

# Pruning the Unsurprising: Efficient LLM Reasoning via First-Token Surprisal

Anonymous ACL submission

## Abstract

Large Reasoning Models (LRMs) have demonstrated remarkable capabilities by scaling up the length of Chain-of-Thought (CoT). However, excessively long reasoning traces pose substantial challenges for training cost and inference latency. While various CoT compression approaches have emerged to address this challenge, they face inherent trade-offs: token-level methods often disrupt syntactic and logical coherence, while step-level methods based on perplexity fail to reliably capture the logically critical reasoning steps because of the dilution of logical information. In this paper, we propose **ASAP** (**A**nchor-guided, **S**urprisal-based **P**runing), a novel coarse-to-fine framework for CoT compression. ASAP first performs anchor-guided pruning to preserve the core reasoning structure, which efficiently reduces the search space for subsequent processing. Leveraging the insight that logical branching choices are concentrated at the onset of reasoning steps, it then enables logic-aware pruning by selecting logically essential reasoning steps based on a novel first-token surprisal metric. Finally, ASAP distills the models to autonomously generate and leverage these concise CoTs at inference time, enabling efficient reasoning. Experiments show that ASAP achieves state-of-the-art accuracy across multiple benchmarks while substantially reducing training and inference costs.

## 1 Introduction

The emergence of Large Reasoning Models, including OpenAI’s o1 (Jaech et al., 2024) and DeepSeek-R1 (Guo et al., 2025), marks a paradigm shift in artificial intelligence. By scaling up Chain-of-Thought (CoT) reasoning (Wei et al., 2022), these models demonstrate emergent capabilities in complex domains such as mathematics (Sun et al., 2025), programming (Shi et al., 2024; Yang et al., 2025b; Hu et al., 2025), and logical reasoning (Liu

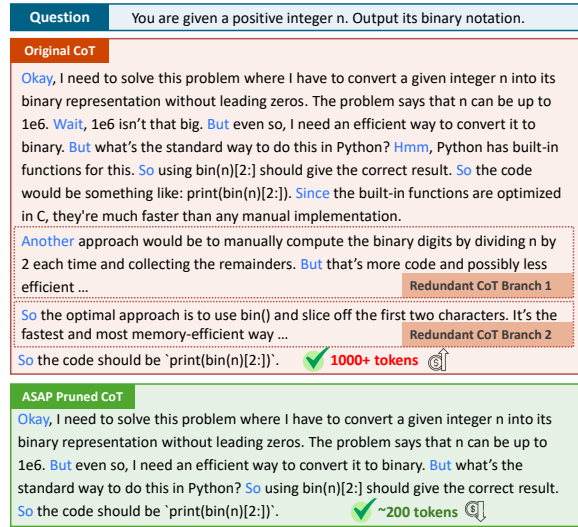


Figure 1: Illustration of CoT pruning by ASAP. The *Original CoT* generated by LRMs exhibits two types of redundancy: (1) **Structural Redundancy**, such as digressive branches (highlighted in red dashed boxes), which are removed by our Stage 1 Anchor-guided pruning; and (2) **Logical Redundancy** within valid paths. ASAP addresses the latter in Stage 2 by computing the surprisal of the first tokens of reasoning steps (marked in blue) to identify and retain only the critical cognitive pivots.

et al., 2025). However, this performance comes at a prohibitive cost: reasoning traces often span thousands of tokens, introducing substantial latency and memory overhead. Crucially, these lengthy traces often contain substantial redundancy, such as over-explaining simple problems or superficially exploring multiple paths for complex ones (Bi et al., 2024; Xie et al., 2025; Wu et al., 2025; Qu et al., 2025). For instance, the *Original CoT* in Figure 1 contains tangential branches (highlighted in red dashed boxes), such as exploring an alternative manual implementation that is subsequently rejected (“*But that’s more code...*”). Furthermore, the reasoning is punctuated by syntactic fillers that contribute little

057 to the core logic. This observation raises a funda- 109  
058 mental question: *Can we identify and retain only* 110  
059 *the “cognitive pivots” of reasoning while discard-* 111  
060 *ing the redundancy?* 112

061 A growing body of research has emerged on 113  
062 CoT compression for efficient reasoning (Qu et al., 114  
063 2025). Token-level methods like TokenSkip (Xia 115  
064 et al., 2025a) adapt general-purpose context com- 116  
065 pressors such as LLMingua-2 (Pan et al., 2024) 117  
066 to prune non-informative tokens. However, indis- 118  
067 criminate token removal risks disrupting the syn- 119  
068 tactic integrity of the reasoning chain. To address 120  
069 this, step-level pruning methods like SPIRIT (Cui 121  
070 et al., 2025) trim entire reasoning steps, thereby pre- 122  
071 serving structural coherence. Nevertheless, these 123  
072 approaches face a fundamental challenge: accu- 124  
073 rately estimating the logical importance of each 125  
074 step. They typically rely on fixed metrics like 126  
075 perplexity (PPL), which measures the overall pre- 127  
076 dictability of a sentence. This holistic measure 128  
077 often dilutes the signal of critical logical leaps with 129  
078 the noise of syntactically predictable but logically 130  
079 trivial content. 131

080 In this work, we ground CoT compression from 132  
081 an information-theoretic perspective. Through an 133  
082 empirical analysis of 10 million reasoning tokens 134  
083 (detailed in Section 2), we find that the logical 135  
084 progression within a CoT sequence is not uni- 136  
085 formly distributed; instead, its information density 137  
086 is highly concentrated at the beginning of each 138  
087 reasoning step—specifically within the first few to- 139  
088 kens (blue-highlighted in Figure 1). These tokens, 140  
089 such as “*But*” (self-correction) or “*So*” (deduction, 141  
090 not continuation), serve as high-entropy **cognitive** 142  
091 **pivots**. By leveraging the surprisal of these initial 143  
092 tokens, we can distinguish between critical logical 144  
093 transitions and predictable elaborations. 145

094 Guided by this insight, we propose **ASAP** 146  
095 (**A**nchor-guided, **S**urpris**A**I-based **P**runing), a 147  
096 coarse-to-fine framework designed to preserve 148  
097 these high-information steps. ASAP in a two-stage 149  
098 cascade that directly addresses the two types of 150  
099 redundancies identified in our study (illustrated in 151  
100 Figure 1): First, it employs **Anchor-guided Prun-** 152  
101 **ing** to remove structural redundancies. By gener- 153  
102 ating a concise step-by-step reasoning trace as a 154  
103 logical backbone, it identifies and prunes the irrel- 155  
104 evant branches (e.g., the red boxes in Figure 1). 156  
105 Second, it performs **Surprisal-based Refining** to 157  
106 eliminate logic-sparse steps. Leveraging our First- 158  
107 Token Surprisal metric, this stage iteratively filters 159  
108 out steps acting as mere fillers while retaining the

high-surprisal cognitive pivots. Finally, we dis- 109  
till these compact, logic-dense CoTs into a target 110  
model, enabling it to generate efficient reasoning 111  
chains. 112

We validate our approach through extensive ex- 113  
periments on the DeepSeek-R1-Distill-Qwen-7B 114  
and DeepSeek-R1-Distill-Llama-8B (Guo et al., 115  
2025) across diverse domains. The results demon- 116  
strate that ASAP establishes a superior Pareto fron- 117  
tier between performance and efficiency. Notably, 118  
on the challenging LiveCodeBench v4\_v5 bench- 119  
mark, ASAP achieves **36.19%** Pass@1 while re- 120  
ducing token generation by **23.5%** and inference 121  
latency by **43.5%** compared to the strongest base- 122  
line. 123

Our main contributions are summarized as fol- 124  
lows: 125

- We present an empirical analysis of the in- 126  
formation concentration of CoTs, uncovering 127  
that the surprisal of the starting token for each 128  
CoT step is a more robust indicator of logical 129  
importance than perplexity. 130
- We propose **ASAP**, a novel CoT compression 131  
framework that combines structural alignment 132  
with information-guided refinement. 133
- Extensive experiments on multiple bench- 134  
marks demonstrate that models fine-tuned on 135  
CoTs pruned by ASAP achieve state-of-the- 136  
art accuracy while substantially reducing com- 137  
putational costs. 138

## 2 Empirical Analysis 139

To investigate the intrinsic distribution of logical in- 140  
formation in CoTs, we conducted a large-scale anal- 141  
ysis on 10 million tokens generated by DeepSeek- 142  
R1-Distill-Qwen-32B (Guo et al., 2025) across 143  
diverse reasoning benchmarks (AIME and Live- 144  
CodeBench (Jain et al., 2024)). 145

**Information Concentration in CoTs.** We ana- 146  
lyze the entropy distribution of the *first token* of 147  
each reasoning step compared to all subsequent 148  
tokens. Entropy, in this context, quantifies the 149  
model’s uncertainty regarding the next state tran- 150  
sition (Shannon, 1948; Malinin and Gales, 2020; 151  
Kuhn et al., 2023; Wang et al., 2025; Cheng et al., 152  
2025). As illustrated in Figure 2(a), a distinct con- 153  
centration is observed. Starting tokens (blue) ex- 154  
hibit a dispersed distribution with a significantly 155

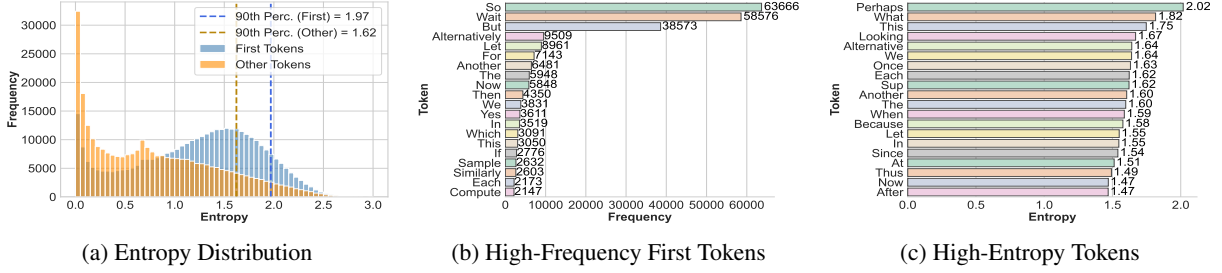


Figure 2: Empirical analysis of 10M tokens from DeepSeek-R1-Distill-Qwen-32B. (a) The entropy distribution reveals a clear information concentration: first tokens (blue) exhibit significantly higher uncertainty (entropy) compared to body tokens (orange), which are highly deterministic. (b) The most frequent first tokens are a mixture of logical operators (e.g., *Wait*) and ubiquitous syntactic connectors (e.g., *So*). (c) High-entropy states filter out predictable fillers like *So* or *Then*, while exclusively highlighting cognitive pivots such as *Perhaps*, *What*, and *Alternative*.

higher 90th percentile. In contrast, other tokens (orange) are heavily concentrated near zero entropy, indicating that once a reasoning step is initiated, its subsequent elaboration is largely deterministic and syntactically driven. This empirical evidence confirms that the logical branching points, where the model actively deliberates on the reasoning path, are structurally concentrated at the beginning of each step.

