

# TokenSkip: Controllable Chain-of-Thought Compression in LLMs

Anonymous ACL submission

## Abstract

Chain-of-Thought (CoT) has been proven effective in enhancing the reasoning capabilities of large language models (LLMs). Recent advancements, such as OpenAI’s o1 and DeepSeek-R1, suggest that scaling up the length of CoT sequences during inference could further boost LLM reasoning performance. However, due to the autoregressive nature of LLM decoding, longer CoT outputs lead to a linear increase in inference latency, adversely affecting user experience, particularly when the CoT exceeds 10,000 tokens. To address this limitation, we analyze the semantic importance of tokens within CoT outputs and reveal that their contributions to reasoning vary. Building on this insight, we propose TokenSkip, a simple yet effective approach that enables LLMs to selectively skip less important tokens, allowing for controllable CoT compression. Extensive experiments across various models and tasks demonstrate the effectiveness of TokenSkip in reducing CoT token usage while preserving strong reasoning performance. Notably, when applied to Qwen2.5-14B-Instruct, TokenSkip reduces reasoning tokens by 40% (from 313 to 181) on GSM8K, with less than a 0.4% performance drop<sup>1</sup>.

## 1 Introduction

Chain-of-Thought (CoT) prompting (Nye et al., 2021; Wei et al., 2022; Kojima et al., 2022) has emerged as a cornerstone strategy for enhancing Large Language Models (LLMs) in complex reasoning tasks. By eliciting step-by-step inference, CoT enables LLMs to decompose intricate problems into manageable subtasks, thereby improving their problem-solving performance (Yao et al., 2023; Wang et al., 2023; Zhou et al., 2023; Shinn et al., 2023). Recent advancements, such as OpenAI’s o1 (OpenAI et al., 2024) and DeepSeek-

<sup>1</sup>All of our codes and checkpoints will be released to facilitate future research.

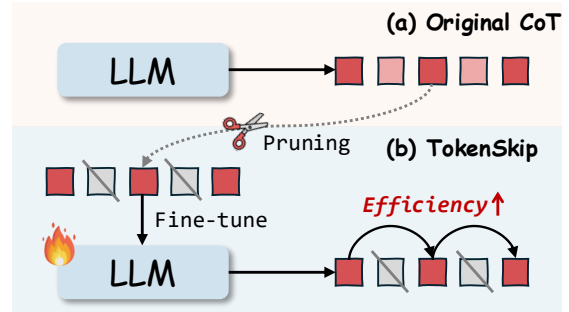


Figure 1: In contrast to vanilla CoT that generates all reasoning tokens sequentially, TokenSkip enables LLMs to skip tokens with less semantic importance (e.g.,  $\square$ ) and learn shortcuts between critical reasoning tokens, facilitating controllable CoT compression.

R1 (DeepSeek-AI et al., 2025), further demonstrate that scaling up CoT lengths from hundreds to thousands of reasoning steps could continuously improve LLM reasoning. These breakthroughs have underscored CoT’s potential to advance LLM capabilities, expanding the boundaries of AI-driven problem-solving.

Despite its effectiveness, the increased length of CoT sequences introduces substantial computational overhead. Due to the autoregressive nature of LLM decoding, longer CoT outputs lead to proportional increases in both inference latency and memory footprints of key-value cache. Additionally, the quadratic computational cost of attention layers further exacerbates this burden. These issues become particularly pronounced when CoT sequences extend into thousands of reasoning steps, resulting in significant computational costs and prolonged response times. While prior research has explored methods for selectively skipping reasoning steps (Ding et al., 2024; Liu et al., 2024), recent findings (Jin et al., 2024; Merrill and Sabharwal, 2024) suggest that such reductions may conflict with test-time scaling (OpenAI, 2024; Snell et al., 2025), ultimately impairing LLM reasoning per-

065 formance. Therefore, striking an optimal balance  
066 between CoT efficiency and reasoning accuracy  
067 remains a critical open challenge.

