

# From Induction to Deduction: Hierarchical Rule Learning for LLM Reasoning Tasks

Anonymous ACL submission

## Abstract

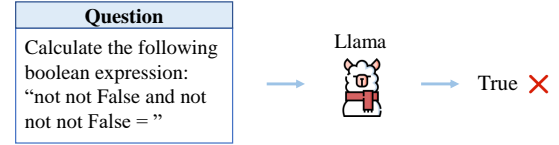
The growing scale of large language models (LLMs) brings powerful reasoning capabilities. Effective instructions can significantly improve LLMs’ reasoning abilities. However, current prompt construction methods mainly focus on enhancing reasoning strategies or knowledge from an inductive perspective, overlooking the necessity of providing models with rules to follow from a deductive viewpoint. This results in the model lacking reliable underlying logic, leading to incorrect answers. This paper proposes an **Induction to Deduction (I2D)** framework to enable LLMs to automatically extract rules from tasks and apply these rules during the actual reasoning process. In the framework, we combine hierarchical clustering and MCTS to extract potential rules in tasks as comprehensively as possible. Experimental results on complex tasks such as GSM8K and Big-Bench Hard demonstrate a superior performance improvement of up to 15% compared to few-shot settings and show the universal applicability of our method.

## 1 Introduction

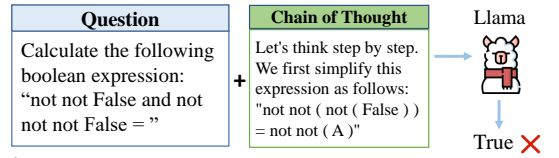
The recent advancements of Large language models (LLMs) have endowed them with strong reasoning capabilities (Brown et al., 2020; Liu et al., 2023c; Wei et al., 2022), where the instruction plays a vital role (Reynolds and McDonell, 2021; Sahoo et al., 2024). It is validated that, by mimicking the thinking patterns of humans, well-written instructions can significantly improve the effectiveness of LLMs in reasoning tasks (Wei et al., 2023; Yao et al., 2023; Dong et al., 2024; Wang et al., 2020).

However, creating high-quality instructions necessitates the expertise of skilled individuals (Wei et al., 2022; Liang et al., 2023; Rae et al., 2022). Even for different datasets of the same task, experts need to make numerous attempts to select better prompts that can cover as many scenarios as possible, which is costly (Yang et al., 2024b;

### a) LLM as Zero-Shot Reasoners



### b) CoT Prompting



### c) From Induction to Deduction

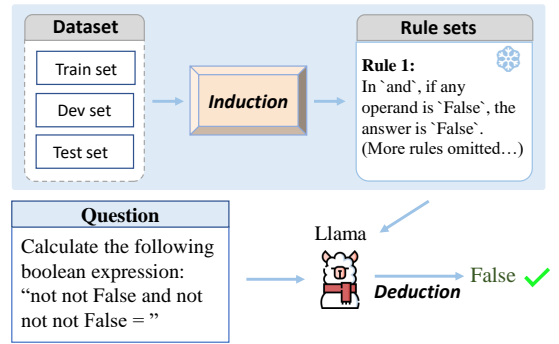


Figure 1: In zero-shot and few-shot settings, the Llama3.1-8b model makes a mistake on the question. Our framework utilizes rule sets pre-built by LLM themselves, which correctly answers this question.

Cao et al., 2023; Zhou et al., 2023b; Wei et al., 2023). Therefore, there is a pressing need for an automated system to generate optimal instructions.

Most prior instructions construction methods for reasoning tasks can be categorized into two main groups: strategy-enhanced reasoning and knowledge-enhanced reasoning (Qiao et al., 2023). The first group focuses on optimizing reasoning strategies by manually or automatically finding the best problem decomposition structure (Zhou et al., 2023a; Imani et al., 2023; Wang et al., 2023a). The second group aims to facilitate the reasoning process by providing more information or similar examples (Trivedi et al., 2023; Yang et al., 2022).



However, these methods listed above primarily target on inductive reasoning, improving LLMs’ reasoning capabilities by observing specific instances. Inspired by various forms of deductive reasoning in humans (Barwise, 1993), we aim to provide the model with more general principles that may be used during the problem-solving process at the task level to help enhance its reasoning abilities (Sun et al., 2024). Some studies (Yang et al., 2023; Zhu et al., 2024) have explored summarizing some rules from a single example and applying them to subsequent samples. But rules induced by these methods follow a fixed format (if..then) and have only been tested on relatively simple tasks such as text classification tasks.

In this work, we introduce an **Induction to Deduction (I2D)** framework, shown in Figure 1. In this boolean expression from Big-Bench Hard (Suzgun et al., 2022), the model, relying solely on in-context learning techniques, yields an incorrect response. However, when provided with explicit rules, the model generates the correct answer. Our I2D framework is divided into two stages. In the first induction stage, we summarize hierarchical sets of rules from the training set. In the deduction stage, we select and refine these rules and apply them to the inference. Specifically, in the induction stage, since hypotheses based on single examples are often insufficiently accurate, we leverage a multi-stage approach to summarize rules from some similar problems, employing the MCTS algorithm to search for the best combination of rules.

For a comprehensive evaluation, we test I2D on 30 challenging reasoning tasks, with difficult levels range from simpler tasks like TweetEval (Barbieri et al., 2020), GSM8K (Cobbe et al., 2021) to Big-Bench Hard (Suzgun et al., 2022) dataset. I2D surpasses the few-shot setting in 27 out of 30 tasks, achieving an absolute improvement of up to 52% in a single task, and an average improvement of up to 15%. Not only does the I2D framework significantly improve performance across various tasks, but it also offers several additional advantages. Firstly, our approach fully automates the process of mining and applying rules, requiring no human intervention, and does not need manual adjustments for different tasks, making it more cost-effective than manually writing prompts. Secondly, compared to parameter optimization, rule sets are more interpretable. Ultimately, our framework can be effortlessly integrated with other prompts or reasoning methods.

Overall, our contributions are three-fold: 1) We propose a hierarchical rule learning framework that enables LLM to automatically summarize hierarchical rules for various reasoning tasks. 2) We apply MCTS in the automated rule mining process to ensure the correctness of the generalization of rules at a lower cost. 3) We evaluate the I2D framework in over 30 challenging tasks, including logical reasoning, mathematics, and code generation, demonstrating its versatility and effectiveness.

## 2 Related Work

**Prompting optimization.** There have been many studies exploring methods for prompt optimization, divided into approaches tuning soft prompts and tuning hard prompts. Methods tuning hard prompts such as Chain-of-Thought (CoT) (Wei et al., 2023), Tree of Thoughts (ToT) (Yao et al., 2023), and Least to Most (Zhou et al., 2023a) rely on some degree of manual crafting. BPO (Cheng et al., 2024a) trains a seq2seq model as the prompt preference optimizer. OPRO (Yang et al., 2024b) demonstrates that LLMs can be used to solve some optimization problems and applies LLMs to optimize prompts. However, these methods require extensive search processes and have poor interpretability. Methods tuning soft prompts (Liu et al., 2023b; Lester et al., 2021; Gu et al., 2022) typically transform prompts into continuous vectors that can be optimized, and then optimize the vectors within this space. These methods often require access to the internal parameters of LLM. Compared to these methods, our work offers better interpretability and transferability.

**Rule following capabilities of LLMs.** LLMs have demonstrated strong reasoning abilities across multiple tasks (Liu et al., 2023a; Hendrycks et al., 2021; Huang et al., 2024), where reasoning according to rules is also an important capability. Several studies have evaluated the rule-following abilities of LLMs. LogicGame (Gui et al., 2024) designs various games that include an initial state and game rules to assess LLMs’ ability to understand and apply game rules to solve problems. ULogic (Wang et al., 2024) constructed a reasoning rule library containing primitive rules and more complex combined rules across five domains, and evaluated LLM’s ability to master reasoning rules on this library. SolverLearner (Cheng et al., 2024b) further differentiated LLM’s abilities in inductive and deductive reasoning, suggesting that they often lack deductive reasoning abilities, especially in tasks