### Identifying Cognitive Pivots with Entropy.

Having identified the informative start tokens, we perform a more in-depth analysis of real logical pivots among the start tokens. We aim to distinguish between superficial syntactic connectors and real logical pivots. Figure 2(b) presents high-frequency start tokens that appear in CoTs, which mix connectors (“*So*”, “*Let*”) with reasoning markers (“*Wait*”, “*But*”). However, when focusing on tokens generated in high-entropy states (Figure 2(c)), a qualitative shift emerges. Predictable connectors like “*So*” and “*Then*” are effectively suppressed due to their low uncertainty. Instead, the distribution is dominated by terms representing cognitive pivots and state transitions, such as: 1) Exploration: “*Alternative*”, “*Another*” (proposing hypotheses or new paths). 2) Causality: “*Because*”, “*Since*” (providing formal justification). 3) Self-Correction: “*Perhaps*”, “*What*” (indicating error detection or logic reversal).

This analysis demonstrates that high entropy is a robust indicator of logical salience. Since the actual next token is known in the given training sequence, we operationalize this insight by using *First-Token Surprisal* as a proxy to identify and preserve these critical reasoning hops in our proposed framework (Fu et al., 2025).

## 3 Methodology

### 3.1 Overall Framework

Formally, we consider a supervised reasoning task defined by a dataset  $\mathcal{D} = \{(Q_i, C_i, A_i)\}_{i=1}^N$ , where  $Q_i$  is the query,  $A_i$  is the predicted answer, and  $C_i$  represents the original CoT generated by the LRMs.  $C_i$  is a sequence of reasoning steps  $C_i = \{s_1, s_2, \dots, s_L\}$ . Our goal is to compress each  $C_i$  into a concise pruned one  $C'_i$  such that  $|C'_i| \ll |C_i|$  while the LRMs maintains the quality of generated reasoning steps and answers when fine-tuned on the dataset  $\mathcal{D}' = \{(Q_i, C'_i, A_i)\}_{i=1}^N$ .

We propose **ASAP**, a coarse-to-fine framework tailored to the redundancy of “Original CoT” ( $C$ ), as illustrated in Figure 3. Stage 1 (Anchor-guided Pruning) reduces *structural redundancy* (e.g., dead ends) by aligning the CoT with a generated logical backbone. The LLM generates a “Direct Thought” ( $\mathcal{P}$ ) from the  $(Q, A)$  pairs.  $\mathcal{P}$  acts as an anchor to prune the  $C$  into a “Coarse-grained Pruned CoT” ( $C_{coarse}$ ). Stage 2 (Surprisal-based Refining) reduces *logical redundancy* (e.g., syntactic fillers) by filtering non-informative steps. We approximate the information of each step in  $C_{coarse}$  with their surprisal of start tokens and prune low-surprisal steps, yielding the final “Fine-grained Pruned CoT” ( $C'$ ). Finally, all  $\{(Q, C', A)\}$  are utilized to fine-tune the target model.

### 3.2 Anchor-guided Pruning

Directly pruning raw CoTs is challenging due to the noise and unstructured digressions inherent in LLM reasoning (Zhou et al., 2024). To address this, we first construct a high-level logical skeleton to narrow the pruning space.

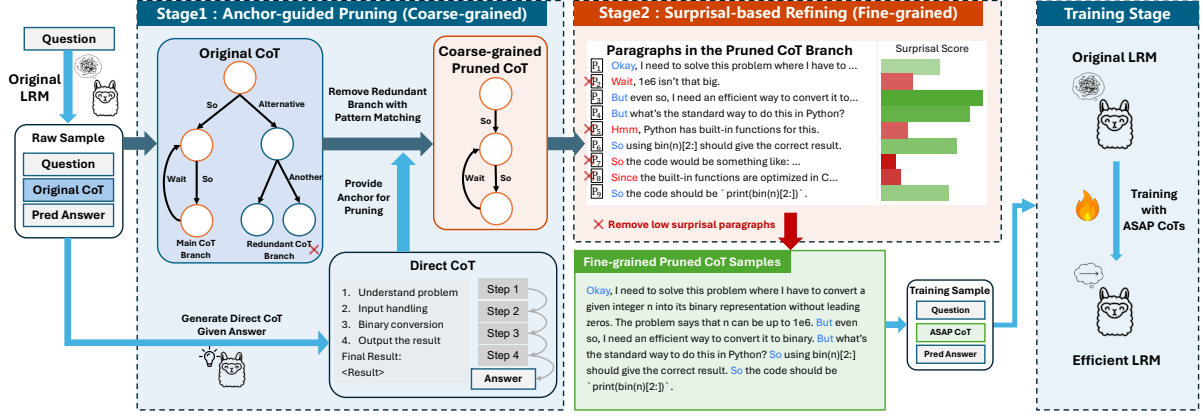


Figure 3: The overall framework of ASAP. The pipeline consists of three phases: (1) **In Stage 1**, the LLM generates a “Direct Thought” ( $\mathcal{P}$ ) from the (Question, Answer) pair.  $\mathcal{P}$  acts as an anchor to prune the “Original CoT” ( $C$ ) into a “Coarse-grained Pruned CoT” ( $C_{coarse}$ ). (2) **In Stage 2**, we compute the *First-Token Surprisal* for each step in  $C_{coarse}$ . High-surprisal steps are retained, while low-surprisal fillers are pruned, yielding the final “Fine-grained Pruned CoT” ( $C'$ ). (3) **In Training Stage**, the data with ASAP pruned CoTs is used to fine-tune the LRM for efficient inference.

**Generate Direct Thoughts.** We prompt the LLM to infer a concise reasoning path called “Direct Thought” ( $\mathcal{P}$ ) based on the ( $Q, A$ ) pair (see Appendix B for prompts). Unlike exploratory CoTs,  $\mathcal{P}$  is generated as a structured, step-by-step explanation that outlines how to derive the answer from the question, exemplified in Appendix B. This  $\mathcal{P}$  acts as a reference anchor, outlining the least reasoning trajectory required to solve the problem.

**Pruning with Pattern Matching.** Guided by the anchor  $\mathcal{P}$ , we prompt the LLM to prune the original CoT  $C$ . Specifically, the LLM is instructed to: 1) remove unnecessary reasoning steps from  $C$ ; 2) retain all key supporting content that aligns with the logic of  $\mathcal{P}$ ; and 3) crucially, preserve the original wording without introducing new information. The prompt used for pruning is shown in Appendix B. The goal is to extract the subsequence of  $C$  that semantically aligns with  $\mathcal{P}$  while discarding irrelevant branches (as shown in the “Original CoT” block of Figure 3).

Crucially, to mitigate LLM hallucination during compression, we enforce an extractive constraint, which validates structural and semantic alignment with  $C$ . Specifically, we design a pattern-matching algorithm that verifies whether each step in  $C_{coarse}$  corresponds to a matching step in  $C$  while preserving their original order. The matching is performed using Gestalt Pattern Matching (Black, 2004) as a text similarity metric. A pruning is considered valid only if all steps in  $C_{coarse}$  achieve a similarity score above a predefined threshold  $\tau$  when matched

against sequential steps in  $C$ . The full pattern-matching algorithm is detailed in Algorithm 1 (see Appendix A). We leverage high-temperature sampling, which provides the necessary diversity to efficiently re-prompt failed cases, ensuring that a valid  $C_{coarse}$  can be eventually generated.

### 3.3 Surprisal-based Refining

Following the coarse-grained pruning, the resulting  $C_{coarse}$  may still contain verbose steps that contribute little to the logic. Grounded in our empirical finding that logical information is concentrated (Section 2), we perform a meticulous, logic-aware refinement in  $C_{coarse}$  to identify more subtle redundancies within the core reasoning path.

#### First-Token Surprisal as Logical Importance.

We introduce *First-Token Surprisal* as a novel metric to precisely quantify the logical importance of each step, enabling us to filter out the least informative ones and produce the final highly condensed CoT. Let a reasoning step  $s$  be a sequence of tokens  $s = (x_1, x_2, \dots, x_T)$ . The informational value of  $s$  within the context of previous steps  $C_{pre}$  is typically estimated by its joint probability. However, our analysis reveals that the *first token*  $x_1$  serves as the “cognitive pivot” carrying the majority of the uncertainty. Therefore, we define the *First-Token Surprisal*  $\mathcal{S}(s)$  as:

$$\mathcal{S}(s | C_{pre}) = -\log P_{\theta}(x_1 | C_{pre}) \quad (1)$$

where  $P_{\theta}$  denotes the probability distribution of the LRM. A high  $\mathcal{S}(s)$  indicates a high-information

transition (e.g., initiating a new deduction or self-correction), whereas a low score suggests a deterministic continuation or syntactic filler.

