
Enhancing Large Language Models through Structured Reasoning

Yubo Dong

Department of Computer Science
Zhejiang University
dongyubo@zju.edu.cn

Hehe Fan

Department of Computer Science
Zhejiang University
hehe.fan.cs@gmail.com

Abstract

Recent Large Language Models (LLMs) have significantly advanced natural language processing and automated decision-making. However, these models still encounter difficulties when performing complex reasoning tasks involving logical deduction and systematic planning, primarily due to their reliance on implicit statistical relationships without structured knowledge representation. Inspired by cognitive science and neurosymbolic AI, we introduce a novel approach to enhance LLMs through explicit structured reasoning. First, we convert unstructured data into structured formats by explicitly annotating reasoning steps. We then employ this structured dataset to train LLMs through Supervised Fine-Tuning (SFT). Additionally, we enhance the structured reasoning capabilities of LLMs using Group Relative Policy Optimization (GRPO), incorporating two innovative algorithms—MAX-Flow and Longest Common Subsequence (LCS)—which notably improve reasoning effectiveness and reduce computational complexity. Experimental results from fine-tuning a DeepSeek-R1-Distill-Qwen-1.5B model demonstrate concise reasoning, robust performance across various scenarios, and improved compatibility with optimization techniques, validating the efficacy of structured reasoning integration in LLMs. Code is available at: <https://github.com/cnsdqd-dyb/Enhancing-Large-Language-Models-through-Structured-Reasoning>

1 Introduction

Recent advancements in Large Language Models (LLMs), exemplified by models such as DeepSeek-R1 [1], OpenAI-o1 [2] and QwQ [3], have impacted various domains including natural language processing, knowledge reasoning, and automated decision-making. Despite their remarkable abilities to generate fluent and contextually coherent text, LLMs usually demonstrate limitations when tackling complex reasoning tasks, such as logical deduction, systematic planning, and sophisticated problems in advanced mathematics and physics. These limitations stem primarily from their reliance on implicit, statistical associations learned from massive unstructured datasets, without explicit mechanisms to systematically represent and manipulate structured knowledge and logical relationships.

Structured reasoning, a fundamental cognitive ability observed in human intelligence ([4–7]), involves explicitly representing problems, systematically manipulating knowledge, and drawing logical conclusions in a structured and interpretable manner. Integrating structured reasoning into LLMs has the potential to significantly improve their performance on tasks that require explicit logic and structured problem solving, bridging the gap between human-like reasoning and statistical text prediction.

This paper proposes a novel approach to enhance LLMs with structured reasoning capabilities, inspired by cognitive science theories and recent advances in neurosymbolic artificial intelligence. Specifically, we introduce mechanisms that explicitly encode structured knowledge representations and reasoning processes in LLMs. Our method aims to leverage both the flexibility of neural networks and the interpretability and precision of symbolic reasoning.

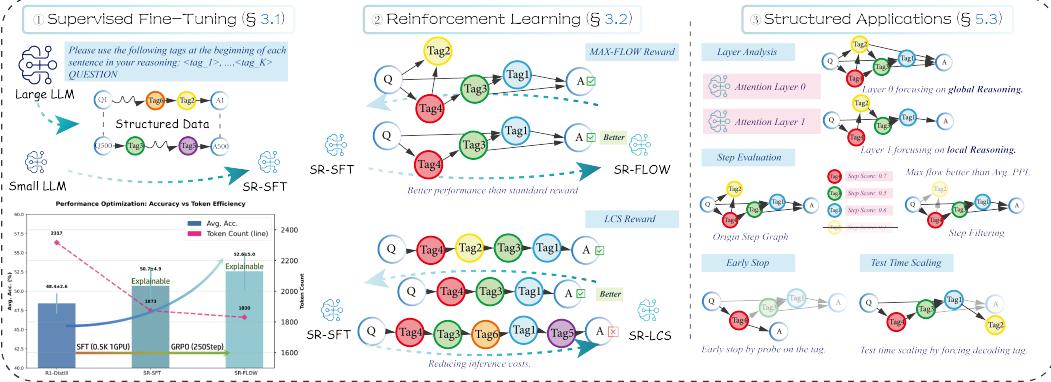


Figure 1: Illustration of enhancing Large Language Models (LLMs) with Structured Reasoning (SR) through (1) Supervised Fine-Tuning (SFT) and (2) Group Relative Policy Optimization (GRPO). In SFT, we construct a structured dataset that is explicitly annotated with reasoning step labels. In GRPO, we apply the MAX-Flow algorithm to emphasize reasoning processes with balanced contributions, improving reasoning accuracy, and utilize the Longest Common Subsequence (LCS) algorithm to prioritize reasoning processes sharing common subsequences, thus reducing computational overhead. As demonstrated in (3), our approach results in more concise reasoning, stable performance across various settings, and improved compatibility with optimization methods.

First, we transform unstructured data into structured data by incorporating explicit reasoning tags or labels that clearly indicate the current step of the reasoning process. In other words, we strengthen LLMs through structured reasoning with explicit step annotations. We utilize this structured dataset to train LLMs through Supervised Fine-Tuning (SFT). Additionally, we employ Group Relative Policy Optimization (GRPO) to further enhance structured reasoning. Specifically, GRPO integrates two key algorithms to promote effective and efficient structured reasoning: (1) *MAX-Flow*, which surpasses traditional perplexity metrics in assessing the effectiveness of reasoning steps; and (2) *Longest Common Subsequence (LCS)*, which preserves reasoning accuracy while substantially decreasing computational complexity. Using only 500 structured reasoning examples and 250 global reinforcement learning, we fine-tuned a DeepSeek-R1-Distill-Qwen-1.5B model. It achieves more concise reasoning, demonstrates stable performance in various scenarios, and offers enhanced compatibility with practical optimization methods. Our contributions are as follows,

- We propose a novel Structured Reasoning approach, that achieves more concise reasoning and stable performance across various scenarios, while enhancing the reasoning capabilities of Large Language Models (LLMs).
- We construct a structured reasoning dataset and demonstrate its effectiveness in enabling structured reasoning within LLMs through Supervised Fine-Tuning (SFT).
- We integrate MAX-Flow and Longest Common Subsequence (LCS) algorithms into Group Relative Policy Optimization (GRPO), reducing computational complexity while maintaining reasoning accuracy.
- Leveraging the interpretability advantages of structured reasoning, we present a novel approach to evaluate reasoning step metrics, while revealing the division of reasoning interests across neural layers, providing valuable insights for future LLM research.

2 Related Work

Reinforcement Learning Helps Efficiency Improvement Recent approaches use RL to improve reasoning efficiency, from basic length penalties [8–10] to more sophisticated methods. L1 [11] embeds length constraints in training instructions, while O1-Pruner [12] balances brevity and accuracy against reference benchmarks. DAST [13] introduces adaptive reasoning through token-length budget, allocating resources based on problem complexity. THINKPRUNE [14] employs a length-aware reward function with progressively tightening constraints, while Think When You Need [15] uses comparative rewards to guide models toward concise yet effective solutions.

Efficient CoT According to Perplexity Several works optimize reasoning chains using perplexity-based methods [16], including stepwise refinement [17], token pruning [18], attack detection [19] and step elimination strategies [20]. Furthermore, [21] proposes exploration based on entropy for multistep reasoning. Our research reveals that perplexity metrics inadequately assess the importance of reasoning steps, demonstrating that our MAX-Flow algorithm outperforms perplexity-based approaches in evaluating the importance of reasoning steps.

Language Model Reasoning (for Math) Since OpenAI-O1 [22], followed by O3 [23] and DeepSeek-R1 [24], researchers have proposed increasingly sophisticated RL algorithms, including LCPO [11], REINFORCE++ [25], DAPO [26], DPO-VP [27], VinePPO [28], CPPO [29], VAPO [30], and GRO [31]. Empirical investigations have explored data scaling [32], curriculum strategies [33, 34], and reward engineering [35–37]. Recent evaluations [38] show many reported improvements fail against properly optimized baselines. Our SR model incorporates structured reasoning through carefully designed SFT and reward approaches, across multiple seeds to ensure reproducibility.

3 Method

3.1 Structured Supervised Fine-Tuning

We transform unstructured data into structured data by incorporating explicit reasoning tags or labels that clearly indicate the current step of the reasoning process. Initially, we establish a structured reasoning format and create a set of reasoning step annotations, followed by supervised fine-tuning.

Structured Reasoning Prompt Given a question, we prompt the model as follows:

```
Please use the following tags at the beginning of each
sentence in your reasoning: <tag_1>, ..., <tag_K> QUESTION.
```

The model shows tagged reasoning, with unrestricted final answers:

```
<think> <tag_i> STEP1 (optional <\tag_i>) <tag_j> STEP2 ...
</think> ANSWER.
```

We first analyze the distribution of the tags used by DeepSeek-R1 671B([1]). After removing low-frequency tags, we retain 23 tags; the full prompts for tasks can be found in A.2.

Supervised Fine-Tuning We fine-tune the model to produce structured reasoning based on our designed prompt. For each Question-Reasoning-Answer triplet (q, r, a) in the chosen dataset $\mathcal{D}_{\text{fine}}$, the question q is combined with our structured reasoning prompt I , which guides the model to employ designated reasoning tags at the beginning of each sentence. Subsequently, the model is trained to create the associated structured reasoning r and provide the answer a . The fine-tuned model, SR_{sft} , is trained as follows:

$$SR_{\text{sft}} = \prod_{(q, r, a) \in \mathcal{D}_{\text{fine}}} P(r, a | q, I), \quad (1)$$

where I denotes our structured reasoning prompt and $\mathcal{D}_{\text{fine}}$ is the set of selected high-quality samples.

3.2 GRPO Rewards for Improved Structured Reasoning

To further enhance the model’s structured reasoning capabilities, we introduce two reasoning rewards to the GRPO (Group Relative Policy Optimization). We also truncated long outputs and shuffled prompt tags during training. The details can be found in A.7.

MAX-Flow Reward To evaluate the importance of the reasoning step, we propose a reward metric based on the max-flow / min-cut theory ([39]).

1. Step-to-Step Attention Matrix. Given the attention tensor $\mathcal{A} \in \mathbb{R}^{H \times L_{\text{seq}} \times L_{\text{seq}}}$ from a certain layer, where H denotes attention heads and L_{seq} the sequence length, we compute the normalized step attention matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ for reasoning steps n . For step i (token range $[s_i^{\text{start}}, s_i^{\text{end}}]$) step j :

$$A_{ij} = \frac{1}{H\mathcal{T}_i} \sum_{h=1}^H \sum_{a \in \mathcal{T}_i} \max_{b \in \mathcal{T}_j} \mathcal{A}_{h,a,b}, \quad (2)$$

where $\mathcal{T}_k \equiv s_k^{\text{end}} - s_k^{\text{start}}$ denotes the token length of step k .

