generalize reasoning to the domain, but has simpler reasoning rules, it is a medium reasoning problem. For real human test questions (e.g., those in LSAT), which require domain knowledge and a lot of logical reasoning to complete, they are complex reasoning problems.

FINAL ANSWER

INPUT:

TITLE: STAIRCASE CLIMBING POSSIBILITIES

Q: IMAGINE YOU ARE CLIMBING A STAIRCASE WITH 17 STEPS TO REACH THE TOP. YOU CAN CLIMB EITHER 1 OR 2 STEPS AT A TIME. HOW MANY DISTINCT WAYS CAN YOU CLIMB TO THE TOP OF THE STAIRCASE?

TARGET:

TO FIND THE NUMBER OF DISTINCT WAYS TO CLIMB A STAIRCASE WITH 17 STEPS USING 1-STEP AND 2-STEP CLIMBS, YOU CAN USE A DYNAMIC PROGRAMMING APPROACH OR RECOGNIZE THAT THIS IS A FIBONACCI-LIKE SEQUENCE PROBLEM.

DEFINE $f(N)$ AS THE NUMBER OF WAYS TO REACH THE N -TH STEP. THE RECURRENCE RELATION IS:

$$f(N) = f(N - 1) + f(N - 2)$$

THIS RELATION HOLDS BECAUSE TO REACH THE N -TH STEP, YOU CAN EITHER COME FROM THE $(N-1)$ -TH STEP (WITH A 1-STEP CLIMB) OR FROM THE $(N-2)$ -TH STEP (WITH A 2-STEP CLIMB).

BASE CASES:

- $f(1) = 1$ (ONLY ONE WAY TO STAND ON THE FIRST STEP)
- $f(2) = 2$ (EITHER TWO 1-STEPS OR ONE 2-STEP)

USING THESE, CALCULATE UP TO $f(17)$:

```
\begin{align*}
f(3) &= f(2) + f(1) = 2 + 1 = 3 \\
f(4) &= f(3) + f(2) = 3 + 2 = 5 \\
f(5) &= f(4) + f(3) = 5 + 3 = 8 \\
f(6) &= f(5) + f(4) = 8 + 5 = 13 \\
f(7) &= f(6) + f(5) = 13 + 8 = 21 \\
f(8) &= f(7) + f(6) = 21 + 13 = 34 \\
f(9) &= f(8) + f(7) = 34 + 21 = 55 \\
f(10) &= f(9) + f(8) = 55 + 34 = 89 \\
f(11) &= f(10) + f(9) = 89 + 55 = 144 \\
f(12) &= f(11) + f(10) = 144 + 89 = 233 \\
f(13) &= f(12) + f(11) = 233 + 144 = 377 \\
f(14) &= f(13) + f(12) = 377 + 233 = 610 \\
f(15) &= f(14) + f(13) = 610 + 377 = 987 \\
f(16) &= f(15) + f(14) = 987 + 610 = 1597 \\
f(17) &= f(16) + f(15) = 1597 + 987 = 2584 \\
\end{align*}
```

THUS, THERE ARE **2584** DISTINCT WAYS TO CLIMB TO THE TOP OF THE STAIRCASE WITH 17 STEPS.



Figure 7: Final Answer

Step 1: Construct Test Sample Inputs

I have an algorithmic problem and its python code, please help me construct thirty different test sample inputs.

1. The constructed test sample inputs need to fulfill the requirements of the algorithmic problem and be compatible with the provided Python code.

2. Please enclose the constructed test sample inputs in the following python format; please enclose each test sample input individually.

```
```python
```

```
Test sample input 1
```

```
Your input here
```

```
```
```

```
```python
```

```
Test sample input 2
```

```
Your input here
```

```
```
```

```
...
```

```
```python
```

```
Test sample input 30
```

```
Your input here
```

```
```
```

3. Ensure that all test samples are unique and as diverse as possible based on the topic and Python code.

4. Consider various aspects of the input type to ensure diversity, such as:

- Range of values: Include small, medium, and large values, as well as edge cases.

- Special cases: Consider cases like empty input, maximum allowed input size, or inputs that might cause edge conditions.

- Pattern variations: If the input is a sequence, vary the sequence patterns (e.g., sorted, reverse-sorted, random order).

- Combining elements: If the input is a composite data structure (e.g., array of strings), combine different types of elements.

5. Generate inputs with varying difficulty levels (low, medium, high) considering the problem statement and the provided Python code:

- Low difficulty: Simple and straightforward inputs that cover basic scenarios.

- Medium difficulty: Moderately complex inputs that include more diverse and realistic scenarios.

- High difficulty: Complex inputs that test edge cases and challenging conditions.

6. Ensure that all test samples adhere to the constraints provided in the problem description.

7. Provide only the input for the test samples, do not include the output.

Algorithmic Questions Title:

```
{algorithmic_problems}
```

python solution:

```
{python_solution}
```

Figure 8: Step 1: Constructing Test Sample Inputs

Step 2: Construct Logical Reasoning Problem

I have an algorithmic question and a corresponding test input; please rewrite the algorithmic question as a text-only logical reasoning question based on the test input.

Instructions:

1. Please incorporate the test input into the description of the algorithm question;
2. Please first give the name of this logical reasoning task; then give the question that contains the test input.

Reference case I:

- algorithmic question: Given a sequence containing only (,), {{, }}, [,], <, >, complete the rest of the sequence, making sure that all the parentheses are properly closed and in the right order.