**Pruning using First-Token Surprisal.** We formulate the fine-grained pruning as a constrained maximization problem. Our goal is to select a subset of steps  $S' \subset C_{coarse}$  that maximizes the total logical information subject to a length budget  $L_{max}$ :

$$\begin{aligned} S^* &= \arg \max_{S' \subseteq C_{coarse}} \sum_{s \in S'} \mathcal{S}(s) \\ \text{s.t.} \quad &\sum_{s \in S'} \text{len}(s) \leq L_{max} \end{aligned} \quad (2)$$

This formulation explicitly prioritizes steps with high information density. To solve this efficiently, we employ a greedy iterative strategy. We calculate the surprisal score for all steps in  $C_{coarse}$  and iteratively remove the step with the lowest  $\mathcal{S}(s)$ , while the relative order of steps in  $S'$  is preserved. The detailed procedure is provided in Algorithm 2 (see Appendix A). This process yields the final fine-grained CoT  $C'$ , which retains the critical ‘‘aha moments’’ (Guo et al., 2025) while meeting efficiency constraints.

### 3.4 Supervised Fine-tuning

Following the pruning, we construct the final training dataset  $\mathcal{D}' = \{(Q_i, C'_i, A_i)\}_{i=1}^N$ . For each instance, we concatenate the pruned CoT ( $C'_i$ ) and the final answer ( $A_i$ ) to form the complete target response  $R_i$ . We then fine-tune the LRM to minimize the standard negative log-likelihood of the target response tokens, conditioned on the input question. Formally, the loss is defined as:

$$\mathcal{L} = - \sum_{i=1}^N \sum_{j=1}^{|R_i|} \log P_{\theta}(r_{i,j} | Q_i, r_{i,<j}) \quad (3)$$

where  $r_{i,j}$  is the  $j$ -th token of the target response  $R_i$ , and  $\theta$  represents the parameters of the model being fine-tuned. This supervised fine-tuning process effectively distills the knowledge from our pruning framework into the model. By training on these compact, logically salient examples, the model learns to internalize efficient reasoning patterns.

## 4 Experiments

### 4.1 Experimental Setup

**Models and Datasets.** All experiments are conducted on the DeepSeek-R1-Distill-Qwen-7B and

DeepSeek-R1-Distill-Llama-8B (Guo et al., 2025), with DeepSeek-R1-Distill-Qwen-7B as the default backbone across all settings. For the code reasoning domains, we use the Python subset of the CodeForces-CoTs (Hugging Face, 2025) dataset. For the math reasoning domain, we adopt the OpenR1-Math (Hugging Face, 2025) dataset and randomly sample 10K instances to match the size of the code subset, ensuring a balanced comparison across domains. The datasets consist of high-quality Chain-of-Thought (CoT) samples generated by DeepSeek-R1, making it particularly suitable for training competitive reasoning tasks. Detailed implementation settings (hyperparameters, hardware, etc.) are provided in Appendix C.

**Benchmarks.** We evaluate our method on a suite of widely used benchmarks that cover both code generation and mathematical reasoning tasks. For code generation, we adopt HumanEval+ (Chen et al., 2021; Liu et al., 2023), LiveCodeBench v1\_v3, LiveCodeBench v4\_v5 (Jain et al., 2024), and LeetCodeDataset (Xia et al., 2025b). For mathematical reasoning, we evaluate on GSM8K (Cobbe et al., 2021), MATH500 (Hendrycks et al., 2021), AIME24, and AIME25.

**Baselines.** We compare our method against a comprehensive set of baselines. **Zero-shot** refers to the original model without any task-specific fine-tuning. **Original** denotes the model fine-tuned on the uncompressed CoTs from the training data. Among compression approaches, Selective Context (Li et al., 2023) prunes redundant lexical units based on self-information; **LLMLingua-2** (Pan et al., 2024) distills GPT-4’s token importance signals into a lightweight Transformer encoder trained as a token classifier; **TokenSkip** (Xia et al., 2025a) learns to skip less informative tokens to achieve controllable compression; and **SPIRIT** (Cui et al., 2025) identifies critical reasoning steps by measuring perplexity shifts. Except for the zero-shot setting, all methods involve fine-tuning on CoTs processed according to their respective compression strategies.

**Metrics.** We evaluate both accuracy and inference efficiency of each approach across three metrics: **Pass@1 (Acc)**, which measures the percentage of problems correctly solved on the first attempt; **Tokens (Tok)**, which denotes the average number of tokens generated by the LRMs; and **La-**

Methods	HE+			LCBv1_v3			LCBv4_v5			LCD		
	Acc ↑	Tok ↓	Lat ↓	Acc ↑	Tok ↓	Lat ↓	Acc ↑	Tok ↓	Lat ↓	Acc ↑	Tok ↓	Lat ↓
Zero-shot	68.29	3051	1.16	42.16	7088	3.59	25.37	8336	5.15	19.74	8680	4.95
Original	<u>75.61</u>	2973	1.12	<u>52.12</u>	6611	3.15	30.97	8289	4.83	<u>25.00</u>	8485	4.72
Selective Context	54.88	2979	1.13	30.23	7025	3.75	16.79	8558	5.35	15.79	8461	4.90
LLMLingua-2	68.29	3075	1.19	38.89	6953	3.60	22.76	8474	5.31	17.54	8513	4.81
TokenSkip	73.78	2823	<u>1.07</u>	32.35	7095	3.85	20.15	8400	5.37	18.42	8503	4.87
SPIRIT	<u>75.61</u>	<u>2764</u>	<u>1.07</u>	50.82	<u>6524</u>	<u>3.09</u>	<u>33.58</u>	<u>7892</u>	<u>4.62</u>	<u>25.00</u>	<u>8186</u>	<u>4.45</u>
ASAP	<b>78.66</b>	<b>2464</b>	<b>0.98</b>	<b>54.74</b>	<b>5177</b>	<b>2.09</b>	<b>36.19</b>	<b>6035</b>	<b>2.61</b>	<b>27.63</b>	<b>7541</b>	<b>3.48</b>

Table 1: Experimental results of different methods on code generation benchmarks with DeepSeek-R1-Distill-Qwen-7B. We report accuracy (Acc), average number of generated tokens (Tok), and average generation latency (Lat) measured in seconds. The best results are highlighted in bold, and the second-best are underlined.

Methods	GSM8K			MATH500			AIME24			AIME25		
	Acc ↑	Tok ↓	Lat ↓	Acc ↑	Tok ↓	Lat ↓	Acc ↑	Tok ↓	Lat ↓	Acc ↑	Tok ↓	Lat ↓
Zero-shot	83.55	1301	0.27	60.40	2629	0.70	<u>36.67</u>	8352	6.76	<u>40.00</u>	8145	6.67
Original	86.35	1250	0.26	63.80	2511	0.66	<b>46.67</b>	8034	6.43	<b>43.33</b>	8026	6.75
Selective Context	75.44	<u>1108</u>	0.24	52.20	2507	0.66	16.67	9610	7.40	10.00	9329	7.40
LLMLingua-2	79.98	1128	0.24	54.60	2802	0.76	<u>36.67</u>	8369	6.60	23.33	8919	7.42
TokenSkip	85.37	1303	0.27	<u>65.60</u>	2483	0.65	<u>36.67</u>	8073	6.61	33.33	8465	7.44
SPIRIT	<u>88.55</u>	1118	<u>0.23</u>	64.20	<u>2144</u>	<u>0.57</u>	<b>46.67</b>	<u>7198</u>	<u>5.78</u>	<b>43.33</b>	<u>7817</u>	<u>6.57</u>
ASAP	<b>90.75</b>	<b>753</b>	<b>0.16</b>	<b>70.80</b>	<b>1649</b>	<b>0.43</b>	<b>46.67</b>	<b>5552</b>	<b>5.04</b>	36.67	<b>5434</b>	<b>5.10</b>

Table 2: Experimental results of different methods on mathematical reasoning benchmarks with DeepSeek-R1-Distill-Qwen-7B. We report accuracy (Acc), average number of generated tokens (Tok), and average generation latency (Lat) measured in seconds. The best results are highlighted in bold, and the second-best are underlined.

**tency (Lat)**, which measures the average time (in seconds) required for the model generation.

## 4.2 Main Results

Tables 1 and 2 present the results of various methods on all benchmarks. The results show that the model fine-tuned on CoTs pruned by ASAP consistently achieves the best trade-off between accuracy and efficiency. It achieves the best accuracy while generating the fewest tokens, leading to the lowest generation latency.

We notice a clear distinction between token-level and step-level pruning strategies. Token-level baselines such as Selective Context, LLMLingua-2, and TokenSkip exhibit a significant performance degradation compared to the original CoTs. This is because the token removal disrupts the syntactic structure and semantic coherence of the original reasoning steps. Consequently, the fine-tuning data becomes fragmented and grammatically unnatural, making it difficult for the model to learn

the intended logical flow of the CoT. Step-level methods, such as SPIRIT, perform significantly better than token-level pruning methods, due to the preservation of sentence-level integrity. While SPIRIT improves efficiency over the Original with comparable accuracy, our method achieves higher efficiency and accuracy at the same time. This improvement is particularly pronounced on the challenging LiveCodeBench v4\_v5 benchmark: ASAP reduces the average number of generated tokens by **23.5%** (from 7892 to 6035) and lowers generation latency by **43.5%** (from 4.62s to 2.61s), while also achieving a **7.8%** improvement in accuracy (Pass@1 increases from 33.58% to 36.19%).