068 In this work, we delve into CoT efficiency and  
069 seek the answer to an important question: “Does  
070 every token in the CoT output contribute equally to  
071 deriving the answer?” We empirically analyze the  
072 semantic importance of tokens within CoT outputs  
073 and reveal that their contributions to the reasoning  
074 performance vary, as depicted in Figure 2. Building  
075 on this insight, we introduce TokenSkip, a simple  
076 yet effective approach that enables LLMs to skip  
077 less important tokens within CoT sequences and  
078 learn shortcuts between critical reasoning tokens,  
079 thereby allowing for controllable CoT compression  
080 with adjustable ratios. Specifically, as shown in  
081 Figure 1, TokenSkip constructs compressed CoT  
082 training data with various compression ratios, by  
083 pruning unimportance tokens from original LLM  
084 CoT trajectories. Then, it conducts a general super-  
085 vised fine-tuning process on target LLMs with this  
086 training data, facilitating LLMs to automatically  
087 trim redundant tokens during reasoning.

088 We conduct extensive experiments across var-  
089 ious models, including LLaMA-3.1-8B-Instruct  
090 and the Qwen2.5-Instruct series, using two widely  
091 recognized math reasoning benchmarks: GSM8K  
092 and MATH-500. The results validate the effec-  
093 tiveness of TokenSkip in compressing CoT out-  
094 puts while maintaining robust reasoning perfor-  
095 mance. Notably, Qwen2.5-14B-Instruct exhibits  
096 almost **NO** performance drop (less than 0.4%) with  
097 a 40% reduction in token usage on GSM8K. On  
098 the challenging MATH-500 dataset, LLaMA-3.1-  
099 8B-Instruct effectively reduces CoT token usage  
100 by 30% with a performance decline of less than  
101 4%, resulting in a 1.4× inference speedup. Further  
102 analysis underscores the coherence of TokenSkip  
103 in specified compression ratios and its potential  
104 scalability with stronger compression techniques.

105 TokenSkip is distinguished by its low training  
106 cost. For Qwen2.5-14B-Instruct, TokenSkip fine-  
107 tunes only 0.2% of the model’s parameters using  
108 LoRA. The size of the compressed CoT training  
109 data is no larger than that of the original training  
110 set, with 7,473 examples in GSM8K and 7,500  
111 in MATH. The training is completed in approxi-  
112 mately 2 hours for the 7B model and 2.5 hours for  
113 the 14B model on two 3090 GPUs. These char-  
114 acteristics make TokenSkip an efficient and repro-  
115 ducible approach, suitable for use in efficient and  
116 cost-effective LLM deployment.

To sum up, our key contributions are:

1. To the best of our knowledge, this work is  
the *first* to investigate the potential of enhanc-  
ing CoT efficiency through *token skipping*,  
inspired by the varying semantic importance  
of tokens in CoT trajectories of LLMs.
2. We introduce TokenSkip, a simple yet effec-  
tive approach that enables LLMs to skip re-  
dundant tokens within CoTs and learn short-  
cuts between critical tokens, facilitating CoT  
compression with adjustable ratios.
3. Our experiments validate the effectiveness of  
TokenSkip. When applied to Qwen2.5-14B-  
Instruct, TokenSkip reduces reasoning tokens  
by 40% (from 313 to 181) on GSM8K, with  
less than a 0.4% performance drop.

## 2 Background and Preliminaries

In this section, we discuss the relevant research  
background and present preliminary studies on to-  
ken efficiency in CoT sequences, exploring its im-  
pact on the reasoning performance of LLMs.

### 2.1 Token Importance

We first investigate a critical research question to  
CoT efficiency: “Does every token in the CoT out-  
put contribute equally to deriving the answer?” In  
other words, we would like to know if there is any  
token redundancy in CoT sequences that could be  
eliminated to improve CoT efficiency.

