

PATH-CONSISTENCY WITH PREFIX ENHANCEMENT FOR EFFICIENT INFERENCE IN LLMs

Anonymous authors

Paper under double-blind review

ABSTRACT

To enhance the reasoning capabilities of large language models (LLMs), self-consistency has become a popular approach, combining multiple samplings with majority voting. However, current methods are computationally expensive and time-consuming due to the need for numerous samplings. To address this, this paper introduces *path-consistency*, which leverages the confidence of earlier-generated answers to identify the most promising prefix and guide the generation of subsequent branches. By dynamically guiding the generation of subsequent branches based on this prefix, *path-consistency* mitigates both the errors and redundancies from random or less useful sampling in self-consistency. This approach reduces errors and redundancies from random sampling, significantly accelerating inference by minimizing token consumption. Our extensive empirical results demonstrate that *path-consistency* improves inference latency by up to 40.5%, while maintaining task accuracy across various tasks, including mathematical reasoning, commonsense reasoning, and symbolic reasoning.

1 INTRODUCTION

The range of tasks that large language models (LLMs) can accomplish is continuously expanding, as the scale and complexity of models continue to grow recently. However, this advancement has not yet endowed LLMs with sufficiently robust reasoning capabilities Rae et al. (2021); Wu et al. (2016); Guo et al. (2018); Chen et al. (2022). To address this shortcoming and further extend the application scope of LLMs, *Chain-of-Thought* (CoT) Wei et al. (2022) prompting has emerged in response. CoT prompting uses reasoning problems and processes as input prompts to guide language models in generating reasoning paths and final answers, aiming to mimic the thought processes humans might use when solving mathematical or logical problems. This enables LLMs to be applied in an increasing number of specific scenarios, such as mathematical reasoning Cobbe et al. (2021); Miao et al. (2020) and logical reasoning Geva et al. (2021). To better accomplish these complex tasks, there are several CoT-based optimization methods Li et al. (2022); Chen et al. (2022). One of the common and effective methods is *self-consistency* Wang et al. (2023), a technique involving multiple sampling and majority voting. In this approach, the model generates multiple reasoning paths for a given input, with the final decision based on the most frequently occurring output among the samples.

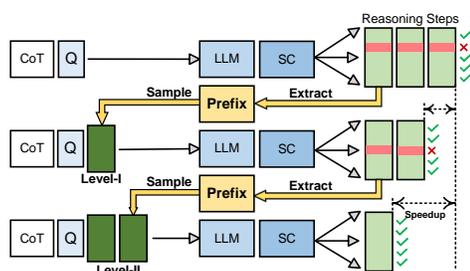


Figure 1: Path-consistency extracts prefixes from earlier generated inference paths to guide the inference of subsequent branches.

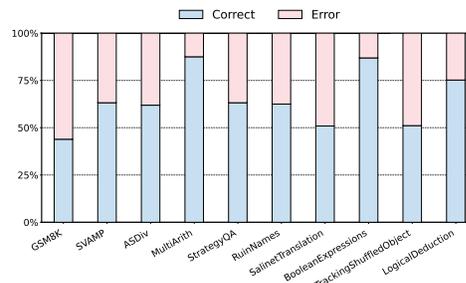


Figure 2: The proportion of tokens generated by self-consistency on correct or incorrect inference paths.

054 Self-consistency significantly enhances the reasoning capabilities of LLMs by sampling a large num-
055 ber of examples during inference. However, it has limitations in practical applications. The basic
056 self-consistency technique frequently invokes the model to generate numerous reasoning paths when
057 solving specific problems Wang et al. (2023). As model size and task complexity increase, the time
058 and computation cost of self-consistency increase sharply, making it a critical issue in practical ap-
059 plications. Motivated by this challenge, this paper proposes *path-consistency*, which is designed to
060 reduce the computation and time cost by leveraging intermediate information from earlier generation
061 results to assist in subsequent generations.

062 As illustrated in fig. 1, the proposed dynamic inference method, path-consistency, continuously
063 extracts appropriate reasoning steps from already generated reasoning paths to serve as “prefixes”.
064 These prefixes are then used to guide and accelerate the generation of subsequent reasoning steps.

065 The proposed *path-consistency* offers several key advantages: (1) It accelerates inference speed and
066 reduces token consumption, while significantly preserving or even improving task accuracy. (2)
067 The method requires no additional computation, fine-tuning, or training, ensuring the generation
068 quality of the model remains intact. (3) It is model-agnostic, making it easy to deploy and apply to
069 various models and tasks in practical scenarios. Furthermore, it integrates seamlessly with existing
070 optimization methods, achieving even better acceleration performance.

071 We evaluated path-consistency using ten different datasets, resulting in an average acceleration of
072 28.7% for mathematical reasoning tasks, 20.9% for commonsense reasoning, and 20.3% for sym-
073 bolic reasoning. Moreover, by continuously extracting potentially correct prefixes, this method has
074 minimal impact on task accuracy and can even enhance it.

076 2 BACKGROUND

077 2.1 REASONING WITH SELF-CONSISTENCY

079 Self-consistency Wang et al. (2023) is a sampling strategy distinct from greedy decoding Vaswani
080 et al. (2017), significantly enhancing the reasoning performance of language models. This approach
081 generates multiple reasoning paths and aggregates the final output through majority voting.

082 Earlier research primarily focused on optimizing reasoning for individual tasks Andor et al. (2019);
083 Ran et al. (2019); Geva et al. (2020); Piekos et al. (2021). More recently, numerous inference strate-
084 gies based on self-consistency have been proposed, particularly using self-evaluation to calibrate
085 LLMs Zhang et al. (2023); Shinn et al. (2023); Madaan et al. (2024); Paul et al. (2024). For exam-
086 ple, self-evaluation guided by random beam search Xie et al. (2024) leverages additional language
087 models to search for the optimal inference path. Other methods, such as Deductive Beam Search
088 Zhu et al. (2024), optimize inference by emphasizing the relationships between each reasoning step.
089 However, these approaches often require increased computational resources and sometimes involve
090 extensive model calls to improve task accuracy in LLMs. This results in additional time costs,
091 always making them inefficient in practical applications and contrary to our original intention.

092 2.2 APPROACHES FOR EFFICIENT REASONING

094 Recent works have attempted to achieve efficient reasoning by using smaller models or increasing
095 inference speed. Distillation Hinton et al. (2015), as an effective model compression technique,
096 has been employed to create smaller models for reasoning tasks Fu et al. (2023); Ho et al. (2023);
097 Magister et al. (2023); Li et al. (2023). Additionally, various approaches have sought to improve
098 inference speed by altering inference strategies, thereby avoiding modifications to the model archi-
099 tecture and retraining Aggarwal et al. (2023); Li et al. (2024). However, these methods often com-
100 promise the diversity inherent in self-consistency, which can negatively impact the quality of the
101 generated outputs. The proposed path-consistency is a model-agnostic method applicable to most
102 common models, including compressed models, to enhance inference efficiency while maintaining
103 task accuracy.

104 3 MOTIVATION

106 Self-consistency in large language models (LLMs) involves generating multiple branches of rea-
107 soning, each potentially leading to different answers. The final answer is determined through ag-

108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161

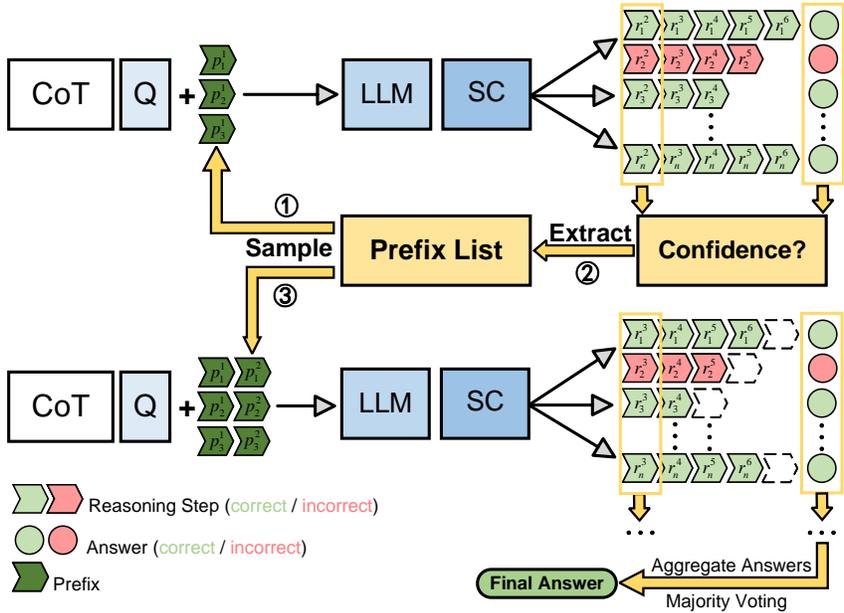


Figure 3: An “extract-and-sample” inference process of the proposed path-consistency. It seeks the “optimal path” in the form of the “prefix”, thereby progressively reducing the number of generated tokens and significantly shortening inference latency.

gregation methods, such as majority voting. Unlike greedy decoding, self-consistency avoids the shortcomings of local optimality and reduces the randomness of single-step sampling. However, the primary drawback of self-consistency is the significant computational redundancy. With N branches, each answer is derived from N similar but independent inference processes, leading to an N -fold increase in computational cost compared to greedy decoding for a single problem. To improve self-consistency, the goal is to achieve similar sampling effects while reducing the time cost of redundant computations.

We propose an intuitive hypothesis: for example, a particular mathematical problem might have five different reasoning paths, p_1 to p_5 . If the model frequently errs on paths p_1 to p_3 while p_4 and p_5 are relatively simpler, then full self-consistency wastes significant computational resources on the problematic paths. As shown in fig. 2, statistics on the number of tokens generated during self-consistency across various datasets reveal that over 25% and sometimes even 50% of tokens are wasted on incorrect branches. By sampling multiple times only on p_4 and p_5 , we could enhance resource utilization and improve output accuracy. Furthermore, storing limited information from paths p_4 and p_5 to guide subsequent branch generation could significantly accelerate inference speed and efficiency.

Additionally, self-consistency involves extensive redundant processes without yielding intermediate results, with the final answer only emerging at the end. If useful information could be identified early in the generation process to guide subsequent branching, outcomes might improve. Intuitively, when tackling complex problems, using simple criteria to preliminarily assess the quality or correctness of the current generation during intermediate stages can enhance the effectiveness of subsequent steps. Our method aims to reduce the time wasted on incorrect branches while increasing the efficiency of generating correct inference paths.