## 4.3 Ablation and Analysis

**Effect of Different Components.** To validate the contribution of each component, we conduct an ablation study on three model variants. 1) *w/o Anchor-guided Pruning*: which skips Stage 1 and applies only surprisal-based pruning to the origi-

Variants	Acc ↑ Tok ↓ Lat ↓
<b>ASAP</b>	<b>36.19 6035 2.61</b>
w/o Anchor-guided Pruning	35.07 7735 4.60
w/o Surprisal-based Refining	31.72 8061 4.83
w/o Both Pruning	30.97 8289 4.83

Table 3: Ablation study of different pruning strategies on LiveCodeBench v4\_v5. We report accuracy (Acc), average number of generated tokens (Tok), and average generation latency (Lat) measured in seconds.

nal CoT. 2) *w/o Surprisal-based Refining*: which omits the surprisal-based refinement stage; and 3) *w/o Both Pruning*: equivalent to the original baseline, where the model is fine-tuned on the full, uncompressed CoT. Table 3 presents the results on LiveCodeBench v4\_v5, which is representative of the consistent trends observed across benchmarks. Additional results are included in Appendix D. The results show that both pruning stages are essential and mutually complementary for optimal accuracy and efficiency. First, removing the anchor-guided pruning leads to a drop in both accuracy and efficiency. While the accuracy decrease is modest, the generation latency increases by a substantial **76.2%** (from 2.61s to 4.60s), underscoring the importance of stage 1. Second, removing the surprisal-based refining results in a significant degradation across all metrics. The accuracy drops by **12.4%** (Pass@1 decreases from 36.19% to 31.72%) relative to the ASAP, and efficiency improvements are largely lost. This highlights that our surprisal-based pruning mechanism is essential to select the most critical steps.

**Generalization to Different Architectures.** To validate the generalizability of ASAP, we replicate our main experiments on the DeepSeek-R1-Distill-Llama-8B. Following the same experimental protocol, we compare ASAP against three strong baselines: Zero-shot, Original, and SPIRIT. We observe consistent trends across all benchmarks. For brevity, we present representative results on two key benchmarks: LiveCodeBench v4\_v5 and AIME24 in Table 4, while reporting the full results in the Appendix E. The results in the Llama3.1 series are highly consistent with our findings in the Qwen2.5 series, confirming the generalization of the ASAP. As shown in Table 4, ASAP achieves the highest accuracy on both benchmarks, and the efficiency improvements are even more pronounced. On LiveCodeBench, for instance, ASAP not only

Methods	LCB			AIME		
	Acc ↑ Tok ↓ Lat ↓	Acc ↑ Tok ↓ Lat ↓	Acc ↑ Tok ↓ Lat ↓	Acc ↑ Tok ↓ Lat ↓	Acc ↑ Tok ↓ Lat ↓	Acc ↑ Tok ↓ Lat ↓
Zero-shot	25.00	8508	8.90	33.33	8445	10.42
Original	31.34	8202	8.60	<b>36.67</b>	8550	10.04
SPIRIT	30.22	7913	8.45	<b>36.67</b>	8788	10.04
<b>ASAP</b>	<b>32.84</b>	<b>4175</b>	<b>2.69</b>	<b>36.67</b>	<b>5314</b>	<b>6.97</b>

Table 4: Experimental results of different methods with DeepSeek-R1-Distill-Llama-8B. We report accuracy (Acc), average number of generated tokens (Tok), and average generation latency (Lat) measured in seconds. The best results are highlighted in bold.

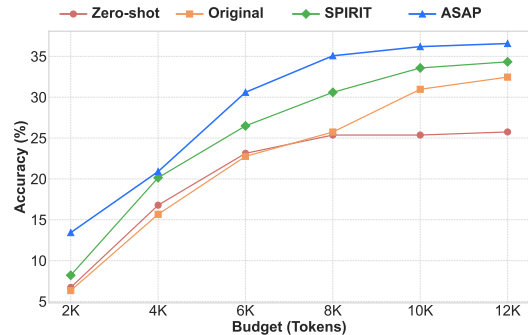


Figure 4: Performance of ASAP on LiveCodeBench v4\_v5 under different token budgets.

surpasses the accuracy of the Original baseline (32.84% vs. 31.34%) but also generates **49.1%** fewer tokens and reduces latency by over **3x** (from 8.60s to 2.69s). This suggests that the ASAP is particularly effective in identifying and distilling the core reasoning patterns, validating its robustness and broad applicability for improving reasoning efficiency across different model families.

**Impact of Token Budget.** To evaluate the scalability and resource sensitivity of our method, we analyze its behavior under varying inference-time token budgets (i.e., the maximum number of tokens to the model). We compare ASAP against the three strong baselines—SPIRIT, Original, and Zero-shot—across all benchmarks, and observe consistent trends. For clarity, we present results on LiveCodeBench v4\_v5 under six budget settings ranging from 2K to 12K tokens. Results for other benchmarks and additional statistics are provided in Appendix F. As shown in Figure 4, ASAP consistently outperforms all baselines across all budget settings. In particular, ASAP exhibits smooth performance scaling with respect to the token budget. We note that ASAP achieves superior performance-efficiency trade-offs. For example, ASAP with just

Methods	Tokens	Time
Original	13023	80.11
Selective Context	6722 (-48.4%)	63.41 (-20.9%)
LLMLingua-2	6919 (-46.9%)	65.25 (-18.6%)
TokenSkip	9813 (-24.6%)	77.27 (-3.6%)
SPIRIT	6082 (-53.3%)	57.45 (-28.3%)
ASAP	<b>3178 (-75.6%)</b>	<b>31.48 (-60.7%)</b>

Table 5: Training efficiency comparison on CodeForces-CoTs dataset. We report the average number of tokens per sample and training time measured in seconds per step. Percentages indicate the reduction relative to the Original baseline.

an 8K token budget achieves higher accuracy than SPIRIT and Original at a much larger 12K budget. These results further validate the practical utility of ASAP in real-world scenarios.

**Training Efficiency.** To quantify the training efficiency gains, we present results of the CodeForces-CoTs dataset in Table 5 and results on other datasets are provided in Appendix G. The results highlight the training efficiency advantage of the ASAP. By generating the most compact yet logically rich CoTs, our approach significantly reduces training overhead. Compared to the uncompressed baseline (Original), our method reduces the number of training tokens by **75.6%** and shortens training time by **60.7%**. These savings substantially exceed those achieved by all other baselines. ASAP enables a more resource-efficient training process, making it a practical and cost-effective solution for real-world deployment.

## 5 Related Work

**Chain-of-Thought and Advanced Reasoning.** Chain-of-Thought (CoT) prompting (Wei et al., 2022) has evolved from heuristic prompting strategies (Yao et al., 2023; Lei et al., 2023; Ling et al., 2023) to the training of specialized Large Reasoning Models (LRMs) like OpenAI’s o1 (Jaech et al., 2024) and DeepSeek-R1 (Guo et al., 2025). These models leverage reinforcement learning to scale test-time compute, generating lengthy reasoning traces to solve complex tasks (Kimi et al., 2025; Yang et al., 2025a; Yu et al., 2025; Wang et al., 2025). Unlike prior works that enhance performance by scaling up CoT length, we focus on pruning redundancy to improve efficiency without compromising reasoning performance.

**Context Compression for LLMs.** To mitigate the computational cost of long contexts, various compression techniques have been proposed. Approaches like Selective Context (Li et al., 2023), LLMLingua series (Jiang et al., 2023; Pan et al., 2024), and LongCodeZip (Shi et al., 2025) employ information-theoretic metrics or small external models to filter redundant tokens. However, these methods typically treat input as unstructured text. Applying them directly to CoT often disrupts the syntactic and logical coherence required for valid reasoning, a limitation that our pruning aims to overcome.

**Efficient Reasoning via Fine-Tuning.** Recent research has explored various efficiency mechanisms (Qu et al., 2025), ranging from compressing thoughts into continuous latent representations (Hao et al., 2024; Cheng and Van Durme, 2024; Shen et al., 2025) to compressing CoTs (Kang et al., 2025; Xia et al., 2025a; Cui et al., 2025). Approaches like TokenSkip (Xia et al., 2025a) and SPIRIT (Cui et al., 2025) reduce length by filtering tokens or steps based on heuristics or perplexity shifts. However, these metrics often struggle to differentiate between syntactic fluency and logical necessity. ASAP differs by combining anchor-guided structural pruning with first-token surprisal, offering a more robust proxy for cognitive pivots.

## 6 Conclusion

In this paper, we address the inefficiency of Large Reasoning Models stemming from the structural and logical redundancies in Chain-of-Thought reasoning. Grounded in an information-theoretic perspective, our large-scale empirical analysis reveals a fundamental property of reasoning traces: Information Concentration, where the logical uncertainty is highly concentrated at the onset of reasoning steps. Guided by this insight, we propose ASAP. This coarse-to-fine framework first aligns the reasoning structure with a logical anchor and then refines it using a novel First-Token Surprisal metric. Extensive experiments across multiple benchmarks demonstrate that ASAP outperforms existing baselines, establishing a new state-of-the-art Pareto frontier between accuracy and efficiency. Our work highlights the potential of using information-theoretic signals for efficient reasoning. Future work will explore applying ASAP to online inference acceleration.

## 571 Limitations

572 While ASAP demonstrates significant improve-  
573 ments in reasoning efficiency, we acknowledge  
574 several limitations. First, our method relies on  
575 the availability of a capable LLM to generate high-  
576 quality “Direct Thoughts” in Stage 1. If the anchor  
577 contains logical errors or hallucinations, it may mis-  
578 guide the subsequent pruning, although our pattern-  
579 matching constraint mitigates this risk. Second,  
580 our experiments primarily focus on code genera-  
581 tion and mathematical reasoning. While we believe  
582 the principle of information concentration applies  
583 broadly, the effectiveness of ASAP on creative writ-  
584 ing or commonsense reasoning tasks remains to be  
585 verified.

## 586 References

587 Zhen Bi, Ningyu Zhang, Yinuo Jiang, Shumin Deng,  
588 Guozhou Zheng, and Huajun Chen. 2024. When  
589 do program-of-thought works for reasoning? In  
590 *Proceedings of the AAAI conference on artificial in-*  
591 *telligence*, volume 38, pages 17691–17699.

592 Paul E Black. 2004. Ratcliff/obershelp pattern recogni-  
593 tion. *Dictionary of algorithms and data structures*,  
594 17.

595 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan,  
596 Henrique Ponde De Oliveira Pinto, Jared Kaplan,  
597 Harri Edwards, Yuri Burda, Nicholas Joseph, Greg  
598 Brockman, and 1 others. 2021. Evaluating large  
599 language models trained on code. *arXiv preprint*  
600 *arXiv:2107.03374*.