- test input: "<> (([[({ { })] [<>]]]"

- text-only logical reasoning question :

Title: Correctly close a Dyck-n word.

Q: Complete the rest of the sequence, making sure that the parentheses are closed properly. Input: <>(([[({ { })] [<>]]]

Reference case II:

- algorithmic question: You are given an integer array 'cards' of length '4'. You have four cards, each containing a number in the range '[1, 9]'. You should arrange the numbers on these cards in a mathematical expression using the operators '[+, '-', '*', /]' and the parentheses '(' and ')' to get the value 24. You are restricted with the following rules: * The division operator '/' represents real division, not integer division.

- test input: "[4, 1, 8, 7]"

- text-only logical reasoning question :

Title: Achieve the Target Value

Q: You are presented with four cards, each bearing a number within the range of 1 to 9. Using the numbers on these cards, form a mathematical expression by arranging them with the operators '+', '-', '*', and '/', as well as parentheses '(' and ')', such that the resulting value of the expression is 24. Note the following rules:

- Division operator '/' represents real division, not integer division.

- Each operation must be performed between two numbers (no unary operations).

- Numbers cannot be concatenated to form multi-digit numbers.

Given the cards with numbers [4, 1, 8, 7], determine if it is possible to form an expression that evaluates to 24.

Can you find such an expression, or prove that it cannot be done?

Refer to the above example of rewriting an algorithmic question into a text-only logical reasoning question based on test input:

- algorithmic question: {algorithmic_question}

- test input: {test_sample_input}

- text-only logical reasoning question:

Figure 9: Step 2: Constructing Logical Reasoning Problems

Step 3: Construct Text Example Code

I have a piece of Python code and a test case input. Please provide the modified code that can directly run this test sample based on the original Python code.

- Please ensure that the generated code can be executed directly.
- Please ensure that after running the code, the output result of the algorithm is returned through the variable 'result'.

Test case input:
{test_sample_input}

python code:
{python_solution}

Figure 10: Step 3-1: Constructing Test Sample Code

Step 4: Obtain Intermediate Variable From Program

Please modify the following code so that it prints out important variables and their detailed descriptions related to the algorithm at appropriate places.

1. Important variables refer to those critical for understanding the algorithm's logic, such as loop counters, function inputs and outputs, key condition judgments, and variables indicating state changes.
2. Ensure that the printed information includes not only the names of the variables but also their roles and meanings within the algorithm, to better understand the execution process of the code.
3. Ensure that the printed information is closely related to the algorithm logic and does not include irrelevant content (such as code errors and exceptions).
4. Ensure that the printed information is detailed enough.

python code:
{test_example_code}

Figure 11: Step 3-2: Rewriting the Code to Print Intermediate Variables

Step 4: Program-Guided Reasoning

There is a logical reasoning question and the intermediate variable output of its code solution. Please answer this logical reasoning question based on the intermediate variable output of the code.

Instructions:

1. Refer to the code's intermediate variable outputs. Use the information provided to help you answer the logical reasoning questions.
2. First, outline your approach to solving the logical reasoning task. Then, provide the exact reasoning process step by step.
3. Do not use code to solve this logical reasoning problem. Instead, use the provided intermediate variable outputs to guide your answer.
4. Do not mention "intermediate variables" in your answer. Focus on solving the logical reasoning question directly.
5. Avoid phrases like "From the intermediate variables" in the answer. Just use them (intermediate variables) to help you answer the logical reasoning question.
6. Please ensure that the logical reasoning questions are answered in a way that reflects each step of the reasoning process in the code.

Reference case:

- Logical reasoning question:

{case_input_question_ref}

- Code intermediate variables:

...

{code_print_ref}

...

- Logical Reasoning Question Answer:

...

{answer_with_print_ref}

...

Refer to the above case to give a solution to a logical reasoning question:

- Logical reasoning question:

{case_input_question}

- Code intermediate variables:

...

{code_print}

...

- Logical Reasoning Question Answer:

Figure 12: Step 4: Constructing the Final Answer

Step 4: Construct The Final Answer For gpt-4

There is a logical reasoning question and the intermediate variable output of its code solution. Please answer this logical reasoning question based on the intermediate variable output of the code.

Instructions:

1. Refer to the code's intermediate variable outputs. Use the information provided to help you answer the logical reasoning questions.
2. First, outline your approach to solving the logical reasoning task. Then, provide the exact reasoning process step by step.
3. Do not use code to solve this logical reasoning problem. Instead, use the provided intermediate variable outputs to guide your answer.
4. Do not mention "intermediate variables" in your answer. Focus on solving the logical reasoning question directly.
5. Avoid phrases like "From the intermediate variables" in the answer. Just use them (intermediate variables) to help you answer the logical reasoning question.
6. Please ensure that the logical reasoning questions are answered in a way that reflects each step of the reasoning process in the code.

Reference case:

- Logical reasoning question:

Title: Balance a Parentheses String

Q: A parentheses string is considered balanced if every left parenthesis '(' has a corresponding two consecutive right parentheses '))' and appears before them. For instance, "(0)", "(0)(0))", and "(0)(0))" are balanced, while ")0", "(0)", and "(0)" are not. You are given a string `s = "))(0)(("` containing only '(' and ')'. Can you determine the minimum number of '(' and ')' characters that need to be inserted at any position in the string to make it balanced?

- Code intermediate variables:

...