2. Graph-Based Flow Analysis. Construct directed graph $G = (V, E)$. Nodes $V = \{1, \dots, n\}$ representing steps (node 1: Question, node n : Answer); Edges $(i, j) \in E$ with capacity A_{ij} when $A_{ij} > \tau$ (threshold $\tau = 0.05$).

3. Importance Scoring. For source $s = 1$ and target $t = n$: Compute max-flow F in G using Ford-Fulkerson algorithm. For each node $k \in V \setminus \{s, t\}$, $\Delta F_k = F - F_{-k}$, where F_{-k} is the max-flow in subgraph G_{-k} (node k removed). The value ΔF_k quantifies how crucial step k is for reaching the conclusion.

4. Robustness Evaluation. The reasoning quality metric $Q \in [0, 1]$ is computed as:

$$Q = 1 - \frac{\sum_{k \in \mathcal{K}_{\text{top}}} \Delta F_k}{\sum_{j=0}^{n-1} \Delta F_j}, \quad (3)$$

where \mathcal{K}_{top} contains the top-25% most important steps. Higher Q indicates more balanced reasoning.

The Longest Common Subsequence Reward Formally, given a set of reasoning completions $\mathcal{R} = \{r_1, \dots, r_n\}$ for a question, let $r_{\text{acc}}(r_i)$ denote the correctness reward for reasoning completion r_i . For each pair (r_i, r_j) , we extract their reasoning steps and compute the longest common subsequence (LCS) of reasoning tags, denoted $\text{LCS}(r_i, r_j)$.

Let L_{lcs} be the total length of the matched steps in the LCS and L_i be the total length of steps in r_i .

To prevent *length hacking* (i.e., artificially increasing the token count of each reasoning step to inflate scores), for each matched step k in the LCS with lengths $\ell_{i,k}$ and $\ell_{j,k}$, we introduce a length suppression factor, defined as $\text{ratio}_k = \frac{\ell_{j,k}}{2\ell_{i,k}}$ if $\ell_{i,k} > \ell_{j,k}$ and $\text{ratio}_k = 1 - \frac{\ell_{i,k}}{2\ell_{j,k}}$ otherwise. Subsequently, the length of the weighted LCS is defined as $L_{\text{lcs}} = \sum_{k \in \text{LCS}(r_i, r_j)} \text{ratio}_k \cdot \ell_{i,k}$.

We define the pairwise LCS score as:

$$\text{Score}_{\text{lcs}}(r_i, r_j) = \begin{cases} \frac{L_{\text{lcs}}}{L_i}, & \text{if both } r_i \text{ and } r_j \text{ are correct,} \\ -\frac{L_{\text{lcs}}}{L_i}, & \text{if both } r_i \text{ and } r_j \text{ are incorrect,} \\ 1 - \frac{L_{\text{lcs}}}{L_i}, & \text{if } r_i \text{ is correct, } r_j \text{ is incorrect,} \\ -1 + \frac{L_{\text{lcs}}}{L_i}, & \text{if } r_i \text{ is incorrect, } r_j \text{ is correct.} \end{cases} \quad (4)$$

Here, a higher LCS ratio is rewarded when compared with correct completions (encouraging consensus on high-quality reasoning), while a lower LCS ratio is rewarded when compared with incorrect completions (encouraging diversity from incorrect reasoning). The length suppression factor ratio_k penalizes unnecessarily long steps and encourages concise reasoning.

Finally, the overall LCS reasoning reward for c_i is averaged over all other completions:

$$r_{\text{lcs}}(c_i) = \frac{1}{n-1} \sum_{j \neq i} \text{Score}_{\text{lcs}}(c_i, c_j). \quad (5)$$

4 Experiments

Datasets. In the first stage, we use the S1 dataset [40], which contains 1,000 high-quality and diverse problems with reference answers and reasoning traces, covering science, technology, engineering and mathematics (STEM) and related domains. For both structured and unstructured reasoning fine-tuning data, we select the 500 most challenging problems (with the lowest accuracy) and generate correct samples from DeepSeek R1. In the second stage, we fine-tune on the DeepScaleR-Preview-Dataset [41], a mathematics dataset containing 40K question-answer pairs drawn from AIME, AMC, Omni-Math [42], and STILL [43].

Evaluation. We evaluate our models on the test sets of nine reasoning datasets: AIME 2024, AIME 2025, AMC, MATH500 [44], Minerva, Olympiad-Bench [45], and out-of-domain benchmarks including GPQA-Diamond [46], LSAT-AR [47], MMLU-ALL-VALID [48]. Following the recommendations of Sober Reasoning [38], we use three different random seeds for sampling and evaluation on large datasets MATH500, Minerva, and Olympiad-Bench, reporting both the mean and variance. For small datasets, we use ten different seeds to ensure statistical reliability.

Models. Our base model is DeepSeek-R1-Distill-Qwen-1.5B [49]. Due to computational constraints, we restrict the maximum context length to 4K tokens during supervised fine-tuning (SFT) using the 500 samples from the S1 dataset, obtaining the **SR-SFT** model.

We then derive three variants through reinforcement learning fine-tuning:

- **SR-ACC**: trained for 250 steps using the standard accuracy reward.
- **SR-LCS**: trained for 250 steps using the LCS reward (Equation 3.2).
- **SR-FLOW**: trained for 250 steps using our MAX-Flow reward (Equation 3.2).

The structured reasoning example can be found in A.5. The experimental setting can be found in A.7.

Baselines. We compare our proposed methods with the following baselines: Qwen2.5-Math-1.5B, Qwen2.5-Math-1.5B-Instruct [50] Qwen2.5-Math-1.5B-Oat-Zero [51] These models are all initialized from Qwen2.5 Math 1.5B and subsequently finetuned via reinforcement learning. DeepSeek-R1-Distill 1.5B [49], FastCurl-1.5B-Preview [52], DeepScaleR-1.5B [41], II-1.5B-Preview, STILL-3-1.5B [53], OpenRS1-1.5B, OpenRS2-1.5B, and OpenRS3-1.5B [54], L1-Qwen-1.5B-Max, L1-Qwen-1.5B-Exact [55]). Note that these models and our model are all initialized from DeepSeek-R1-Distill-1.5B and subsequently fine-tuned via reinforcement learning (e.g. GRPO [56]) to enhance mathematical reasoning capabilities.

5 Results

5.1 Main Experiment Results.

In this subsection, we evaluate the effectiveness of our proposed methods by reporting Pass@1 accuracy (mean \pm standard deviation) across six math benchmarks using standardized evaluation protocols. For AIME24, AIME25 and AMC23, we perform evaluations in 10 seeds each, while MATH500, Minerva and OlympiadBench are evaluated in 3 seeds each. We report two aggregate metrics: the average score across all six benchmarks (Avg. score) and more stable large average score (Large Avg.) computed from the three benchmarks with lower variance (MATH500, Minerva, and OlympiadBench). We perform comprehensive comparisons between the base model, our SR model, and the standard SFT model across both in-domain and out-of-domain datasets (MATH500, Minerva, OlympiadBench, GPQA-Diamond, LSAT-AR, MMLU-ALL-VALID) using 8K context length. Additionally, for correct completions of MATH500, we analyze the average token length and the average number of reasoning steps to characterize the reasoning process.

Table 1: Training Details. To ensure consistency in counting training steps, we standardized the batch size to 128. This means that two steps with a batch size of 64 are considered equivalent to one step with a batch size of 128. For SR, we performed 5 epochs of SFT followed by 250 steps of GRPO.

Model	Training Steps	Training Stages	Number of GPUs Used in Each Stage
SR	SFT+RL(~ 23)	2	1, 4
FastCuRL	RL(~ 860)	4	8, 8, 8, 8
DeepScaleR	RL($\sim 1,750$)	3	8, 16, 32

Competitive Performance at Minimal Cost. After fine-tuning, our model achieves a high performance as shown in Table 2, improving the average score from 48.4 to 50.7 and the large benchmark average from 55.9 to 56.7. Building upon our MAX Flow Reward method (Section 3.2), we further trained the model using GRPO, achieving additional performance gains, increasing the average score from 50.7 to 52.6 and the large benchmark average from 56.7 to 58.1. Our experiments validate that structured reasoning LLMs can rapidly approach the performance of current state-of-the-art models (trained with multistage 8k-64k long-context large-scale datasets; see Table 1) when trained efficiently with: (1) limited amounts of high-quality structured data and (2) reward functions specifically designed for structured reasoning.

Table 2: Benchmark Results (Pass@1 Accuracy). All results are reported as mean \pm standard deviation. Avg. score calculates the average across all six benchmarks, while Large Avg. focuses on the more stable MATH500, Minerva, and Olympiad benchmarks. Top-3 models in each category are highlighted with increasing gray intensity.