601 Daixuan Cheng, Shaohan Huang, Xuekai Zhu, Bo Dai,  
602 Wayne Xin Zhao, Zhenliang Zhang, and Furu Wei.  
603 2025. Reasoning with exploration: An entropy per-  
604 spective. *arXiv preprint arXiv:2506.14758*.

605 Jeffrey Cheng and Benjamin Van Durme. 2024. Com-  
606 pressed chain of thought: Efficient reasoning  
607 through dense representations. *arXiv preprint*  
608 *arXiv:2412.13171*.

609 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,  
610 Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias  
611 Plappert, Jerry Tworek, Jacob Hilton, Reiichiro  
612 Nakano, and 1 others. 2021. Training verifiers  
613 to solve math word problems. *arXiv preprint*  
614 *arXiv:2110.14168*.

615 Yingqian Cui, Pengfei He, Jingying Zeng, Hui Liu,  
616 Xianfeng Tang, Zhenwei Dai, Yan Han, Chen Luo,  
617 Jing Huang, Zhen Li, and 1 others. 2025. Stepwise  
618 perplexity-guided refinement for efficient chain-of-  
619 thought reasoning in large language models. *arXiv*  
620 *preprint arXiv:2502.13260*.

Yichao Fu, Xuewei Wang, Yuandong Tian, and Jiawei  
Zhao. 2025. Deep think with confidence. *arXiv*  
*preprint arXiv:2508.15260*. 621  
622  
623

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao  
Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shi-  
rong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025.  
Deepseek-r1: Incentivizing reasoning capability in  
llms via reinforcement learning. *arXiv preprint*  
*arXiv:2501.12948*. 624  
625  
626  
627  
628  
629

Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li,  
Zhiting Hu, Jason Weston, and Yuandong Tian. 2024.  
Training large language models to reason in a contin-  
uous latent space. *arXiv preprint arXiv:2412.06769*. 630  
631  
632  
633

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul  
Arora, Steven Basart, Eric Tang, Dawn Song, and  
Jacob Steinhardt. 2021. Measuring mathematical  
problem solving with the math dataset. *NeurIPS*. 634  
635  
636  
637

Minghao Hu, Junzhe Wang, Weisen Zhao, Qiang Zeng,  
and Lannan Luo. 2025. Flowmaltrans: Unsuper-  
vised binary code translation for malware detec-  
tion using flow-adapter architecture. *arXiv preprint*  
*arXiv:2508.20212*. 638  
639  
640  
641  
642

Hugging Face. 2025. Open r1: A fully open reproduc-  
tion of deepseek-r1. 643  
644

Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richard-  
son, Ahmed El-Kishky, Aiden Low, Alec Helyar,  
Aleksander Madry, Alex Beutel, Alex Carney, and 1  
others. 2024. Openai o1 system card. *arXiv preprint*  
*arXiv:2412.16720*. 645  
646  
647  
648  
649

Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia  
Yan, Tianjun Zhang, Sida Wang, Armando Solar-  
Lezama, Koushik Sen, and Ion Stoica. 2024. Live-  
codebench: Holistic and contamination free eval-  
uation of large language models for code. *arXiv*  
*preprint arXiv:2403.07974*. 650  
651  
652  
653  
654  
655

Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing  
Yang, and Lili Qiu. 2023. Llmlingua: Compressing  
prompts for accelerated inference of large language  
models. In *Proceedings of the 2023 Conference on*  
*Empirical Methods in Natural Language Processing*,  
pages 13358–13376. 656  
657  
658  
659  
660  
661

Yu Kang, Xianghui Sun, Liangyu Chen, and Wei Zou.  
2025. C3ot: Generating shorter chain-of-thought  
without compromising effectiveness. In *Proceedings*  
*of the AAAI Conference on Artificial Intelligence*,  
volume 39, pages 24312–24320. 662  
663  
664  
665  
666

Kimi, Angang Du, Bofei Gao, Bowei Xing, Changjiu  
Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chen-  
zhuang Du, Chonghua Liao, and 1 others. 2025.  
Kimi k1. 5: Scaling reinforcement learning with llms.  
*arXiv preprint arXiv:2501.12599*. 667  
668  
669  
670  
671

Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023.  
Semantic uncertainty: Linguistic invariances for un-  
certainty estimation in natural language generation.  
*arXiv preprint arXiv:2302.09664*. 672  
673  
674  
675

676	Bin Lei, Chunhua Liao, Caiwen Ding, and 1 others.	Jiankai Sun, Chuanyang Zheng, Enze Xie, Zhengying	730
677	2023. Boosting logical reasoning in large language	Liu, Ruihang Chu, Jianing Qiu, Jiaqi Xu, Mingyu	731
678	models through a new framework: The graph of	Ding, Hongyang Li, Mengzhe Geng, and 1 others.	732
679	thought. <i>arXiv preprint arXiv:2308.08614</i> .	2025. A survey of reasoning with foundation mod-	733
680	Yucheng Li, Bo Dong, Frank Guerin, and Chenghua Lin.	els: Concepts, methodologies, and outlook. <i>ACM</i>	734
681	2023. Compressing context to enhance inference	<i>Computing Surveys</i> , 57(11):1–43.	735
682	efficiency of large language models. In <i>Proceedings</i>	Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shix- uan Liu, Rui Lu, Kai Dang, Xionghui Chen, Jianxin	736
683	<i>of the 2023 Conference on Empirical Methods in</i>	Yang, Zhenru Zhang, and 1 others. 2025. Beyond	737
684	<i>Natural Language Processing</i> , pages 6342–6353.	the 80/20 rule: High-entropy minority tokens drive	738
685	Zhan Ling, Yunhao Fang, Xuanlin Li, Zhiao Huang,	effective reinforcement learning for llm reasoning.	739
686	Mingu Lee, Roland Memisevic, and Hao Su. 2023.	<i>arXiv preprint arXiv:2506.01939</i> .	740
687	Deductive verification of chain-of-thought reasoning.	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	742
688	<i>Advances in Neural Information Processing Systems</i> ,	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	743
689	36:36407–36433.	and 1 others. 2022. Chain-of-thought prompting elic- its reasoning in large language models. <i>Advances</i>	744
690	Hanmeng Liu, Zhizhang Fu, Mengru Ding, Ruoxi Ning,	<i>in neural information processing systems</i> , 35:24824–	745
691	Chaoli Zhang, Xiaozhang Liu, and Yue Zhang. 2025.	24837.	746
692	Logical reasoning in large language models: A sur- vey. <i>arXiv preprint arXiv:2502.09100</i> .	Yuyang Wu, Yifei Wang, Ziyu Ye, Tianqi Du, Stefanie	747
693	694	Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Ling- ming Zhang. 2023. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. <i>Advances in Neural</i>	748
694	695	Information Processing Systems, 36:21558–21572.	749
696	697	Andrey Malinin and Mark Gales. 2020. Uncertainty esti- mation in autoregressive structured prediction. <i>arXiv</i>	750
698	699	<i>preprint arXiv:2002.07650</i> .	751
700	701	Zhuoshi Pan, Qianhui Wu, Huiqiang Jiang, Menglin Xia, Xufang Luo, Jue Zhang, Qingwei Lin, Victor Rühle, Yuqing Yang, Chin-Yew Lin, and 1 others.	752
702	703	2024. Llm-lingua-2: Data distillation for efficient and faithful task-agnostic prompt compression. In <i>Find- ings of the Association for Computational Linguistics</i>	753
704	705	<i>ACL 2024</i> , pages 963–981.	754
706	707	Xiaoye Qu, Yafu Li, Zhaochen Su, Weigao Sun, Jianhao Yan, Dongrui Liu, Ganqu Cui, Daizong Liu, Shuxian	755
708	709	Liang, Junxian He, and 1 others. 2025. A survey of efficient reasoning for large reasoning models: Lan- guage, multimodality, and beyond. <i>arXiv preprint</i>	756
710	711	<i>arXiv:2503.21614</i> .	757
712	713	Claude E Shannon. 1948. A mathematical theory of communication. <i>The Bell system technical journal</i> ,	758
714	715	27(3):379–423.	759
716	717	Zhenyi Shen, Hanqi Yan, Linhai Zhang, Zhanghao Hu, Yali Du, and Yulan He. 2025. Codi: Compress- ing chain-of-thought into continuous space via self- distillation. <i>arXiv preprint arXiv:2502.21074</i> .	760
718	719	720	761
721	722	723	762
722	723	724	763
724	725	726	764
726	727	728	765
728	729	730	766
730	731	732	767
732	733	734	768
734	735	736	769
736	737	738	770
738	739	740	771
740	741	742	772
742	743	744	773
744	745	746	774
746	747	748	775
748	749	750	776
750	751	752	777
752	753	754	778
754	755	756	779
756	757	758	780
758	759	760	781
760	761	762	782
762	763	764	783
764	765	766	784
766	767	768	785
768	769	770	
770	771	772	
772	773	774	
774	775	776	
776	777	778	
778	779	780	
780	781	782	
782	783	784	
784	785		

Zhanke Zhou, Rong Tao, Jianing Zhu, Yiwen Luo, Zeng-mao Wang, and Bo Han. 2024. Can language models perform robust reasoning in chain-of-thought prompting with noisy rationales? *Advances in Neural Information Processing Systems*, 37:123846–123910.