Testing with input string: `))(0)((`

Input string: `s = "))(0)(("`

Variable 'left' initialized to 0, representing the number of open parentheses encountered so far.

Variable 'res' initialized to 0, representing the minimum number of insertions required.

Variable 'n' set to the length of the input string, which is 8

Loop counter 'i' initialized to 0, iterating through the input string.

Current iteration: `i = 0`

Encountered a close parenthesis without a matching open parenthesis, incrementing 'res' to 1

Current iteration: `i = 1`

Encountered a close parenthesis without a matching open parenthesis, incrementing 'res' to 2

Reached the end of the string or encountered a close parenthesis without a matching open parenthesis, incrementing 'res' to 3

Current iteration: `i = 3`

Encountered an open parenthesis, incrementing 'left' to 1

Current iteration: `i = 4`

Encountered a close parenthesis with a matching open parenthesis, decrementing 'left' to 0

Current iteration: `i = 5`

Encountered a close parenthesis without a matching open parenthesis, incrementing 'res' to 4

Reached the end of the string or encountered a close parenthesis without a matching open parenthesis, incrementing 'res' to 5

Current iteration: `i = 7`

Encountered an open parenthesis, incrementing 'left' to 1

Final result: `res = 5` plus 1 unmatched open parentheses, totaling 7

Result: 7

...

- Logical Reasoning Question Answer:

To determine the minimum number of '(' and ')' characters that need to be inserted at any position in the string `))(0)(("` to make it balanced, we can follow these steps:

1. **Initialization**:

- Initialize two variables: 'left' to track the number of unmatched '(' characters and 'res' to track the number of insertions required.

- Set 'left' to 0 and 'res' to 0.

Figure 13: Step 4: Constructing the Final Answer for gpt4-1

Step 4: Construct The Final Answer For gpt-4

```
...
2. Iterate through the String:
- Loop through each character in the string.
- For each character:
- If it is a '(', increment `left` by 1.
- If it is a ')', check if there is an unmatched '(' (i.e., `left > 0`):
- If there is, decrement `left` by 1.
- If there isn't, increment `res` by 1 because an insertion of '(' is needed before this ')'.

3. Final Adjustments:
- After the loop, any remaining unmatched '(' (i.e., `left > 0`) need to be balanced by inserting ')' for each unmatched '('.
- Therefore, add `2 * left` to `res`.

Step-by-Step Reasoning:

1. Initialization:
- Input string: `s = ")()())("`
- `left = 0`
- `res = 0`
- Length of the string, `n = 8`

2. Iteration:
- Iteration 1 (i = 0):
- Encountered ')', no matching '(', increment `res` to 1.
- Iteration 2 (i = 1):
- Encountered ')', no matching '(', increment `res` to 2.
- Iteration 3 (i = 2):
- Encountered '(', increment `left` to 1.
- Iteration 4 (i = 3):
- Encountered '(', increment `left` to 2.
- Iteration 5 (i = 4):
- Encountered ')', matching '(', decrement `left` to 1.
- Iteration 6 (i = 5):
- Encountered ')', matching '(', decrement `left` to 0.
- Iteration 7 (i = 6):
- Encountered '(', increment `left` to 1.
- Iteration 8 (i = 7):
- Encountered '(', increment `left` to 2.

3. Final Adjustments:
- After the loop, `left = 2` (unmatched '(').
- Each unmatched '(' needs ')' to balance: `2 * left = 4`.
- Total insertions required: `res + 2 * left = 2 + 4 = 6`.

Therefore, the minimum number of '(' and ')' characters that need to be inserted to make the string balanced is 6.
Refer to the above case to give a solution to a logical reasoning question:
- Logical reasoning question:
{case_input_question}
- Code intermediate variables:
...
{code_print}
...
- Logical Reasoning Question Answer:
```