Token redundancy has been recognized as a  
longstanding and fundamental issue in LLM ef-  
ficiency (Hou et al., 2022; Zhang et al., 2023; Lin  
et al., 2024; Chen et al., 2024). Recently, it has gar-  
nered intensive research attention in prompt com-  
pression (Li et al., 2023; Jiang et al., 2023; Pan  
et al., 2024), which focuses on removing redundant  
tokens from input prompt to reduce API token us-  
age. To address this issue, Selective Context (Li  
et al., 2023) proposed to measure the importance  
of tokens in a piece of text based on the semantic  
confidence of LLMs:

$$I_1(x_i) = -\log P(x_i | \mathbf{x}_{<i}; \theta_{\mathcal{M}_L}), \quad (1)$$

where  $\mathbf{x} = \{x_i\}_{i=1}^n$  is the given text,  $x_i$  denotes  
a token, and  $\mathcal{M}_L$  denotes the LLM used to com-  
pute the confidence of each token. Intuitively, such  
measurement could be seamlessly applied to CoT  
tokens generated by LLMs. We show an example  
of this measurement in Figure 2.

**Problem:** Marcus is half of Leo’s age and five years younger than Deanna. Deanna is 26. How old is Leo?

**Chain-of-Thought:** Let’s break it down step by step: 1. Deanna is 26 years old. 2. Marcus is five years younger than Deanna, so Marcus is  $26 - 5 = 21$  years old. 3. Marcus is half of Leo’s age, so Leo’s age is twice Marcus’s age. 4. Since Marcus is 21, Leo’s age is  $2 \times 21 = 42$ . (Selective Context)

**Chain-of-Thought:** Let’s break it down step by step: 1. Deanna is 26 years old. 2. Marcus is five years younger than Deanna, so Marcus is  $26 - 5 = 21$  years old. 3. Marcus is half of Leo’s age, so Leo’s age is twice Marcus’s age. 4. Since Marcus is 21, Leo’s age is  $2 \times 21 = 42$ . (LLMLingua-2)

**Final Answer:** 42.

Figure 2: Visualization of token importance within a CoT sequence, with darker colors indicating higher values. This figure compares two token importance measurements: Selective Context and LLMLingua-2.

Despite its simplicity, LLMLingua-2 (Pan et al., 2024) argued that there exist two major limitations in the aforementioned measurement that hinder the compression performance. Firstly, as shown in Figure 2, the intrinsic nature of LLM perplexity leads to lower importance measures (i.e., higher confidence) for tokens at the end of the sentence. Such position dependency impacts the factual importance measurement of each token. Furthermore, the unidirectional attention mechanism in causal LMs may fail to capture all essential information needed for token importance within the text.

To tackle these limitations, LLMLingua-2 introduced utilizing a bidirectional BERT-like LM (Devlin et al., 2019) for token importance measurement. It utilizes GPT-4 (OpenAI, 2023) to label each token as “important” or not and trains the bidirectional LM with a token classification objective. The token importance is measured by the predicted probability of each token:

$$I_2(x_i) = P(x_i | \mathbf{x}_{\leq n}; \theta_{\mathcal{M}_B}), \quad (2)$$

where  $\mathcal{M}_B$  denotes the bidirectional LM.

In this study, we apply LLMLingua-2 as the token importance measurement to LLM CoT outputs. Similar to plain text, we observe that the semantic importance of tokens within CoT outputs varies, as shown in Figure 2. For instance, mathematical equations tend to have a greater contribution to the final answer, consistent with recent research (Ma et al., 2024). In contrast, semantic connectors such as “so” and “since” generally contribute less. These findings highlight the token redundancy in LLM

#### Recovering the Compressed Chain-of-Thought

**Compressed CoT:** break down Deanna 26 Marcus five younger 26 - 5 21 Marcus half Leo’s age twice Marcus Marcus 21, Leo’s age  $2 \times 21 = 42$ .

**Recovered Compressed CoT:** Let’s break it down step by step. Deanna is 26 years old. Marcus is five years younger than Deanna:  $M = D - 5$ . Marcus’s age:  $M = 26 - 5 = 21$ . Marcus is half of Leo’s age:  $M = L / 2$ . Leo is twice Marcus’s age:  $L = 2M$ . Leo’s age:  $L = 2 \times 21 = 42$ .