4 METHODOLOGY

4.1 PATH-CONSISTENCY

Based on the internal mechanisms of self-consistency, we propose an automated dynamic reasoning approach *path-consistency* that continually seeks the “optimal path” in the form of “prefix”. This will progressively reduce the number of generated tokens and significantly shorten inference latency. The methodology is shown in fig. 3 and table 1, and can be described as the following “extract-and-sample” process:

- Determine the maximum number of branches and the highest prefix level for the specific task, and use these to divide all branches into multiple windows of equal length. Begin by generating a small number of branches for the first window.
- Then, use a confidence assessment metric¹ to assess the confidence of the answers generated in the current window. If the confidence exceeds the set threshold, extract shorter prefixes from these optimal paths (e.g., the first step of the current optimal reasoning paths) to guide subsequent generation; if the confidence is low, continue generating branches for the subsequent windows at the current prefix level to ensure a safer prefix selection.
- Randomly sample from the extracted prefixes as part of the prompt to guide subsequent generation. After the branches for the next window are generated, continue using the confidence assessment metric to find the optimal inference path and extract the prefix for the next level (e.g., the first two steps of the current optimal reasoning paths).
- Repeat the above steps iteratively, extending the prefix length until the optimal reasoning path is identified and aggregated to produce the final answer.

Table 1 illustrates the reasoning process of path-consistency with an example. We have taken a reasoning path from different stages of the branching process as samples in table 1. In the initial stage, the model generates the full reasoning paths and answers normally. When a particular answer appears frequently and has a higher confidence level compared to other answers within the current branch, the first step of the corresponding reasoning path is extracted and used as the input prompt for the next stage. As shown in table 1, as the number of branches increases, the prefix length grows while the newly generated portion of the path subsequently shortens.

Table 1: The generated parts are indicated in blue in the table, with “Level- X prefix” denoting the action of extracting the prefix for the X^{th} time. After three prefix extractions, the model only needs to generate the final sentence of the conclusion based on the prefix.

CoT-prompting:

Q: A robe takes 2 bolts of blue fiber and half that much white fiber. How many bolts in total does it take?

A: To make a robe, we need 2 bolts of blue fiber and half that much white fiber. We can write this as $2 + 1/2 * 2$. Now $1/2 * 2 = 1$, so the equation becomes $2 + 1 = 3$. The answer is 3.

Level-I prefix:

Q: A robe takes 2 bolts of blue fiber and half that much white fiber. How many bolts in total does it take?

A: To make a robe, we need 2 bolts of blue fiber and half that much white fiber. $2 * 1/2 = 1$ bolt of white fiber. $2 + 1 = 3$ bolts in total. The answer is 3.

Level-II prefix:

Q: A robe takes 2 bolts of blue fiber and half that much white fiber. How many bolts in total does it take?

A: To make a robe, we need 2 bolts of blue fiber and half that much white fiber. $2 * 1/2 = 1$ bolt of white fiber. $1 + 2 = 3$ bolts in total. The answer is 3.

Level-III prefix:

Q: A robe takes 2 bolts of blue fiber and half that much white fiber. How many bolts in total does it take?

A: To make a robe, we need 2 bolts of blue fiber and half that much white fiber. $2 * 1/2 = 1$ bolt of white fiber. $1 + 2 = 3$ bolts in total. The answer is 3.

When solving complex problems, LLM takes a question, represented as q , and generates a distribution of answers, denoted as $P(a | q)$, after producing a reasoning path. During inference, the model generates a multi-step reasoning path $R = [r^1, r^2, \dots, r^T] = r^{1:T}$ under the guidance of CoT. This process can also be expressed as

$$P(a | q) = \mathbb{E}_{R \sim P(R|q)} P(a | q, R). \quad (1)$$

In basic self-consistency, the generation process is mechanically repeated numerous times without variation, whereas path-consistency attempts to gather advantageous intermediate information from the early stages of the generation process for subsequent generations. Path-consistency gradually identifies sufficiently confident partial reasoning steps $R_{\text{prefix}} = [r^1, r^2, \dots, r^t]$ and incorporates

¹The proposed path-consistency allows a flexible choice of confidence assessment metrics (i.e., stopping criteria), based on intended objective and requirements. The confidence assessment metric used in the evaluation is the beta confidence criteria Aggarwal et al. (2023).

216 them as part of the input, thereby better guiding the subsequent reasoning steps and the generation
 217 of the answer a . Thus, this process can be expressed as

$$218 P(a | q, R_{\text{prefix}}) = \mathbb{E}_{R \sim P(R|q, R_{\text{prefix}})} P(a | q, R). \quad (2)$$

221 Combining the above expressions, we obtain the following formula via the law of total probability:
 222

$$223 P(a | q) = \sum_{R_{\text{prefix}}} P(R_{\text{prefix}} | q) P(a | q, R_{\text{prefix}}). \quad (3)$$

226 Path-consistency utilizes basic majority voting or other lightweight confidence metrics to model the
 227 distribution $P(R_{\text{prefix}} | q)$, facilitating the transformation from $P(a | q)$ to $P(a | q, R_{\text{prefix}})$. This
 228 approach enhances efficiency while maintaining the model’s performance across various tasks.
 229

230 4.2 PROBLEMS WITH “TRUTH IS IN HANDS OF A FEW” 231

232 In the basic self-consistency method, when generating answers for particularly challenging ques-
 233 tions, the problems with “Truth Is in the Hands of a Few” may still occur despite the repeated gen-
 234 eration of multiple branches. This means that the correct answer may not be the most frequently oc-
 235 ccurring one, leading to an incorrect final answer in majority voting. In the proposed path-consistency
 236 method, during the continuous exploration for the optimal path, if the problems with “Truth Is in the
 237 Hands of a Few” are encountered at a specific prefix selection, there’s a concern that this undesirable
 238 phenomenon may be exacerbated in subsequent branches. We will use the following to demonstrate
 239 that the proposed approach does not worsen this problem.

240 To present the analysis, we make the following assumptions: the probability of generating the correct
 241 answer is p_0 , and apart from the correct answer, the model generates only one unique incorrect
 242 answer. The total number of branches is set to N , and prefix selection based on majority voting
 243 is performed only once at $\frac{N}{2}$. If the correct answer is the majority at the time of prefix selection,
 244 the correct prefix will be selected to guide subsequent generation, increasing the probability of
 245 generating the correct answer in the remaining branches to $p_1 = p_0 + \Delta p$; Conversely, if the incorrect
 246 answer is the majority at this point, the probability of generating the correct answer in subsequent
 247 branches decreases to $p_2 = p_0 - \Delta p$. Using the binomial distribution formula, the probability of the
 248 correct answer being in the majority during the vote is given by:

$$249 P_{\text{vote}} = \sum_{k=\lceil \frac{N}{4} \rceil}^{\frac{N}{2}} \binom{\frac{N}{2}}{k} p_0^k (1 - p_0)^{\frac{N}{2} - k}. \quad (4)$$

254 After prefix selection, the probability of obtaining the correct answer in the subsequent $\frac{N}{2}$ branches
 255 increases or decreases to

$$\begin{aligned} 256 P_{\text{mc}} &= \sum_{k=\lceil \frac{N}{4} \rceil}^{\frac{N}{2}} \binom{\frac{N}{2}}{k} p_1^k (1 - p_1)^{\frac{N}{2} - k} \\ 257 &= \sum_{k=\lceil \frac{N}{4} \rceil}^{\frac{N}{2}} \binom{\frac{N}{2}}{k} (p_0 + \Delta p)^k (1 - p_0 - \Delta p)^{\frac{N}{2} - k} \\ 258 &= \sum_{k=\lceil \frac{N}{4} \rceil}^{\frac{N}{2}} \binom{\frac{N}{2}}{k} p_0^k \left(1 + \frac{\Delta p}{p_0}\right)^k (1 - p_0)^{\frac{N}{2} - k} \left(1 - \frac{\Delta p}{1 - p_0}\right)^{\frac{N}{2} - k} \\ 259 &\approx \sum_{k=\lceil \frac{N}{4} \rceil}^{\frac{N}{2}} \binom{\frac{N}{2}}{k} p_0^k (1 - p_0)^{\frac{N}{2} - k} \left(1 + \frac{k}{p_0} \Delta p\right) \left(1 - \frac{\frac{N}{2} - k}{1 - p_0} \Delta p\right) \\ 260 &\approx \sum_{k=\lceil \frac{N}{4} \rceil}^{\frac{N}{2}} \binom{\frac{N}{2}}{k} p_0^k (1 - p_0)^{\frac{N}{2} - k} \left(1 + \frac{k - \frac{N}{2} p_0}{p_0(1 - p_0)} \Delta p\right), \end{aligned} \quad (5)$$

and

$$\begin{aligned}
 P_{\text{dec}} &= \sum_{k=\lceil \frac{N}{4} \rceil}^{\frac{N}{2}} \binom{\frac{N}{2}}{k} p_2^k (1-p_2)^{\frac{N}{2}-k} \\
 &= \sum_{k=\lceil \frac{N}{4} \rceil}^{\frac{N}{2}} \binom{\frac{N}{2}}{k} (p_0 - \Delta p)^k (1-p_0 + \Delta p)^{\frac{N}{2}-k} \\
 &\approx \sum_{k=\lceil \frac{N}{4} \rceil}^{\frac{N}{2}} \binom{\frac{N}{2}}{k} p_0^k (1-p_0)^{\frac{N}{2}-k} \left(1 - \frac{k - \frac{N}{2} p_0}{p_0(1-p_0)} \Delta p \right).
 \end{aligned} \tag{6}$$

The two equations above use the first-order Taylor expansion, neglecting higher-order terms. The relationship between the two values is

$$P_{\text{inc}} + P_{\text{dec}} = 2 \cdot P_{\text{vote}}. \tag{7}$$

Meanwhile, the final probability of obtaining the correct answer after prefix selection is

$$P'_{\text{correct}} = P_{\text{vote}} \cdot P_{\text{inc}} + (1 - P_{\text{vote}}) \cdot P_{\text{dec}}. \tag{8}$$

If no prefix selection is performed, the probability of generating the correct answer in subsequent steps remains unchanged and is equal to the probability at $\frac{N}{2}$ with majority voting: $P_{\text{correct}} = P_{\text{vote}}$. To ensure that accuracy is not adversely affected, we require $P'_{\text{correct}} \geq P_{\text{correct}}$, which can be simplified to:

$$P_{\text{vote}} \geq \frac{P_{\text{vote}} - P_{\text{dec}}}{P_{\text{inc}} - P_{\text{dec}}}. \tag{9}$$

By combining eq. (7), we can obtain $P_{\text{vote}} \geq 0.5$. This proves that if the model’s initial performance ensures $P_{\text{vote}} \geq 0.5$, path-consistency will be sufficiently reliable and will not harm accuracy.