Figure 14: Step 4: Constructing the Final Answer for gpt4-2

| Base Model | SFT Data | BBH ²⁷ | GSM8K | HellaSwag | Logiqa | Reclor | RTE | Average |
|-------------|----------------------------------|-------------------|-------|-----------|--------|--------|------|-------------|
| GPT-4 | - | 74.9 | 94.2 | 83.0 | 69.0 | 88.5 | 77.5 | 76.1 |
| ChatGPT | - | 50.3 | 65.3 | 73.0 | 44.5 | 61.0 | 90.0 | 52.8 |
| Qwen1.5-7B | General-100k | 45.0 | 65.5 | 69.5.0 | 50.0 | 62.5 | 82.5 | 48.3 |
| | Gen-Logic | 44.8 | 65.2 | 71.0 | 54.0 | 63.5 | 84.0 | 48.3 |
| | + Code | 45.6 | 65.6 | 74.5 | 55.0 | 61.5 | 84.0 | 49.1 |
| | + LogicPro_{ours} | 47.9 | 66.2 | 71.0 | 55.5 | 64.0 | 84.0 | 51.1 |
| Deepseek-7B | General-100k | 35.7 | 21.5 | 33.5 | 39.5 | 45.0 | 53.0 | 36.2 |
| | Gen-Logic | 35.7 | 23.0 | 36.3 | 40.0 | 45.0 | 54.0 | 36.3 |
| | + Code | 35.9 | 22.9 | 36.3 | 38.5 | 43.0 | 54.5 | 36.4 |
| | + LogicPro_{ours} | 38.0 | 27.3 | 39.3 | 44.0 | 48.5 | 52.5 | 38.7 |
| Llama-2-7B | General-100k | 34.9 | 23.4 | 29.5 | 40.0 | 40.0 | 52.5 | 35.2 |
| | Gen-Logic | 35.0 | 24.6 | 40.0 | 37.5 | 45.0 | 52.5 | 35.8 |
| | + Code | 34.0 | 25.2 | 35.0 | 38.5 | 44.5 | 52.5 | 34.8 |
| | + LogicPro_{ours} | 35.8 | 25.7 | 35.0 | 40.0 | 49.5 | 53.0 | 36.6 |
| Llama-3-8B | General-100k | 43.3 | 63.8 | 38.8 | 57.5 | 60.0 | 76.0 | 45.9 |
| | Gen-Logic | 44.7 | 63.3 | 41.5 | 59.0 | 64.0 | 77.5 | 47.2 |
| | + Code | 45.2 | 63.7 | 41.3 | 57.0 | 67.0 | 77.5 | 47.7 |
| | + LogicPro_{ours} | 49.3 | 65.7 | 44.3 | 59.0 | 67.5 | 78.0 | 51.4 |
| Yi-1.5-9B | General-100k | 43.3 | 63.8 | 38.8 | 57.5 | 60.0 | 76.0 | 45.9 |
| | Gen-Logic | 52.2 | 71.0 | 52.0 | 60.5 | 69.5 | 84.0 | 54.5 |
| | + Code | 52.4 | 71.6 | 51.5 | 61.5 | 68.0 | 85.5 | 54.8 |
| | + LogicPro_{ours} | 53.8 | 74.8 | 50.0 | 62.0 | 70.0 | 85.0 | 56.1 |
| Llama-2-13B | General-100k | 40.9 | 34.2 | 35.5 | 44.0 | 56.0 | 75.5 | 42.2 |
| | Gen-Logic | 38.8 | 35.4 | 38.3 | 43.0 | 59.5 | 77.0 | 40.6 |
| | + Code | 38.8 | 35.9 | 35.8 | 45.0 | 57.5 | 76.5 | 40.6 |
| | + LogicPro_{ours} | 40.2 | 44.0 | 37.3 | 46.5 | 60.0 | 81.0 | 42.4 |
| Qwen1.5-14B | General-100k | 50.8 | 69.8 | 47.0 | 62.0 | 68.0 | 84.0 | 53.2 |
| | Gen-Logic | 51.2 | 71.7 | 63.5 | 65.0 | 69.0 | 85.0 | 54.2 |
| | + Code | 51.4 | 70.3 | 64.0 | 64.0 | 68.0 | 86.0 | 54.4 |
| | + LogicPro_{ours} | 51.8 | 71.8 | 66.0 | 65.0 | 66.5 | 86.0 | 54.8 |

Table 6: The results of LogicPro data on different models. General-100k and Gen-Logic serve as baseline results (Section Train Datasets). + Code refers to training on Gen-Logic combined with 2,360 original algorithm problems. **+ LogicPro_{ours}** is trained on Gen-Logic combined with our LogicPro data.