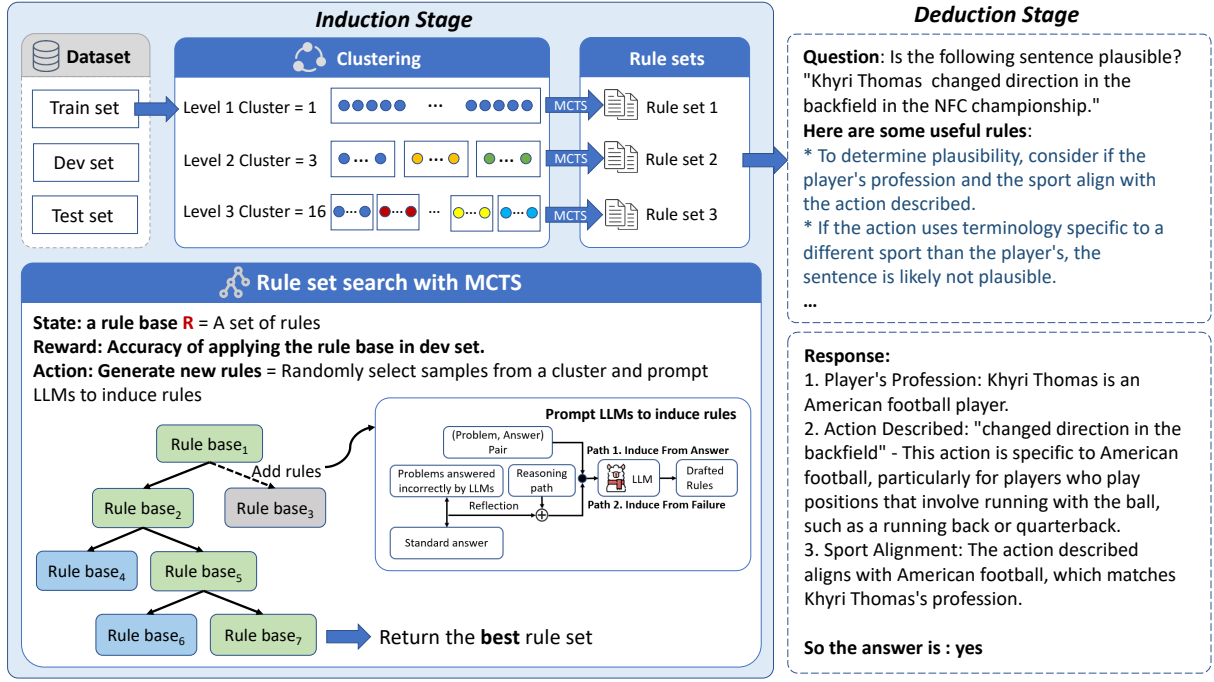


Figure 2: Overview of the I2D framework. We perform hierarchical clustering on the train set, and for each layer, MCTS algorithm is used to search for an optimal rule set. When generating rules, both positive and negative samples are used to guide the derivation of rules LLMs, which is described in detail in Section 3.3.1.

involving "counterfactual" reasoning. These works partially demonstrate the feasibility of our work.

**Rule induction.** Traditional rule induction (Agrawal and Srikant, 1994; Galárraga et al., 2013) aims to discover interesting relationships between variables in large databases, while our work targets more general rules. Recently, some work has applied rule learning to natural language, HtT (Zhu et al., 2024) extracts rules from LLMs' inference process using regular expressions and tracks the usage of rules to determine their effectiveness. TRAN (Yang et al., 2023) proposes a Tuning-free Rule Accumulation (TRAN) framework, which guides LLMs in improving their performance by learning from previous mistakes. Self-Discover (Zhou et al., 2024) allows LLMs to self-assemble atomic reasoning modules into an effective reasoning structure, which is then applied to subsequent tasks, sharing similarities with human reasoning patterns. DEER (Yang et al., 2024d) creates a dataset containing 1.2k rule-fact pairs. Using this dataset, the author analyzed how well LLMs can induce natural language rules from facts. Yang et al. proposes a rule distillation algorithm that explicitly encodes knowledge into the parameters of LLMs by learning from the in-context signals generated within the model. Compared to these works, our work provides a superior rule mining

and application framework, achieving significant performance improvements.

### 3 Method

Summarizing rules plays a significant role in problem-solving, as it helps us quickly identify and generalize the key elements of a problem, thereby enhancing the efficiency and accuracy of problem resolution. For example, in the boolean expression "A or B", if the value of A is true, the result of the entire expression is necessarily true, so there is no need to calculate B. Using this rule can avoid unnecessary calculations and if we provide it to LLMs, it may prevent LLMs from outputting unnecessary errors (Sun et al., 2024).

#### 3.1 Task Definition

As previously discussed, our task is to summarize rules to enhance LLM's reasoning ability. Formally, we divide the entire dataset containing questions and answers  $\{(x_i, y_i)\}_{i=1}^N$  into train set  $S_{train}$ , validation set  $S_{val}$ , and test set  $S_{test}$ . Our goal is for LLM to automatically construct a rule set  $R = \{r_1, r_2, r_3, \dots, r_n\}$  in natural language based on the data in  $S_{train}$ , such that using rules in  $R$  can improve LLM's performance on  $S_{test}$ .

Rules are conclusions derived from certain facts that can be used to guide or simplify subsequent



reasoning processes, such as "*Rule of Logical Negation: Double negation returns the original value.*" A rule  $r$  is expressed in natural language but is not limited to the if...then structure.

## 3.2 Overview

In previous work, LLMs are prompted to summarize rules from a single example each time. However, in real life, it is easier to propose a correct rule if we observe some samples with similar characteristics. Therefore, we need to cluster the samples, sampling some similar ones from the same cluster each time, and prompt LLMs to summarize rules from these samples.

In addition, we believe that rules are hierarchical and can be divided into multiple levels from general to specific. The most general rules can be applied to all situations under a certain task, such as the need to calculate boolean expressions from left to right, while specific rules only take effect in certain specific situations, such as the short-circuit principle. To enable the model to induce rules hierarchically, we choose to use hierarchical clustering methods. Because more clusters mean that the problems within each cluster are closer, we construct clusters of different sizes at different levels to allow LLMs to summarize specific rules from more similar samples and general rules from more diverse samples.

Based on the above ideas, we propose the I2D framework. The overall flow of our framework is shown in the Figure 2. Our approach is divided into the Induction Stage and the Deduction Stage.

## 3.3 Induction Stage

In previous works, they mostly provide LLMs samples directly and prompt for a hypothesis rule which is not well-designed. DEER (Yang et al., 2024d) uses PLM to construct a multi-step rule summarizer. We go further by constructing a multi-step, multi-method induction stage, divided into three parts as shown in Figure 2.

### 3.3.1 Rule Summarization

For a given sample set, we aim to extract rules using LLM. We first test the samples on the target LLM and collect the LLM's thought process. The samples are then divided into correct and incorrect parts based on whether they are answer correctly, and each part is processed separately. For the correct part, we want LLM to learn from the answer. We concatenate all samples and thoughts into LLM

and ask LLM to summarize the rules or patterns. For the incorrect part, we want LLM to learn from failure, identifying commonalities in errors to prevent future mistakes. We will have LLM reflect on each incorrect sample's thought process against the correct answer, concatenate all thoughts and reflection results into LLM, and ask LLM to summarize rules to prevent future errors. Finally, we input the samples and outputs of both parts into LLM, asking it to concisely summarize all correct rules. The full prompts are shown in Appendix C. Through learning on both positive and negative samples, more comprehensive rules are summarized.

### 3.3.2 Sample Clustering

However, due to the limited input length of the model, we cannot input all the data into the model at once. Additionally, as we mentioned earlier, rules also have different levels. At a higher level, we need to extract some rules with a broader scope; at a lower level, we need to extract more specific rules. If we divide all the problems into more sets, the problems within each set will be more similar.

Therefore, we consider first performing hierarchical clustering on all samples. Specifically, at each level, we apply the k-means clustering algorithm to all samples, setting a different number of cluster centers at different levels, with fewer centers at higher levels. As a result, the clustering at different levels is uncorrelated, and two samples that belong to the same class at a lower level may be assigned to different classes at a higher level.

After clustering, the questions within the same cluster are relatively similar, with varying degrees of similarity at different levels.

### 3.3.3 Rule Search

Since LLM is a probabilistic model, the rules it summarizes each time will not be the same, and we want to find the best rules. At each clustering level, we will perform an MCTS using the clustering results of that level to obtain the optimal rule set for that level. The MCTS algorithm has 4 main steps: **Selection**, **Expansion**, **Simulation**, and **Backpropagation**. At each node in the search tree, we store the rule set corresponding to that node. Initially, the tree has a root node with the rule set corresponding to the best rule set from the previous layer.

**Selection.** Before expanding a node, we first need to determine which node should be expanded. We start from the root node and, at each step,



choose the most promising node among its children based on the Upper Confidence bounds applied to Trees (UCT) formula (Kocsis and Szepesvári, 2006), until we reach a leaf node. The UCT value expression is  $UCT_i = q_i + c \cdot \sqrt{\frac{\ln N_i}{n_i}}$  where  $q_i$  represents the score of the  $i$ -th node,  $n_i$  denotes the number of times the  $i$ -th thought has been visited, and  $N_i$  is the total number of visits to the parent of the  $i$ -th node.

**Expansion.** When expanding nodes, we randomly select a set from the clustering results of that layer, randomly choose 10 samples for rule summarization, then merge the newly generated rules with the existing rule set of the node, and call LLM to determine if there are any duplicate rules, finally adding the non-duplicate rules to the rule set of the newly expanded node.