Model	AIME'24	AIME'25	AMC'23	MATH500	Minerva	Olympiad	Avg.	Large Avg.
Based on: Qwen2.5-Math-1.5B (RL)								
Math	11.3 \pm 3.6	5.7 \pm 2.7	44.0 \pm 4.9	51.7 \pm 5.5	11.3 \pm 2.2	26.0 \pm 0.6	25.0 \pm 3.3	29.7 \pm 2.8
Oat-Zero	16.0 \pm 3.2	6.7 \pm 3.4	52.5 \pm 2.9	73.5 \pm 1.7	26.3 \pm 0.8	37.2 \pm 1.3	32.0 \pm 2.2	45.7 \pm 1.3
Math	12.0 \pm 1.7	11.7 \pm 5.7	54.8 \pm 5.3	74.7 \pm 0.5	26.7 \pm 1.8	37.9 \pm 0.2	36.3 \pm 2.5	46.4 \pm 0.8
Based on: Deepseek-R1-Distill-Qwen-1.5B (RL)								
R1-Distill	28.7 \pm 4.8	22.3 \pm 5.2	71.5 \pm 3.9	84.9 \pm 0.3	30.5 \pm 1.0	52.4 \pm 0.4	48.4 \pm 2.6	55.9 \pm 0.6
L1-Exact	24.4 \pm 3.3	22.3 \pm 4.2	70.5 \pm 3.7	86.6 \pm 0.8	31.5 \pm 1.7	52.5 \pm 1.3	47.9 \pm 2.5	56.9 \pm 1.3
L1-Max	27.7 \pm 4.2	21.0 \pm 5.0	73.2 \pm 6.0	84.7 \pm 0.1	33.3 \pm 0.9	52.3 \pm 0.6	48.7 \pm 2.8	56.8 \pm 0.5
Open-RS1	28.9 \pm 6.0	21.3 \pm 4.2	75.0 \pm 3.3	85.1 \pm 0.8	30.4 \pm 0.2	53.2 \pm 1.9	49.0 \pm 2.7	56.2 \pm 1.0
Open-RS2	31.3 \pm 7.7	22.7 \pm 5.6	73.0 \pm 5.7	84.1 \pm 0.2	29.2 \pm 1.1	53.7 \pm 0.6	49.0 \pm 3.5	55.7 \pm 0.6
Open-RS3	29.7 \pm 4.6	24.7 \pm 6.5	69.2 \pm 5.5	84.2 \pm 1.1	28.6 \pm 2.3	51.8 \pm 0.8	48.0 \pm 3.5	54.9 \pm 1.4
STILL-3	34.7 \pm 5.5	24.0 \pm 6.4	72.5 \pm 5.4	86.6 \pm 1.9	30.0 \pm 0.6	53.9 \pm 1.5	50.3 \pm 3.6	56.8 \pm 1.3
II-Thought	32.0 \pm 5.9	24.0 \pm 4.1	79.5 \pm 5.1	86.6 \pm 0.6	31.7 \pm 0.6	54.9 \pm 0.4	51.5 \pm 2.8	57.7 \pm 0.5
FastCuRL	36.3 \pm 4.3	27.0 \pm 3.7	78.8 \pm 4.1	87.9 \pm 1.2	30.8 \pm 1.4	56.5 \pm 0.6	52.9 \pm 2.6	58.4 \pm 1.1
DeepScaleR	37.0 \pm 6.6	30.3 \pm 4.3	76.2 \pm 4.6	87.8 \pm 1.0	31.0 \pm 1.5	55.5 \pm 1.1	53.0 \pm 3.2	58.1 \pm 1.2
Ours Based on: Deepseek-R1-Distill-Qwen-1.5B (RL)								
SR-SFT	32.7 \pm 8.7	25.3 \pm 8.0	75.8 \pm 6.7	85.6 \pm 1.6	31.3 \pm 2.8	53.3 \pm 1.9	50.7 \pm 4.9	56.7 \pm 2.1
SR-FLOW	36.7 \pm 8.9	27.0 \pm 8.2	77.8 \pm 6.6	85.3 \pm 1.6	34.2 \pm 2.9	54.9 \pm 1.9	52.6 \pm 5.0	58.1 \pm 2.1

Reward Analysis During Training. Figure 2 and Table 6 presents the learning dynamics of our models during RL fine-tuning, with evaluations conducted at 50-step intervals. The curves demonstrate the evolution of: (1) SR model performance (both Avg. and Large Avg. scores), (2) the average number of reasoning steps in correct completions on MATH500, and (3) the average token length of solutions on MATH500.

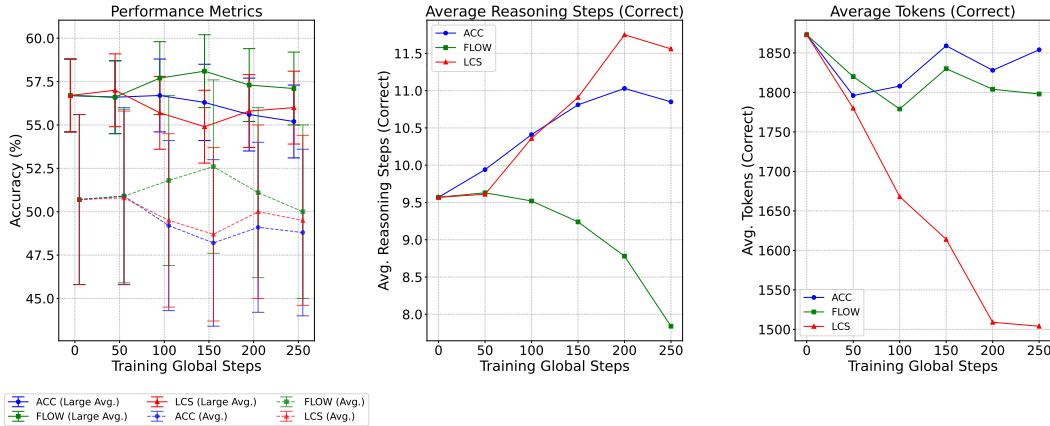


Figure 2: Comparison across Global Steps: Pass@1 Performance (Left), Reasoning Steps (Middle), and Average Tokens (Right) for **SR-ACC**(standard reward), **SR-FLOW**, and **SR-LCS**.

Reasoning Efficiency Improvements. Our proposed MAX-Flow reward and LCS reward offer different advantages in the trade-off between higher performance and better reasoning efficiency.

As shown in Figure 2 and Table 6, SR-FLOW achieves **Higher Accuracy**: +1.9% average gain on the "Larger Avg." metric at 150 steps compared to SR-ACC. **More Concise Reasoning**: Steps reduced from 9.57 to 7.84 while maintaining 84.7% accuracy on MATH500. We observed that although we did not explicitly require the model to reduce the reasoning steps or token count, SR-FLOW spontaneously began to perform the fusion of the reasoning steps.

SR-LCS achieves **Token Efficiency**: 1504 tokens vs. 1873 (baseline) with a slight drop in Large Avg. from 56.7 to 56.0. Figure 7 shows more responses in the 256-1024 token range and fewer that exceed 8192 tokens, indicating more efficient reasoning. With accuracy rates of 84.4%, 84.3%, and 84.5% for LCS-50, LCS-150, and LCS-250 respectively, the model maintains performance while improving efficiency under 8k context length constraints.

Table 3: Comprehensive Comparison on In-domain and Out-of-domain Data with 8k Context Length.

Method	MATH500	Minerva	OlympiaBench	GPQA	LSAT-AIR	MMLU
R1-Distill-1.5B (Base)	80.33% \pm 1.78%	31.00% \pm 2.81%	44.49% \pm 1.91%	35.02% \pm 3.40%	26.26% \pm 2.90%	50.60% \pm 1.28%
SFT (Standard)	82.00% \pm 1.72%	28.68% \pm 2.75%	48.84% \pm 1.93%	33.66% \pm 3.34%	25.87% \pm 2.89%	44.83% \pm 1.27%
SR-SFT	84.53% \pm 1.62%	32.11% \pm 2.84%	49.53% \pm 1.93%	36.20% \pm 3.39%	26.00% \pm 2.89%	46.66% \pm 1.28%
SR-Flow (150 steps)	84.74% \pm 1.58% +0.25%	33.52% \pm 2.88% +4.39%	50.37% \pm 1.92% +1.70%	36.70% \pm 3.43% +1.38%	26.78% \pm 2.92% +3.00%	46.81% \pm 1.28% +0.32%

We also conducted a comprehensive evaluation on both in-domain and out-of-domain benchmarks (Table 3). During the reinforcement learning stage, our SR-Flow outperforms SR-SFT across all benchmarks (ranging from +0.25% to +4.39%), validating the universal advantage of our approach in enhancing structured reasoning capabilities. As expected, the SFT model shows improvements on in-domain tasks but degradation on out-of-domain benchmarks due to domain shifting.

5.2 What Are the Advantages of Structured Reasoning?

Structured Reasoning Naturally Produces More Concise Token Length Distributions. For large models, Table 4 summarizes the average token lengths and accuracies across several benchmarks, including MATH500, GPQA-Diamond, MMLU-ALL-VALID, AMC23, and AIME24. Notably, the structured reasoning model achieves similar or higher accuracy with much shorter answers, e.g., on MATH500, the average reasoning token length drops from 2945 (base) to 1577 (SR-Prompt), while accuracy remains above 92% (See Appendix A.8 for the fill in the middle prompting strategy).

Table 4: Token length and accuracy comparison for DeepSeek-R1 671B on several benchmarks.

Model	MATH500 Acc./Len.	GPQA Acc./Len.	MMLU Acc./Len.	AMC23 Acc./Len.
Base	92.9%/2945	70.3%/6537	88.8%/989	100%/1716
SR-Prompt	93.0%/1577 -46.5%	71.1%/4028 -38.4%	89.6%/512 -48.2%	100%/2053 +19.6%

For small models, we constructed structured and unstructured reasoning datasets using the same problems from s1k for fine-tuning and analyzed the model outputs truncated at 2k, 4k, and 8k tokens. Figure 3a shows that models fine-tuned on structured data (SR-SFT) have better performance across all token limits (2k, 4k, and 8k) compared to standard SFT. Specifically, at the 8k limit, the unstructured model had 17.0% of responses truncated with an average correct reasoning length of 2058 tokens, while the standard SFT model had only 10% truncated with an average length of 1735 tokens, demonstrating a natural tendency to produce more concise answers.

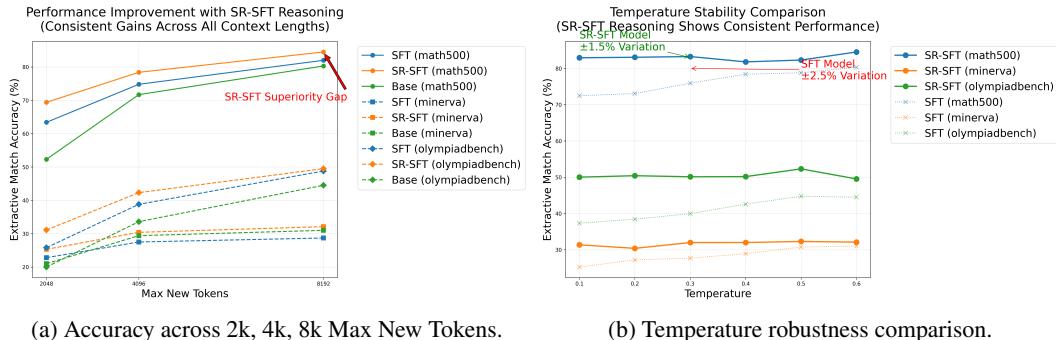


Figure 3: Performance Comparison of 1.5B Models across Maximum New Tokens (Left) and Temperature Settings (Right), averaged over three seeds.

Structured Reasoning Models Produce More Stable Outputs. As shown in Figure 3b, we observed an interesting phenomenon across experiments conducted with 3 seeds per task using an 8k context length. When the sampling temperature increases from 0 to 0.6 (0.6 is the recommended temperature from DeepSeek), the unstructured SFT model's accuracy rises significantly on all datasets: MATH500 (from 72.47 to 80.33), Minerva (from 25.25 to 31.00), and OlympiadBench (from 37.33 to 44.49). In contrast, our SR-SFT model maintains more consistent accuracy across temperature values: MATH500 (from 82.93 to 84.53), Minerva (from 31.37 to 32.11), and OlympiadBench (from 50.02 to 49.53). This temperature sensitivity is much more pronounced in unstructured reasoning models, while structured reasoning models demonstrate stability across the temperature range.