Figure 3: Recovering the compressed CoT for GSM8K math word problem using LLaMA-3.1-8B-Instruct.

CoT outputs and the substantial potential to enhance CoT efficiency by trimming this redundancy.

## 2.2 CoT Recovery

We further explore the following research question: “Are LLMs capable of restoring the CoT process from compressed outputs?” The answer is yes. As shown in Figure 3 and detailed in Appendix A, examples restored from compressed CoTs using LLaMA-3.1-8B-Instruct demonstrate that LLMs could effectively comprehend the semantic information encoded in the compressed CoT and restore the CoT process. This capability ensures that the interpretability of compressed CoTs is maintained. Additionally, when required by users, the complete CoT process can be recovered and presented.

In summary, the empirical analysis above underscores the potential of trimming redundant tokens to enhance CoT efficiency, as well as the ability of LLMs to restore CoT from compressed outputs. However, enabling LLMs to autonomously skip redundant CoT tokens and identify shortcuts between critical reasoning tokens presents a non-trivial challenge. To the best of our knowledge, this work is the *first* to explore CoT compression through *token skipping*. In the following sections, we present our proposed methodology in detail.

## 3 TokenSkip

We introduce TokenSkip, a simple yet effective approach that enables LLMs to skip less important tokens, enabling controllable CoT compression with adjustable ratios. This section demonstrates the details of our methodology, including token pruning (§3.1), training (§3.2), and inference (§3.3).

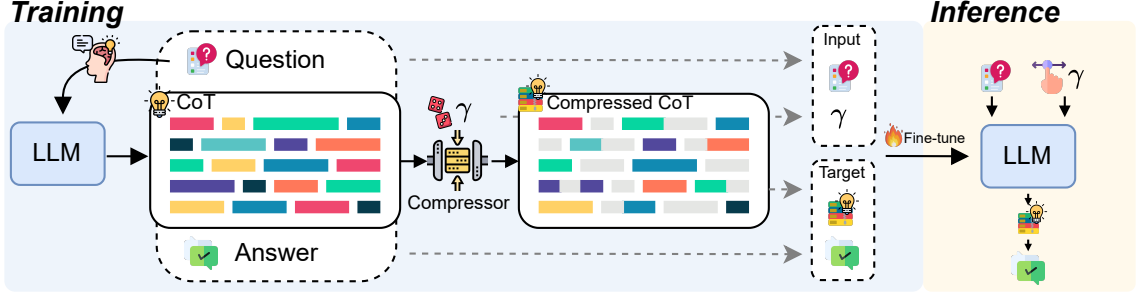


Figure 4: Illustration of TokenSkip. During the training phase, TokenSkip first generates CoT trajectories from the target LLM. These CoTs are then compressed to a specified ratio,  $\gamma$ , based on the semantic importance of tokens. TokenSkip fine-tunes the target LLM using compressed CoTs, enabling controllable CoT inference at the desired  $\gamma$ .

### 3.1 Token Pruning

The key insight behind TokenSkip is that “each reasoning token contributes differently to deriving the answer.” To enhance CoT efficiency, we propose to trim redundant tokens from LLM CoT outputs and fine-tune LLMs using these trimmed CoT trajectories. The token pruning process is guided by the concept of *token importance*, as detailed in Section 2.1.