**Simulation.** After expanding a new node, we use the accuracy on the validation set as the reward. For each question in the validation set, we concatenate the rule set of the new node into the prompt, then send it to the evaluation LLM. We check if the output is correct, and finally calculate the average accuracy on the validation set as the reward.

**Backpropagation.** The calculated reward will be propagated back to each parent node of the new node. For each parent node, their score will be updated to the larger value between the original score and the reward calculated this time.

### 3.4 Deduction Stage

When applying rules, we will concatenate all rules in the order from general to specific, using this as the rules section. We will then combine it with the other parts in the sequence of {examples} + "rules:" + {rules} + "Q:" + {question} to input into the evaluation LLM. The specific prompt design is detailed in Appendix C. When we want the model to produce a reasoning chain, we will include the phrase "think step by step." It is noteworthy that although some special rules only take effect in certain situations, we found that directly providing all rules to LLMs yields the best results in experiments.

## 4 Experiments

### 4.1 Experiment Setup

**Datasets.** We selected a series of reasoning tasks of varying difficulty, starting from the easiest to the most challenging: Tweeteval (Barbieri et al., 2020), Big-Bench Hard (BBH) (Suzgun et al., 2022), GSM8K (Cobbe et al., 2021) and KQA

Pro (Cao et al., 2022). Tweeteval is a dataset for text classification in social scenarios, consisting of 7 tasks. Similar to (Yang et al., 2023), we chose the "Offensive" sub-task out of these tasks. Big-Bench Hard is composed of 23 challenging tasks extracted from Big-Bench, encompassing a variety of tasks, most of which require LLM multi-step reasoning capabilities. GSM8K is a dataset of 8.5K high-quality, linguistically diverse elementary math word problems created by human question writers. We also tested on the most challenging datasets KQA Pro. KQA Pro is a semantic parsing task that requires LLMs to generate KoPL query statements from natural language.

We use accuracy to measure the performance of the model on Tweeteval, Big-Bench Hard, KQA Pro, and GSM8K. We provide detailed settings of each task in Appendix A.

**Models.** In the process of summarizing rules, we attempted to use GPT-4o from OpenAI (OpenAI, 2023), Grok-beta (xAI, 2024), and DeepSeek-R1-Distill-Llama-70B (DeepSeek-AI et al., 2025) to summarize the rules. For evaluation LLM, there are open-source models like Llama-3.1-8b, Llama-3.1-70b, Qwen2-7b (Yang et al., 2024a) and Mistral-7b (Jiang et al., 2023), as well as closed-source models like GPT3.5 and Deepseek-V2.5 (DeepSeek-AI, 2024).

**Baselines.** We compared our method with other prompting methods for reasoning tasks. (1) **Few-shot:** We provide the model with three input-output examples and ask the model to output the answer without showing intermediate steps. We demonstrate few-shot results under various settings with sample numbers ranging from 3 to 32. (2) **CoT** (Wei et al., 2023): We provide the model with three manual chain-of-thought examples and then prompt the model to generate a reasoning process to arrive at the final result by thinking step by step. (3) **Self-Consistency (SC)** (Wang et al., 2023b): For each question, we sample  $k=5$  reasoning chains and take the most frequently occurring answer as the final answer. (4) **HtT** (Zhu et al., 2024): Before the testing phase, we extract 5 test samples, prompting LLMs to summarize rules from each sample, and then directly concatenate these rules to add them to the prompt during the testing phase.

**Cost.** We analyze the usage of LLMs to summarize the rule set for each dataset, averaging 156 API calls and approximately 250k tokens with rollouts = 10. Detailed data can be found in Appendix D. During the rule summarization process, approxi-



Method	Tweeteval	BBH	GSM8K	KQAPro
Direct	76.16	45.77	81.73	29.70
3-shot	73.26	46.82	81.50	39.42
8-shot	69.88	47.75	81.50	47.66
16-shot	69.53	47.27	82.71	<u>52.34</u>
32-shot	70.58	45.68	80.89	<b>56.46</b>
CoT	-	60.93	84.00	-
SC	73.37	65.64	<b>87.79</b>	-
HtT	72.44	60.96	83.09	39.36
I2D (Grok)	77.09	51.25	80.89	33.08
+Few-shot	78.26	59.64	82.03	42.62
+CoT	-	63.38	84.31	-
I2D (R1-70b)	80.93	54.88	81.58	33.08
+Few-shot	<u>81.86</u>	58.70	81.58	42.62
+CoT	-	<u>65.66</u>	82.79	-
I2D (GPT-4o)	81.05	53.19	81.27	31.60
+Few-shot	<b>82.09</b>	61.05	82.26	43.66
+CoT	-	<b>66.05</b>	<u>84.58</u>	-

Table 1: Comparison of accuracy on four datasets under both zero-shot and few-shot settings. We mark the best and the second best performance in **bold** and underline. LLM which induces the rule is marked in brackets, and the evaluation model is Llama3.1-8B. In Few-shot, SC, HtT, and all CoT methods, we provided 3 examples in the prompt.

mately 1k repeatable examples were accessed and formed 17.6 rules per task.

**Implementation Details.** For the KQA Pro dataset, we set rollouts = 16, while for other datasets, rollouts = 10. For rule generation LLM, we set the temperature = 1.0. For rule validation LLM, we use the vLLM (Kwon et al., 2023) framework to accelerate inference, setting the temperature=0.0 to ensure consistent output each time. For all datasets, we set the questions to three layers. In the first and second layer, the number of clusters is 1 and 8 respectively, allowing the model to summarize three rules based on 10 samples. In the third layer, the number of clusters is set to an integer that makes the average number of questions per cluster equal to 10, and the model summarizes one rule based on 5 samples.

## 4.2 Main Results

The I2D framework enhances the performance of LLMs on a variety of tasks. Table 1 shows the overall results on Tweeteval, Big-Bench Hard (BBH), GSM8K and KQA Pro. Compared to each corresponding setting, our method results in a certain performance improvement. Adding more samples to the prompt does not always improve accuracy; in BBH, even the best-case improvement is only 2%.

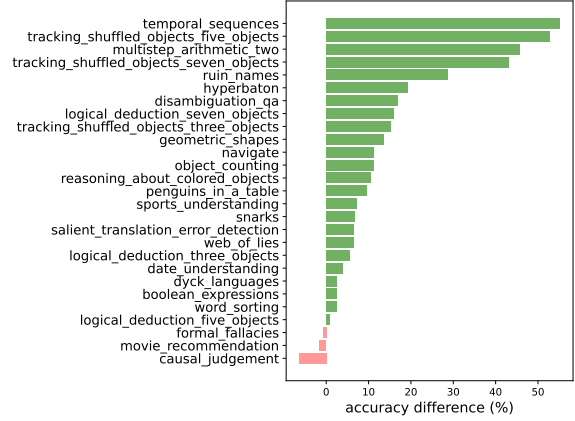


Figure 3: Accuracy improvement between applying rules and few-shot setting on 27 BBH tasks. Detailed descriptions of each task are provided in Appendix A.

Compared to self-consistency method, which has 5 times the computational complexity, our method still achieves up to a 7% improvement in accuracy, with overall performance being comparable.

On TweetEval, our method improves the accuracy by 5% with the rule set induced by GPT-4o and 1% by Grok. In this task, we found few-shot does not work as well as summarized rules and rules induced by stronger LLM are more useful.

On Big-Bench Hard, our method improves 14.2% and 5% upon few-shot and CoT respectively. Figure 3 shows the gap on each task separately and we provide detailed per-task accuracy in Appendix B. For the rule set summarized by GPT-4o, the biggest improvements were made on three separate tasks. The improvement is not obvious in some tasks, which we believe is because these tasks focus more on basic rules that LLMs have already mastered. We also observed a performance drop on certain tasks, which we believe is due to the model learning incorrect rules. We will analyze this in the following section.

For more complex tasks GSM8K and KQA Pro, the improvement was relatively modest, with a 2~4% gain on KQA Pro and a 0.5% gain on GSM8K. On KQA Pro, the rules summarized by the model primarily concern the correspondence between grammar and functions with natural language. For example, “If a question involves entities (people, places, organizations, etc.), use the ‘Find()’ function to locate each entity in the database”. Through error analysis, we discover they can effectively reduce errors in function usage, but are less effective for formatting issues.



Model	GSM8K		BBH		Tweeteval	
	FS	Ours	FS	Ours	FS	Ours
Llama3.1-8b	81.50	+0.76	46.82	+14.23	76.16	+4.89
Llama3.1-70b	93.18	+0.60	65.10	+1.22	77.79	+5.93
Deepseek-70b	88.55	+5.54	83.40	+2.00	77.56	+3.25
Qwen2-7b	82.34	+2.65	50.41	+0.71	67.79	+11.74
Mistral-7b	14.33	+19.71	44.77	+0.68	71.28	+8.49
GLM4-9b	13.80	+6.90	51.19	+2.22	78.49	+1.51
GPT-3.5	34.57	+2.05	40.65	+4.73	60.00	+8.14
Deepseek-V2.5	45.19	+0.07	66.45	+1.36	75.35	+6.39

Table 2: Generalization testing of the rule set. **FS** stands for few-shot setting. **Deepseek-70b** refers to DeepSeek-R1-Distill-Llama-70b model. The rule set demonstrates good generalization, indicating that our method is effective for a wide range of models.