Table 5: Comparing Trigger Counts and Distances to First Correct Answer Across Different Methods.

Trigger Type	Settings	Avg. Trigger Count ↓	Avg. Distance to First Correct Answer (tokens) ↓
Top tags	"verify", "summarize", etc.	2.02	78.01
Token chunks	128-token intervals	3.93	131.05
Keywords	"but", "wait", "however", etc.	2.69	139.97

Improved Compatibility with Test-time Scaling and Early Stopping. For Test-time Scaling, existing work extends model outputs by injecting prompt tokens at thought-stopping points. Our method simplifies this by guiding outputs through the most likely next tag at stopping points. For early stopping, our tag-based approach outperforms traditional methods. In our experiment with 705 correct MATH500 reasoning completions (Table 7), we compared interval-based (128-token), keyword-based ("but", "wait", "however", etc.), and tag-based ("verify", "summarize", etc.) detection strategies. As Table 5 shows, our structured approach reduces average Probe-In-Middle interventions to just 2.02 while maintaining closest proximity to correct answers (78.01 tokens).

5.3 What Are the Gains from Structured Reasoning Models?

Structured Analysis Helps Identify Redundant Reasoning Steps. We designed an experiment (Interference Injection and Selective Removal, IISR) (A.6) to assess whether we can better analyze the importance of reasoning steps. Since existing datasets rarely provide direct importance annotations for reasoning steps, and LLM-based scoring is noisy, we injected obviously irrelevant reasoning steps into existing chains. While we cannot confirm the relative importance of original steps, we are certain about the irrelevance of the injected ones. We compared our max-flow algorithm (3.2), top-p/top-k backtracking (A.3), average step perplexity (A.3), and random selection for their Error Filtering Efficiency (A.4) when removing 1-11 steps from mixed reasoning chains.

This experiment used 70 correctly reasoned examples from S1k covering science, technology, engineering, and mathematics (STEM) and related domains with longer trace lengths and more uniform reasoning steps. We defined four types of interference prompts: (1) **Redundant** - statements like “Let’s summarize what we’ve done so far, our previous work is correct” that add no value to reasoning; (2) **Distracted** - comments indicating distraction such as “This reminds me of another problem”; (3) **Harmful** - randomly injected reasoning steps from other problems; and (4) **Confused** - copies of current reasoning steps randomly injected at incorrect positions in the reasoning chain.

Through extensive randomized experiments, we found that as more reasoning steps were removed, our proposed methods based on step-matrix (See Section 3.2) (top-k, top-p, and max-flow) significantly outperformed random removal. The specific example can be found in Appendix A.6.

Takeaway 1

Using average perplexity of reasoning steps as a reasoning step evaluation metric is **unreliable**.

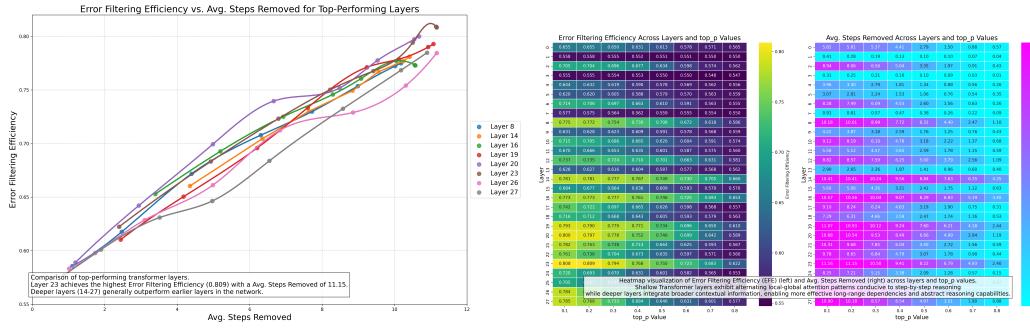
Additionally, in our comparison with perplexity-based algorithms, we found that removing steps with the lowest PPL (PPL Bottom) performed similarly (though slightly worse) to our methods when dealing with redundant but harmless information, as such information typically has low information content and low perplexity. Interestingly, for logically confused interference, removing steps with the highest PPL (PPL Top) performed slightly better, as steps appearing in inappropriate positions caused significantly increased perplexity. This shows that PPL primarily reflects information quantity and cannot distinguish valuable reasoning from disruptive content. Our step-matrix-based methods significantly outperformed PPL-based approaches across all tasks.

Reasoning Models may Exhibit Distinct Layer-wise Division of Reasoning Roles. Based on our Top-P reasoning step filtering algorithm, which showed excellent performance in the IISR experiment, we further compared the changes in Error Filtering Efficiency across 0-28 layers.

Takeaway 2

During the reasoning process, the earlier layers of small models **alternate** between focusing on local and global reasoning, while later layers predominantly attend to global reasoning.

As shown in the Figure 4b, when using the step matrix of layer 0 to calculate the Top-P algorithm with a retention rate of 0.1, an average of 5.82 steps can be removed, indicating that the attention span of this layer spans at least more than 5 steps. However, under the same conditions, layer 1 can only remove 0.41 reasoning steps, suggesting that this layer's attention primarily focuses on the previous step. This alternating pattern consistently appears up to layer 13. When using the step matrix from layers 14-27 for the top-P algorithm with a retention rate of 0.1, the number of steps removed is consistently greater than 8, reaching a maximum of 11.06. This shows that later layers are interested in information from a wider range of reasoning steps. Similarly, we also found that these layers with broader attention spans can better describe the importance of reasoning steps (Figure 4a). These findings may inspire future work on pruning reasoning models.



(a) Comparison of Error Filtering Efficiency across different model layers, demonstrating how later layers better capture reasoning step importance. (b) Attention heatmap visualization showing the alternating local-global focus pattern in early layers (0-13) versus consistent global attention in later layers (14-27).

Figure 4: Layer-wise Analysis Using Step Matrix: Error Filtering Efficiency across Different Layers (Left) and Reasoning Step Intervals (Right) Based on Top-P Filtering Algorithm.

Analysis of Model Reasoning Patterns. In this section, we select 1k SR-FLOW model reasoning completions and construct a graph of directed edge relationship between the reasoning steps (Figure 9). We ignore edges with low frequency for clarity of presentation.

Takeaway 3

During the problem-solving phase, various reasoning **loops** of different sizes form around steps like <verify>, depending on the difficulty of the problem.

We find that when reasoning, SR-FLOW typically starts from problem definition (assumption), then proceeds through problem decomposition (decompose), moves to formalize, and finally reaches a conclusion through verify consequence. This represents a common trajectory. The paths involving assumption and summarize components are quite stable. However, we discovered multiple loop structures centered around verify: for example, verify-alternative, verify-caseAnalysis, and other complex loops. The model's different answers are often formed and corrected within these loops.

Figure 8 shows the distribution of the relative positions of each tag within the reasoning process.

6 Conclusion

In this paper, we present a novel approach to enhance LLMs' structured reasoning capabilities. Our method combines cognitive science principles with neurosymbolic AI advances by explicitly incorporating structured knowledge representations and reasoning processes. The key contributions include: (1) a structured data transformation method with explicit reasoning step annotations, (2) two reinforcement learning algorithms, MAX-Flow for reasoning effectiveness evaluation, and LCS for computational efficiency, and (3) a lightweight fine-tuning strategy using only 500 examples and 250 RL steps. Our experiments on the DeepSeek-R1-Distill-Qwen-1.5B model demonstrate that this approach achieves more concise reasoning while maintaining robust performance across various scenarios. These results suggest that incorporating explicit structure into LLM reasoning processes can lead to more efficient and reliable AI systems.

References

- [1] DeepSeek-AI, D. Guo, D. Yang, et al. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning, 2025. ArXiv:2501.12948 [cs].
- [2] OpenAI. OpenAI o1, 2024.
- [3] QwQ. QwQ-32B: Embracing the Power of Reinforcement Learning | Qwen, 2025.
- [4] Bronkhorst, H., G. Roorda, C. Suhre, et al. Students' use of formalisations for improved logical reasoning. *Research in Mathematics Education*, 2022.
- [5] Forstmann, B. U., R. Ratcliff, E.-J. Wagenmakers. Sequential sampling models in cognitive neuroscience. *Annual review of psychology*, 67:641–666, 2016.
- [6] Evans, J. S. B. Dual-process theories. In *The Routledge international handbook of thinking and reasoning*, pages 157–174. Routledge, 2018.
- [7] Miller, E. K., J. D. Cohen. An integrative theory of prefrontal cortex function. *Annual review of neuroscience*, 24(1):167–202, 2001.
- [8] Team, K., A. Du, B. Gao, et al. Kimi k1. 5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599*, 2025.
- [9] Li, C., N. Liu, K. Yang. Adaptive group policy optimization: Towards stable training and token-efficient reasoning. *arXiv preprint arXiv:2503.15952*, 2025.
- [10] Arora, D., A. Zanette. Training language models to reason efficiently. *arXiv preprint arXiv:2502.04463*, 2025.
- [11] Aggarwal, P., S. Welleck. L1: Controlling how long a reasoning model thinks with reinforcement learning. *arXiv preprint arXiv:2503.04697*, 2025.
- [12] Luo, H., L. Shen, H. He, et al. O1-pruner: Length-harmonizing fine-tuning for o1-like reasoning pruning. *arXiv preprint arXiv:2501.12570*, 2025.
- [13] Shen, Y., J. Zhang, J. Huang, et al. Dast: Difficulty-adaptive slow-thinking for large reasoning models. *arXiv preprint arXiv:2503.04472*, 2025.
- [14] Hou, B., Y. Zhang, J. Ji, et al. Thinkprune: Pruning long chain-of-thought of llms via reinforcement learning. *arXiv preprint arXiv:2504.01296*, 2025.
- [15] Yang, J., K. Lin, X. Yu. Think when you need: Self-adaptive chain-of-thought learning. *arXiv preprint arXiv:2504.03234*, 2025.
- [16] Jelinek, F., R. L. Mercer, L. R. Bahl, et al. Perplexity—a measure of the difficulty of speech recognition tasks. *The Journal of the Acoustical Society of America*, 62(S1):S63–S63, 2005.
- [17] Cui, Y., P. He, J. Zeng, et al. Stepwise perplexity-guided refinement for efficient chain-of-thought reasoning in large language models, 2025.
- [18] Xia, H., Y. Li, C. T. Leong, et al. Tokenskip: Controllable chain-of-thought compression in llms, 2025.
- [19] Alon, G., M. Kamfonas. Detecting language model attacks with perplexity, 2023.
- [20] Liu, T., Q. Guo, X. Hu, et al. Can language models learn to skip steps?, 2024.
- [21] Zhang, J., X. Wang, F. Mo, et al. Entropy-based exploration conduction for multi-step reasoning, 2025.
- [22] Jaech, A., A. Kalai, A. Lerer, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
- [23] OpenAI. OpenAI o3-mini System Card, 2025.
- [24] DeepSeek-AI. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning, 2025.
- [25] Hu, J. Reinforce++: A simple and efficient approach for aligning large language models. *arXiv preprint arXiv:2501.03262*, 2025.
- [26] Yu, Q., Z. Zhang, R. Zhu, et al. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*, 2025.