Specifically, given a target LLM  $\mathcal{M}$ , one of its CoT trajectories  $c = \{c_i\}_{i=1}^m$ , and a desired compression ratio  $\gamma \in [0, 1]$ , TokenSkip first calculates the semantic importance of each CoT token  $I(c)$ , as defined in Eq (2). The tokens are then ranked in descending order based on their importance values. Next, the  $\gamma$ -th percentile of these importance values is computed, representing the threshold for token pruning:

$$I_\gamma = \text{np.percentile}([I(c_1), \dots, I(c_m)], \gamma). \quad (3)$$

Finally, CoT tokens with an importance value greater than or equal to  $I_\gamma$  are retained in the compressed CoT trajectory:

$$\tilde{c} = \{c_i \mid I(c_i) \geq I_\gamma, 1 \leq i \leq m\}. \quad (4)$$

### 3.2 Training

Given a training dataset  $\mathcal{D}$  with  $N$  samples and a target LLM  $\mathcal{M}$ , we first obtain  $N$  CoT trajectories with  $\mathcal{M}$ . Then, we filter out trajectories with incorrect answers to ensure the high quality of training data. For the remaining CoT trajectories, we prune each CoT with a randomly selected compression ratio  $\gamma$ , as demonstrated in Section 3.1. For each  $\langle \text{question, compressed CoT, answer} \rangle$ , we inserted the compression ratio  $\gamma$  after the question. Finally, each training sample is formatted as follows:

$$\mathcal{Q} [\text{EOS}] \gamma [\text{EOS}] \text{Compressed CoT } \mathcal{A},$$

where  $\langle \mathcal{Q}, \mathcal{A} \rangle$  indicates the  $\langle \text{question, answer} \rangle$  pair. Formally, given a question  $\mathbf{x}$ , compression ratio  $\gamma$ , and the output sequence  $\mathbf{y} = \{y_i\}_{i=1}^l$ , which includes the compressed CoT  $\tilde{c}$  and the answer  $\mathbf{a}$ , we fine-tune the target LLM  $\mathcal{M}$ , enabling it to perform chain-of-thought in a compressed pattern by minimizing

$$\mathcal{L} = \sum_{i=1}^l \log P(y_i \mid \mathbf{x}, \gamma, \mathbf{y}_{<i}; \theta_{\mathcal{M}}), \quad (5)$$

where  $\mathbf{y} = \{\tilde{c}_1, \dots, \tilde{c}_{m'}, a_1, \dots, a_t\}$ . Note that the compression is performed solely on CoT sequences, and we keep the answer  $\mathbf{a} = \{a_i\}_{i=1}^t$  unchanged. To preserve LLMs’ reasoning capabilities, we also include a portion of the original CoT trajectories in the training data, with  $\gamma$  set to 1.

### 3.3 Inference

The inference of TokenSkip follows autoregressive decoding. Compared to original CoT outputs that may contain redundancy, TokenSkip facilitates LLMs to skip *unimportant* tokens during the chain-of-thought process, thereby enhancing reasoning efficiency. Formally, given a question  $\mathbf{x}$  and the compression ratio  $\gamma$ , the input prompt of TokenSkip follows the same format adopted in fine-tuning, which is  $\mathcal{Q} [\text{EOS}] \gamma [\text{EOS}]$ . The LLM  $\mathcal{M}$  sequentially predicts the output sequence  $\hat{\mathbf{y}}$ :

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}^*} \sum_{j=1}^{l'} \log P(y_j \mid \mathbf{x}, \gamma, \mathbf{y}_{<j}; \theta_{\mathcal{M}}),$$

where  $\hat{\mathbf{y}} = \{\hat{c}_1, \dots, \hat{c}_{m''}, \hat{a}_1, \dots, \hat{a}_{t'}\}$  denotes the output sequence, which includes CoT tokens  $\hat{c}$  and the answer  $\hat{a}$ . We illustrate the training and inference process of TokenSkip in Figure 4.



## 4 Experiments

### 4.1 Experimental Setup

**Models and Datasets** We primarily evaluate our method using LLaMA-3.1-8B-Instruct (Dubey et al., 2024) and Qwen2.5-Instruct series (Yang et al., 2024). The evaluation leverages two widely-used math reasoning benchmarks: GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021). For training, we use the respective training sets from both datasets. Regarding the MATH dataset, due to the computation cost, we assess our method on a subset, MATH-500, which is identical to the test set used in Lightman et al. (2024). The subset comprises 500 representative problems, and we find that its evaluation yields results comparable to those from the full dataset.