### 4.3 General Effectiveness

To verify the general effectiveness of our approach, we create a set of rules using GPT-4o as the inducer and Llama-3.1-8b as the evaluation model. Then we apply these rules to various other models. Table 2 shows our results.

Since our rules are represented by natural language, they should be effective for various LLMs. It can be found that all models showed improvement, indicating that the rules summarized by our method exhibit strong generalization capabilities. However, the degree of improvement varied across models. For instance, the GSM8K task running on Mistral-7b achieve nearly a 20% gain, while the Tweeteval task running on Qwen2-7b saw a 12% improvement. On the powerful reasoning model DeepSeek-R1-Distill-Llama-70b, our method remains effective, improving by 5.5% and 2% on GSM8K and BBH respectively. Overall, models with lower baseline performance on a given task tends to experience more significant gains after applying the rules.

### 4.4 Error Analysis

To analyze how many of the rules proposed by the model are incorrect, we iterate through each rule and calculate the accuracy when used it alone on the validation set. Then we compares it with the accuracy without the rule set, the difference between the two is used as a metric to measure the value of the rule. We consider the rule to be incorrect when this value is negative. When these two are within the margin of error, we consider the rule to be redundant. The remaining rules are considered correct. We have executed this step on several tasks from various accuracy segments.

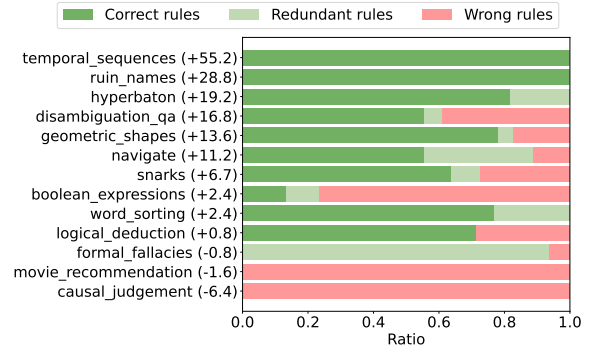


Figure 4: The proportion of correct and incorrect rules in the rule set generated by GPT-4o. The accuracy increase under few-shot setting is shown in parentheses.

Based on this result, we conduct manual checks.

As shown in the Figure 4, for tasks where accuracy has improved, most rules are correct. Except for the task of calculating Boolean expressions, which we consider to be relatively basic with fewer underlying rules, the remaining rules are overfitting products, thus showing errors when used alone. Most incorrect rules have loose conditions and only work correctly when combined with others.

### 4.5 Ablation Studies

**The effect of clustering.** We divide the rule set into several layers to generate rules with different generalization. We compared the following options: 1) using hierarchical clustering which we proposed, 2) not clustering, sampling from the entire training set each time.

Method	Tweeteval	BBH	GSM8K
<b>I2D</b>	<b>81.05</b>	<b>61.05</b>	<b>82.26</b>
-clustering	80.01	58.47	81.79

Table 3: Ablation on clustering across different tasks.

Table 3 shows that our method with hierarchical clustering achieves better final accuracies in all datasets. It proved that hierarchical clustering can make LLMs to summarize rules that are more beneficial for solving problems.

**The effect of MCTS.** In our approach, we employ the MCTS method to search for the optimal combination of rules, as we believe that merging different rules into a rule set will yield varying effects. If only the best-performing rules are retained, it may lead to local optima. We compared three rule search methods: 1) MCTS, using the default parameters mentioned above; 2) Replacing



BBH-Ruin names		BBH-Hyperbaton	
Example: Which of the following is a humorous edit of this artist or movie name: 'rain man'?		Example: Which sentence has the correct adjective order:	
(A) ruin man (B) rains man (C) rain men (D) rainmman		(A) midsize old grey Brazilian sweater (B) midsize grey Brazilian old sweater	
Value	Rule	Value	Rule
27.2	If a word in the option sounds similar to a word in the original name but has a different meaning that introduces humor, then select that option.	10.4	The standard order of adjectives is: Opinion, Size or Quantity, Age, Shape, Color, National Origin, Material or Composition, Purpose or Function.
20.0	Avoid options that are merely typographical errors or simple misspellings without any wordplay or clever substitution.	5.6	Non-adjective words disrupt the adjective order and typically make the sentence incorrect.
BBH-Boolean Expressions		GSM8K	
Example: not True or False or ( False ) is		Example: If there are 10 eggs in a basket, and there are twice as many eggs in a second basket, how many eggs are in both baskets put together?	
Value	Rule	Value	Rule
3.2	If the expression is a combination of "False" and "or" without "and" and contains "True", then the answer is True.	0.31	If a problem involves calculating totals, then break the problem into smaller parts, perform necessary arithmetic operations, and combine results to find the answer.
3.2	If the expression contains "False and" followed by any nested False expressions, then the result is False.	0.16	If the problem involves multiple components contributing to a total (e.g., costs, points, time), then calculate each component separately and sum them up.

Figure 5: Examples of I2D framework induced rule set. For each task, we show the two rules with the highest value, where the value of a rule is the increase in accuracy on the task’s test set after using only the rule.

MCTS with a greedy approach, where only one rule set is maintained at each layer. In each iteration, if adding a rule improves accuracy, the newly generated rule is added to the rule set. 3) We simulated the method in Tuning-free Rule Accumulation (TRAN) framework (Yang et al., 2023), where we retained only one layer, summarized rules only from error samples, and used the greedy approach in 2) to construct the rule set.

Task	Method		
	1	2	3
Ruin names	<b>88.80</b>	82.40	79.20
Hyperbaton	<b>74.40</b>	67.20	72.40
Snarks	<b>64.04</b>	63.29	62.92
Causal judgement	50.00	53.19	<b>57.45</b>
Movie recommendation	<b>73.60</b>	68.80	68.00
Word sorting	48.00	<b>49.60</b>	47.20
Boolean expressions	<b>76.00</b>	74.40	75.20
<b>Average</b>	<b>67.83</b>	65.55	66.05

Table 4: Ablation on search method across tasks.

The experimental results, shown in Table 4, indicate that MCTS performs the best at the same cost. The results suggest that using the MCTS algorithm to find the optimal combination of rules is closer to the global optimum, maximizing the ability of

LLMs to induce rules.

## 4.6 Case Study

We present examples of tasks and the best rules corresponding to these tasks in Figure 5.

In Big-Bench Hard tasks, LLMs generate detailed and effective rules, clearly outlining scenarios and providing accurate descriptions. The rules in Boolean Expressions highlight the advantages of hierarchical clustering, with LLMs summarizing secondary conclusions to supplement implicit knowledge, further enhancing performance.

In GSM8K, the rules are quite general. The best two rules primarily describe the use of divide-and-conquer algorithms to solve partial problems. There are no special solutions tailored to more specific types of problems as we might expect. However, these rules still improve the performance of LLMs to some extent.

## 5 Conclusion

We proposed I2D framework to utilize LLMs to build rule sets for reasoning tasks. We have observed significant improvements of up to 15% in challenging reasoning tasks across various LLMs. The ablation study of I2D demonstrates that the hierarchical clustering and MCTS methods we utilize are both very effective. We hope to explore more to stimulate LLMs to solve more complex tasks.



## 6 Limitation

Although I2D has demonstrated effectiveness and great potential in inference tasks, there are still some shortcomings that need to be explored in future work.

**Powerful LLM is needed to summarize the rules.** In the paper, we use GPT-4o to summarize the rules. In the initial phase of our work, we attempted to use a 7b model to summarize the rules, and the final result was that only half of the task performance was improved. Fortunately, we have found that using the Deepseek-distilled Llama-70b model can also achieve a level close to GPT-4o, which means we no longer need to submit data to closed-source models, thereby reducing the risk of leakage. We hope to summarize the rules by fine-tuning smaller models, not only improving the ability to summarize rules but also reducing computational overhead.

**Overfitting has occurred.** During testing on the GSM8K dataset, we noticed that while the accuracy on the validation set increased, the accuracy on the test set remained unchanged. We attempted to penalize rule sets with too many rules in the MCTS reward, but this did not yield satisfactory results. We hope to select more appropriate optimization methods to enhance the ability to summarize rules in complex tasks.