- [27] Tu, S., J. Lin, X. Tian, et al. Enhancing llm reasoning with iterative dpo: A comprehensive empirical investigation. *arXiv preprint arXiv:2503.12854*, 2025.
- [28] Kazemnejad, A., M. Aghajohari, E. Portelance, et al. Vineppo: Unlocking rl potential for llm reasoning through refined credit assignment. *arXiv preprint arXiv:2410.01679*, 2024.
- [29] Lin, Z., M. Lin, Y. Xie, et al. Cppo: Accelerating the training of group relative policy optimization-based reasoning models. *arXiv preprint arXiv:2503.22342*, 2025.
- [30] Yue, Y., Y. Yuan, Q. Yu, et al. Vapo: Efficient and reliable reinforcement learning for advanced reasoning tasks. *arXiv preprint arXiv:2504.05118*, 2025.
- [31] Cai, X. One framework to rule them all: Unifying rl-based and rl-free methods in rlhf. *arXiv preprint arXiv:2503.19523*, 2025.
- [32] Shen, W., G. Liu, Z. Wu, et al. Exploring data scaling trends and effects in reinforcement learning from human feedback. *arXiv preprint arXiv:2503.22230*, 2025.
- [33] Wen, L., Y. Cai, F. Xiao, et al. Light-r1: Curriculum sft, dpo and rl for long cot from scratch and beyond. *arXiv preprint arXiv:2503.10460*, 2025.
- [34] Roux, N. L., M. G. Bellemare, J. Lebensold, et al. Tapered off-policy reinforce: Stable and efficient reinforcement learning for llms. *arXiv preprint arXiv:2503.14286*, 2025.
- [35] Gao, J., S. Xu, W. Ye, et al. On designing effective rl reward at training time for llm reasoning. *arXiv preprint arXiv:2410.15115*, 2024.
- [36] Cui, G., L. Yuan, Z. Wang, et al. Process reinforcement through implicit rewards. *arXiv preprint arXiv:2502.01456*, 2025.
- [37] Ma, Y. J., W. Liang, G. Wang, et al. Eureka: Human-level reward design via coding large language models. *arXiv preprint arXiv:2310.12931*, 2023.
- [38] Hochlehnert, A., H. Bhatnagar, V. Udandarao, et al. A sober look at progress in language model reasoning: Pitfalls and paths to reproducibility, 2025.
- [39] Ford, L. R., D. R. Fulkerson. Maximal flow through a network. *Canadian Journal of Mathematics*, 8:399–404, 1956.
- [40] Muennighoff, N., Z. Yang, W. Shi, et al. s1: Simple test-time scaling, 2025.
- [41] Luo, M., S. Tan, J. Wong, et al. Deepscaler: Surpassing o1-preview with a 1.5b model by scaling rl. <https://pretty-radio-b75.notion.site/DeepScaleR-Surpassing-01-Preview-with-a-1-5B-Model-by-Scaling-RL-19681902c1468005bed8ca303013a4e2>, 2025. Notion Blog.
- [42] Gao, B., F. Song, Z. Yang, et al. Omni-math: A universal olympiad level mathematic benchmark for large language models, 2024.
- [43] Song, H., J. Jiang, Y. Min, et al. R1-searcher: Stimulating the search capability of llm from zero via reinforcement learning. 2025.
- [44] Hendrycks, D., C. Burns, S. Kadavath, et al. Measuring mathematical problem solving with the math dataset. *NeurIPS*, 2021.
- [45] He, C., R. Luo, Y. Bai, et al. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems, 2024.
- [46] Rein, D., B. L. Hou, A. C. Stickland, et al. Gpqa: A graduate-level google-proof q&a benchmark, 2023.
- [47] Zhong, W., R. Cui, Y. Guo, et al. Agieval: A human-centric benchmark for evaluating foundation models, 2023.
- [48] Hendrycks, D., C. Burns, S. Basart, et al. Measuring massive multitask language understanding, 2021.
- [49] DeepSeek-AI, D. Guo, D. Yang, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025.
- [50] Yang, A., B. Zhang, B. Hui, et al. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement, 2024.

- [51] Liu, Z., C. Chen, W. Li, et al. There may not be aha moment in r1-zero-like training — a pilot study. <https://oatllm.notion.site/oat-zero>, 2025. Notion Blog.
- [52] Song, M., M. Zheng, Z. Li, et al. Fastcurl: Curriculum reinforcement learning with progressive context extension for efficient training r1-like reasoning models, 2025.
- [53] Min, Y., Z. Chen, J. Jiang, et al. Imitate, explore, and self-improve: A reproduction report on slow-thinking reasoning systems, 2024.
- [54] Dang, Q.-A., C. Ngo. Reinforcement learning for reasoning in small llms: What works and what doesn't, 2025.
- [55] Aggarwal, P., S. Welleck. L1: Controlling how long a reasoning model thinks with reinforcement learning, 2025.
- [56] Shao, Z., P. Wang, Q. Zhu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024.

A Appendix

A.1 Part of Figures and Tables

For better layout and presentation, we have placed some figures and tables in a unified location in the Appendix.

Table 6 presents a comprehensive comparison of three Small Structure Reasoning (SR) methods across various mathematical benchmarks. SR-FLOW demonstrates superior performance, achieving the highest average accuracy (58.1%) while requiring fewer reasoning steps. SR-LCS offers the most token-efficient approach, using approximately 20% fewer tokens while maintaining competitive accuracy. Highlighted cells indicate top performances for each benchmark and method, showing that different reasoning approaches excel in different problem domains.

Table 6: Comparative Analysis of Reasoning Methods (Pass@1 Accuracy & Efficiency)

Method	Accuracy (%)							Steps Avg.	Tokens Avg.
	AIME24	AIME25	AMC23	MATH500	Minerva	Olympiad	Large Avg.		
SR-ACC	32.7±8.7	25.3±8.0	75.8±6.7	85.6±1.6	31.3±2.8	53.3±1.9	56.7±2.1	9.57	1873
	36.7±8.9	26.7±8.2	72.0±7.1	85.5±1.6	31.1±2.8	53.3±1.9	56.6±2.1	9.94	1796
	30.0±8.5	21.3±7.6	74.0±7.0	84.7±1.6	32.4±2.8	52.9±1.9	56.7±2.1	10.41	1808
	24.0±7.8	22.5±7.7	74.0±7.0	83.7±1.7	33.0±2.9	52.2±1.9	56.3±2.2	10.81	1859
	30.3±8.7	24.2±7.8	73.5±6.6	84.2±1.6	31.7±2.8	50.8±1.9	55.6±2.1	11.03	1828
	36.7±7.7	19.5±7.2	70.6±7.4	83.7±1.7	31.1±2.8	50.9±1.9	55.2±2.1	10.85	1854
SR-FLOW	32.7±8.7	25.3±8.0	75.8±6.7	85.6±1.6	31.3±2.8	53.3±1.9	56.7±2.1	9.57	1873
	33.0±8.5	26.0±8.1	76.5±6.8	84.9±1.6	31.3±2.8	53.6±1.9	56.6±2.1	9.63	1820
	34.7±8.3	26.3±7.8	76.8±6.8	85.6±1.6	34.1±2.9	53.5±1.9	57.7±2.1	9.52	1779
	36.7±8.9	27.0±8.2	77.8±6.6	85.3±1.6	34.2±2.9	54.8±1.9	58.1±2.1	9.24	1830
	33.5±8.2	25.7±7.6	75.3±6.9	85.0±1.6	33.2±2.9	53.7±1.9	57.3±2.1	8.78	1804
	30.3±8.7	24.2±7.8	74.0±7.0	84.7±1.6	32.3±2.8	54.3±1.9	57.1±2.1	7.84	1798
SR-LCS	32.7±8.7	25.3±8.0	75.8±6.7	85.6±1.6	31.3±2.8	53.3±1.9	56.7±2.1	9.57	1873
	34.0±8.7	24.7±7.9	75.0±6.9	84.9±1.6	32.0±2.8	54.1±1.9	57.0±2.1	9.61	1780
	33.3±8.7	22.5±7.8	74.0±7.0	84.4±1.6	30.9±2.8	51.8±1.9	55.7±2.1	10.36	1668
	30.3±8.5	23.3±7.8	74.0±7.0	83.3±1.7	29.9±2.8	51.6±1.9	54.9±2.1	10.91	1614
	33.7±8.7	23.7±7.9	75.0±6.9	84.9±1.6	31.8±2.8	50.6±1.9	55.8±2.1	11.75	1509
	31.0±8.5	23.3±7.7	74.8±7.1	84.8±1.6	30.5±2.8	52.7±1.9	56.0±2.1	11.56	1504

This table demonstrates the evolution of various metrics during GRPO training from 0 to 250 global steps under three reward functions: ACC, FLOW, and LCS. The metrics are tracked throughout the training process to show how different reward mechanisms influence model performance.

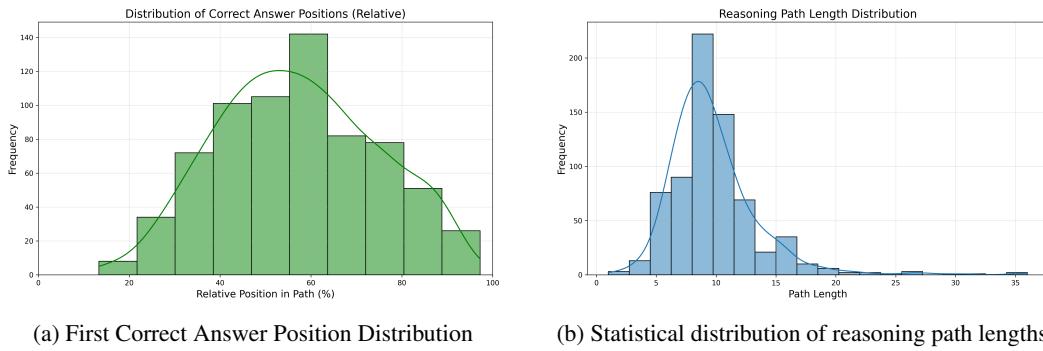


Figure 5: Analysis of Model Reasoning Patterns: Distribution of First Correct Answers (Left) and Reasoning Path Lengths (Right).