During prefix selection, employing confidence-based criteria can make this process more reliable. For instance, using the beta confidence criteria Aggarwal et al. (2023)

$$\int_0^{0.5} p^{\frac{N}{2}-v_m} \cdot (1-p)^{v_m} dp, \tag{10}$$

where v_m represents the number of major elements. A prefix is considered sufficiently reliable for selection only if this value exceeds a predefined confidence threshold $C_{\text{threshold}}$.

While this calculation is a rough estimate and the actual scenario is likely more complex, it provides insight into the effect of path-consistency on problems where “Truth Is in the Hands of a Few”. Essentially, path-consistency tends to enhance the accuracy of self-consistency: when self-consistency performs well for a particular input, path-consistency may perform even better; when self-consistency performs poorly, the accuracy of path-consistency might decrease, but not significantly. Moreover, prefix selection guided by a confidence threshold provides a safer alternative to direct majority-vote selection.

5 EVALUATION

5.1 EXPERIMENTAL SETTING

Benchmarks: We evaluated the performance and efficiency improvement of path-consistency across the following types of tasks: (1) Arithmetic Reasoning, including GSM8K Cobbe et al. (2021), SVAMP Patel et al. (2021), ASDiv Miao et al. (2020), and MultiArith Roy & Roth (2015); (2) Commonsense Reasoning, including StrategyQA Geva et al. (2021), Ruin Names, and Salient Translation; (3) Symbolic Reasoning, including Boolean Expressions, Tracking Shuffled Objects, as well as Logical Deduction Srivastava et al. (2023).

Model: We used the open-source model Llama3-8B Dubey et al. (2024) as our backbone model. During input, we employed prompts similar to those used in CoT Wei et al. (2022).

Hyperparameters: For LLM inference with self-consistency and path-consistency, we employed nucleus sampling with a temperature of 0.6 and top-p of 0.9, generating 20 paths for each example.

According to Wang et al. (2023), 20 sampled reasoning paths almost achieve convergence in task accuracy, with only marginal efficiency gains observed when increasing to 40 paths. Therefore, 20 paths are more suitable for evaluating the optimization capability of path-consistency in terms of efficiency. Meanwhile, we used the beta confidence criteria in eq. (10) for confidence calculation in path-consistency Aggarwal et al. (2023), with threshold values ranging from 0.5 to 1.0.

Metrics: We employed multiple metrics to comprehensively compare the performance of path-consistency against baselines, including reasoning accuracy, inference latency, and token consumption during reasoning. Since inference latency is highly dependent on the operating environment and hardware configuration, considering the absolute value of inference latency is not meaningful. Instead, we calculated speedup under identical conditions as a standardized evaluation metric. All experiments were conducted on a single NVIDIA GeForce RTX 3090 GPU.

5.2 EVALUATION OF PATH-CONSISTENCY

5.2.1 A CASE STUDY ON GSM8K

Table 2: Comparison of inference latency speedup and average token consumption reduction under different prefix levels, demonstrating the effect of path-consistency on GSM8K.

Method	Accuracy (%)	Speedup (%)				Tokens (#)	Decrease (%)			
		Total	Level-1	Level-2	Level-3		Total	Level-1	Level-2	Level-3
SC	64.1				96.60					
PC (C=0)	66.6 (+2.5)	+28.9	+18.8	+42.6	+78.9	73.37	-24.0	-16.9	-32.7	-47.4
PC (C=0.7)	67.1 (+3.0)	+18.0	+10.0	+25.3	+46.5	81.09	-16.1	-9.4	-21.3	-33.8
PC (C=0.8)	67.8 (+3.7)	+17.4	+9.8	+25.0	+43.9	80.91	-16.2	-9.8	-21.8	-33.3
PC (C=0.9)	66.6 (+2.5)	+10.2	+4.2	+13.3	+27.0	87.01	-9.9	-4.2	-12.7	-23.0

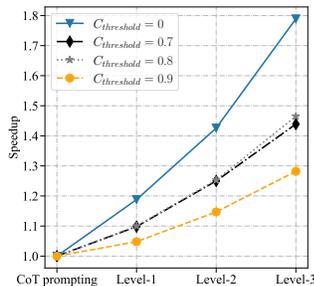


Figure 4: Speedup of inference.

Table 2 presents the performance and efficiency of path-consistency on the GSM8K dataset across different confidence thresholds and prefix levels. The experiments compare path-consistency with the basic self-consistency method in terms of task accuracy, inference latency, and token consumption. The results show that, across all confidence criteria, path-consistency not only enhances task accuracy but also provides a maximum acceleration of 28.9%. Observing the acceleration process of path-consistency, it is evident that as the prefix level increases, the speed improvement becomes more significant, reaching up to 78.9% at level-3 while maintaining at least 27.0%, as shown in fig. 4.

5.2.2 PERFORMANCE ON DIFFERENT TASKS

Arithmetic Reasoning: Table 3, rows 1-4, shows the task performance on arithmetic reasoning datasets, along with three additional datasets. Path-consistency ensures task performance that is almost comparable to or even better than the baseline under various confidence settings. On GSM8K, SVAMP, ASDiv, and MultiArith, accuracy improvements of up to 3.8%, 0.6%, 0.2%, and 2.2% are

Table 3: The evaluation results of path-consistency across various tasks, including (i) Arithmetic reasoning: GSM8K, SVAMP, ASDiv, MultiArith (MA); (ii) Commonsense reasoning: StrategyQA (SQA), RuinNames (RN), SalientTranslation (ST); (iii) Symbolic reasoning: BooleanExpression (BE), TrackingShuffledObjects (TSO), LogicalDeduction (LD). C=0 indicates that no confidence threshold is applied; the prefix corresponding to the majority element is extracted directly.

Dataset	Metrics (%)	SC	PC (C=0)	PC (C=0.7)	PC (C=0.8)	PC (C=0.9)	Dataset	Metrics (%)	SC	PC (C=0)	PC (C=0.7)	PC (C=0.8)	PC (C=0.9)
GSM8K	Acc.	64.1	66.6 (+2.6)	67.1 (+3.1)	67.8 (+3.8)	66.6 (+2.6)	RN	Acc.	76.0	71.2 (-4.8)	74.4 (-1.6)	72 (-4.0)	71.6 (-4.4)
	Speedup	-	+28.9	+18.0	+17.4	+10.2		Speedup	-	+10.4	+7.8	+7.4	+4.7
	Decrease	-	-24.0	-16.1	-16.2	-9.9		Decrease	-	-9.8	-7.2	-6.6	-4.7
SVAMP	Acc.	79.3	78.8 (-0.5)	78.5 (-0.8)	79.2 (-0.1)	79.9 (+0.6)	ST	Acc.	54.8	58.4 (+3.6)	57.2 (+2.4)	56.4 (+1.6)	55.2 (+0.4)
	Speedup	-	+48.3	+36.6	+35.7	+24.8		Speedup	-	+30.5	+19.6	+17.3	+10.2
	Decrease	-	-36.6	-29.4	-29.7	-22.3		Decrease	-	-25.8	-18.3	-16.4	-10.6
ASDiv	Acc.	81.0	80.1 (-0.9)	81.2 (+0.2)	80 (-1.0)	80.6 (-0.4)	BE	Acc.	88.4	90 (+1.6)	89.6 (+1.2)	88.8 (+0.4)	89.6 (+1.2)
	Speedup	-	+42.0	+32.1	+31.1	+24.9		Speedup	-	+26.6	+24.3	+25.2	+22.0
	Decrease	-	-32.7	-27.2	-27.1	-22.6		Decrease	-	-21.3	-19.8	-20.3	-18.3
MA	Acc.	96.7	97.8 (+1.1)	97.2 (+0.5)	98.9 (+2.2)	98.3 (+1.6)	TSO	Acc.	54.0	58 (+4.0)	56 (+2.0)	59.2 (+5.2)	58.4 (+4.4)
	Speedup	-	+42.7	+40.8	+40.5	+33.8		Speedup	-	+18.4	+12.9	+11.3	+6.8
	Decrease	-	-34.8	-33.3	-33.1	-27.2		Decrease	-	-16.7	-12.4	-10.8	-6.8
SQA	Acc.	71.4	70.8 (-0.6)	70.6 (-0.8)	71.3 (-0.1)	71.1 (-0.3)	LD	Acc.	78.8	76.8 (-2.0)	78.8 (+0.0)	77.2 (-1.6)	76 (-2.8)
	Speedup	-	+34.2	+29.4	+24.3	+16.4		Speedup	-	+31.6	+23.1	+30.1	+19.1
	Decrease	-	-27.8	-24.0	-22.1	-16.0		Decrease	-	-26.1	-20.5	-24.8	-17.5

observed, respectively. Due to the varying difficulty levels of the datasets for the LLM, the optimal confidence threshold differs for each dataset. For instance, higher or lower confidence thresholds may enhance performance on SVAMP and ASDiv.

The results also highlight the efficiency improvements on arithmetic reasoning datasets. Path-consistency achieves acceleration rates of up to 28.9%, 48.3%, 42.0%, and 42.7% on GSM8K, SVAMP, ASDiv, and MultiArith, respectively, along with a corresponding decrease in token consumption. Comparing these results reveals that for a specific dataset, there exists a proper confidence threshold at which path-consistency simultaneously enhances both task performance and efficiency.

Commonsense Reasoning: Table 3, rows 5-7, reports the performance of path-consistency on commonsense reasoning datasets. In terms of task accuracy, path-consistency slightly underperforms compared to the baseline on the StrategyQA and Ruins Names tasks. However, this decline can be mitigated by properly adjusting the confidence threshold. On the Salient Translation task, path-consistency still maintains better task performance compared to the baseline. Regarding acceleration, path-consistency provides up to 34.2%, 10.4%, and 30.5% speed improvements on StrategyQA, Ruins Names, and Salient Translation datasets, respectively.

According to the results, it can be found that the performance of path-consistency on the Ruins Names task is less pronounced compared to other datasets. This is because the Ruins Names task is more difficult, with longer input content and longer required prompts, making the benefits of the prefixing behavior less noticeable.