## References

- Rakesh Agrawal and Ramakrishnan Srikant. 1994. Fast algorithms for mining association rules in large databases. In *Proceedings of the 20th International Conference on Very Large Data Bases, VLDB '94*, page 487–499, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Francesco Barbieri, Jose Camacho-Collados, Leonardo Neves, and Luis Espinosa-Anke. 2020. [Tweeteval: Unified benchmark and comparative evaluation for tweet classification](#). *Preprint*, arXiv:2010.12421.
- Jon Barwise. 1993. [Everyday reasoning and logical inference](#). *Behavioral and Brain Sciences*, 16(2):337–338.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario

- Amodei. 2020. [Language models are few-shot learners](#). *Preprint*, arXiv:2005.14165.
- Shulin Cao, Jiaxin Shi, Liangming Pan, Lunyiu Nie, Yutong Xiang, Lei Hou, Juanzi Li, Bin He, and Hanwang Zhang. 2022. [KQA pro: A dataset with explicit compositional programs for complex question answering over knowledge base](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6101–6119, Dublin, Ireland. Association for Computational Linguistics.
- Shulin Cao, Jiajie Zhang, Jiaxin Shi, Xin Lv, Zijun Yao, Qi Tian, Juanzi Li, and Lei Hou. 2023. [Probabilistic tree-of-thought reasoning for answering knowledge-intensive complex questions](#). *Preprint*, arXiv:2311.13982.
- Jiale Cheng, Xiao Liu, Kehan Zheng, Pei Ke, Hongning Wang, Yuxiao Dong, Jie Tang, and Minlie Huang. 2024a. [Black-box prompt optimization: Aligning large language models without model training](#). *Preprint*, arXiv:2311.04155.
- Kewei Cheng, Jingfeng Yang, Haoming Jiang, Zhengyang Wang, Binxuan Huang, Ruirui Li, Shiyang Li, Zheng Li, Yifan Gao, Xian Li, Bing Yin, and Yizhou Sun. 2024b. [Inductive or deductive? rethinking the fundamental reasoning abilities of llms](#). *Preprint*, arXiv:2408.00114.
- 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](#). *Preprint*, arXiv:2110.14168.
- DeepSeek-AI. 2024. [Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model](#). *Preprint*, arXiv:2405.04434.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojuan Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan,



687	Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen,	Xiao, Yiyuan Li, Fan Zhou, Steffi Chern, Yiwei	746
688	Shanghao Lu, Shangyan Zhou, Shanhuang Chen,	Qin, Yan Ma, Jiadi Su, Yixiu Liu, Yuxiang Zheng,	747
689	Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng	Shaoting Zhang, Dahua Lin, Yu Qiao, and Pengfei	748
690	Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing	Liu. 2024. <a href="#">Olympicarena: Benchmarking multi-</a>	749
691	Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun,	<a href="#">discipline cognitive reasoning for superintelligent ai.</a>	750
692	T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu,	<i>Preprint</i> , arXiv:2406.12753.	751
693	Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao		
694	Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan	Shima Imani, Liang Du, and Harsh Shrivastava. 2023.	752
695	Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin	<a href="#">Mathprompter: Mathematical reasoning using large</a>	753
696	Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li,	<a href="#">language models.</a> <i>Preprint</i> , arXiv:2303.05398.	754
697	Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin,		
698	Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxi-	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-	755
699	ang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang,	sch, Chris Bamford, Devendra Singh Chaplot, Diego	756
700	Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	757
701	Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng	laume Lample, Lucile Saulnier, L��lio Renard Lavaud,	758
702	Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi,	Marie-Anne Lachaux, Pierre Stock, Teven Le Scao,	759
703	Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang,	Thibaut Lavril, Thomas Wang, Timoth��e Lacroix,	760
704	Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo,	and William El Sayed. 2023. <a href="#">Mistral 7b.</a> <i>Preprint</i> ,	761
705	Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yu-	arXiv:2310.06825.	762
706	jia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You,		
707	Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu,	Levente Kocsis and Csaba Szepesv��ri. 2006. Bandit	763
708	Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu,	based monte-carlo planning. In <i>Machine Learning:</i>	764
709	Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan,	<i>ECML 2006</i> , pages 282–293, Berlin, Heidelberg.	765
710	Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean	Springer Berlin Heidelberg.	766
711	Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao,		
712	Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zi-	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying	767
713	jia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song,	Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.	768
714	Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu	Gonzalez, Hao Zhang, and Ion Stoica. 2023. Effi-	769
715	Zhang, and Zhen Zhang. 2025. <a href="#">Deepseek-r1: Incent-</a>	cient memory management for large language model	770
716	<a href="#">tivizing reasoning capability in llms via reinforce-</a>	serving with pagedattention. In <i>Proceedings of the</i>	771
717	<a href="#">ment learning.</a> <i>Preprint</i> , arXiv:2501.12948.	<i>ACM SIGOPS 29th Symposium on Operating Systems</i>	772
		<i>Principles.</i>	773
718	Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan		
719	Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu,	Brian Lester, Rami Al-Rfou, and Noah Constant. 2021.	774
720	Tianyu Liu, Baobao Chang, Xu Sun, Lei Li, and	<a href="#">The power of scale for parameter-efficient prompt</a>	775
721	Zhifang Sui. 2024. <a href="#">A survey on in-context learning.</a>	<a href="#">tuning.</a> <i>Preprint</i> , arXiv:2104.08691.	776
722	<i>Preprint</i> , arXiv:2301.00234.		
723	Luis Antonio Gal��rraga, Christina Teflioudi, Katja Hose,	Percy Liang, Rishi Bommasani, Tony Lee, Dimitris	777
724	and Fabian Suchanek. 2013. <a href="#">Amie: association rule</a>	Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian	778
725	<a href="#">mining under incomplete evidence in ontological</a>	Zhang, Deepak Narayanan, Yuhuai Wu, Ananya	779
726	<a href="#">knowledge bases.</a> In <i>Proceedings of the 22nd Inter-</i>	Kumar, Benjamin Newman, Binhang Yuan, Bobby	780
727	<i>national Conference on World Wide Web, WWW</i>	Yan, Ce Zhang, Christian Cosgrove, Christopher D.	781
728	’13, page 413–422, New York, NY, USA. Association	Manning, Christopher R��, Diana Acosta-Navas,	782
729	for Computing Machinery.	Drew A. Hudson, Eric Zelikman, Esin Durmus,	783
		Faisal Ladhak, Frieda Rong, Hongyu Ren, Huaxiu	784
730	Yuxian Gu, Xu Han, Zhiyuan Liu, and Minlie Huang.	Yao, Jue Wang, Keshav Santhanam, Laurel Orr,	785
731	2022. <a href="#">Ppt: Pre-trained prompt tuning for few-shot</a>	Lucia Zheng, Mert Yuksekgonul, Mirac Suzgun,	786
732	<a href="#">learning.</a> <i>Preprint</i> , arXiv:2109.04332.	Nathan Kim, Neel Guha, Niladri Chatterji, Omar	787
		Khattab, Peter Henderson, Qian Huang, Ryan Chi,	788
733	Jiayi Gui, Yiming Liu, Jiale Cheng, Xiaotao Gu, Xiao	Sang Michael Xie, Shibani Santurkar, Surya Gan-	789
734	Liu, Hongning Wang, Yuxiao Dong, Jie Tang, and	guli, Tatsunori Hashimoto, Thomas Icard, Tianyi	790
735	Minlie Huang. 2024. <a href="#">Logicgame: Benchmarking</a>	Zhang, Vishrav Chaudhary, William Wang, Xuechen	791
736	<a href="#">rule-based reasoning abilities of large language mod-</a>	Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. 2023.	792
737	<a href="#">els.</a> <i>Preprint</i> , arXiv:2408.15778.	<a href="#">Holistic evaluation of language models.</a> <i>Preprint</i> ,	793
		arXiv:2211.09110.	794
738	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,		
739	Mantas Mazeika, Dawn Song, and Jacob Steinhardt.	Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xu-	795
740	2021. <a href="#">Measuring massive multitask language under-</a>	anyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding,	796
741	<a href="#">standing.</a> <i>Preprint</i> , arXiv:2009.03300.	Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang	797
		Deng, Aohan Zeng, Zhengxiao Du, Chenhui Zhang,	798
742	Zhen Huang, Zengzhi Wang, Shijie Xia, Xuefeng Li,	Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun,	799
743	Haoyang Zou, Ruijie Xu, Run-Ze Fan, Lyumanshan	Minlie Huang, Yuxiao Dong, and Jie Tang. 2023a.	800
744	Ye, Ethan Chern, Yixin Ye, Yikai Zhang, Yuqing	<a href="#">Agentbench: Evaluating llms as agents.</a> <i>Preprint</i> ,	801
745	Yang, Ting Wu, Binjie Wang, Shichao Sun, Yang	arXiv:2308.03688.	802