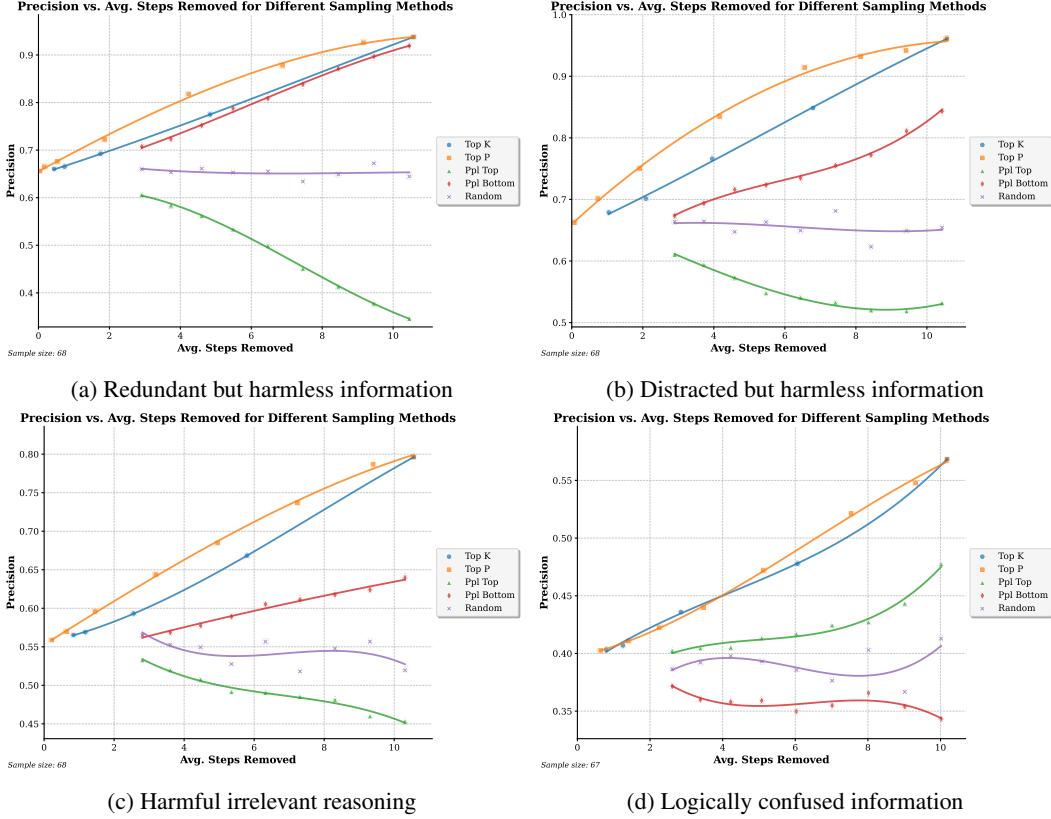


Figure 6: IISR Results: Error Filtering Efficiency (Precision) of different algorithms when removing 1-11 steps under four types of information interference.

Table 7: Early Stopping Detection Parameters and Sample Statistics

Metric	Value
Number of valid samples	705
Top five frequent tags	verify, summarize, equivalent, formalize, consequence
Token interval	128
Top five useful words	but, wait, however, check, alternatively

A.2 Full Prompts

Mathematical Problem Solving Template

Please use the following tags at the beginning of each sentence in your reasoning: <rephrase>, <inference>, <analogy>, <equivalent>, <association>, <reverse>, <summarize>, <verify>, <complete>, <decompose>, <counterexample>, <assumption>, <constraint>, <case_analysis>, <contradiction>, <abstraction>, <formalize>, <generalize>, <specialize>, <critique>, <alternative>, <consequence>, <intuition>.

{Question}

Please reason step by step, and put your final answer within boxed{ }.

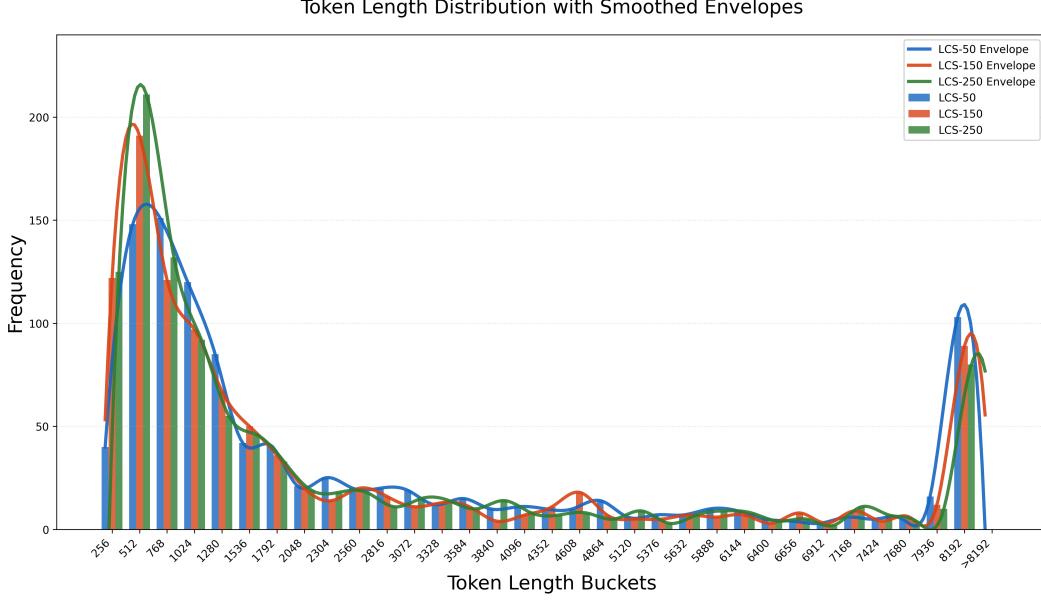


Figure 7: Token length distribution of LCS models under different training stages. The smoothed envelopes show how reward training shifts the distribution towards optimal reasoning lengths (512–1024 tokens) while maintaining performance.

Multiple Choice Problem Template

Please use the following tags at the beginning of each sentence in your reasoning: <rephrase>, <inference>, <analogy>, <equivalent>, <association>, <reverse>, <summarize>, <verify>, <complete>, <decompose>, <counterexample>, <assumption>, <constraint>, <case_analysis>, <contradiction>, <abstraction>, <formalize>, <generalize>, <specialize>, <critique>, <alternative>, <consequence>, <intuition>.

{Question}

- A) {A}
- B) {B}
- C) {C}
- D) {D}

Please reason step by step, and answer the following multiple choice question. The last line of your response should be of the following format: 'Answer: \$LETTER' (without quotes) where LETTER is one of ABCD.

A.3 Other Step Importance Evaluation Algorithm Implementation

Top-P and Top-K Selection. Based on the step matrix computed from different layers (See Section 3.2), we implement backtracking selection methods:

$$\text{SelectSteps}(A, k, p) = \{s_i\}_{i=0}^m, \quad (6)$$

where $A \in \mathbb{R}^{n \times n}$ is the step attention matrix, and we select steps starting from the last step s_{n-1} by either: Top-K: For each step s_i , select up to k preceding steps with highest attention scores. Top-P: Select preceding steps with cumulative normalized attention exceeding threshold p .

The algorithm traverses backward from the final step, adding important preceding steps to a visited set based on attention weights, ensuring all critical reasoning dependencies are captured. **Average Perplexity.** For each step, we compute token-level perplexity:

$$\text{Perplexity}(t_i) = \frac{1}{P(t_i|x, t_1, \dots, t_{i-1})}, \quad (7)$$

where $P(t_i|x, t_1, \dots, t_{i-1})$ is the probability of token t_i given the prompt x and all preceding tokens, derived from the softmax of logits:

$$P(t_i|x, t_1, \dots, t_{i-1}) = \frac{\exp(\text{logits}_i)}{\sum_j \exp(\text{logits}_j)}. \quad (8)$$

The average perplexity for a step s containing tokens $\{t_1, t_2, \dots, t_m\}$ is:

$$\text{AvgPerplexity}(s) = \exp \left(-\frac{1}{m} \sum_{i=1}^m \log P(t_i|x, t_1, \dots, t_{i-1}) \right). \quad (9)$$

Random Selection. A baseline approach where steps are selected randomly without leveraging attention patterns or perplexity metrics.

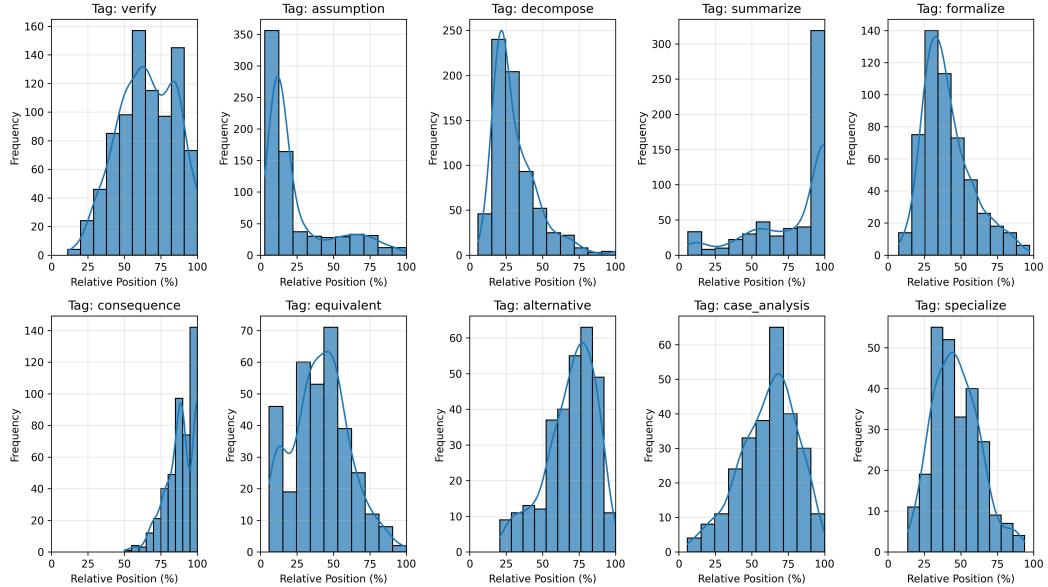


Figure 8: Distribution of the relative positions of each tag within the reasoning process.