**Implementation Details** We utilize LLMingua-2 (Pan et al., 2024) as the token importance metric to generate our compressed CoT training data. The compression ratio  $\gamma$  is randomly selected from  $\{0.5, 0.6, 0.7, 0.8, 0.9, 1.0\}$  for each training sample. We adopt LoRA (Hu et al., 2022), an efficient and reproducible approach that has been widely verified as effective in LLM fine-tuning, to train our models. The rank  $r$  is set to 8, and the scaling parameter  $\alpha$  is set to 16. TokenSkip is characterized by its low training cost, with training taking  $\sim 2$  hours for the 7B model and  $\sim 2.5$  hours for the 14B model on 3090 GPUs. During inference, the maximum number of tokens `max_len` is set to 512 for GSM8K and 1024 for MATH<sup>2</sup>. All experiments are conducted using Pytorch 2.1.0 on  $2 \times$  NVIDIA GeForce RTX 3090 GPU (24GB) with CUDA 12.1, and an Intel(R) Xeon(R) Platinum 8370C CPU with 32 cores. We include more implementation details in Appendix B.1.

**Baselines** In our main experiments, we compare TokenSkip to two commonly used length control baselines: **1) Prompt-based Reduction.** In this approach, we instruct the LLM to reduce a fixed proportion of output tokens in the CoT process. Specifically, we append a prompt such as “Please reduce 50% of the words in your Chain-of-Thought process.” to the input instruction. **2) Truncation.** This method involves brute-force length truncation, where the maximum number of output tokens is restricted, compressing the CoT output to a fixed

<sup>2</sup>Since many samples reach the maximum length when testing TokenSkip on MATH-500, we adjust its length budget to `max_len` $\times\gamma$ , with no adjustment for GSM8K.

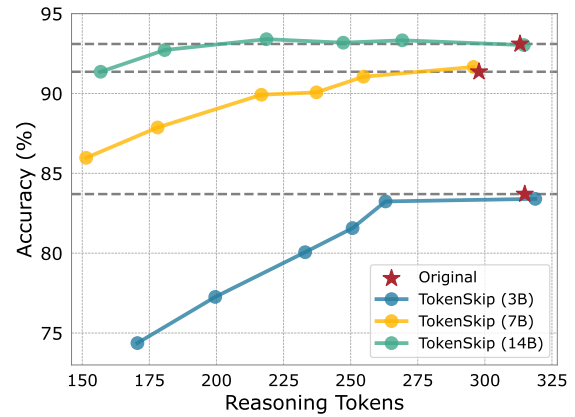


Figure 5: Compression performance of TokenSkip on Qwen2.5-Instruct models. Qwen2.5-14B-Instruct shows almost **no** performance drop with 40% token trimming.

ratio. These baselines are referred to as Prompt and Truncation in Table 1, respectively.

**Evaluation Metrics** We evaluate TokenSkip using three widely used metrics: accuracy, the number of CoT tokens, and inference latency per sample. Model performance is assessed using scripts from DeepSeek-Math<sup>3</sup>. Greedy decoding is employed to generate the outputs from the target LLM. Inference latency is measured on a single NVIDIA 3090 GPU with a batch size of 1. In addition to these metrics, we report the actual compression ratio of the CoTs to assess whether the compression aligns with the specified ratio.

### 4.2 Main Results

The performance of TokenSkip on GSM8K using the Qwen2.5-Instruct series<sup>4</sup> is illustrated in Figure 5. As the model scale increases, there is less performance degradation at higher compression ratios, indicating that larger LLMs are better at identifying shortcuts between critical reasoning tokens, enabling more efficient CoT generation. Notably, Qwen2.5-14B-Instruct exhibits almost **NO** performance drop (less than 0.4%) with 40% token trimming. Even at a compression ratio of 0.5, the model maintains strong reasoning capabilities, with only 2% performance degradation. These results highlight the substantial potential of TokenSkip to reduce CoT token usage and accelerate reasoning in large-scale LLMs. Due to computational constraints, experiments with larger models are not conducted and are left for future exploration.