Symbolic Reasoning: Table 3, rows 8-10, compares the performance of different methods on the symbolic reasoning datasets. Path-consistency performs exceptionally well in terms of task accuracy, especially in the Tracking Shuffled Objects task, where it achieves the highest improvement of 5.2% among all datasets. Additionally, compared to the baseline, it delivers approximately a 20% speedup in inference latency across tasks such as Boolean Expressions, Tracking Shuffled Objects, and Logical Deduction.

In general, lower confidence thresholds indicate a more aggressive prefix selection strategy, often resulting in more significant efficiency improvements. However, in the Logical Deduction task, the efficiency improvement is actually higher at a confidence threshold of 0.8 compared to 0.7. This is because a more aggressive prefix selection strategy is more likely to choose suboptimal prefixes, which makes the subsequent branching generation less effective. Specifically, we found that with a confidence threshold of 0.7, the efficiency improvement at prefix level-2 is more significant than at prefix level-3, as prefix level-3 generates many erroneous paths, leading to this phenomenon.

5.2.3 RESULT ANALYSIS

Reducing Error and Redundancy: In section 3, we characterized the challenges of self-consistency. On one hand, it wastes computational resources on incorrect branches. On the other hand, it repeats the same computations without obtaining useful intermediate information. Figure 5 shows the changes in the proportion of tokens generated by path-consistency on correct or incorrect reasoning paths. The results show a decrease in the number of tokens wasted on incorrect branches. Additionally, for correct branches, there is a significant reduction in redundant tokens, which greatly improves efficiency while maintaining accuracy.

Hyperparameters Analysis: We explored the impact of prefix extraction frequency on path-consistency. With the setting of 20 branches, if the highest prefix level is set to level-3, the pre-

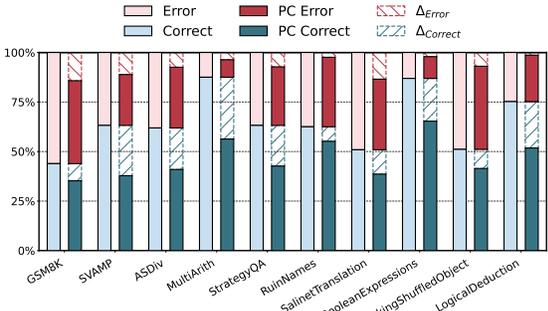


Figure 5: The change in the proportion of tokens generated by path-consistency on correct or incorrect paths.

Table 4: Impacts of maximum path levels of path-consistency on GSM8K.

Method	Level	Accuracy (%)	Speedup (%)	Decrease (%)
PC (C=0.7)	3	67.1	+18.0	-16.1
	4	66.3	+20.5	-20.0
PC (C=0.8)	3	67.8	+17.4	-16.2
	4	67.8	+18.0	-17.8

fix is extracted every 5 branches; if the highest prefix level is level-4, the prefix is extracted every 4 branches. Table 4 indicates that maintaining the highest prefix level at level-3 preserves higher task accuracy. Increasing the highest level enables a more aggressive strategy, leading to a more noticeable acceleration.

5.3 COMPARISON WITH SIMILAR WORK

We compared path-consistency (PC) to the following methods: (1) *Self-consistency* (SC) Wang et al. (2023) is the baseline that combines multiple samplings with majority voting. (2) *Adaptive-consistency* (AC) Aggarwal et al. (2023) introduces the concept of confidence for the first time. If a highly confident answer is identified, the sampling process is terminated, and the answer is selected as the final output, preventing any further branch sampling and inference. (3) *Early-stopping self-consistency* (ESC) Li et al. (2024) divides branches into equally sized windows. When all the answers generated within a given window are identical, then exit directly.

As shown in the experimental data in table 5, the path-consistency achieves task accuracy comparable to other optimization methods while consuming fewer tokens. The diversity of the reasoning paths is the key to a better performance in self-consistency Wang et al. (2023). AC and ESC achieve significant acceleration by sacrificing diversity through early exit strategies under large branching. In contrast, path-consistency strives to preserve the diversity of reasoning paths during the sampling process from the prefix set. It makes more efficient use of available information, achieving task accuracy comparable to other methods with fewer branches and tokens.

Table 5: The performance of various inference techniques on GSM8K, Svamp and StrategyQA.

Method	GSM8K		Svamp		StrategyQA	
	Accuracy (%)	Tokens (# per problem)	Accuracy (%)	Tokens (# per problem)	Accuracy (%)	Tokens (# per problem)
SC Wang et al. (2023)	68.4	3,860	80.4	2,321	71.7	1,805
AC Aggarwal et al. (2023)	67.2	2,248 (-41.7%)	79.7	935 (-59.7%)	71.2	706 (-60.9%)
ESC Li et al. (2024)	67.2	2,710 (-29.8%)	79.2	1,154 (-50.3%)	70.8	782 (-56.7%)
PC	67.8	1,630 (-57.8%)	79.2	816 (-64.8%)	71.3	700 (-61.2%)

Table 6: Evaluation results of different inference techniques on DeepSeek-V3.

Method	GSM8K		Svamp		StrategyQA	
	Accuracy (%)	Tokens (# per problem)	Accuracy (%)	Tokens (# per problem)	Accuracy (%)	Tokens (# per problem)
PC	91.8	786	94.7	263	84.5	230
AC Aggarwal et al. (2023)	91.6	828 (+5.34%)	94.7	363 (+38.0%)	84.5	326 (+41.7%)
ESC Li et al. (2024)	91.6	868 (+10.4%)	94.7	358 (+36.1%)	84.2	315 (+36.9%)

5.4 DISCUSSION ON ROBUSTNESS AND SCALABILITY

Path-consistency doesn’t alter the model generation process. It is a lightweight, model-agnostic approach that requires no additional computation. As demonstrated in table 6, it maintains strong robustness even when applied to another open-source model, DeepSeek-V3 Liu et al. (2024).

We can also observe the scalability of path-consistency with larger model sizes. Due to the inherently strong reasoning capabilities of larger models, the impact of the three optimization methods on accuracy remains relatively stable. However, path-consistency still achieves a substantial reduction in token consumption. This demonstrates that path-consistency exhibits excellent robustness and scalability, consistently providing efficiency improvements across different models and at larger scales. Furthermore, as shown in the Appendix, we also conducted experiments on a much smaller model.

6 CONCLUSION

This paper proposes path-consistency, which achieves internal optimization of self-consistency, extracting information from early branches in the form of “prefixes”, guiding the generation of subsequent branches more efficiently. The effectiveness of path-consistency across a broad range of tasks, including arithmetic, commonsense, and symbolic reasoning, demonstrates its robustness in various application areas. Additionally, path-consistency is a lightweight, model-agnostic method that requires minimal additional computational resources and scales effectively to larger models.

REFERENCES

- 486
487
488 Pranjali Aggarwal, Aman Madaan, Yiming Yang, et al. Let’s sample step by step: Adaptive-
489 consistency for efficient reasoning and coding with llms. *arXiv preprint arXiv:2305.11860*, 2023.
- 490 Daniel Andor, Luheng He, Kenton Lee, and Emily Pitler. Giving bert a calculator: Finding oper-
491 ations and arguments with reading comprehension. In *Proceedings of the 2019 Conference on*
492 *Empirical Methods in Natural Language Processing and the 9th International Joint Conference*
493 *on Natural Language Processing (EMNLP-IJCNLP)*, pp. 5947–5952, 2019.
- 494 Wenhui Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. Program of thoughts prompt-
495 ing: Disentangling computation from reasoning for numerical reasoning tasks. *arXiv preprint*
496 *arXiv:2211.12588*, 2022.
- 497
498 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
499 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to
500 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 501 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
502 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
503 *arXiv preprint arXiv:2407.21783*, 2024.
- 504
505 Yao Fu, Hao Peng, Litu Ou, Ashish Sabharwal, and Tushar Khot. Specializing smaller language
506 models towards multi-step reasoning. In *International Conference on Machine Learning*, pp.
507 10421–10430. PMLR, 2023.
- 508 Mor Geva, Ankit Gupta, and Jonathan Berant. Injecting numerical reasoning skills into language
509 models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Lin-*
510 *guistics*, pp. 946–958, 2020.
- 511 Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. Did aristotle
512 use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions of*
513 *the Association for Computational Linguistics*, 9:346–361, 2021.
- 514
515 Jiaxian Guo, Sidi Lu, Han Cai, Weinan Zhang, Yong Yu, and Jun Wang. Long text generation via
516 adversarial training with leaked information. In *Proceedings of the AAAI conference on artificial*
517 *intelligence*, volume 32, 2018.
- 518 Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *stat*,
519 1050:9, 2015.
- 520
521 Namgyu Ho, Laura Schmid, and Seyoung Yun. Large language models are reasoning teachers. In
522 *61st Annual Meeting of the Association for Computational Linguistics, ACL 2023*, pp. 14852–
523 14882. Association for Computational Linguistics (ACL), 2023.
- 524 Liunian Harold Li, Jack Hessel, Youngjae Yu, Xiang Ren, Kai Wei Chang, and Yejin Choi. Symbolic
525 chain-of-thought distillation: Small models can also “think” step-by-step. In *61st Annual Meeting*
526 *of the Association for Computational Linguistics, ACL 2023*, pp. 2665–2679. Association for
527 Computational Linguistics (ACL), 2023.
- 528
529 Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. On the
530 advance of making language models better reasoners. *arXiv preprint arXiv:2206.02336*, 2022.
- 531 Yiwei Li, Peiwen Yuan, Shaoxiong Feng, Boyuan Pan, Xinglin Wang, Bin Sun, Heda Wang, and
532 Kan Li. Escape sky-high cost: Early-stopping self-consistency for multi-step reasoning. In *The*
533 *Twelfth International Conference on Learning Representations*, 2024.
- 534
535 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,
536 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint*
537 *arXiv:2412.19437*, 2024.
- 538 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri
539 Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement
with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2024.