## A Algorithms

---

### Algorithm 1 Pattern Matching

---

**Require:** Original CoT  $C$ , Coarse-grained Pruned CoT  $C_{coarse}$ , Threshold  $\tau$

**Ensure:** *True* if  $C_{coarse}$  is valid, *False* otherwise.

```

1: function PATTERNMATCH( $C, C_{coarse}, \tau$ )
2:    $S_{origin} \leftarrow \text{SplitStepsByBlankLine}(C)$ 
3:    $S_{coarse} \leftarrow \text{SplitStepsByBlankLine}(C_{coarse})$ 
4:    $origin\_idx \leftarrow 0$ 
5:   for each step  $s_{coarse}$  in  $S_{coarse}$  do
6:      $found\_match \leftarrow \text{False}$ 
7:     while  $origin\_idx < \text{Length}(S_{origin})$  do
8:        $s_{origin} \leftarrow S_{origin}[origin\_idx]$ 
9:        $score \leftarrow \text{GestaltSimilarity}(s_{origin}, s_{coarse})$ 
10:      if  $score \geq \tau$  then
11:         $found\_match \leftarrow \text{True}$ 
12:         $origin\_idx \leftarrow origin\_idx + 1$ 
13:      break
14:    end if
15:     $origin\_idx \leftarrow origin\_idx + 1$ 
16:  end while
17:  if not  $found\_match$  then
18:    return False
19:  end if
20: end for
21: return True
22: end function

```

---



---

### Algorithm 2 Iterative Pruning via First-Token Surprisal

---

**Require:** Coarse-grained Pruned CoT  $C_{coarse}$ , Max Tokens  $L_{max}$ , Model  $M$ , Tokenizer  $T$

**Ensure:** Fine-grained Pruned CoT  $C'$

```

1: function FINEGRAINEDPRUNE( $C_{coarse}, L_{max}, M, T$ )
2:   if  $\text{Length}(T(C_{coarse})) \leq L_{max}$  then
3:     return  $C_{coarse}$ 
4:   end if
5:    $S \leftarrow \text{SplitStepsByBlankLine}(C_{coarse})$ 
6:    $SurprisalScores \leftarrow \text{CalculateAll}(S, M, T)$ 
7:    $StepsToPrune \leftarrow \text{SortByScore}(S, SurprisalScores)$ 
8:    $S_{current} \leftarrow S$ 
9:   for each step  $s_{prune}$  in  $StepsToPrune$  do
10:     $S_{temp} \leftarrow S_{current} \setminus \{s_{prune}\}$ 
11:     $C_{temp} \leftarrow \text{Join}(S_{temp})$ 
12:    if  $\text{Length}(T(C_{temp})) \leq L_{max}$  then
13:       $S_{current} \leftarrow S_{temp}$ 
14:    break
15:  end if
16:   $S_{current} \leftarrow S_{temp}$ 
17: end for
18:   $C' \leftarrow \text{Join}(S_{current})$ 
19:  return  $C'$ 
20: end function

```

---

## B Prompt Templates and Generated Examples

We prompt the LLM to produce a direct thought: a concise, structured, step-by-step explanation that outlines how to derive the answer from the given question. A direct thought always consists of a small number of clear logical steps that directly contribute to the solution, and terminates with a single, explicit final answer. We use the prompt below to elicit such direct thoughts (the model is instructed to return only a detailed step-by-step solution containing only “Step-by-Step Solution” and “Final Answer”).

### Prompt for Generating the Direct CoT $\mathcal{P}$

Given a question, please tell me how to get this answer step by step.

Question:

```

{question}

```

Answer:

```python

{answer}

```

Only return a detailed step-by-step solution (containing only “Step-by-Step Solution” and “Final Answer”). The detailed step-by-step solution is:

The following example illustrates a direct thought for a simple programming task: converting an integer into its binary representation. This example demonstrates the defining characteristics of direct CoTs, namely, concise step-by-step reasoning that directly derives the solution, without redundant exploration or verbose commentary.

### Example of a Direct CoT

#### Step-by-Step Solution

1. Understand the Problem: The task is to convert a given positive integer  $n$  into its binary representation without leading zeros. The constraints are lenient ( $n \leq 10^6$ ), so even a straightforward approach is efficient.

2. Input Handling: Read the integer  $n$  from standard input. Since the input consists of a single integer, we can directly use `input()` and convert it to an integer.

3. Binary Conversion: Python's built-in `bin()` function converts an integer to its binary string representation (prefixed with `0b`). To remove the prefix, slice the string from the third character onward (`[2:]`).

4. Output the Result: Print the resulting binary string. This ensures no leading zeros are present.

5. Efficiency: The approach is optimal because conversion and slicing both run in  $O(\log n)$  time, and memory usage is minimal.

#### Final Code:

```
```python
n = int(input())
print(bin(n)[2:])
```
```

### Prompt for Coarse-grained Pruning

Compress the given thinking by referring to the provided solution. The goal is to remove irrelevant reasoning paths while retaining all content along the core reasoning path. Compression must be based on thinking, ensuring that the original wording and structure are preserved as much as possible. Follow these strict rules:

1. Use thinking as the foundation: Do not rewrite or replace its content with solution—only use solution to determine which parts are relevant.

2. Remove unnecessary reasoning: Aggressively remove alternative paths that are not part of the core reasoning path.

3. Retain key supporting content: Keep examples, reflections, and tests that help illustrate, verify, or analyze the core reasoning path.

4. Preserve original words: Do not paraphrase, reorder, or change any words.

5. Do not add new words: Do not introduce new concepts, symbols, or abbreviations.

If you understand, compress the following thinking based on the given solution.

Solution:

```
```
{solution}
```
```

Thinking:

```
```
{think}
```
```

The compressed thinking is:

Given the direct thought as a reference, this stage instructs the model to prune redundant reasoning from the original chain of thought while preserving the essential logical structure. Specifically, the model is required to 1) remove unnecessary reasoning steps from the original CoT, 2) retain all key supporting content that aligns with the logic of direct CoT, and 3) strictly preserve the original wording and sentence order without introducing new information. This ensures that the compressed reasoning remains faithful to the original thought process while aligning with the concise, goal-oriented structure of the direct CoT. The following detailed prompt is used to elicit such coarse-grained pruning behavior.

## C Implementation Details

**Software and Hardware.** For fine-tuning, we utilized the `unsloth` library<sup>1</sup> for its memory-efficient optimizations. For inference, we employed the `vLLM` engine<sup>2</sup> to maximize throughput and efficiency. All experiments were conducted on NVIDIA H20 GPUs and Intel Xeon Platinum 8480+ CPUs.

**Fine-tuning Configuration.** We performed full-parameter fine-tuning for all models in our experiments. Key hyperparameters included precision set to `bf16`, `num_train_epochs` set to 10, and a `max_seq_length` of 16384. We used a `per_device_train_batch_size` of 1 with `gradient_accumulation_steps` set to 16, resulting in an effective batch size of 16. For the optimizer, we used AdamW with a `cosine_with_min_lr` learning rate scheduler. The `warmup_ratio` was set to 0.03, and the scheduler's `min_lr_rate` was 0.1 of the peak learning rate. To stabilize training, we applied gradient clipping with a `max_grad_norm` of 0.2. Based on pre-

<sup>1</sup><https://pypi.org/project/unsloth/2025.5.6/>

<sup>2</sup><https://pypi.org/project/vllm/0.8.4/>

liminary experiments, we set the peak learning rate to  $4 \times 10^{-5}$  for the DeepSeek-R1-Distill-Qwen-7B and  $2 \times 10^{-5}$  for the DeepSeek-R1-Distill-Llama-8B. Due to the high computational cost of full-parameter fine-tuning, the model is fine-tuned by a single run with a fixed random seed 42.

**Inference and Evaluation Protocol.** All inference benchmarks were run using the vLLM engine with dtype set to bfloat16 and gpu\_memory\_utilization set to 0.9. To ensure deterministic and reproducible outputs, we set the sampling temperature to 0.0 and set enable\_prefix\_caching to False. The default token budget for generation is adjusted based on the task difficulty. Specifically, it is 2K for GSM8K, 4K for MATH500, 6K for HumanEval+, and 10K for AIME24, AIME25, LiveCodeBench, and LeetCodeDataset. Results with other token budget settings are shown in Appendix F.

**Baseline Details.** Following established practices, we used a consistent scoring model; as our primary model is DeepSeek-R1-Distill checkpoints, we employed DeepSeek-R1-Distill-Qwen-7B for all model-scoring tasks. To ensure a fair comparison, we standardize the input format across all methods by preserving the original question and final answer, and applying compression only to the CoT reasoning steps. To balance compression ratio and content retention, we set the target compression ratio to 0.5 for all baseline methods, except for TokenSkip, where we follow its original design that allows a controllable compression ratio between 0.5 and 1.0. Additionally, since the original SPIRIT method is computationally expensive when applied to extremely long CoTs, we adopt a modified version to ensure fair comparison: specifically, we compute perplexity once per reasoning step and iteratively remove steps until the target ratio is met. This variant retains the core idea of SPIRIT while improving scalability in our evaluation setting.

**Hyperparameters for Our Method.** Our method involves several stages. For the LLM-guided Coarse-grained Pruning stage, we employed DeepSeek-V3 for economic reasons. When generating the direct thought  $\mathcal{P}$ , we used a deterministic setting (temperature=0.0, top\_p=1.0), while for making the final pruning result, we increased exploration (temperature=1.0, top\_p=1.0). For Pattern Matching, the similarity threshold  $\tau$  was set to 0.6. Finally, during Surprisal-based Fine-grained

Pruning, the maximum token budget was set to 4096 to ensure a deep level of compression.