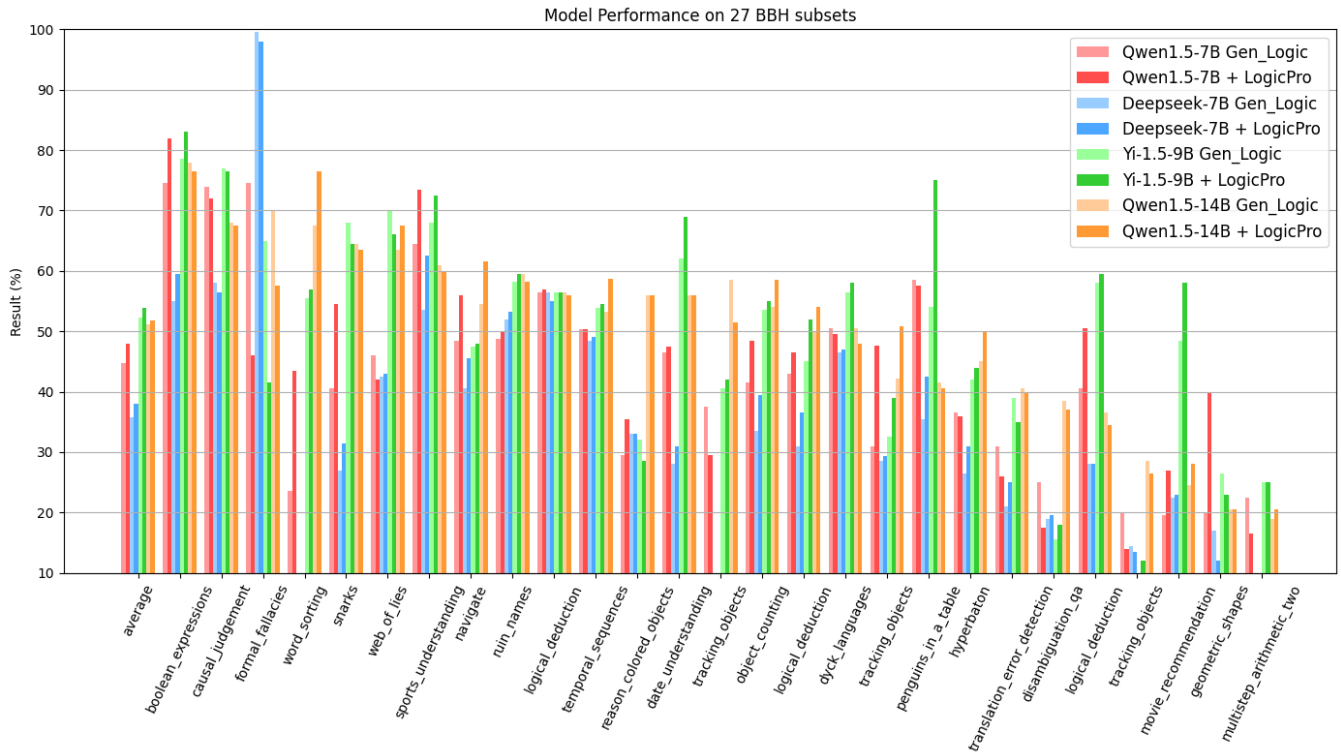


Figure 15: Comparison of LogicPro-Train results on different types of base models (deepseek-Yi-Qwen). The dashed line indicates the result of baseline, and the solid line indicates the results after LogicPro-Train training.

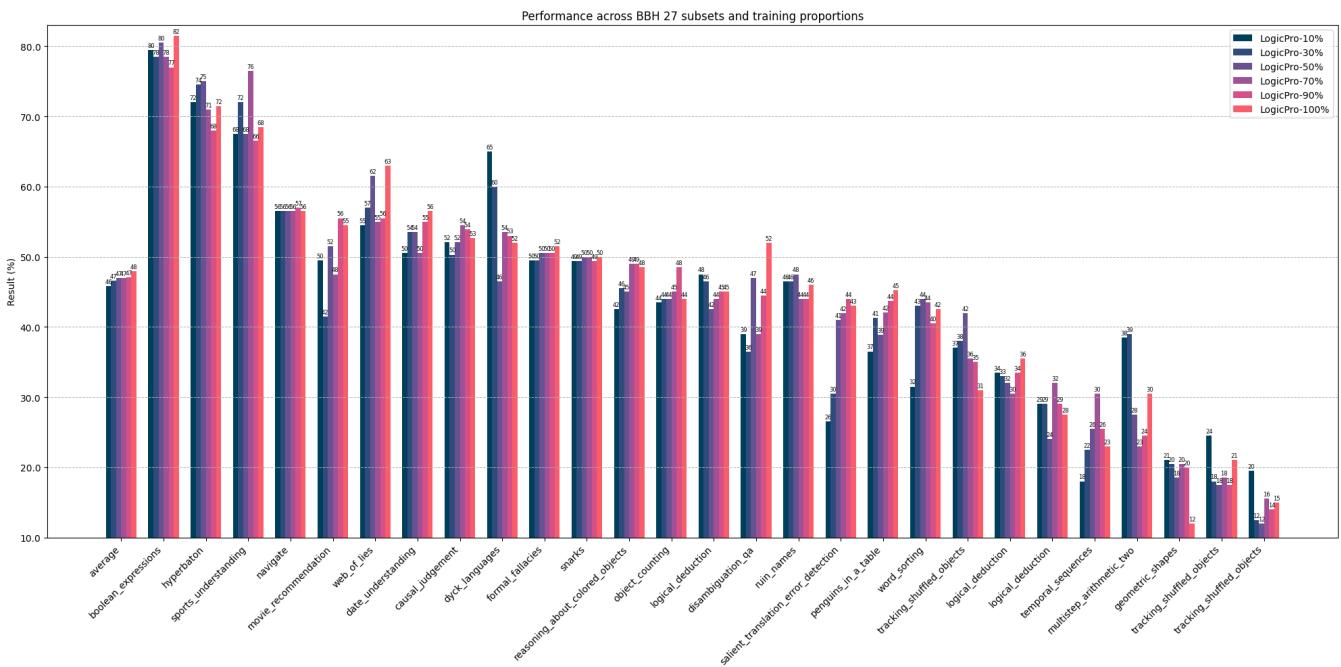



Figure 16: Results of Different Proportions of LogicPro- Train Data on BBH 27 subsets. (Qwen1.5-7B).

Question:

Here is a table where the first line is a header and each subsequent line is a penguin: name, age, height (cm), weight (kg) Louis, 7, 50, 11 Bernard, 5, 80, 13 Vincent, 9, 60, 11 Gwen, 8, 70, 15 For example: the age of Louis is 7, the weight of Gwen is 15 kg, the height of Bernard is 80 cm. We now add a penguin to the table:
James, 12, 90, 12


And here is a similar table, but listing giraffes:
name, age, height (cm), weight (kg)
Jody, 5, 430, 620
Gladys, 10, 420, 590
Marian, 2, 310, 410
Donna, 9, 440, 650
What is the name of the last animal?

Options:
(A) Bernard
(B) Donna
(C) Gladys
(D) Gwen
(E) James

Answer (Qwen1.5-7B w. LogicPro) : 

1. The table of penguins:
- The last penguin added is James, with the attributes: name = James, age = 12, height = 90 cm, weight = 12 kg.

2. The table of giraffes:
- The last giraffe added is Donna, with the attributes: name = Donna, age = 9, height = 440 cm, weight = 650 kg.


Answer (Qwen1.5-7B Gen_Logic) : 

The name of the last animal in the table of penguins is:
(E) James

Figure 17: Success Case for LogicPro: bbh penguins in a table


Question:

Is the following sentence plausible? "Nathan MacKinnon bricked the three pointer in the Eastern Conference Finals."