We further compare TokenSkip with two widely

<sup>3</sup><https://github.com/deepseek-ai/DeepSeek-Math>

<sup>4</sup>For detailed results, please refer to Appendix B.2.

Methods	Ratio	GSM8K				MATH-500			
		Accuracy $\uparrow$	Tokens $\downarrow$	Latency (s) $\downarrow$	<i>ActRatio</i>	Accuracy $\uparrow$	Tokens $\downarrow$	Latency (s) $\downarrow$	<i>ActRatio</i>
Original	-	86.2 <sub>(0.0<math>\downarrow</math>)</sub>	213.17	5.96 <sub>1.0<math>\times</math></sub>	-	48.6 <sub>(0.0<math>\downarrow</math>)</sub>	502.60	16.37 <sub>1.0<math>\times</math></sub>	-
Prompt	0.9	84.1 <sub>(2.1<math>\downarrow</math>)</sub>	226.37	6.12 <sub>1.0<math>\times</math></sub>	1.06	48.6 <sub>(0.0<math>\downarrow</math>)</sub>	468.04	15.39 <sub>1.1<math>\times</math></sub>	0.93
	0.7	84.9 <sub>(1.3<math>\downarrow</math>)</sub>	209.39	5.51 <sub>1.1<math>\times</math></sub>	0.98	48.4 <sub>(0.4<math>\downarrow</math>)</sub>	472.13	15.55 <sub>1.1<math>\times</math></sub>	0.94
	0.5	83.7 <sub>(2.5<math>\downarrow</math>)</sub>	188.82	4.97 <sub>1.2<math>\times</math></sub>	0.89	47.8 <sub>(0.4<math>\downarrow</math>)</sub>	471.11	15.48 <sub>1.1<math>\times</math></sub>	0.94
Truncation	0.9	70.2 <sub>(26.0<math>\downarrow</math>)</sub>	202.06	5.29 <sub>1.1<math>\times</math></sub>	0.95	47.8 <sub>(0.8<math>\downarrow</math>)</sub>	440.33	14.56 <sub>1.1<math>\times</math></sub>	0.88
	0.7	25.9 <sub>(60.3<math>\downarrow</math>)</sub>	149.99	3.97 <sub>1.5<math>\times</math></sub>	0.70	45.0 <sub>(3.6<math>\downarrow</math>)</sub>	386.89	12.85 <sub>1.3<math>\times</math></sub>	0.77
	0.5	7.0 <sub>(79.2<math>\downarrow</math>)</sub>	103.69	2.95 <sub>2.0<math>\times</math></sub>	0.49	27.4 <sub>(21.2<math>\downarrow</math>)</sub>	283.70	9.40 <sub>1.7<math>\times</math></sub>	0.56
TokenSkip	1.0	86.7 <sub>(0.5<math>\uparrow</math>)</sub>	213.60	5.98 <sub>1.0<math>\times</math></sub>	1.00	48.2 <sub>(0.4<math>\downarrow</math>)</sub>	504.79	16.43 <sub>1.0<math>\times</math></sub>	1.00
	0.9	86.1 <sub>(0.1<math>\downarrow</math>)</sub>	198.01	5.65 <sub>1.1<math>\times</math></sub>	0.93	47.8 <sub>(0.8<math>\downarrow</math>)</sub>	448.31	15.26 <sub>1.1<math>\times</math></sub>	0.89
	0.8	84.3 <sub>(1.9<math>\downarrow</math>)</sub>	169.89	5.13 <sub>1.2<math>\times</math></sub>	0.80	47.3 <sub>(1.3<math>\downarrow</math>)</sub>	398.94	13.39 <sub>1.2<math>\times</math></sub>	0.79
	0.7	82.5 <sub>(3.7<math>\downarrow</math>)</sub>	150.12	4.36 <sub>1.4<math>\times</math></sub>	0.70	46.7 <sub>(1.9<math>\downarrow</math>)</sub>	349.13	11.55 <sub>1.4<math>\times</math></sub>	0.69
	0.6	81.1 <sub>(5.1<math>\downarrow</math>)</sub>	129.38	3.81 <sub>1.6<math>\times</math></sub>	0.61	42.0 <sub>(6.6<math>\downarrow</math>)</sub>	318.36	10.58 <sub>1.6<math>\times</math></sub>	0.63
	0.5	78.2 <sub>(8.0<math>\downarrow</math>)</sub>	113.05	3.40 <sub>1.8<math>\times</math></sub>	0.53	40.2 <sub>(8.4<math>\downarrow</math>)</sub>	292.17	9.67 <sub>1.7<math>\times</math></sub>	0.58