- 540 Lucie Charlotte Magister, Jonathan Mallinson, Jakub Adamek, Eric Malmi, and Aliaksei Severyn.
541 Teaching small language models to reason. In *Proceedings of the 61st Annual Meeting of the*
542 *Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 1773–1781, 2023.
- 543 Shen-Yun Miao, Chao-Chun Liang, and Keh-Yih Su. A diverse corpus for evaluating and developing
544 english math word problem solvers. In *Proceedings of the 58th Annual Meeting of the Association*
545 *for Computational Linguistics*, pp. 975–984, 2020.
- 546 Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are nlp models really able to solve simple
547 math word problems? In *Proceedings of the 2021 Conference of the North American Chapter of*
548 *the Association for Computational Linguistics: Human Language Technologies*, pp. 2080–2094,
549 2021.
- 550 Debjit Paul, Mete Ismayilzada, Maxime Peyrard, Beatriz Borges, Antoine Bosselut, Robert West,
551 and Boi Faltings. Refiner: Reasoning feedback on intermediate representations. In *Proceedings*
552 *of the 18th Conference of the European Chapter of the Association for Computational Linguistics*
553 *(Volume 1: Long Papers)*, pp. 1100–1126, 2024.
- 554 Piotr Piekos, Mateusz Malinowski, and Henryk Michalewski. Measuring and improving bert’s math-
555 ematical abilities by predicting the order of reasoning. In *Proceedings of the 59th Annual Meeting*
556 *of the Association for Computational Linguistics and the 11th International Joint Conference on*
557 *Natural Language Processing (Volume 2: Short Papers)*, pp. 383–394, 2021.
- 558 Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John
559 Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. Scaling language models:
560 Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*, 2021.
- 561 Qiu Ran, Yankai Lin, Peng Li, Jie Zhou, and Zhiyuan Liu. Numnet: Machine reading comprehen-
562 sion with numerical reasoning. In *Proceedings of the 2019 Conference on Empirical Methods in*
563 *Natural Language Processing and the 9th International Joint Conference on Natural Language*
564 *Processing (EMNLP-IJCNLP)*, pp. 2474–2484, 2019.
- 565 Subhro Roy and Dan Roth. Solving general arithmetic word problems. In *Proceedings of the 2015*
566 *Conference on Empirical Methods in Natural Language Processing*, pp. 1743–1752, 2015.
- 567 Noah Shinn, Beck Labash, and Ashwin Gopinath. Reflexion: an autonomous agent with dynamic
568 memory and self-reflection. *arXiv preprint arXiv:2303.11366*, 2(5):9, 2023.
- 569 Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam
570 Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. Beyond the
571 imitation game: Quantifying and extrapolating the capabilities of language models. *Transactions*
572 *on Machine Learning Research*, 2023.
- 573 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
574 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-*
575 *tion processing systems*, 30, 2017.
- 576 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha
577 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language
578 models. In *The Eleventh International Conference on Learning Representations*, 2023.
- 579 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
580 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*
581 *neural information processing systems*, 35:24824–24837, 2022.
- 582 Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey,
583 Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. Google’s neural machine trans-
584 lation system: Bridging the gap between human and machine translation. *arXiv preprint*
585 *arXiv:1609.08144*, 2016.
- 586 Yuxi Xie, Kenji Kawaguchi, Yiran Zhao, James Xu Zhao, Min-Yen Kan, Junxian He, and Michael
587 Xie. Self-evaluation guided beam search for reasoning. *Advances in Neural Information Process-*
588 *ing Systems*, 36, 2024.

Tianyi Zhang, Tao Yu, Tatsunori B Hashimoto, Mike Lewis, Wen-tau Yih, Daniel Fried, and Sida I Wang. Coder reviewer reranking for code generation. In *Proceedings of the 40th International Conference on Machine Learning*, pp. 41832–41846, 2023.

Tinghui Zhu, Kai Zhang, Jian Xie, and Yu Su. Deductive beam search: Decoding deducible rationale for chain-of-thought reasoning. *arXiv preprint arXiv:2401.17686*, 2024.

A APPENDIX

A.1 FURTHER ANALYSIS

A.1.1 PERFORMANCE IN SMALL MODEL

As presented in the table 7, path-consistency substantially enhances the inference efficiency of small models on reasoning tasks. Nevertheless, due to the inherent accuracy limitations of the base models, the performance of path-consistency is inevitably influenced as well.

Table 7: Evaluation results of different inference techniques on Llama-3.2-1B-Instruct.

Method	GSM8K		Svamp		StrategyQA	
	Accuracy (%)	Tokens (# per problem)	Accuracy (%)	Tokens (# per problem)	Accuracy (%)	Tokens (# per problem)
PC	47.3	1430	62.0	711	63.1	799
AC Aggarwal et al. (2023)	48.1	2748	61.1	800	63.9	1401
ESC Li et al. (2024)	47.9	2324	61.9	924	63.7	1136

A.1.2 ROBUSTNESS TO CONFIDENCE THRESHOLD

As shown in Figure 6, we examined the relationship between the confidence of the majority element in the tenth branch of the self-consistency baseline and the final answer across all datasets. It was observed that a confidence level around 0.8 serves as a critical threshold for determining the correctness of the final answer. Therefore, selecting confidence levels in the range of [0.7, 0.8, 0.9] is a more appropriate choice.

A.1.3 PROTECTION OF HIGH-CONFIDENCE ANSWERS

As illustrated in the Figure 7, the distribution of confidence for the final answers across all samples is shown. The implementation of path-consistency has led to a significant increase in the number of high-confidence answers, while reducing the number of samples with confidence levels around 0.5. This suggests that path-consistency effectively resolves many ambiguous cases. Consequently, this approach not only avoids compromising task accuracy but may even enhance it.

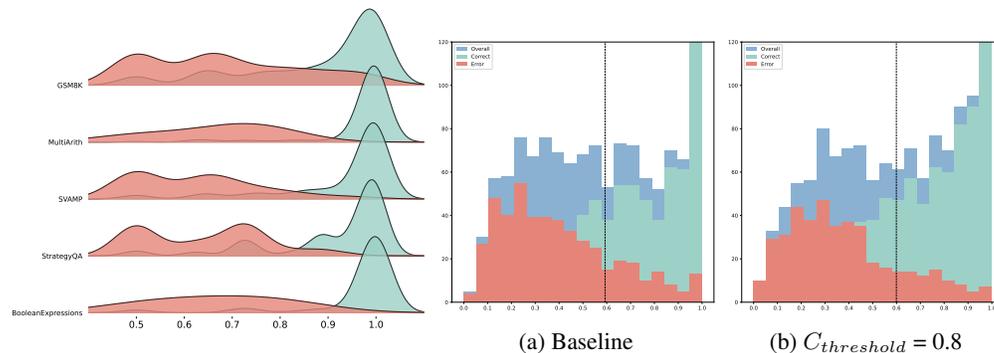


Figure 6: The distribution of the majority element’s self-confidence across all samples when self-consistency reaches the 10^{th} branch. Figure 7: The distribution of confidence for the final answers across all samples on GSM8K.

We analyzed the per-branch behavior of PC and found that prefix correctness does not need to be perfect in early iterations. Let’s see this case:

```

648 Input:
649 {"input": "Two girls each got 16 of the 24 liters of water. Then a boy got 6 liters of water. How many
650 liters of water were left?", "target": 10}
651
652 Output:
653 {"id": 80, "result": 1, "target": 10, "answer": 10.0,
654 "answers": [8.0, 10.0, 14.0, 14.0, 14.0, 2.0, 14.0, 14.0, 10.0, 16.0, 10.0, 10.0, 10.0, 10.0, 12.0, 16.0, 14.0,
655 2.0, 10.0, 14.0],
656 "tokens_gen": [70, 90, 63, 76, 84, 52, 57, 59, 54, 61, 15, 15, 15, 15, 29, 68, 65, 76, 94, 77],
657 "prefix_level": [0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 2, 2, 2, 2, 0, 0, 0, 0, ...]}

```

In this GSM8K instance using the Llama-3.2-1B-Instruct model, the first several sampled branches do not converge to the correct answer, and the initially extracted prefix is itself incorrect. However, as the sampling branches expands, the distribution shifts, the confidence of the current prefix falls below the Beta-gated threshold, and PC automatically discards the erroneous prefix. This “prefix correction” step allows later branches to align toward the correct reasoning (final answer: 10 liters), despite early inconsistencies.

This illustrates that iterative prefix selection does not increase the risk of propagating early minority errors. Instead, incorrect prefixes are statistically unlikely to survive multiple levels of gating: each Beta gate independently requires a majority concentration above the promotion threshold. As shown in the appendix, the probability that a wrong prefix is promoted across multiple levels decays exponentially with depth. Thus, even in cases where early reasoning is noisy or inconsistent, the PC mechanism naturally self-corrects, preventing error amplification and improving stability relative to standard SC.

A.1.4 EVALUATION ON LONG-CHAIN REASONING

Path-consistency is primarily designed for short-form chain-of-thought scenarios where solution paths consist of semantically discrete reasoning steps. In such tasks, path-consistency effectively reduces unnecessary token generation and accelerates inference.

As shown in Table 8, to evaluate the performance on longer reasoning chains, we have added experiments on the AIME2025 dataset using DeepSeek-V3, which produces long CoT sequences. The results show that while the speedup effect remains present on long chains, it is less pronounced, achieving a 4% reduction in token consumption. Meanwhile, there is a substantial improvement in accuracy when applying path-consistency, because path-consistency filter out early unstable reasoning branches and enforce prefix agreement, which helps long CoT sequences converge to more coherent and correct solutions. Thank you for the helpful suggestion.

Table 8: Evaluation results on AIME2025.

Method	Accuracy (%)	Decrease (%)
Self-consistency	60.0	-
Path-consistency	76.7	-4.1%

A.1.5 COMPUTATIONAL OVERHEAD

As shown in Table 9, we empirically evaluated the runtime cost of path-consistency on GSM8K using Llama3-8B with 20-branch sampling per problem. The additional overhead of prefix extraction and Beta confidence evaluation was negligible compared to the model sampling time.

Table 9: Relative Time Overhead of Beta Evaluation and Prefix Extraction

Component	Time (s)	Ratio to Generation Time
Model generation	2.2	100%
Beta confidence evaluation	0.000096	0.0044%
Prefix extraction	0.00013	0.0059%

702 A.1.6 EVALUATION ON MISTRAL-7B

703
704 Path-consistency also generalizes to other models. Table 10 reports results on Mistral-7B evaluated
705 on GSM8K. The trend remains consistent with our main experiments: PC provides substantial effi-
706 ciency gains (up to +22% speedup and -20.3% token reduction) while maintaining accuracy close to
707 or slightly above the SC baseline. Notably, with a moderate confidence threshold ($C_{threshold} = 0.8$),
708 PC even improves accuracy (+0.6%) while reducing computation by more than 11%. These results
709 demonstrate that path-consistency is not tied to a specific architecture and remains effective on
710 stronger open-weight models such as Mistral-7B.