Answer (Qwen1.5-7B w. LogicPro) : 

1. The sentence mentions Nathan MacKinnon, who is a professional ice hockey player.
2. It also mentions a three-pointer, which is a term used in basketball, not ice hockey.
3. The sentence refers to the "Eastern Conference Finals," which is a term used in basketball, not ice hockey.

Based on these points, the sentence is not plausible because it combines elements from two different sports (ice hockey and basketball) and uses terms that are specific to basketball.

Answer (Qwen1.5-7B Gen_Logic) : 


1. Nathan MacKinnon is a professional ice hockey player.
2. Bricking a three pointer is a term used to describe a missed shot.
3. The Eastern Conference Finals is a stage in the NHL playoffs.

Based on these points, the sentence is is plausible, Nathan MacKinnon is a professional ice hockey player, and it is possible for him to miss a shot during a game. The Eastern Conference Finals is a stage in the NHL playoffs, so it is plausible for him to participate in such a game.

Figure 18: Success Case for LogicPro: bbh sports understanding

Question:

It is 4/19/1969 today. What is the date 24 hours later in MM/DD/YYYY?
Options:
(A) 06/21/1969
(B) 01/18/1969
(C) 04/16/1969
(D) 05/18/1969
(E) 04/20/1969
(F) 05/11/1969

Answer (Qwen1.5-7B w. LogicPro) : 

To find the date 24 hours later from 4/19/1969, we can follow these steps:


1. Add 24 hours to the day of the week (4/19/1969).
2. Check if the resulting day is a valid day of the week.
3. If the day is valid, add the month and year to the result.

Step 1: Add 24 hours to the day of the week (4/19/1969).
- 4/19/1969 is a Thursday.
- Adding 24 hours to Thursday results in Friday, 4/20/1969.

Step 2: Check if the resulting day is a valid day of the week.
- 4/20/1969 is a valid day of the week.

Step 3: Add the month and year to the result.
- The result is 4/20/1969.

Therefore, the date 24 hours later from 4/19/1969 is **(E) 04/20/1969**.

Answer (Qwen1.5-7B Gen_Logic) : 

To find the date 24 hours later, we need to add 24 hours to the given date, 4/19/1969.

First, let's convert the given date to a day of the year (DOY) format. April has 30 days, so the DOY for 4/19/1969 is 19.


Now, let's add 24 hours to the DOY ($19 + 24 = 43$). Since April has 30 days, we need to adjust the DOY to the next month. In this case, it will be May.

So, the date 24 hours later is 5/18/1969. Therefore, the correct answer is:
(D) 05/18/1969

Figure 19: Success Case for LogicPro: bbh date understanding


Question:

Sentence: The chief told the counselor that he took the day off.
Options:
(A) The chief took the day off
(B) The counselor took the day off
(C) Ambiguous

Answer (Qwen1.5-7B w. LogicPro) : 

In the given sentence, the pronoun "he" is used to refer to the person who took the day off. The antecedent of the pronoun "he" is the person who is being referred to as the one who took the day off.

Based on the context, the most likely antecedent for the pronoun "he" is the chief. Therefore, the correct answer is:
(A) The chief took the day off

Answer (Qwen1.5-7B Gen_Logic) : 

The antecedent of the pronoun "he" in the sentence "The chief told the counselor that he took the day off" is **ambiguous**. It could refer to either the chief or the counselor, as both are possible antecedents for the pronoun "he."

Figure 20: Success Case for LogicPro: bbh disambiguation qa

Question:

The following paragraphs each describe a set of seven objects arranged in a fixed order. The statements are logically consistent within each paragraph. In a golf tournament, there were seven golfers: Eve, Rob, Dan, Mel, Ana, Eli, and Ada. Ada finished above Rob. Eve finished below Rob. Mel finished above Eli. Ada finished below Dan. Ana finished third. Eli finished second.

Options:

- (A) Eve finished third-to-last
- (B) Rob finished third-to-last
- (C) Dan finished third-to-last
- (D) Mel finished third-to-last
- (E) Ana finished third-to-last
- (F) Eli finished third-to-last
- (G) Ada finished third-to-last

Answer (Qwen1.5-7B w. LogicPro) :



Let's analyze the given information and order of the golfers:

- 1. Ada finished above Rob.
- 2. Eve finished below Rob.
- 3. Mel finished above Eli.
- 4. Ada finished below Dan.
- 5. Ana finished third.
- 6. Eli finished second.

Now, let's determine the order of the golfers from first to last:

- 1. Ana (third)
- 2. Eli (second)
- 3. Mel (above Eli)
- 4. Dan (below Ada)
- 5. Ada (above Rob)
- 6. Rob (below Eve)
- 7. Eve (below Rob)

So, the order from first to last is: Ana, Eli, Mel, Dan, Ada, Rob, Eve.

Now, let's find the third-to-last golfer:

- 1. Ada finished above Rob.
- 2. Rob finished below Ada.
- 3. Ada finished third-to-last.

Therefore, the correct answer is **(G) Ada finished third-to-last.**

Answer (Qwen1.5-7B Gen_Logic) :



We know that Ana finished third, so she is in the third position. We also know that Eli finished second, so he is in the second position. Since there are seven golfers, the remaining positions are filled by the other golfers in order from first to last. We can use the information given about the positions of Ada, Rob, Dan, and Mel to determine the positions of the other golfers.