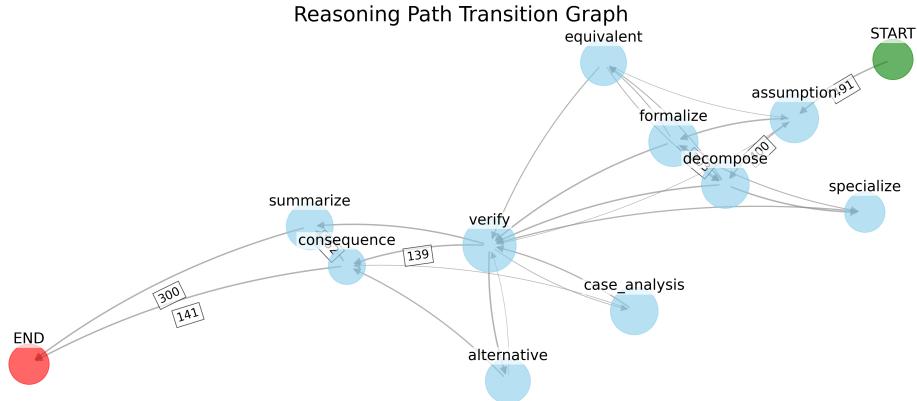


Figure 9: Illustration of Reasoning Path Transition Graph.

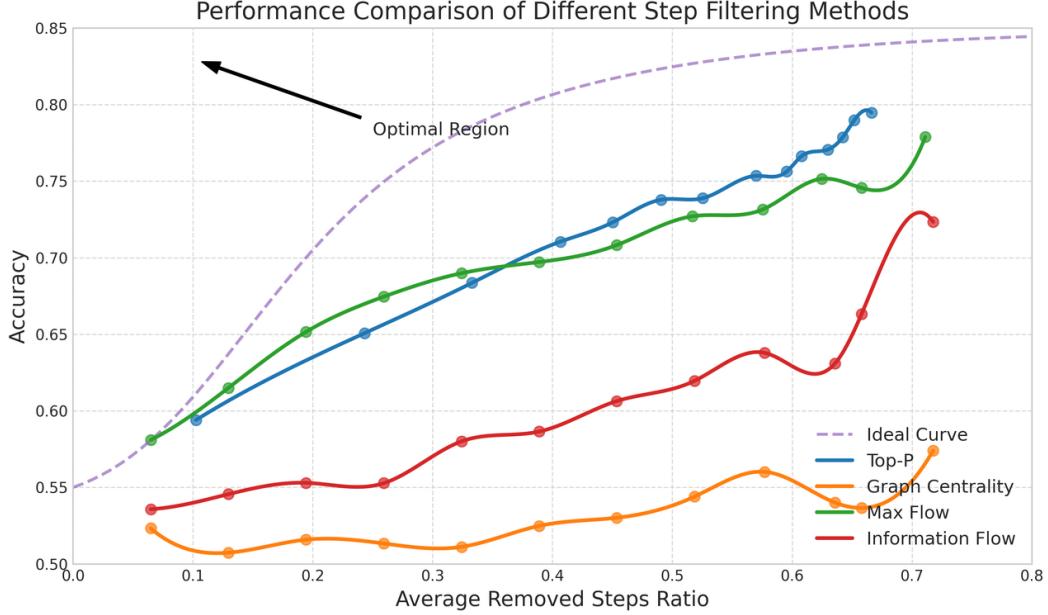


Figure 10: Comparison of Algorithms for Error Filtering Efficiency Averaged Across Four Tasks.

A.4 Error Filtering Efficiency (EFE) Evaluation Formula

For the IISR experiment, where we randomly inject N interference steps into an M -step reasoning process, the Error Filtering Efficiency is calculated as:

$$EFE = 1 - \frac{|\text{RetainedIrrelevantSteps}|}{|\text{IrrelevantSteps}|}, \quad (10)$$

where $|\text{IrrelevantSteps}|$ is the total number of interference steps injected (N), $|\text{RetainedIrrelevantSteps}|$ is the number of interference steps that were incorrectly retained after filtering

EFE measures the algorithm's ability to identify and remove irrelevant steps, with a value of 1.0 indicating perfect filtering (all interference steps removed) and 0.0 indicating no filtering capability.

As shown in Figure 6, we first compared Top K, Top P, Ppl Top (where higher perplexity indicates higher step importance), Ppl Bottom (the opposite), and Random. We evaluated the Error Filtering Efficiency under four different interference injection methods. The results show that Ppl-based methods exhibit unstable performance across different tasks.

As illustrated in Figure 10, we further compared the better performing methods: Top-P, Max-Flow, and Information-Flow. We found that the Max Flow method demonstrates a superior ability in evaluating reasoning steps, particularly when removing a small number of steps.

A.5 Structure Reasoning Process Demonstration

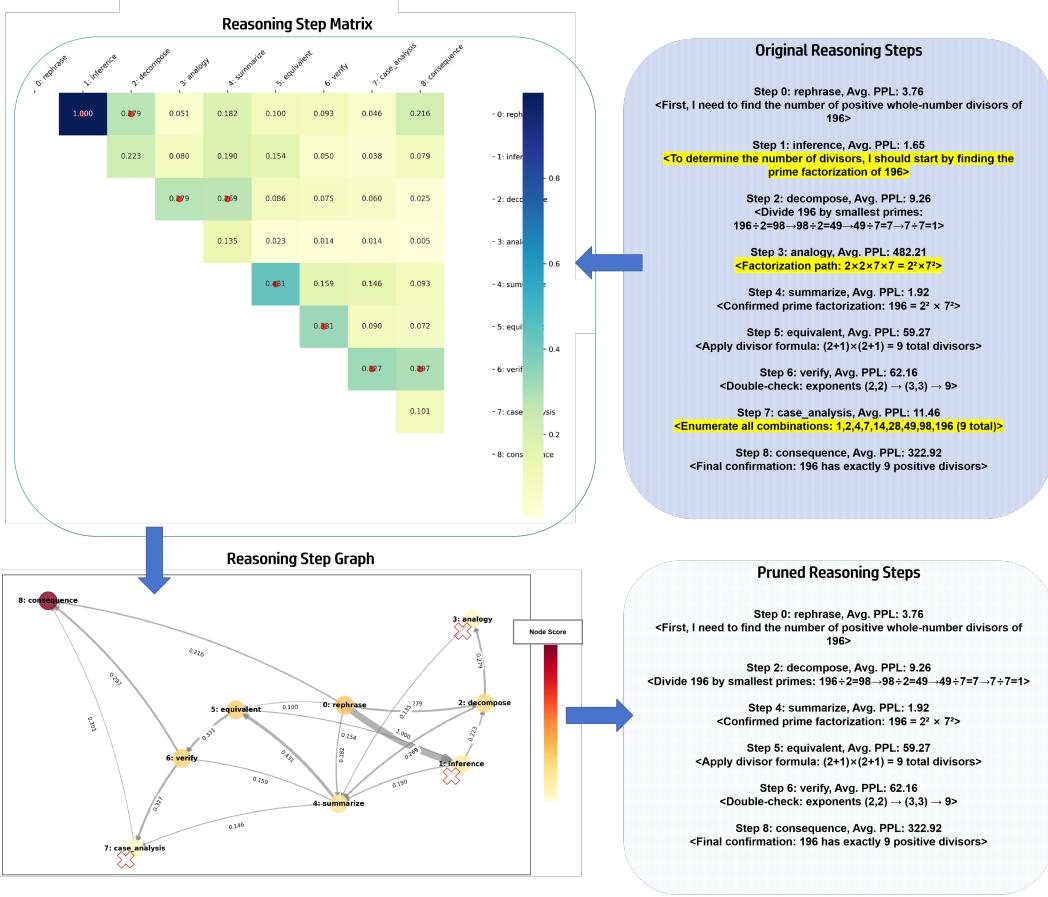


Figure 11: Using Max-flow Method for Filtering Steps during IISR Experiment.

Input Prompt:

Please use the following tags at the beginning of each sentence in your reasoning: <rephrase>, <inference>, <analogy>, <equivalent>, <association>, <reverse>, <summarize>, <verify>, <complete>, <decompose>, <counterexample>, <assumption>, <constraint>, <case_analysis>, <contradiction>, <abstraction>, <formalize>, <generalize>, <specialize>, <critique>, <alternative>, <consequence>, <intuition>

Convert the point (0, 3) in rectangular coordinates to polar coordinates. Enter your answer in the form (r, θ) , where $r > 0$ and $0 \leq \theta < 2\pi$.

Please reason step by step, and put your final answer within

<summarize> The task is to convert the rectangular coordinates (0, 3) to polar coordinates (r, θ) with $r > 0$ and $0 \leq \theta < 2\pi$.

<assumption> In rectangular (Cartesian) coordinates, a point is represented as (x, y) . To convert to polar coordinates (r, θ) , we need to find:

- r : distance from origin
- θ : angle with positive x-axis

<formalize> Conversion formulas:

$$r = \sqrt{x^2 + y^2}$$
$$\theta = \arctan(y/x) \quad (\text{when } x \neq 0)$$

<decompose> Given point (0, 3):

- $x = 0$
- $y = 3$

<specialize> Calculate r :

$$r = \sqrt{0^2 + 3^2} = 3$$

<reverse> For θ : $\arctan(3/0)$ is undefined. Observation:

- $x = 0$ and $y > 0 \Rightarrow$ point lies on positive y-axis

<case_analysis> Special case when $x = 0$:

$$\theta = \begin{cases} \pi/2 & \text{if } y > 0 \\ 3\pi/2 & \text{if } y < 0 \end{cases}$$

<verify> Check θ range: $\pi/2 \in [0, 2\pi)$ is valid.

<complete> Final polar coordinates:

$$\boxed{\left(3, \frac{\pi}{2}\right)}$$

Output:

The polar coordinates are $\boxed{\left(3, \frac{\pi}{2}\right)}$.

Input Prompt:

Please use the following tags at the beginning of each sentence in your reasoning: **<rephrase>**, **<inference>**, **<analogy>**, **<equivalent>**, **<association>**, **<reverse>**, **<summarize>**, **<verify>**, **<complete>**, **<decompose>**, **<counterexample>**, **<assumption>**, **<constraint>**, **<case_analysis>**, **<contradiction>**, **<abstraction>**, **<formalize>**, **<generalize>**, **<specialize>**, **<critique>**, **<alternative>**, **<consequence>**, **<intuition>**

What is $\sqrt{53}$ in simplest radical form?

Please reason step by step, and put your final answer within \square .

<assumption> To find the simplest radical form of $\sqrt{53}$, I need to check if 53 has any perfect square factors.

<decompose> First, let's factorize 53 to see its prime components. Since 53 is a prime number, its only factors are 1 and 53.