## D Effect of Different Components.

To validate the contribution and necessity of each component in our two-stage pruning framework, we conduct a detailed ablation study. Specifically, we evaluate the following three variants: *ASAP w/o Coarse-grained Pruning*, *ASAP w/o Fine-grained Pruning*, and *ASAP w/o Any Pruning*. We present results on the HumanEval+, LiveCodeBench v1\_v3, and LeetCodeDatasets benchmarks in Table 6, Table 7, and Table 8.

| Variants                   | Acc $\uparrow$ | Tok $\downarrow$ | Lat $\downarrow$ |
|----------------------------|----------------|------------------|------------------|
| <b>ASAP</b>                | <b>78.66</b>   | <b>2464</b>      | <b>0.98</b>      |
| w/o Coarse-grained Pruning | 78.05          | 2839             | 1.10             |
| w/o Fine-grained Pruning   | 67.07          | 2897             | 1.10             |
| w/o Any Pruning            | 75.61          | 2973             | 1.12             |

Table 6: Ablation study of different pruning strategies for ASAP on HumanEval+. We report accuracy (Acc), average number of generated tokens (Tok), and average generation latency (Lat) measured in seconds.

| Variants                   | Acc $\uparrow$ | Tok $\downarrow$ | Lat $\downarrow$ |
|----------------------------|----------------|------------------|------------------|
| <b>ASAP</b>                | <b>54.74</b>   | <b>5177</b>      | <b>2.09</b>      |
| w/o Coarse-grained Pruning | 53.92          | 6107             | 2.77             |
| w/o Fine-grained Pruning   | 51.14          | 6599             | 3.20             |
| w/o Any Pruning            | 52.12          | 6611             | 3.15             |

Table 7: Ablation study of different pruning strategies for ASAP on LiveCodeBench v1\_v3. We report accuracy (Acc), average number of generated tokens (Tok), and average generation latency (Lat) measured in seconds.

| Variants                   | Acc $\uparrow$ | Tok $\downarrow$ | Lat $\downarrow$ |
|----------------------------|----------------|------------------|------------------|
| <b>ASAP</b>                | <b>27.63</b>   | <b>7541</b>      | <b>3.48</b>      |
| w/o Coarse-grained Pruning | 24.12          | 7954             | 3.75             |
| w/o Fine-grained Pruning   | 25.44          | 8326             | 4.77             |
| w/o Any Pruning            | 25.00          | 8485             | 4.72             |

Table 8: Ablation study of different pruning strategies for ASAP on LeetCodeDataset. We report accuracy (Acc), average number of generated tokens (Tok), and average generation latency (Lat) measured in seconds.

## E Generalization to Different Architectures

To evaluate the generalizability of ASAP, we replicate our main experiments on the DeepSeek-R1-Distill-Llama-8B. Following the same experimental protocol, we compare ASAP against three baselines: Zero-shot, Original, and SPIRIT. The results of the code generation task on the HumanEval+, LiveCodeBench v1\_v3, LiveCodeBench v4\_v5, and LeetCodeDataset benchmarks are shown in Table 9. The results of the mathematical reasoning task on the GSM8K, MATH500, AIME24, and AIME25 benchmarks are shown in Table 10.

## F Performance under Different Token Budgets

To evaluate the performance scalability and resource sensitivity of our method, we analyze its behavior under varying inference-time token budgets (i.e., the maximum number of tokens the model is allowed to generate). We compare ASAP with three strong baselines—SPIRIT, Original, and Zero-shot—on HumanEval+, LiveCodeBench v1\_v3, LiveCodeBench v4\_v5, LeetCodeDataset, GSM8K, MATH500, AIME24, and AIME25. For simpler benchmarks (including HumanEval+, GSM8K, and MATH500), we evaluate the performance under four budget settings, ranging from 1K to 6K tokens. For more complex benchmarks (including LiveCodeBench v1\_v3, LiveCodeBench v4\_v5, LeetCodeDataset, AIME24, and AIME25), we evaluate the performance under six budget settings, ranging from 2K to 12K tokens. Results are shown in Table 11, Table 12, Table 13, Table 14, Table 15, Table 16, Table 17, and Table 18.

## G Training Efficiency

To quantify the training efficiency gains, we present results of the CodeForces-CoTs dataset in Table 5 and results of the OpenR1-Math dataset in Table 19. We report two key metrics: the *average number of tokens* per sample and the *average training time* measured in seconds per step.

| Methods           | Tokens               | Time                  |
|-------------------|----------------------|-----------------------|
| Original          | 5807                 | 47.82                 |
| Selective Context | 3149 (-45.8%)        | 25.85 (-45.9%)        |
| LLMLingua-2       | 3478 (-40.1%)        | 28.75 (-39.9%)        |
| TokenSkip         | 4728 (-18.6%)        | 39.20 (-18.0%)        |
| SPIRIT            | 2858 (-50.8%)        | 23.67 (-50.5%)        |
| ASAP              | <b>1834 (-68.4%)</b> | <b>15.36 (-67.9%)</b> |

Table 19: Training efficiency comparison on OpenR1-Math dataset. We report the average number of tokens per sample and training time measured in seconds per step. Percentages indicate the reduction relative to the Original baseline.

| Methods   | HE+   |       |       | LCBv1_v3 |       |       | LCBv4_v5 |       |       | LCD   |       |       |
|-----------|-------|-------|-------|----------|-------|-------|----------|-------|-------|-------|-------|-------|
|           | Acc ↑ | Tok ↓ | Lat ↓ | Acc ↑    | Tok ↓ | Lat ↓ | Acc ↑    | Tok ↓ | Lat ↓ | Acc ↑ | Tok ↓ | Lat ↓ |
| Zero-shot | 64.02 | 3334  | 1.86  | 44.12    | 7162  | 6.92  | 25.00    | 8508  | 8.90  | 27.19 | 8358  | 8.65  |
| Original  | 76.22 | 2978  | 1.63  | 52.61    | 6614  | 6.16  | 31.34    | 8202  | 8.60  | 26.32 | 8413  | 8.85  |
| SPIRIT    | 72.56 | 3159  | 1.74  | 52.61    | 6280  | 5.84  | 30.22    | 7913  | 8.45  | 26.75 | 8449  | 8.73  |
| ASAP      | 76.83 | 2494  | 1.30  | 48.86    | 3605  | 2.18  | 32.84    | 4175  | 2.69  | 27.63 | 3792  | 2.42  |

Table 9: Experimental results of different methods on code generation benchmarks with DeepSeek-R1-Distill-Llama-8B. We report accuracy (Acc), average number of generated tokens (Tok), and average generation latency (Lat) measured in seconds.

| Methods   | GSM8K |       |       | MATH500 |       |       | AIME24 |       |       | AIME25 |       |       |
|-----------|-------|-------|-------|---------|-------|-------|--------|-------|-------|--------|-------|-------|
|           | Acc ↑ | Tok ↓ | Lat ↓ | Acc ↑   | Tok ↓ | Lat ↓ | Acc ↑  | Tok ↓ | Lat ↓ | Acc ↑  | Tok ↓ | Lat ↓ |
| Zero-shot | 79.15 | 1262  | 0.36  | 57.20   | 2612  | 1.08  | 33.33  | 8445  | 10.42 | 26.67  | 8597  | 10.54 |
| Original  | 84.91 | 1310  | 0.37  | 63.00   | 2534  | 1.01  | 36.67  | 8550  | 10.04 | 30.00  | 8268  | 10.05 |
| SPIRIT    | 85.67 | 1256  | 0.35  | 62.60   | 2533  | 1.01  | 36.67  | 8788  | 10.04 | 36.67  | 8094  | 9.57  |
| ASAP      | 87.34 | 768   | 0.20  | 66.00   | 1734  | 0.65  | 36.67  | 5314  | 6.97  | 33.33  | 5348  | 7.05  |

Table 10: Experimental results of different methods on mathematical reasoning benchmarks with DeepSeek-R1-Distill-Llama-8B. We report accuracy (Acc), average number of generated tokens (Tok), and average generation latency (Lat) measured in seconds.

| Budget    | Zero-shot |       |       | Original |       |       | SPIRIT |       |       | ASAP  |       |       |
|-----------|-----------|-------|-------|----------|-------|-------|--------|-------|-------|-------|-------|-------|
|           | Acc ↑     | Tok ↓ | Lat ↓ | Acc ↑    | Tok ↓ | Lat ↓ | Acc ↑  | Tok ↓ | Lat ↓ | Acc ↑ | Tok ↓ | Lat ↓ |
| <b>1K</b> | 9.76      | 1007  | 0.28  | 14.63    | 983   | 0.28  | 10.98  | 995   | 0.28  | 23.78 | 946   | 0.27  |
| <b>2K</b> | 42.68     | 1813  | 0.53  | 43.29    | 1702  | 0.49  | 47.56  | 1690  | 0.49  | 54.88 | 1502  | 0.44  |
| <b>4K</b> | 66.46     | 2561  | 0.85  | 65.85    | 2511  | 0.82  | 69.51  | 2401  | 0.80  | 71.34 | 2116  | 0.72  |
| <b>6K</b> | 68.29     | 3051  | 1.16  | 75.61    | 2973  | 1.12  | 75.61  | 2764  | 1.07  | 78.66 | 2464  | 0.98  |

Table 11: Results of different methods under different budgets on HumanEval+. We report accuracy (Acc), average number of generated tokens (Tok), and average generation latency (Lat) measured in seconds.

| Budget     | Zero-shot |       |       | Original |       |       | SPIRIT |       |       | ASAP  |       |       |
|------------|-----------|-------|-------|----------|-------|-------|--------|-------|-------|-------|-------|-------|
|            | Acc ↑     | Tok ↓ | Lat ↓ | Acc ↑    | Tok ↓ | Lat ↓ | Acc ↑  | Tok ↓ | Lat ↓ | Acc ↑ | Tok ↓ | Lat ↓ |
| <b>2K</b>  | 16.50     | 1966  | 0.52  | 17.16    | 1920  | 0.51  | 18.95  | 1908  | 0.51  | 21.57 | 1833  | 0.49  |
| <b>4K</b>  | 32.68     | 3499  | 1.06  | 30.72    | 3432  | 1.05  | 34.80  | 3370  | 1.03  | 34.97 | 3244  | 1.00  |
| <b>6K</b>  | 39.05     | 4806  | 1.70  | 42.65    | 4673  | 1.67  | 43.14  | 4605  | 1.64  | 46.24 | 4358  | 1.54  |
| <b>8K</b>  | 44.28     | 5903  | 2.46  | 47.71    | 5723  | 2.43  | 51.80  | 5515  | 2.27  | 52.61 | 4919  | 1.90  |
| <b>10K</b> | 42.16     | 7088  | 3.59  | 52.12    | 6611  | 3.15  | 50.82  | 6524  | 3.09  | 54.74 | 5177  | 2.09  |
| <b>12K</b> | 43.95     | 7988  | 5.10  | 54.41    | 7473  | 4.22  | 51.63  | 7362  | 4.09  | 55.56 | 5322  | 2.27  |

Table 12: Results of different methods under different budgets on LiveCodeBench v1\_v3. We report accuracy (Acc), average number of generated tokens (Tok), and average generation latency (Lat) measured in seconds.