711 Table 10: Evaluation results on Mistral-7B with path-consistency on GSM8K.

712 Method	713 $C_{threshold}$	714 Accuracy (%)	715 Speedup (%)	716 Decrease (%)
717 SC	718 -	719 59.7	720 -	721 -
722 PC	723 0	724 58.5	725 +22.0	726 -20.3
	727 0.8	728 60.3	729 +11.1	730 -11.4
	731 0.9	732 60.4	733 +6.2	734 -6.6

735 A.2 CODE ANALYSIS

736 Section 4.1 introduces the “extract-and-sample” process of path-consistency. Following is the code
737 implementation of the key parts:

738 A.2.1 SAMPLE

739 Randomly sample a prefix from a list of prefixes.

740 Listing 1: Sample

```
741 1 def sample_prefix(prefix_list : list[str]) -> str:
742 2     if len(prefix_list) > 0:
743 3         prefix_index = random.randint(0, len(prefix_list)-1)
744 4         return prefix_list[prefix_index]
745 5     return ""
```

746 A.2.2 EXTRACT

747 Generate a list of potential prefixes for future inference based on the confidence of previous answers.

748 Listing 2: Extract

```
749 1 def get_prefix(answers, conf_thresh, gens, prefix, prefix_level, frq:int
750 = 5):
751 2     confidence = confidence_criteria(answers, conf_thresh)
752 3     prefix_list = []
753 4     if confidence['conf'] == False:
754 5         return prefix, prefix_level
755 6     for i, gen in enumerate(gens[-frq:]):
756 7         index = len(gens) - frq + i
757 8         parts = gen.split("answer is")
758 9         if len(parts) > 1:
759 10            pre_gen = parts[0]
760 11            if answers[index] == confidence['most_common']:
761 12                parts = re.split(r'(\.s|\.\n|\n)', pre_gen)
762 13                sentences = []
763 14                for i in range(0, len(parts) - 1, 2):
764 15                    sentences.append(parts[i] + parts[i + 1])
765 16            sentence = ""
766 17            for j in range(prefix_level + 1):
```

```

756 18         if j < len(sentences) - 1:
757 19             sentence += sentences[j]
758 20             # sentence += ". "
759 21             prefix_list.append(sentence)
760 22     prefix_level += 1
761 23     return prefix_list, prefix_level

```

763 A.2.3 MAIN INFERENCE

764 A class to perform path-consistency inference using a model. It recursively samples and integrates
765 answers based on confidence thresholds, ultimately returning a final answer.

767 Listing 3: Inference

```

768
769 1 def inference(self, prompt: str, **kwargs):
770 2     answers = []
771 3     prefix_level = 0
772 4     prefix_list = []
773 5     reasoning = []
774 6     times = []
775 7     generations = []
776 8     for branch_id in range(self.max_branch):
777 9         prefix = sample_prefix(prefix_list)
778 10        prompt_plus_prefix = prompt + prefix
779 11        start_time = time.time()
780 12        generation =
781 13            self.model.completion_function(prompt_plus_prefix, **kwargs)
782 14        end_time = time.time()
783 15        times.append(end_time - start_time)
784 16        reasoning.append(prefix + generation)
785 17        answer = extract_answer(generation, self.ans_type)
786 18        answers.append(answer)
787 19        generations.append(generation)
788 20        if (branch_id + 1) % (self.max_branch / (self.max_level + 1)) ==
789 21            0:
790 22            prefix_list, prefix_level = get_prefix(answers,
791 23                self.confidence_thres, reasoning, prefix_list,
792                prefix_level)
793        final_answer = integrate_answer(answers)
794        info = {'answer' : final_answer, 'answers' : answers, 'latency' :
795            times, 'generations' : generations}
796        return info

```

794 A.2.4 WRAP THE MODEL

795 Abstract base class for completion models defines the interface for completion models. All sub-
796 classes must implement the “completion_function” method. Concrete implementation of Comple-
797 tionModel using the Llama model uses a text generation model (like Llama) to generate completions
798 based on a given prompt.

800 Listing 4: Inference

```

801 1 class CompletionModel(ABC):
802 2     @abstractmethod
803 3     def completion_function(self, prompt: str) -> str:
804 4         pass
805 5
806 6 class LlamaModel(CompletionModel):
807 7     def __init__(self, generator):
808 8         self.generator = generator
809 9
810 10    def completion_function(self, prompt: str, **kwargs) -> str:
811 11        prompt_list = [prompt]

```

```
810
811 12
812 13     results = self.generator.text_completion(
813 14         prompt_list,
814 15         **kwargs
815 16     )
816 17     return results[-1]['generation']
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
```

A.3 PROOF OF PROMOTION SOUNDNESS

Theorem 1 (Single-Level Soundness of Beta-Gated Prefix Selection). *Consider a fixed level of prefix extraction with the following setup:*

- A denotes a finite answer set.
- The algorithm performs promotion checks at a sequence of increasing sample counts:

$$0 < n^{(1)} < n^{(2)} < \dots < n^{(K)} = n_{\max}.$$

At each check k , we draw the first $n^{(k)}$ i.i.d. samples:

$$X_1, \dots, X_{n^{(k)}} \in \mathcal{A}.$$

- For each $a \in \mathcal{A}$, let $\pi(a) := \Pr(X_i = a)$ denote the model’s sampling probability under the prefix. The correct answer is a^* with

$$q := \pi(a^*), \quad R := 1 - q = \sum_{a \neq a^*} \pi(a).$$

- For each check k , define the observed counts:

$$N_a^{(k)} := \sum_{i=1}^{n^{(k)}} \mathbf{1}\{X_i = a\}, \quad \hat{a}^{(k)} := \arg \max_{a \in \mathcal{A}} N_a^{(k)}, \quad s^{(k)} := N_{\hat{a}^{(k)}}^{(k)}.$$

Assume a Beta-gated promotion threshold $\tau \geq 0.5$, where promotion occurs at check k if

$$T(s^{(k)}, n^{(k)}) := \Pr(\theta > 1/2 \mid \text{Beta}(s^{(k)} + 1, n^{(k)} - s^{(k)} + 1)) > \tau.$$

Then, if $q \geq 1/2 + \Delta$ for some $\Delta > 0$, the probability that the correct prefix fails to be promoted at check k is bounded by

$$\Pr(A_k) \leq \exp(-2n^{(k)}\Delta^2),$$

where $A_k := \{\text{not promoted at check } k\}$.

Proof. Promotion at check k requires that $\hat{a}^{(k)} = a^*$ and $T(s^{(k)}, n^{(k)}) > \tau$. Thus failure at check k implies

$$A_k \subseteq \{N_{a^*}^{(k)} < k_{\min}(n^{(k)}, \tau)\},$$

where

$$k_{\min}(n, \tau) := \min\{r \in \{0, \dots, n\} : T(r, n) > \tau\}.$$

Since $\tau \geq 0.5$, symmetry of the Beta posterior implies

$$k_{\min}(n^{(k)}, \tau) > n^{(k)}/2.$$

Hence

$$A_k \subseteq \{N_{a^*}^{(k)} \leq n^{(k)}/2\}.$$

But $N_{a^*}^{(k)}$ is a sum of $n^{(k)}$ independent Bernoulli(q) variables with mean at least $n^{(k)}(1/2 + \Delta)$. By Hoeffding’s inequality,

$$\Pr\left(N_{a^*}^{(k)} \leq n^{(k)}/2\right) = \Pr\left(N_{a^*}^{(k)} \leq \mathbb{E}[N_{a^*}^{(k)}] - n^{(k)}\Delta\right) \leq \exp(-2n^{(k)}\Delta^2).$$

□

Corollary 1 (Multi-Level Promotion Soundness). *Consider m consecutive prefix-extraction levels indexed by ℓ_1, \dots, ℓ_m . At each level ℓ_j , promotion checks occur at sample sizes $n_{\ell_j}^{(1)} < \dots < n_{\ell_j}^{(K_j)}$, with correct-answer probabilities $q_{\ell_j} \geq 1/2 + \Delta_{\ell_j}$ and Beta thresholds $\tau_{\ell_j} \geq 0.5$. Let $A_{\ell_j}^{(k)}$ denote the failure-to-promote event at the k -th check of level ℓ_j . Then, under conditional independence across levels,*

$$\Pr(\text{failure along full promotion path}) \leq \prod_{j=1}^m \left(\sum_{k=1}^{K_j} \exp(-2n_{\ell_j}^{(k)}\Delta_{\ell_j}^2) \right).$$

918 *Proof.* For each level ℓ_j , we assume the true success probability remains $q_{\ell_j} \geq 1/2 + \Delta_{\ell_j}$. This
 919 assumption is valid due to the one-time selection guarantee proved in Section 4.2, which ensures
 920 that the prefix extraction does not hurt accuracy. Thus, prefix quality is not degraded by previous
 921 selection decisions; it preserves a positive advantage Δ_{ℓ_j} at every level. At each promotion check k
 922 of level ℓ_j , requiring posterior probability

$$923 \Pr(\theta_{\ell_j} > 1/2) > \tau_{\ell_j} \geq 0.5,$$

924
 925 Hoeffding's bound implies that the probability that a truly good prefix fails to promote at this check
 926 satisfies

$$927 \Pr(A_{\ell_j}^{(k)}) \leq \exp(-2n_{\ell_j}^{(k)} \Delta_{\ell_j}^2).$$

928 Since promotion only needs to succeed once within a level, the probability of failure across all
 929 checks of level ℓ_j is bounded by

$$930 \Pr\left(\bigcap_k A_{\ell_j}^{(k)}\right) \leq \sum_{k=1}^{K_j} \exp(-2n_{\ell_j}^{(k)} \Delta_{\ell_j}^2).$$

931
 932 Finally, assuming conditional independence across levels,

$$933 \Pr(\text{failure along full promotion path}) = \Pr\left(\bigcap_{j=1}^m \bigcap_{k=1}^{K_j} A_{\ell_j}^{(k)}\right) \leq \prod_{j=1}^m \left(\sum_{k=1}^{K_j} \exp(-2n_{\ell_j}^{(k)} \Delta_{\ell_j}^2)\right).$$

934
 935 Thus the overall probability of failing to promote a prefix with true advantage $\Delta_{\ell_j} > 0$ decays
 936 exponentially with the total effective sample budget across levels. \square

937
 938 **Corollary 2** (Worst-Case Wrong Answer Distribution). *For fixed total wrong-mass $R = 1 - q$, the
 939 probability that a wrong answer is promoted is maximized when all wrong probability is concen-
 940 trated on a single alternative. In that case, the probability of wrong promotion is*

$$941 \Pr(N_{a_{\text{bad}}} \geq k_{\min}(n, \tau)), \quad a_{\text{bad}} \in \mathcal{A} \setminus \{a^*\}.$$

942
 943 *Proof.* Let $T := \sum_{a \neq a^*} N_a$ denote the total count of wrong answers. Then

$$944 \{\exists a \neq a^* : N_a \geq k_{\min}\} \subseteq \{T \geq k_{\min}\}.$$

945
 946 When the entire wrong-mass R is assigned to a single a_{bad} , $T = N_{a_{\text{bad}}}$, and the inclusion becomes
 947 equality. Hence the probability of wrong promotion is maximized in this configuration. \square

A.4 FULL PROMPTS

Table 11: Few-shot exemplars for arithmetic reasoning.