<verify> Since 53 is prime, there are no smaller factors other than 1 and itself. Therefore, $\sqrt{53}$ cannot be simplified further because there are no perfect squares greater than 1 that divide 53.

<conclusion> Thus, the simplest radical form of $\sqrt{53}$ is $\sqrt{53}$ itself.

Output:

The simplest radical form is $\boxed{\sqrt{53}}$.

A.6 Example Demonstration of Interference Injection and Selective Removal (IISR)

Table 8: Steps retained by different filtering methods

Method	Retained Steps
Top-P (p=0.7)	[1, 3, 5, 10, 11]
Maximum Flow	[1, 3, 5, 7, 10, 11]
Graph Centrality	[3, 4, 6, 8, 9, 10, 11]
Information Flow	[1, 2, 3, 4, 5, 6, 7, 8, 10, 11]

- Steps marked in red: [2, 4, 6, 9]
- Interference types:
 - Irrelevant pauses (Step 2)
 - Off-topic reasoning (Steps 4)
 - Disordered repetition (Steps 6)
 - Redundant Reasoning (Step 9)

Input Prompt:

Please use the following tags at the beginning of each sentence in your reasoning: <rephrase>, <inference>, <analogy>, <equivalent>, <association>, <reverse>, <summarize>, <verify>, <complete>, <decompose>, <counterexample>, <assumption>, <constraint>, <case_analysis>, <contradiction>, <abstraction>, <formalize>, <generalize>, <specialize>, <critique>, <alternative>, <consequence>, <intuition>

What is the smallest positive perfect cube that can be written as the sum of three?

Please reason step by step, and put your final answer within

<**rephrase**> The problem is to find the smallest positive perfect cube expressible as the sum of three consecutive integers.

<**inference**> [Interference Step] Pausing to consider the implications of the next step.

<**inference**> First, recall that perfect cubes are numbers like 1, 8, 27. The sum of three consecutive integers can be represented as $3n$.

<**inference**> [Off-topic Interference] Gravity cancels out, so mass and distances are the key variables...

<**analogy**> Let the integers be $(n - 1, n, n + 1)$, then sum is $3n$. The smallest cube divisible by 3 is 27.

<**inference**> [Disordered Reasoning] Smaller cubes fail: $1 = 1^3$ gives non-integer $n = \frac{1}{3}$, $8 = 2^3$ gives $n = \frac{8}{3}$.

<**verify**> Check 27: $8 + 9 + 10 = 27$. Is there a smaller cube? No, since 1 and 8 aren't divisible by 3.

<**assumption**> Alternative representation: $(n, n + 1, n + 2)$ sums to $3(n + 1)$. Still requires divisibility by 3.

<**reverse**> [Redundant Reasoning] Alternatively, maybe the three consecutive integers could be...

<**consequence**> Smaller cubes fail: $1 = 1^3$ gives non-integer $n = \frac{1}{3}$, $8 = 2^3$ gives $n = \frac{8}{3}$.

<**summarize**> Therefore, the smallest positive perfect cube meeting the criteria is:

A.7 Experimental Details

Truncation Robustness.

Our observations reveal that truncated long outputs can induce notable gradient fluctuations and result in unstable training processes (Figure 12). To mitigate this issue, we mask the truncated completions to disregard their reward values and gradient updates. This approach effectively stabilizes optimization by omitting samples surpassing a predefined length limit.

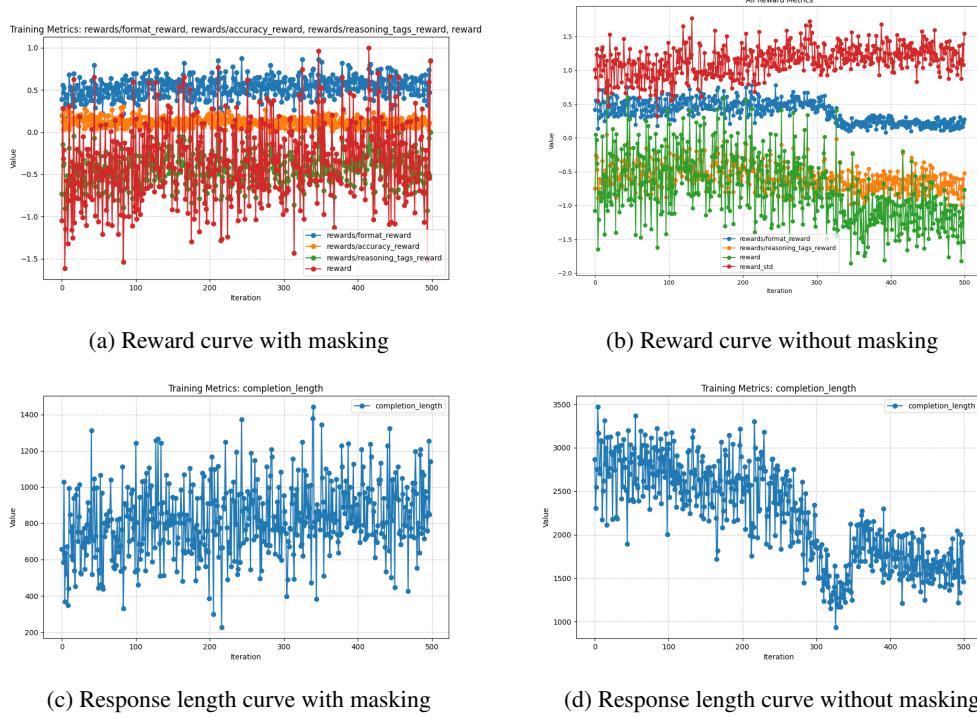


Figure 12: Impact of Truncated Completion Masking on Training Stability

Tag Randomization for Robustness

Inspired by DeepSeek-R1’s reasoning completions, we introduce randomization in the order of reasoning tags in the prompt during training. Specifically, for each question, we retain the top 5 tags and randomly sample 0–5 additional tags from the remaining set, shuffling their order in the prompt. This approach reduces overfitting to fixed reasoning patterns and encourages the model to generalize reasoning strategies.

Supervised Fine-Tuning (SFT)

We conducted supervised fine-tuning using the DeepSeek-R1-Distill-Qwen-1.5B model as our base architecture. The model was trained for 5 epochs with a learning rate of 1e-5 using a cosine scheduler with no minimum learning rate. We applied a weight decay of 1e-4 for regularization. The training was performed on a single NVIDIA A6000 GPU with a micro-batch size of 1 and no gradient accumulation. We used a context window of 4096 tokens and trained on the "S1 Structured 500 Completions" dataset. The model was trained using mixed precision (bf16) with fully sharded data parallelism for memory efficiency. The final model was saved to the "SR-SFT" directory without pushing to the Hugging Face Hub.

Reinforcement Learning (RL)

For the reinforcement learning phase, we used our fine-tuned SR-SFT model as the starting point. We implemented GRPO (Generalized Reinforcement from Policy Optimization) training on the DeepScaleR-Preview-Dataset. The model was trained with bf16 precision with Flash Attention 2 for improved efficiency. We used VLLM for inference with 70% GPU memory utilization and a maximum sequence length of 4096 tokens. Training was configured with a batch size of 6 per device, gradient accumulation of 4 steps, and gradient checkpointing to manage memory usage. We used a learning rate of 1e-6 with a cosine scheduler having a minimum learning rate at 10% of the initial

value. For each training example, we generated 6 completions with a temperature of 0.6. The reward function combines format scoring and max flow scoring (fixed weights 1.0 and 2.0). Format scoring serves as a fundamental component integrated into all reward functions. The training was conducted on 4 NVIDIA A100 GPUs. The final model was saved to the "SR-FLOW".

A.8 Structured reasoning through Fill In the Middle API

The full message template is structured the same as the prompt in Appendix A.2.

The API call is implemented as:

```
messages = [
    {"role": "user", "content": full_message},
    {
        "role": "assistant",
        "reasoning_content": "<rephrase>\nOkay, I will organize my thoughts
process in a hierarchical manner.\n</rephrase>\n",
        "content": "",
        "prefix": True
    }
]
response = await self.client.chat.completions.create(
    model=model,
    messages=messages
)
```

In particular, we initialize the reasoning process by injecting a <rephrase> tag and a metacognitive statement. This approach is effective in guiding DeepSeek-R1 to perform structured reasoning in a zero-shot setting, leading to more stable and organized reasoning patterns without additional training.

A.9 Limitation

While our structured approach has demonstrated effectiveness across mathematical reasoning, natural sciences, and logical reasoning tasks, its impact on non-reasoning tasks such as story generation and agent API calls remains unclear. Although our MAX-Flow algorithm provides a more effective way than Perplexity to evaluate the importance of reasoning steps for answers, it requires researchers to process and implement model attention mechanisms more meticulously to obtain the reasoning step matrix. The training process also slightly increases memory overhead. These functionalities are not yet optimized by existing large model libraries and inference acceleration frameworks.

A.10 Broader Impact

Positive Impact Our research further validates the effectiveness of high-quality data fine-tuning combined with reinforcement learning training, providing examples for research in other tasks. We discovered that LLMs, similar to humans in cognitive science, can improve reasoning effectiveness through structured reasoning. By encouraging spontaneous labeling of structured reasoning processes, we can simplify reasoning while maintaining performance. This raises several considerations: should we develop deeper reward-based reasoning processes? Are there more complex reasoning structures (like CoT, XML, or code formats) that could benefit models across multiple tasks? Additionally, our structured reasoning approach can be applied to early exit strategies, test-time scaling, or Monte Carlo tree search phases, benefiting other tasks. Our further analysis reveals that different model layers show varying attention spans to reasoning steps - some layers primarily focus on the immediate previous step, while others attend to a broader range of steps. This insight opens possibilities for customized model optimization, model KV cache reduction, and inference pruning optimization (we found that reducing cache for layers focusing only on the previous step doesn't significantly impact performance, while the opposite can lead to question forgetting). Furthermore, structured reasoning enables clearer definition of reasoning steps, rather than relying on logical transition words like "but" or line breaks for parsing, making it more conducive to reasoning compression and analysis tasks.

Negative Impact Our method for evaluating reasoning effectiveness involves injecting irrelevant or harmful reasoning steps for accurate assessment. However, this approach could potentially be misused for reasoning attacks or harmful reasoning bias manipulation. Through layer-wise analysis, we discovered different attention interests across layers, which might make models vulnerable to attacks based on parameter modifications that simply redirect model attention patterns. This could lead to significant performance degradation, resulting in issues like forgetfulness (forgetting the input question after reasoning for a while) or repetitive outputs (entering text output loops that cannot exit, causing substantial computational resource waste). These issues require further attention.