| Budget     | Zero-shot |       |       | Original |       |       | SPIRIT |       |       | ASAP  |       |       |
|------------|-----------|-------|-------|----------|-------|-------|--------|-------|-------|-------|-------|-------|
|            | Acc ↑     | Tok ↓ | Lat ↓ | Acc ↑    | Tok ↓ | Lat ↓ | Acc ↑  | Tok ↓ | Lat ↓ | Acc ↑ | Tok ↓ | Lat ↓ |
| <b>2K</b>  | 6.72      | 2021  | 0.59  | 6.34     | 1999  | 0.57  | 8.21   | 1993  | 0.56  | 13.43 | 1930  | 0.54  |
| <b>4K</b>  | 16.79     | 3820  | 1.22  | 15.67    | 3799  | 1.20  | 20.15  | 3712  | 1.18  | 20.90 | 3594  | 1.15  |
| <b>6K</b>  | 23.13     | 5444  | 2.07  | 22.76    | 5397  | 2.00  | 26.49  | 5237  | 1.93  | 30.60 | 4988  | 1.85  |
| <b>8K</b>  | 25.37     | 6927  | 3.27  | 25.74    | 6882  | 3.24  | 30.60  | 6634  | 3.09  | 35.07 | 5793  | 2.38  |
| <b>10K</b> | 25.37     | 8336  | 5.15  | 30.97    | 8289  | 4.83  | 33.58  | 7892  | 4.62  | 36.19 | 6035  | 2.61  |
| <b>12K</b> | 25.75     | 9706  | 7.44  | 32.46    | 9567  | 7.10  | 34.33  | 8987  | 6.73  | 36.57 | 6128  | 2.76  |

Table 13: Results of different methods under different budgets on LiveCodeBench v4\_v5. We report accuracy (Acc), average number of generated tokens (Tok), and average generation latency (Lat) measured in seconds.

| Budget     | Zero-shot |       |       | Original |       |       | SPIRIT |       |       | ASAP  |       |       |
|------------|-----------|-------|-------|----------|-------|-------|--------|-------|-------|-------|-------|-------|
|            | Acc ↑     | Tok ↓ | Lat ↓ | Acc ↑    | Tok ↓ | Lat ↓ | Acc ↑  | Tok ↓ | Lat ↓ | Acc ↑ | Tok ↓ | Lat ↓ |
| <b>2K</b>  | 7.02      | 2028  | 0.53  | 6.14     | 2020  | 0.53  | 7.02   | 2001  | 0.53  | 10.09 | 1965  | 0.53  |
| <b>4K</b>  | 13.16     | 3848  | 1.21  | 13.16    | 3854  | 1.21  | 16.23  | 3789  | 1.19  | 15.79 | 3758  | 1.19  |
| <b>6K</b>  | 16.23     | 5553  | 2.04  | 16.67    | 5548  | 2.04  | 18.86  | 5407  | 2.00  | 19.30 | 5387  | 2.00  |
| <b>8K</b>  | 19.30     | 7165  | 3.27  | 22.37    | 7104  | 3.18  | 22.37  | 6882  | 3.04  | 23.25 | 6722  | 2.88  |
| <b>10K</b> | 19.74     | 8680  | 4.95  | 25.00    | 8485  | 4.72  | 25.00  | 8186  | 4.45  | 27.63 | 7541  | 3.48  |
| <b>12K</b> | 21.49     | 10142 | 7.58  | 28.07    | 9717  | 7.09  | 26.32  | 9354  | 6.86  | 27.63 | 7902  | 3.83  |

Table 14: Results of different methods under different budgets on LeetCodeDataset. We report accuracy (Acc), average number of generated tokens (Tok), and average generation latency (Lat) measured in seconds.

| Budget    | Zero-shot |       |       | Original |       |       | SPIRIT |       |       | ASAP  |       |       |
|-----------|-----------|-------|-------|----------|-------|-------|--------|-------|-------|-------|-------|-------|
|           | Acc ↑     | Tok ↓ | Lat ↓ | Acc ↑    | Tok ↓ | Lat ↓ | Acc ↑  | Tok ↓ | Lat ↓ | Acc ↑ | Tok ↓ | Lat ↓ |
| <b>1K</b> | 48.75     | 963   | 0.19  | 59.29    | 942   | 0.19  | 54.66  | 925   | 0.18  | 83.93 | 693   | 0.14  |
| <b>2K</b> | 83.55     | 1301  | 0.27  | 86.35    | 1250  | 0.26  | 88.55  | 1118  | 0.23  | 90.75 | 753   | 0.16  |
| <b>4K</b> | 88.65     | 1553  | 0.37  | 90.37    | 1432  | 0.34  | 90.52  | 1227  | 0.28  | 91.28 | 778   | 0.18  |
| <b>6K</b> | 89.23     | 1714  | 0.46  | 91.05    | 1513  | 0.39  | 91.28  | 1297  | 0.33  | 91.81 | 790   | 0.20  |

Table 15: Results of different methods under different budgets on GSM8K. We report accuracy (Acc), average number of generated tokens (Tok), and average generation latency (Lat) measured in seconds.

| Budget    | Zero-shot |       |       | Original |       |       | SPIRIT |       |       | ASAP  |       |       |
|-----------|-----------|-------|-------|----------|-------|-------|--------|-------|-------|-------|-------|-------|
|           | Acc ↑     | Tok ↓ | Lat ↓ | Acc ↑    | Tok ↓ | Lat ↓ | Acc ↑  | Tok ↓ | Lat ↓ | Acc ↑ | Tok ↓ | Lat ↓ |
| <b>1K</b> | 19.00     | 1020  | 0.19  | 28.00    | 1017  | 0.19  | 18.40  | 1012  | 0.19  | 36.40 | 935   | 0.19  |
| <b>2K</b> | 42.20     | 1804  | 0.39  | 52.00    | 1767  | 0.39  | 54.40  | 1592  | 0.36  | 59.80 | 1347  | 0.31  |
| <b>4K</b> | 60.40     | 2629  | 0.70  | 63.80    | 2511  | 0.66  | 64.20  | 2144  | 0.57  | 70.80 | 1649  | 0.43  |
| <b>6K</b> | 66.60     | 3100  | 0.94  | 70.60    | 2843  | 0.84  | 69.60  | 2460  | 0.74  | 71.00 | 1758  | 0.52  |

Table 16: Results of different methods under different budgets on MATH500. We report accuracy (Acc), average number of generated tokens (Tok), and average generation latency (Lat) measured in seconds.

| Budget     | Zero-shot |       |       | Original |       |       | SPIRIT |       |       | ASAP  |       |       |
|------------|-----------|-------|-------|----------|-------|-------|--------|-------|-------|-------|-------|-------|
|            | Acc ↑     | Tok ↓ | Lat ↓ | Acc ↑    | Tok ↓ | Lat ↓ | Acc ↑  | Tok ↓ | Lat ↓ | Acc ↑ | Tok ↓ | Lat ↓ |
| <b>2K</b>  | 0.00      | 2048  | 0.93  | 16.67    | 2048  | 0.93  | 6.67   | 2048  | 0.93  | 10.00 | 1978  | 0.93  |
| <b>4K</b>  | 20.00     | 4003  | 2.12  | 20.00    | 3984  | 2.11  | 23.33  | 3877  | 2.10  | 36.67 | 3415  | 1.99  |
| <b>6K</b>  | 33.33     | 5682  | 3.58  | 36.67    | 5695  | 3.55  | 40.00  | 5216  | 3.32  | 36.67 | 4410  | 3.03  |
| <b>8K</b>  | 30.00     | 7073  | 5.30  | 40.00    | 7093  | 5.14  | 46.67  | 6243  | 4.47  | 40.00 | 5159  | 4.10  |
| <b>10K</b> | 36.67     | 8352  | 6.76  | 46.67    | 8034  | 6.43  | 46.67  | 7198  | 5.78  | 46.67 | 5552  | 5.04  |
| <b>12K</b> | 40.00     | 9318  | 7.84  | 46.67    | 8990  | 7.75  | 46.67  | 8363  | 7.21  | 46.67 | 5767  | 5.88  |

Table 17: Results of different methods under different budgets on AIME24. We report accuracy (Acc), average number of generated tokens (Tok), and average generation latency (Lat) measured in seconds.

| Budget     | Zero-shot |       |       | Original |       |       | SPIRIT |       |       | ASAP  |       |       |
|------------|-----------|-------|-------|----------|-------|-------|--------|-------|-------|-------|-------|-------|
|            | Acc ↑     | Tok ↓ | Lat ↓ | Acc ↑    | Tok ↓ | Lat ↓ | Acc ↑  | Tok ↓ | Lat ↓ | Acc ↑ | Tok ↓ | Lat ↓ |
| <b>2K</b>  | 6.67      | 2046  | 0.95  | 13.33    | 2048  | 0.94  | 3.33   | 2044  | 0.94  | 10.00 | 2020  | 0.95  |
| <b>4K</b>  | 20.00     | 3851  | 2.17  | 26.67    | 3834  | 2.16  | 16.67  | 3792  | 2.15  | 20.00 | 3511  | 2.15  |
| <b>6K</b>  | 30.00     | 5369  | 3.62  | 33.33    | 5452  | 3.63  | 36.67  | 5360  | 3.63  | 30.00 | 4484  | 3.14  |
| <b>8K</b>  | 36.67     | 6848  | 5.30  | 36.67    | 6798  | 5.32  | 36.67  | 6611  | 5.16  | 33.33 | 5002  | 4.10  |
| <b>10K</b> | 40.00     | 8145  | 6.67  | 43.33    | 8026  | 6.75  | 43.33  | 7817  | 6.57  | 36.67 | 5434  | 5.10  |
| <b>12K</b> | 36.67     | 9442  | 8.25  | 40.00    | 9461  | 8.32  | 46.67  | 8598  | 7.78  | 36.67 | 5720  | 6.06  |

Table 18: Results of different methods under different budgets on AIME25. We report accuracy (Acc), average number of generated tokens (Tok), and average generation latency (Lat) measured in seconds.