803	Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding,	rule following of large language models. <i>Preprint</i> ,	862
804	Yujie Qian, Zhilin Yang, and Jie Tang. 2023b. <a href="#">Gpt</a>	arXiv:2407.08440.	863
805	<a href="#">understands, too</a> . <i>Preprint</i> , arXiv:2103.10385.		
806	Yiheng Liu, Tianle Han, Siyuan Ma, Jiayue Zhang,	Mirac Suzgun, Nathan Scales, Nathanael Schärli, Se-	864
807	Yuanyuan Yang, Jiaming Tian, Hao He, Antong Li,	bastian Gehrmann, Yi Tay, Hyung Won Chung,	865
808	Mengshen He, Zhengliang Liu, Zihao Wu, Lin Zhao,	Aakanksha Chowdhery, Quoc V. Le, Ed H. Chi,	866
809	Dajiang Zhu, Xiang Li, Ning Qiang, Dingang Shen,	Denny Zhou, and Jason Wei. 2022. <a href="#">Challenging</a>	867
810	Tianming Liu, and Bao Ge. 2023c. <a href="#">Summary of</a>	<a href="#">big-bench tasks and whether chain-of-thought can</a>	868
811	<a href="#">chatgpt-related research and perspective towards the</a>	<a href="#">solve them</a> . <i>Preprint</i> , arXiv:2210.09261.	869
812	<a href="#">future of large language models</a> . <i>Meta-Radiology</i> ,		
813	1(2):100017.	Harsh Trivedi, Niranjana Balasubramanian, Tushar	870
814	OpenAI. 2023. <a href="#">Gpt-4 technical report</a> . <i>Preprint</i> ,	Khot, and Ashish Sabharwal. 2023. <a href="#">Interleav-</a>	871
815	arXiv:2303.08774.	<a href="#">ing retrieval with chain-of-thought reasoning for</a>	872
		<a href="#">knowledge-intensive multi-step questions</a> . <i>Preprint</i> ,	873
		arXiv:2212.10509.	874
816	Shuofei Qiao, Yixin Ou, Ningyu Zhang, Xiang Chen,	Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu,	875
817	Yunzhi Yao, Shumin Deng, Chuanqi Tan, Fei Huang,	Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim.	876
818	and Huajun Chen. 2023. <a href="#">Reasoning with lan-</a>	2023a. <a href="#">Plan-and-solve prompting: Improving zero-</a>	877
819	<a href="#">guage model prompting: A survey</a> . <i>Preprint</i> ,	<a href="#">shot chain-of-thought reasoning by large language</a>	878
820	arXiv:2212.09597.	<a href="#">models</a> . <i>Preprint</i> , arXiv:2305.04091.	879
821	Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie	Siyuan Wang, Zhongyu Wei, Yejin Choi, and Xiang Ren.	880
822	Millican, Jordan Hoffmann, Francis Song, John	2024. <a href="#">Can llms reason with rules? logic scaffold-</a>	881
823	Aslanides, Sarah Henderson, Roman Ring, Susan-	<a href="#">ing for stress-testing and improving llms</a> . <i>Preprint</i> ,	882
824	nah Young, Eliza Rutherford, Tom Hennigan, Ja-	arXiv:2402.11442.	883
825	cob Menick, Albin Cassirer, Richard Powell, George	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V.	884
826	van den Driessche, Lisa Anne Hendricks, Mari-	Le, Ed H. Chi, Sharan Narang, Aakanksha Chowd-	885
827	beth Rauh, Po-Sen Huang, Amelia Glaese, Jo-	hery, and Denny Zhou. 2023b. <a href="#">Self-consistency</a>	886
828	hannes Welbl, Sumanth Dathathri, Saffron Huang,	<a href="#">improves chain of thought reasoning in language</a>	887
829	Jonathan Uesato, John Mellor, Irina Higgins, Anto-	<a href="#">models</a> . In <i>The Eleventh International Conference</i>	888
830	nia Creswell, Nat McAleese, Amy Wu, Erich Elsen,	<a href="#">on Learning Representations, ICLR 2023, Kigali,</a>	889
831	Siddhant Jayakumar, Elena Buchatskaya, David Bud-	<a href="#">Rwanda, May 1-5, 2023</a> . OpenReview.net.	890
832	den, Esme Sutherland, Karen Simonyan, Michela Pa-	Yaqing Wang, Quanming Yao, James Kwok, and Li-	891
833	ganini, Laurent Sifre, Lena Martens, Xiang Lorraine	onel M. Ni. 2020. <a href="#">Generalizing from a few ex-</a>	892
834	Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena	<a href="#">amples: A survey on few-shot learning</a> . <i>Preprint</i> ,	893
835	Gribovskaya, Domenic Donato, Angeliki Lazaridou,	arXiv:1904.05046.	894
836	Arthur Mensch, Jean-Baptiste Lespiau, Maria Tsim-	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel,	895
837	poukelli, Nikolai Grigorev, Doug Fritz, Thibault Sot-	Barret Zoph, Sebastian Borgeaud, Dani Yogatama,	896
838	tiaux, Mantas Pajarskas, Toby Pohlen, Zhitao Gong,	Maarten Bosma, Denny Zhou, Donald Metzler, Ed H.	897
839	Daniel Toyama, Cyprien de Masson d’Autume, Yujia	Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy	898
840	Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin,	Liang, Jeff Dean, and William Fedus. 2022. <a href="#">Emer-</a>	899
841	Aidan Clark, Diego de Las Casas, Aurelia Guy,	<a href="#">gent abilities of large language models</a> . <i>Preprint</i> ,	900
842	Chris Jones, James Bradbury, Matthew Johnson,	arXiv:2206.07682.	901
843	Blake Hechtman, Laura Weidinger, Iason Gabriel,	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	902
844	William Isaac, Ed Lockhart, Simon Osindero, Laura	Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and	903
845	Rimell, Chris Dyer, Oriol Vinyals, Kareem Ayoub,	Denny Zhou. 2023. <a href="#">Chain-of-thought prompting elic-</a>	904
846	Jeff Stanway, Lorraine Bennett, Demis Hassabis, Ko-	<a href="#">its reasoning in large language models</a> . <i>Preprint</i> ,	905
847	ray Kavukcuoglu, and Geoffrey Irving. 2022. <a href="#">Scaling</a>	arXiv:2201.11903.	906
848	<a href="#">language models: Methods, analysis &amp; insights from</a>	xAI. 2024. <a href="#">xai grok</a> .	907
849	<a href="#">training gopher</a> . <i>Preprint</i> , arXiv:2112.11446.	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng,	908
850	Laria Reynolds and Kyle McDonell. 2021. <a href="#">Prompt</a>	Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan	909
851	<a href="#">programming for large language models: Beyond the</a>	Li, Dayiheng Liu, Fei Huang, Guanting Dong, Hao-	910
852	<a href="#">few-shot paradigm</a> . <i>Preprint</i> , arXiv:2102.07350.	ran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian	911
853	Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha,	Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin	912
854	Vinija Jain, Samrat Mondal, and Aman Chadha.	Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang	913
855	2024. <a href="#">A systematic survey of prompt engineering in</a>	Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang,	914
856	<a href="#">large language models: Techniques and applications</a> .	Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng	915
857	<i>Preprint</i> , arXiv:2402.07927.	Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin,	916
858	Wangtao Sun, Chenxiang Zhang, Xueyou Zhang, Xu-	Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu,	917
859	anqing Yu, Ziyang Huang, Pei Chen, Haotian Xu,		
860	Shizhu He, Jun Zhao, and Kang Liu. 2024. <a href="#">Be-</a>		
861	<a href="#">yond instruction following: Evaluating inferential</a>		



Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. 2024a. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.

Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V. Le, Denny Zhou, and Xinyun Chen. 2024b. [Large language models as optimizers](#). *Preprint*, arXiv:2309.03409.

Wenkai Yang, Yankai Lin, Jie Zhou, and Ji-Rong Wen. 2024c. [Distilling rule-based knowledge into large language models](#). *Preprint*, arXiv:2311.08883.

Zeyuan Yang, Peng Li, and Yang Liu. 2023. [Failures pave the way: Enhancing large language models through tuning-free rule accumulation](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1751–1777, Singapore. Association for Computational Linguistics.

Zhicheng Yang, Jinghui Qin, Jiaqi Chen, Liang Lin, and Xiaodan Liang. 2022. [Logicsolver: Towards interpretable math word problem solving with logical prompt-enhanced learning](#). *Preprint*, arXiv:2205.08232.

Zonglin Yang, Li Dong, Xinya Du, Hao Cheng, Erik Cambria, Xiaodong Liu, Jianfeng Gao, and Furu Wei. 2024d. [Language models as inductive reasoners](#). *Preprint*, arXiv:2212.10923.

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](#). *Preprint*, arXiv:2305.10601.

Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, and Ed Chi. 2023a. [Least-to-most prompting enables complex reasoning in large language models](#). *Preprint*, arXiv:2205.10625.

Pei Zhou, Jay Pujara, Xiang Ren, Xinyun Chen, Heng-Tze Cheng, Quoc V. Le, Ed H. Chi, Denny Zhou, Swaroop Mishra, and Huaixiu Steven Zheng. 2024. [Self-discover: Large language models self-compose reasoning structures](#). *Preprint*, arXiv:2402.03620.

Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. 2023b. [Large language models are human-level prompt engineers](#). *Preprint*, arXiv:2211.01910.