- 1. Ana finished third.
- 2. Eli finished second.
- 3. Ada finished above Rob, so she is in the fourth position.
- 4. Rob finished below Eve, so he is in the fifth position.
- 5. Dan finished above Ada, so he is in the sixth position.
- 6. Mel finished above Eli, so she is in the seventh position.

Now, we can determine the positions of the other golfers:

- 1. Ana finished third.
- 2. Eli finished second.
- 3. Ada finished fourth.
- 4. Rob finished fifth.
- 5. Dan finished sixth.
- 6. Mel finished seventh.

Since Eve finished below Rob, she must be in the first position. Therefore, the correct answer is:

(A) Eve finished third-to-last.

Figure 21: Failure Case for LogicPro: bbh tracking shuffled objects

LeetCode 420: Strong Password Checker

A password is considered strong if the below conditions are all met:

- * It has at least `6` characters and at most `20` characters.
- * It contains at least **one lowercase** letter, at least **one uppercase** letter, and at least **one digit**.
- * It does not contain three repeating characters in a row (i.e., `B**aaa**bb0` is weak, but `B**aa**b**a**0` is strong).

Given a string `password`, return `_` the minimum number of steps required to make `password` strong. if `password` is already strong, return `0`.

In one step, you can:

- * Insert one character to `password`,
- * Delete one character from `password`, or
- * Replace one character of `password` with another character.

Example 1:

Input: password = "a"

Output: 5

Example 2:

Input: password = "aA1"

Output: 3

Example 3:

Input: password = "1337C0d3"

Output: 0

Constraints:

- * `1 <= password.length <= 50`
- * `password` consists of letters, digits, dot `'.'` or exclamation mark `'!'`.

Test Sample Inputs:

```
# Test sample input 1
password = "a"

# Test sample input 2
password = "A1"

# Test sample input 3
password = "1337C0d3"

# Test sample input 4
password = "abcdef"

# Test sample input 5
password = "ABCDEF"

# Test sample input 6
password = "123456"

# Test sample input 7
password = "aA1."

# Test sample input 8
password = "aA1!"

# Test sample input 9
password = "aA1abc"

# Test sample input 10
password = "aA1abcABC123"
...
```

Figure 22: Test Sample Inputs: LeetCode 420: Strong Password Checker

LeetCode 524: Longest Word in Dictionary

Given an array of strings `words` representing an English Dictionary, return the longest word in `words` that can be built one character at a time by other words in `words`.

If there is more than one possible answer, return the longest word with the smallest lexicographical order. If there is no answer, return the empty string.

Note that the word should be built from left to right with each additional character being added to the end of a previous word.

Example 1:

Input: words = ["w", "wo", "wor", "worl", "world"]

Output: "world"

Explanation: The word "world" can be built one character at a time by "w", "wo", "wor", and "worl".

Example 2:

Input: words = ["a", "banana", "app", "appl", "ap", "apply", "apple"]

Output: "apple"

Explanation: Both "apply" and "apple" can be built from other words in the dictionary. However, "apple" is lexicographically smaller than "apply".

Constraints:

* 1 <= words.length <= 1000

* 1 <= words[i].length <= 30

* words[i] consists of lowercase English letters.

Test Sample Inputs:

Test sample input 1
words = ["a", "b", "c", "d", "e"]

Test sample input 2
words = ["apple", "app", "appl", "apply"]

Test sample input 3
words = ["world", "worl", "wor", "wo", "w"]

Test sample input 4
words = ["banana", "ban", "ba", "b"]

Test sample input 5
words = ["abc", "abcd", "abcde", "abcdef"]

Test sample input 6
words = ["x", "xy", "xyz", "xyza", "xyzab"]

Test sample input 7
words = ["hello", "hell", "he", "h"]

Test sample input 8
words = ["abc", "abcd", "abcde", "abcdefg"]

Test sample input 9
words = ["a", "ab", "abc", "abcd", "abcde", "abcdef"]

Test sample input 10
words = ["xyz", "xy", "x"]

...

Figure 23: Test Sample Inputs: LeetCode 524: Longest Word in Dictionary

LeetCode 2231: Largest Number After Digit Swaps by Parity

You are given a positive integer `num`. You may swap any two digits of `num` that have the same **parity** (i.e. both odd digits or both even digits).
Return the **largest** possible value of `num` after **any** number of swaps.

Example 1:
Input: num = 1234
Output: 3412
Explanation: Swap the digit 3 with the digit 1, this results in the number 3214. Swap the digit 2 with the digit 4, this results in the number 3412. Note that there may be other sequences of swaps but it can be shown that 3412 is the largest possible number. Also note that we may not swap the digit 4 with the digit 1 since they are of different parities.

Example 2:
Input: num = 65875
Output: 87655
Explanation: Swap the digit 8 with the digit 6, this results in the number 85675. Swap the first digit 5 with the digit 7, this results in the number 87655. Note that there may be other sequences of swaps but it can be shown that 87655 is the largest possible number.

Constraints:
 * $1 \leq \text{num} \leq 10^9$

Test Sample Inputs:

| | |
|--|--|
| # Test sample input 1
num = 135792 | # Test sample input 6
num = 2468109 |
| # Test sample input 2
num = 468135 | # Test sample input 7
num = 1357924 |
| # Test sample input 3
num = 753192 | # Test sample input 8
num = 4681357 |
| # Test sample input 4
num = 1234567 | # Test sample input 9
num = 7531928 |
| # Test sample input 5
num = 9876543 | # Test sample input 10
num = 12345678 |
| | ... |

Figure 24: Test Sample Inputs: LeetCode 2231: Largest Number After Digit Swaps by Parity

LeetCode 2208: Minimum Operations to Halve Array Sum

You are given an array `nums` of positive integers. In one operation, you can choose **any** number from `nums` and reduce it to **exactly** half the number. (Note that you may choose this reduced number in future operations.)
Return the **minimum** number of operations to reduce the sum of `nums` by **at least** half.