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

A: We start with 15 trees. Later we have 21 trees. The difference must be the number of trees they planted. So, they must have planted $21 - 15 = 6$ trees. The answer is 6.

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

A: There are 3 cars in the parking lot already. 2 more arrive. Now there are $3 + 2 = 5$ cars. The answer is 5.

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

A: Leah had 32 chocolates and Leah's sister had 42. That means there were originally $32 + 42 = 74$ chocolates. 35 have been eaten. So in total they still have $74 - 35 = 39$ chocolates. The answer is 39.

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

A: Jason had 20 lollipops. Since he only has 12 now, he must have given the rest to Denny. The number of lollipops he has given to Denny must have been $20 - 12 = 8$ lollipops. The answer is 8.

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

A: He has 5 toys. He got 2 from mom, so after that he has $5 + 2 = 7$ toys. Then he got 2 more from dad, so in total he has $7 + 2 = 9$ toys. The answer is 9.

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

A: There are 4 days from monday to thursday. 5 computers were added each day. That means in total $4 * 5 = 20$ computers were added. There were 9 computers in the beginning, so now there are $9 + 20 = 29$ computers. The answer is 29.

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

A: Michael initially had 58 balls. He lost 23 on Tuesday, so after that he has $58 - 23 = 35$ balls. On Wednesday he lost 2 more so now he has $35 - 2 = 33$ balls. The answer is 33.

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

A: She bought 5 bagels for \$3 each. This means she spent $5 * \$3 = \15 on the bagels. She had \$23 in beginning, so now she has $\$23 - \$15 = \$8$. The answer is 8.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

Table 12: Few-shot exemplars for StrategyQA.

Q: Do hamsters provide food for any animals?
A: Hamsters are prey animals. Prey are food for predators. Thus, hamsters provide food for some animals. The answer is yes.

Q: Yes or no: Hydrogen’s atomic number squared exceeds number of Spice Girls?
A: Hydrogen has an atomic number of 1. 1 squared is 1. There are 5 Spice Girls. Thus, Hydrogen’s atomic number squared is less than 5. The answer is no.

Q: Yes or no: Would a pear sink in water?
A: The density of a pear is about 0.6 g/cm^3 , which is less than water. Objects less dense than water float. Thus, a pear would float. The answer is no.

Q: Yes or no: Is it common to see frost during some college commencements?
A: College commencement ceremonies can happen in December, May, and June. December is in the winter, so there can be frost. Thus, there could be frost at some commencements. The answer is yes.

Q: Yes or no: Could a llama birth twice during War in Vietnam (1945-46)?
A: The War in Vietnam was 6 months. The gestation period for a llama is 11 months, which is more than 6 months. Thus, a llama could not give birth twice during the War in Vietnam. The answer is no.

Q: Could Brooke Shields succeed at University of Pennsylvania?
A: Brooke Shields went to Princeton University. Princeton University is about as academically rigorous as the University of Pennsylvania. Thus, Brooke Shields could also succeed at the University of Pennsylvania. The answer is yes.

Table 13: Few-shot exemplars for RuinNames.

1080
1081
1082 **Q:** Which of the following is a humorous edit of this artist or movie name: 'whitesnake'?

1083 **Options:** (A) whitesnape, (B) whitesnapke, (C) whitesnuake, (D) mwhitesnake

1084 **A:** The original name is "whitesnake". This is the name of an old English hard rock band. It is a
1085 compound word, formed by the words "white" and "snake".

1086 (A) "whitesnape": It is formed by the combination of "white" and "snake"; therefore, "snake" has
1087 been changed to "snape". Snape makes a reference to the fictional character Severus Snape in the
1088 Harry Potter series, so (A) is indeed a meaningful and funny edit.

1089 (B) "whitesnapke": It is formed by the combination of "white" and "snapke", but "snapke" is not an
1090 actual word; therefore, "whitesnapke" is not humorous.

1091 (C) "whitesnuake": It is formed by the combination of "white" and "snuake", but "snuake" is not an
1092 actual word; therefore, "whitesnuake" is not humorous.

1093 (D) "mwhitesnake": It is formed by the combination of "m", "white", and "snake", but the prefix
1094 "m" seems arbitrary; therefore, "mwhitesnake" is not meaningful or humorous.

1095 Above the above, the only humorous edit is (A). So the answer is (A).

1096
1097 **Q:** Which of the following is a humorous edit of this artist or movie name: 'one of our dinosaurs is
1098 missing'?

1099 **Options:** (A) ofne of our dinosaurs is missing, (B) one af our dinosaurs is missing, (C) one of our
1100 dinosaurs is pissing, (D) one of our dinosaur is missing

1101 **A:** The original name is "one of our dinosaurs is missing". This is the name of an old British movie.
1102 (A) "ofne of our dinosaurs is missing": Here "one of" is changed to "ofne", but the word "ofne" is
1103 not an actual word.

1104 (B) "one af our dinosaurs is missing": Here the word "of" is changed to "af", but the word "af" is
1105 not an actual word.

1106 (C) "one of our dinosaurs is pissing": Here the word "missing" is changed to "pissing", and "one of
1107 our dinosaurs is pissing" is indeed a very whimsical and mischievous edit. This change truly ruins
1108 the original title of the movie.

1109 (D) "one of our dinosaur is missing": Here the word "dinosaurs" is changed to "dinosaur", but
1110 "dinosaur" is singular but should be plural in the title; this change therefore feels arbitrary and not
1111 humorous.

1112 Above the above, the only humorous edit is (C). So the answer is (C).

1113
1114 **Q:** Which of the following is a humorous edit of this artist or movie name: 'counting crows'?

1115 **Options:** (A) countingy crows, (B) counting cows, (C) courting crows, (D) countnting crows

1116 **A:** The original name is "counting crows". This is the name of an American rock band. Historically,
1117 the band name comes from the British nursery rhyme "One for Sorrow", which is about counting of
1118 magpies.

1119 (A) "countingy crows": Here the word "counting" is changed to "countingy", but the word "count-
1120 ingy" is not an actual word.

1121 (B) "counting cows": Here the word "crows" is changed to "cows", and this is indeed a playful and
1122 meaningful edit that ruins the original name of the band.

1123 (C) "courting crows": Here the word "counting" is changed to "courting", and "courting" is an
1124 actual word; however, "courting crows" does not sound as humorous as "counting crows".

1125 (D) "countnting crows": Here the word "counting" is changed to "countnting", but the word "countnt-
1126 ing" is not an actual word.

1127 Above the above, the only humorous edit is (B). So the answer is (B).

1128
1129
1130
1131
1132
1133

Table 14: Few-shot exemplars for SalientTranslation.

1134
1135
1136 **Q:** The following translations from German to English contain a particular error. That error will be one of the
1137 following types: Named Entities: An entity (names, places, locations, etc.) is changed to a different entity.
1138 Numerical Values: Numerical values (ordinals or cardinals), dates, and/or units are changed. Modifiers or
1139 Adjectives: The modifiers and adjectives pertaining to a noun are changed. Negation or Antonyms: Introduce
1140 or remove a negation or change comparatives to their antonyms. Facts: Trivial factual errors not pertaining to
1141 the above classes are introduced in the translations. Dropped Content: A significant clause in the translation is
1142 removed. Please identify that error.
1143 **Source:** In der Liste der Baudenkmale in Lenzen (Elbe) sind alle Baudenkmale der brandenburgischen Stadt
1144 Lenzen (Elbe) und ihrer Ortsteile aufgelistet.
1145 **Translation:** In the list of architectural monuments in Lenzen all architectural monuments of the Brandenburg
1146 city of Lenzen and its districts are listed.
1147 **The translation contains an error pertaining to:**
1148 **Options:** (A) Modifiers or Adjectives, (B) Numerical Values, (C) Negation or Antonyms, (D) Named Entities,
1149 (E) Dropped Content, (F) Facts
1150 **A:** Let's think step by step. We solve this question by first translating the source sentence to English and then by
1151 comparing our translation with the provided translation. According to Google Translate, the correct translation
1152 of the source sentence from German to English is "The list of monuments in Lenzen (Elbe) includes all the
1153 monuments in the Brandenburg town of Lenzen (Elbe) and its districts." On the other hand, the provided trans-
1154 lation is "In the list of architectural monuments in Lenzen all architectural monuments of the Brandenburg city
1155 of Lenzen and its districts are listed." Note that Lenzen (Elbe) is changed to Lenzen in the original translation;
1156 so, there is a named entity error. Because an entity in the original source sentence is changed to a different
1157 entity in the translation, the translation contains an error pertaining to Named Entities. So the answer is (D).

1158 **Q:** The following translations from German to English contain a particular error. That error will be one of the
1159 following types: Named Entities: An entity (names, places, locations, etc.) is changed to a different entity.
1160 Numerical Values: Numerical values (ordinals or cardinals), dates, and/or units are changed. Modifiers or
1161 Adjectives: The modifiers and adjectives pertaining to a noun are changed. Negation or Antonyms: Introduce
1162 or remove a negation or change comparatives to their antonyms. Facts: Trivial factual errors not pertaining to
1163 the above classes are introduced in the translations. Dropped Content: A significant clause in the translation is
1164 removed. Please identify that error.
1165 **Source:** Auf dieser Seite sind die Baudenkmaler der oberbayerischen Groben Kreisstadt Landsberg am Lech
1166 zusammengestellt.
1167 **Translation:** On this page are compiled the architectural monuments of the town of Landsberg am Lech.
1168 **The translation contains an error pertaining to:**
1169 **Options:** (A) Modifiers or Adjectives, (B) Numerical Values, (C) Negation or Antonyms, (D) Named Entities,
1170 (E) Dropped Content, (F) Facts
1171 **A:** Let's think step by step. We solve this question by first translating the source sentence to English and then by
1172 comparing our translation with the provided translation. According to Google Translate, the correct translation
1173 of the source sentence from German to English is "The monuments of the Upper Bavarian district town of
1174 Landsberg am Lech are compiled on this page." On the other hand, the provided translation is "On this page
1175 are compiled the architectural monuments of the town of Landsberg am Lech." Note that an important detail
1176 about the location of Landsberg am Lech is omitted in the original translation: The translation should have
1177 said "Upper Bavarian district town of Landsberg am Lech". Because a significant clause in the translation was
1178 removed, the translation contains an error pertaining to Dropped Content. So the answer is (E).