Zhaocheng Zhu, Yuan Xue, Xinyun Chen, Denny Zhou, Jian Tang, Dale Schuurmans, and Hanjun Dai. 2024. [Large language models can learn rules](#). *Preprint*, arXiv:2310.07064.



## A Dataset Description

In this section, we introduce the details of the datasets used.

We evaluated our method on the "Offensive" sub-task of **Tweeteval** (Barbieri et al., 2020), using the original train and test set splits from the dataset, and set the validation set to the first 200 samples of the training set.

For each task in **Big-Bench Hard** (Suzgun et al., 2022), we used the first half of the dataset as the training and validation set, and the second half as the test set. Detailed description of Big-Bench Hard tasks are shown in Table 7.

For the **GSM8K** (Cobbe et al., 2021) dataset, we use the first 4000 samples as the training set, samples 4001 to 4500 as the validation set, and the original test set as the test set.

**KQA Pro** is a dataset for complex question answering over knowledge base. The main task is to transfer natural language questions to a compositional and interpretable programming language KoPL. Table 5 Shows some examples from the KQA Pro dataset. In **KQA Pro** (Cao et al., 2022), we use the first 5000 samples as the training set, samples 5001 to 5500 as the validation set, and the first 5000 samples of the original test set as the test set. The specific dataset sizes are provided in Table 6.

## B Big-Bench Hard Per Task Performance

Per-task performance on Big-Bench Hard with rule set induced by GPT-4o and Grok-beta are shown in Table 8 and Table 9, respectively.

## C Prompt Design

During the process of summarizing the rules, we will interact with LLM multiple times, and the prompts for this part are shown in Table 12. {task\_description} refers to a simple description of each task provided in Table 10. {training\_examples} refers to all samples provided for this rule summarization process, concatenated as "Q:"+{question}+"A:"+{answer}, where {answer} includes the model's step-by-step reasoning process. {rules} refers to the rules summarized in the previous step. {error\_examples} refers to all samples provided for this rule summarization process, concatenated as "Q:"+{question}+"A:"+{answer}+{error\_reason}.

Dataset	Description
Tweeteval	classification task
Big-bench Hard	do a specific task
GSM8K	solve the following mathematical questions step by step
KQAPro	convert a question into a query program

Table 10: The description for all datasets.

When validating rules on a certain dataset, we designed corresponding prompts for each dataset as provided in Table 13. There are some options in the prompt template. When applying the rule set, we concatenate the rule into {rule}, and when we want the model to output a reasoning chain, add the phrase "think step by step." The examples provided for the model will be concatenated in format "Q:"+{question}+"A:"+{answer} and placed in {training\_examples}.

## D Cost of rule summarization

Table 14 shows the detailed usage information during the rule summarization process.

## E Error analysis in KQA Pro

Table 11 displays the error distribution when running corresponding KoPL program on the KQA Pro dataset before and after adding rules. Syntax errors include mismatched brackets, type errors refer to parameter mismatches.

Category	Few-shot	I2D
Correct	1970	2183
Wrong answer	714	633
Wrong Format (type error)	922	919
Wrong Format (syntax error)	521	472
Wrong Format (key error)	42	46
Wrong Format (others)	831	747

Table 11: Statistics on Results of Different Error Types



<b>Question</b>	Who received an Oscar nomination for Best Writing, Adapted Screenplay for A Simple Plan?
<b>Kopl Program</b>	QueryRelationQualifier(Find('A Simple Plan'),Find('Academy Award for Best Writing, Adapted Screenplay'),'nominated for', 'nominee')
<b>Question</b>	How many Pennsylvania counties have population of 4000000, or population above 60000?
<b>Kopl Program</b>	Count(Or(FilterConcept(FilterNum(FindAll()),'population', '4000000', '='),'county of Pennsylvania'),FilterConcept(FilterNum(FindAll()),'population', '60000', '>'),'county of Pennsylvania'))
<b>Question</b>	How many music genres were created after 1945?
<b>Kopl Program</b>	Count(FilterConcept(FilterYear(FindAll()),'inception', '1945', '>'),'music genre'))

Table 5: Examples of KQA Pro dataset.

Dataset	Task	Train	Dev	Test
Tweeteval	Offensive	11,916	200	860
	Hyperbaton	125	125	125
	Penguins In A Table	73	73	73
	Formal Fallacies	125	125	125
	Logical Deduction Five Objects	125	125	125
	Disambiguation QQ	125	125	125
	Object Counting	125	125	125
	Logical Deduction Seven Objects	125	125	125
	Tracking Shuffled Objects Three Objects	125	125	125
	Word Sorting	125	125	125
	Navigate	125	125	125
	Ruin Names	125	125	125
	Causal Judgement	93	93	94
	Temporal Sequences	125	125	125
	Salient Translation Error Detection	125	125	125
	Web Of Lies	125	125	125
	Boolean Expressions	125	125	125
	Dyck Languages	125	125	125
	Reasoning About Colored Objects	125	125	125
	Logical Deduction Three Objects	125	125	125
BBH	Multistep Arithmetic Two	125	125	125
	Movie Recommendation	125	125	125
	Tracking Shuffled Objects Five Objects	125	125	125
	Date Understanding	125	125	125
	Geometric Shapes	125	125	125
	Snarks	89	89	89
	Sports Understanding	125	125	125
	Tracking Shuffled Objects Seven Objects	125	125	125
	GSM8K	4,000	500	1,319
	KQAPro	5,000	500	5,000

Table 6: The statistics of the datasets.



---

<b>Boolean Expressions.</b>	Evaluate the truth value of a random Boolean expression consisting of Boolean constants (True, False) and basic Boolean operators (and, or, and not).
<b>Causal Judgment.</b>	Given a short story (involving moral, intentional, or counterfactual analysis), determine how a typical person would answer a causal question about the story.
<b>Date Understanding.</b>	Given a small set of sentences about a particular date, answer the provided question (e.g., “The concert was scheduled to be on 06/01/1943, but was delayed by one day to today. What is the date yesterday in MM/DD/YYYY?”).
<b>Disambiguation QA.</b>	Given a sentence with an “ambiguous” pronoun, either determine whether the sentence is inherently ambiguous (i.e., the thing that the pronoun refers to cannot be inferred by given information) or, if the pronoun can be implicitly deduced, state the antecedent of the pronoun (i.e., the noun to which the pronoun refers).
<b>Dyck Languages.</b>	Predict the sequence of the closing parentheses of a Dyck-4 word without its last few closing parentheses.
<b>Formal Fallacies Syllogisms Negation.</b>	Given a context involving a set of statements (generated by one of the argument schemes), determine whether an argument—presented informally—can be logically deduced from the provided context.
<b>Geometric Shapes.</b>	Given a full SVG path element containing multiple commands, determine the geometric shape that would be generated if one were to execute the full path element.
<b>Hyperbaton (Adjective Ordering).</b>	Given two English-language sentences, determine the one with the correct adjective order.
<b>Logical Deduction.</b>	Deduce the order of a sequence of objects based on the clues and information about their spacial relationships and placements.
<b>Movie Recommendation.</b>	Given a list of movies a user might have watched and liked, recommend a new, relevant movie to the user out of the four potential choices user might have.
<b>Multi-Step Arithmetic.</b>	Solve multi-step equations involving basic arithmetic operations (addition, subtraction, multiplication, and division).
<b>Navigate.</b>	Given a series of navigation steps to an agent, determine whether the agent would end up back at its initial starting point.
<b>Object Counting.</b>	Given a collection of possessions that a person has along with their quantities (e.g., three pianos, two strawberries, one table, and two watermelons), determine the number of a certain object/item class (e.g., fruits).
<b>Penguins in a Table.</b>	Given a unique table of penguins (and sometimes some new information), answer a question about the attributes of the penguins.
<b>Reasoning about Colored Objects.</b>	Given a context, answer a simple question about the color of an object on a surface.
<b>Ruin Names.</b>	Given an artist, band, or movie name, identify a one-character edit to the name that changes the meaning of the input and makes it humorous.
<b>Salient Translation Error Detection.</b>	Given a source sentence written in German and its translation in English, determine the type of translation error that the translated sentence contains.
<b>Snarks.</b>	Given two nearly-identical sentences, determine which one is sarcastic.
<b>Sports Understanding.</b>	Determine whether a factitious sentence related to sports is plausible.
<b>Temporal Sequences.</b>	Given a series of events and activities a person has completed in the course of a day, determine what time, during the day, they might have been free to perform another activity.
<b>Tracking Shuffled Objects.</b>	Given the initial positions of a set of objects and a series of transformations (namely, pairwise swaps) applied to them, determine the final positions of the objects.
<b>Web of Lies.</b>	Evaluate the truth value of a random Boolean function expressed as a natural-language word problem.
<b>Word Sorting.</b>	Given a list of words, sort them lexicographically.