Example 1:
Input: nums = [5,19,8,1]
Output: 3

Example 2:
Input: nums = [3,8,20]
Output: 3

Constraints:
 * $1 \leq \text{nums.length} \leq 10^5$
 * $1 \leq \text{nums}[i] \leq 10^7$

Test Sample Inputs:

| | |
|---|---|
| # Test sample input 1
nums = [5, 19, 8, 1] | # Test sample input 6
nums = [10, 20, 30, 40, 50, 60, 70, 80, 90, 100] |
| # Test sample input 2
nums = [3, 8, 20] | # Test sample input 7
nums = [100, 200, 300, 400, 500, 600, 700, 800, 900, 1000] |
| # Test sample input 3
nums = [1, 1, 1, 1, 1] | # Test sample input 8
nums = [1, 1, 1, 1, 1, 1, 1, 1, 1] |
| # Test sample input 4
nums = [1000000, 1000000, 1000000] | # Test sample input 9
nums = [2, 4, 6, 8, 10, 12, 14, 16, 18, 20] |
| # Test sample input 5
nums = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10] | # Test sample input 10
nums = [1000, 1000, 1000, 1000, 1000] |

Figure 25: Test Sample Inputs: LeetCode 2208: Minimum Operations to Halve Array Sum

leetcode 524: Longest Word in Dictionary through Deleting

INPUT:

Title: Form the Longest Word

Q: You are given a string "abcdefghij" and a list of words ["abcdefghij", "abcdefgh", "abcdef", "abcd", "abc"]. Determine the longest word from the list that can be formed by deleting some characters from the given string. If there are multiple possible results, find the longest word with the smallest alphabetical order. If no such word exists, conclude that the result is an empty string.

TARGET:

To determine the longest word in the list that can be formed by deleting some characters from the given string "abcdefghij", follow these steps:

1. **Initialization**:

- Begin with an empty string to represent the longest word found so far.

2. **Check Each Word in the List**:

- Iterate through the list of words, starting from the first word.
- For each word, match it with the given string "abcdefghij" by sequentially checking characters.
- If the current word can be formed by deleting some characters from the given string, update the longest word found.

3. **Choose the Longest Word**:

- If multiple words of the same length are found, the one found first (having a higher position in the list) will be considered the answer as it will be lexicographically smaller, assuming the list is sorted.

Step-by-Step Reasoning:

1. **Initialization**:

- Start with an empty string as the longest word found.

2. **Word Checking**:

- Begin with the first word in the list, "abc". Sequentially match characters with the given string "abcdefghij", and since it can be formed by deleting characters, update the longest word to "abc".
- Proceed to the next word, "abcd". Match characters sequentially with the given string. Since it can be formed by deleting characters, update the longest word to "abcd".
- For the next word, "abcdef", follow the same process. Since it can be formed, update the longest word to "abcdef".
- Moving on to "abcdefgh", it can also be formed by deleting characters, so update the longest word to "abcdefgh".
- Finally, check "abcdefghij". Since it can be formed by not deleting any characters, update the longest word to "abcdefghij".

3. **Final Check and Result**:

- Since "abcdefghij" is the longest word found and can be fully matched with the given string, no further checks are necessary.

Therefore, the longest word that can be formed from the given string by deleting some characters, while also being lexicographically smallest if there are multiple results of the same length, is "abcdefghij".

Figure 26: LogicPro: leetcode 524 Longest Word in Dictionary through Deleting

leetcode 557: Reverse Words in a String III

INPUT:

Title: Reverse Word Characters

Q: You are given a sentence "Special characters @#\$. Reverse the order of characters in each word within the sentence, while preserving the whitespace and initial word order. Can you determine the resulting sentence?

TARGET:

To determine the resulting sentence after reversing the order of characters in each word while preserving the whitespace and initial word order in the sentence "Special characters @#\$", follow these steps:

1. **Identify the Words**:

- Break down the sentence into individual words, recognizing that a word is a sequence of characters separated by spaces.

2. **Reverse Each Word**:

- For each identified word, reverse the sequence of characters.

3. **Maintain White Space and Word Order**:

- After reversing each word, assemble the sentence by placing the reversed words back into their original positions, ensuring the whitespace between words and the original word order are preserved.

Step-by-Step Reasoning:

1. **Identification of Words**:

- The sentence given is: "Special characters @# \$"

- The words identified are: "Special", "characters", "@# \$"

2. **Reversal of Each Word**:

- Word 1: "Special"

- Reverse of "Special" is "laicepS"

- Word 2: "characters"

- Reverse of "characters" is "sretcarahc"

- Word 3: "@# \$"

- Reverse of "@# \$" is "\$#@"

3. **Reconstructing the Sentence**:

- Place each reversed word back in its original position with original spacing:

- "laicepS sretcarahc \$#@"

The final result after reversing the order of characters in each word while preserving the whitespace and initial word order in the sentence "Special characters @# \$" is "laicepS sretcarahc \$#@".

Figure 27: LogicPro: leetcode 557: Reverse Words in a String III