1179 **Q:** The following translations from German to English contain a particular error. That error will be one of the
1180 following types: Named Entities: An entity (names, places, locations, etc.) is changed to a different entity.
1181 Numerical Values: Numerical values (ordinals or cardinals), dates, and/or units are changed. Modifiers or
1182 Adjectives: The modifiers and adjectives pertaining to a noun are changed. Negation or Antonyms: Introduce
1183 or remove a negation or change comparatives to their antonyms. Facts: Trivial factual errors not pertaining to
1184 the above classes are introduced in the translations. Dropped Content: A significant clause in the translation is
1185 removed. Please identify that error.
1186 **Source:** Leba ist eine Kleinstadt und ein Badeort im Powiat Leborski der polnischen Woiwodschaft Pommern.
1187 **Translation:** Eba is not a small town and seaside resort in the Powiat Leborski county of the Pomeranian
1188 Voivodeship of Poland.
1189 **The translation contains an error pertaining to:**
1190 **Options:** (A) Modifiers or Adjectives, (B) Numerical Values, (C) Negation or Antonyms, (D) Named Entities,
1191 (E) Dropped Content, (F) Facts
1192 **A:** Let's think step by step. We solve this question by first translating the source sentence to English and
1193 then by comparing our translation with the provided translation. According to Google Translate, the correct
1194 translation of the source sentence from German to English is "Leba is a small town and seaside resort in
1195 the Powiat Leborski of the Polish Pomeranian Voivodeship." On the other hand, the provided translation is
1196 "Leba is not a small town and seaside resort in the Powiat Leborski county of the Pomeranian Voivodeship of
1197 Poland." Note that the provided sentence says, "Leba is not a small town ..." However, the translation should
1198 have been "Leba is a small town ..." Because a negation is introduced at the beginning of the sentence and
1199 has fundamentally changed the meaning of the original source, the translation contains an error pertaining to
1200 Negation or Antonyms. So the answer is (C).

Table 15: Few-shot exemplars for BooleanExpressions.

1188 **Q:** not ((not not True)) is
1189 **A:** Let’s think step by step.
1190 Remember that (i) expressions inside brackets are always evaluated first and that (ii) the order of
1191 operations from highest priority to lowest priority is ”not”, ”and”, ”or”, respectively.
1192 We first simplify this expression ”Z” as follows: ”Z = not ((not not True)) = not ((A))” where
1193 ”A = not not True”.
1194 Let’s evaluate A: A = not not True = not (not True) = not False = True.
1195 Plugging in A, we get: Z = not ((A)) = not ((True)) = not True = False. So the answer is False.
1196
1197
1198 **Q:** True and False and not True and True is
1199 **A:** Let’s think step by step.
1200 Remember that (i) expressions inside brackets are always evaluated first and that (ii) the order of
1201 operations from highest priority to lowest priority is ”not”, ”and”, ”or”, respectively.
1202 We first simplify this expression ”Z” as follows: ”Z = True and False and not True and True = A and
1203 B” where ”A = True and False” and ”B = not True and True”.
1204 Let’s evaluate A: A = True and False = False.
1205 Let’s evaluate B: B = not True and True = not (True and True) = not (True) = False.
1206 Plugging in A and B, we get: Z = A and B = False and False = False. So the answer is False.
1207
1208
1209 **Q:** not not (not (False)) is
1210 **A:** Let’s think step by step.
1211 Remember that (i) expressions inside brackets are always evaluated first and that (ii) the order of
1212 operations from highest priority to lowest priority is ”not”, ”and”, ”or”, respectively.
1213 We first simplify this expression ”Z” as follows: ”Z = not not (not (False)) = not not (A)” where
1214 ”A = not (False)”.
1215 Let’s evaluate A: A = not (False) = not False = True.
1216 Plugging in A, we get: Z = not not (A) = not not (True) = not not False = True. So the answer is
1217 True.

1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

Table 16: Few-shot exemplars for TrackingShuffledObject.

1242
1243
1244 **Q:** Alice, Bob, and Claire are playing a game. At the start of the game, they are each holding a ball:
1245 Alice has a yellow ball, Bob has a blue ball, and Claire has a pink ball. As the game progresses,
1246 pairs of players trade balls. First, Claire and Alice swap balls. Then, Alice and Bob swap balls.
1247 Finally, Claire and Bob swap balls. At the end of the game, Bob has the
1248 **Options:** (A) yellow ball, (B) blue ball, (C) pink ball
1249 **A:** Let’s think step by step.
1250 (0) At the start: Alice: yellow, Bob: blue, Claire: pink.
1251 (1) Claire and Alice swap balls: Alice: pink, Bob: blue, Claire: yellow.
1252 (2) Alice and Bob swap balls: Alice: blue, Bob: pink, Claire: yellow.
1253 (3) Claire and Bob swap balls: Alice: blue, Bob: yellow, Claire: pink.
1254 At the end of the game, Bob has the yellow ball. So the answer is (A).

1255 **Q:** Alice, Bob, and Claire are playing a game. At the start of the game, they are each holding a ball:
1256 Alice has a white ball, Bob has a purple ball, and Claire has a pink ball. As the game progresses,
1257 pairs of players trade balls. First, Bob and Alice swap balls. Then, Bob and Claire swap balls.
1258 Finally, Bob and Alice swap balls. At the end of the game, Alice has the
1259 **Options:** (A) white ball, (B) purple ball, (C) pink ball
1260 **A:** Let’s think step by step.
1261 (0) At the start: Alice: white, Bob: purple, Claire: pink.
1262 (1) Bob and Alice swap balls: Alice: purple, Bob: white, Claire: pink.
1263 (2) Bob and Claire swap balls: Alice: purple, Bob: pink, Claire: white.
1264 (3) Bob and Alice swap balls: Alice: pink, Bob: purple, Claire: white.
1265 At the end of the game, Alice has the pink ball. So the answer is (C).

1266 **Q:** Alice, Bob, and Claire are dancers at a square dance. At the start of a song, they each have a
1267 partner: Alice is dancing with Lola, Bob is dancing with Rodrigo, and Claire is dancing with Patrick.
1268 Throughout the song, the dancers often trade partners. First, Alice and Bob switch partners. Then,
1269 Claire and Bob switch partners. Finally, Bob and Alice switch partners. At the end of the dance,
1270 Alice is dancing with
1271 **Options:** (A) Lola, (B) Rodrigo, (C) Patrick
1272 **A:** Let’s think step by step.
1273 (0) At the start: Alice: Lola, Bob: Rodrigo, Claire: Patrick.
1274 (1) Alice and Bob switch partners: Alice: Rodrigo, Bob: Lola, Claire: Patrick.
1275 (2) Claire and Bob switch partners: Alice: Rodrigo, Bob: Patrick, Claire: Lola.
1276 (3) Bob and Alice switch partners: Alice: Patrick, Bob: Rodrigo, Claire: Lola.
1277 At the end of the dance, Alice is dancing with Patrick. So the answer is (C).

1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

Table 17: Few-shot exemplars for LogicalDeduction.

1296	
1297	
1298	Q: The following paragraphs each describe a set of three objects arranged in a fixed order. The
1299	statements are logically consistent within each paragraph. In a golf tournament, there were three
1300	golfers: Amy, Eli, and Eve. Eve finished above Amy. Eli finished below Amy.
1301	Options: (A) Amy finished last, (B) Eli finished last, (C) Eve finished last
1302	A: Let’s think step by step.
1303	(1) Eve finished above Amy: ”(above) ? Eve ? Amy ? (below)”.
1304	(2) Eli finished below Amy: ”(above) ? Amy ? Eli ? (below)”.
1305	(3) Combining (1) and (2) we get the following ordering: ”(above) Eve Amy Eli (below)”.
1306	According to this ordering, the person who finished last (the one at the bottom of this list) is Eli.
1307	Eli finished last. So the answer is (B).
1308	
1309	Q: The following paragraphs each describe a set of three objects arranged in a fixed order. The
1310	statements are logically consistent within each paragraph. On a shelf, there are three books: a white
1311	book, a green book, and an orange book. The green book is to the right of the white book. The
1312	orange book is the rightmost.
1313	Options: (A) The white book is the leftmost, (B) The green book is the leftmost, (C) The orange
1314	book is the leftmost
1315	A: Let’s think step by step.
1316	(1) The green book is to the right of the white book: ”(left) ? white ? green ? (right)”.
1317	(2) The orange book is the rightmost: ”(left) ? white ? green orange (right)”.
1318	(3) Combining (1) and (2) we get the following ordering: ”(left) white green orange (right)”.
1319	According to this ordering, the leftmost book is the white book.
1320	The white book is the leftmost. So the answer is (A).
1321	
1322	Q: The following paragraphs each describe a set of three objects arranged in a fixed order. The
1323	statements are logically consistent within each paragraph. On a shelf, there are three books: a red
1324	book, a gray book, and a white book. The white book is to the left of the gray book. The red book
1325	is the second from the left.
1326	Options: (A) The red book is the leftmost, (B) The gray book is the leftmost, (C) The white book
1327	is the leftmost
1328	A: Let’s think step by step.
1329	(1) The white book is to the left of the gray book: ”(left) ? white ? gray ? (right)”.
1330	(2) The red book is the second from the left: ”(left) ? white red gray ? (right)”.
1331	(3) Combining (1) and (2) we get the following ordering: ”(left) white red gray (right)”.
1332	According to this ordering, the leftmost book is the white book.
1333	The white book is the leftmost. So the answer is (C).
1334	
1335	
1336	
1337	
1338	
1339	
1340	
1341	
1342	
1343	
1344	
1345	
1346	
1347	
1348	
1349	