---

Table 7: Detailed description of Big-Bench Hard tasks. For convenience, the original descriptions from [Suzgun et al. \(2022\)](#) is duplicated here



Task	Direct	Few-shot	CoT	I2D	I2D + few-shot	I2D + CoT
boolean_expressions	64	73.6	86.4	71.2	76	91.2
causal_judgement	54.26	56.38	57.45	54.26	50	56.38
date_understanding	27.2	36	64	19.2	40	72
disambiguation_qa	63.2	46.4	56	66.4	63.2	52
dyck_languages	47.2	44	45.6	44	46.4	36.8
formal_fallacies	92	100	73.6	100	99.2	99.2
geometric_shapes	24	26.4	38.4	23.2	40	57.6
hyperbaton	56.8	55.2	56	76.8	74.4	64.8
logical_deduction_five_objects	42.4	32	41.6	31.2	32.8	42.4
logical_deduction_seven_objects	32.8	23.2	24.8	40	39.2	28
logical_deduction_three_objects	54.4	52.8	70.4	60	58.4	68
movie_recommendation	72.8	75.2	56.8	72.8	73.6	52
multistep_arithmetic_two	7.2	14.4	58.4	42.4	60	68.8
navigate	67.2	66.4	81.6	72	77.6	70.4
object_counting	31.2	27.2	54.4	28	38.4	69.6
penguins_in_a_table	24.66	52.05	71.23	41.1	61.64	82.19
reasoning_about_colored_objects	38.4	41.6	78.4	41.6	52	77.6
ruin_names	68.8	60	65.6	87.2	88.8	84.8
salient_translation_error_detection	32.8	35.2	48.8	38.4	41.6	43.2
snarks	58.43	57.3	65.17	58.43	64.04	55.06
sports_understanding	65.6	73.6	80	54.4	80.8	93.6
temporal_sequences	46.4	39.2	80.8	92	94.4	93.6
tracking_shuffled_objects_five_objects	18.4	17.6	44	35.2	70.4	62.4
tracking_shuffled_objects_seven_objects	11.2	17.6	42.4	44	60.8	46.4
tracking_shuffled_objects_three_objects	38.4	40.8	62.4	53.6	56	62.4
web_of_lies	55.2	54.4	88	52.8	60.8	88.8
word_sorting	40.8	45.6	52.8	36	48	64

Table 8: Big Bench-Hard per-task performance with rule set induced by GPT-4o



Task	Direct	Few-shot	CoT	I2D	I2D + few-shot	I2D + CoT
boolean_expressions	64	73.6	86.4	58.4	72	88.8
causal_judgement	54.26	56.38	57.45	46.81	54.26	56.38
date_understanding	27.2	36	64	25.6	53.6	60.8
disambiguation_qa	63.2	46.4	56	53.6	56	36
dyck_languages	47.2	44	45.6	46.4	52	40.8
formal_fallacies	92	100	73.6	98.4	98.4	100
geometric_shapes	24	26.4	38.4	35.2	50.4	57.6
hyperbaton	56.8	55.2	56	78.4	78.4	66.4
logical_deduction_five_objects	42.4	32	41.6	43.2	43.2	48
logical_deduction_seven_objects	32.8	23.2	24.8	38.4	31.2	29.6
logical_deduction_three_objects	54.4	52.8	70.4	59.2	56.8	56
movie_recommendation	72.8	75.2	56.8	76.8	73.6	61.6
multistep_arithmetic_two	7.2	14.4	58.4	24	38.4	69.6
navigate	67.2	66.4	81.6	65.6	70.4	83.2
object_counting	31.2	27.2	54.4	32.8	32	64
penguins_in_a_table	24.66	52.05	71.23	20.55	52.05	78.08
reasoning_about_colored_objects	38.4	41.6	78.4	40.8	56.8	77.6
ruin_names	68.8	60	65.6	86.4	85.6	83.2
salient_translation_error_detection	32.8	35.2	48.8	43.2	47.2	44.8
snarks	58.43	57.3	65.17	53.93	64.04	60.67
sports_understanding	65.6	73.6	80	60	74.4	91.2
temporal_sequences	46.4	39.2	80.8	95.2	99.2	95.2
tracking_shuffled_objects_five_objects	18.4	17.6	44	42.4	52	63.2
tracking_shuffled_objects_seven_objects	11.2	17.6	42.4	24	44.8	36.8
tracking_shuffled_objects_three_objects	38.4	40.8	62.4	45.6	64.8	60.8
web_of_lies	55.2	54.4	88	47.2	59.2	78.4
word_sorting	40.8	45.6	52.8	41.6	49.6	22.4

Table 9: Big Bench-Hard per-task performance with rule set induced by Grok



---

**Path 1 Step 1 - Draft Rule**

The following data is about to {task\_description}. Please analyze the relationship between the problem and the corresponding answer, as well as the similarities between the answer. Then summarize {number\_of\_rules} rules you discovered. You can make the rules more precise and detailed, for example, by including an example.

Training Examples:

{training\_examples}

---

**Path 1 Step 2 - Rule Formatting**

The following data is about to {task\_description}. I will give you some rules and comments, please help me to find useful rules.

Training Examples:

{training\_examples}

Rules:

{rules}

ONLY output the rules or patterns, each rule or pattern must be written in one line, add a \* before each rule.

---

**Path 2 Step 1 - Error Reason**

The following data is about to {task\_description}. I will provide you with your last answer and the standard answer. If your answer does not match the standard answer, please analyze the reason for the error and how to fix these errors based on the question and the standard answer.

Question:{question}

Your last answer:{answer}

Standard answer:{standard\_answer}

---

**Path 2 Step 2 - Draft Rule**

The following data is about to {task\_description}. I will give you some error samples, please help me find the similarities between the errors, then formulate {number\_of\_rules} rules or patterns you find in errors to help perform this task next time. It should not be suggestions. Your rules need to be more specific rather than general. You can make the rules more precise and detailed, for example, by including an example.

Error Examples:

{error\_examples}

---

**Path 2 Step 3 - Rule Formatting**

Same as "Path 1 Step 2 - Rule Formatting".

---

**Consistency Check**

I will give you two rules. Please help me classify whether the contents of these two rules are describing the same thing.

You are only allowed to give me the answer, selecting from "same" and "not same"

Rule1:

{rule1}

Rule2:

{rule2}

---

Table 12: The prompt design for all tasks.



---

**Tweeteval**

Help me perform a {task\_description}. I will give you a review and you should help me by figuring whether this review is semantically offensive.

Rules:

{rules}

You are only allowed to give me the answer, selecting from "offensive" and "not offensive".

Review: {question}

Answer:

---

**Big-bench Hard**

Your have to follow below instruction to {task\_description}.

Training Examples:

{training\_examples}

Rules:

{rules}

- Don't forget to put your final answer after "The answer is:".

Q: {question}

A: (Let's think step by step.)

---

**GSM8K**

As an expert problem solver {task\_description}.

Training Examples:

{training\_examples}

Rules: {rules}

- Don't forget to put your final answer after "The answer is:".

Q: {question}

A: (Let's think step by step.)

---

**KQA Pro**

Your have to follow below instruction to {task\_description}.

- Use the training examples to understand the task.

- You can use the rules to help you generate the program.

(- Please think step by step.)

Training Examples:

{training\_examples}

Rules:

{rules}

Test Examples:

Question: {question}

- Output your final answer in the last line without quotation mark or prefix like "KoPL program".

Your answer:

---

Table 13: The evaluation prompt for all tasks.



Task	Tokens	API Calls	Rules
BBH_boolean_expressions	133,140	116	30
BBH_causal_judgement	352,464	154	7
BBH_date_understanding	344,771	156	21
BBH_disambiguation_qa	256,521	158	18
BBH_dyck_languages	275,346	154	11
BBH_formal_fallacies	235,303	88	16
BBH_geometric_shapes	399,712	158	23
BBH_hyperbaton	216,257	158	11
BBH_logical_deduction_five_objects	348,370	156	14
BBH_logical_deduction_seven_objects	386,725	144	21
BBH_logical_deduction_three_objects	213,325	144	12
BBH_movie_recommendation	314,784	160	7
BBH_multistep_arithmetic_two	208,835	134	7
BBH_navigate	243,360	154	9
BBH_object_counting	174,149	134	8
BBH_penguins_in_a_table	296,659	154	14
BBH_reasoning_about_colored_objects	212,080	142	14
BBH_ruin_names	211,984	134	4
BBH_salient_translation_error_detection	339,666	126	14
BBH_snarks	239,571	158	11
BBH_sports_understanding	172,665	130	9
BBH_temporal_sequences	320,183	126	6
BBH_tracking_shuffled_objects_five_objects	303,542	144	24
BBH_tracking_shuffled_objects_seven_objects	349,765	144	24
BBH_tracking_shuffled_objects_three_objects	247,217	154	32
BBH_web_of_lies	180,597	128	13
BBH_word_sorting	390,045	136	13
Average	272,853	142	14.6
GSM8K	278,978	182	12
KQA Pro	298,458	313	18
Tweeteval	150,365	145	50

Table 14: The token usage and number of API calls of various tasks.