Table 1: Experimental results of TokenSkip on LLaMA-3.1-8B-Instruct. We report accuracy, average CoT token count (Tokens), average latency per sample, and actual compression ratio (*ActRatio*) for comparison.

used length control baselines — prompt-based reduction and truncation. The experimental results are presented in Table 1. As shown, prompt-based reduction fails to achieve the specified compression ratio, with the actual ratio exceeding 0.89 even when the target is set to 0.5. While truncation adheres to the specified ratio, it results in significant degradation in reasoning performance. Specifically, at a compression ratio of 0.5, truncation causes a 79% accuracy drop on GSM8K and a 21% drop on MATH-500. In contrast, TokenSkip ensures adherence to the specified compression ratio (see Figure 6) while preserving strong reasoning capabilities. Notably, TokenSkip achieves an actual compression ratio of **0.53** on GSM8K with only a 10% performance drop, resulting in a **1.8 $\times$**  speedup in average latency. On the challenging MATH-500 dataset, TokenSkip effectively reduces CoT token usage by **30%** with a performance drop of less than 4%. These results validate the effectiveness of TokenSkip.

### 4.3 Analysis

**Compression Ratio** In our main results, we focus on compression ratios greater than 0.5. To further investigate the performance of TokenSkip at lower compression ratios, we train an additional variant, denoted as *More Ratio*, with extra compression ratios of 0.3 and 0.4. As shown in Figure 6, the ratio adherence of models largely degrades at these lower ratios. We attribute this decline to the excessive trimming of reasoning tokens, which likely causes a loss of critical information in the completions, hindering the effective training of

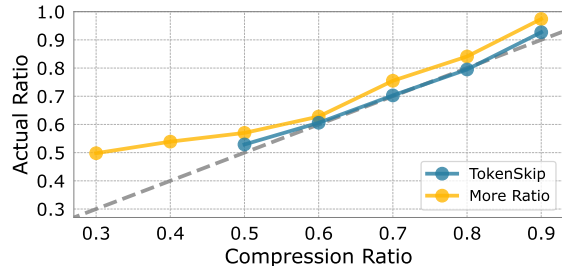


Figure 6: Comparison of ratio adherence across different compression ratio settings. The experimental results are obtained with LLaMA-3.1-8B-Instruct on GSM8K.

LLMs to learn CoT compression. Furthermore, we observe that the overall adherence of *More Ratio* is not as good as TokenSkip with the default settings, which further supports our hypothesis.

**Importance Metric** Figure 7 presents a performance comparison of TokenSkip across different token importance metrics. In addition to the metrics discussed in Section 2.1, we include GPT-4o<sup>5</sup> as a strong token importance metric for comparison. Specifically, for a given CoT trajectory, we prompt GPT-4o to trim redundant tokens according to a specified compression ratio, without adding any additional tokens. Additionally, we ask GPT-4o to suggest the *optimal* compression format of the CoT trajectory, referred to as GPT-4o-Optimal in Figure 7. We incorporate all training data generated by GPT-4o to train a variant of TokenSkip. We use the “[optimal]” token to prompt the model, obtaining the results of GPT-4o-Optimal.

As illustrated in Figure 7, TokenSkip utilizing

<sup>5</sup>We use the gpt-4o-2024-08-06 version for experiments.

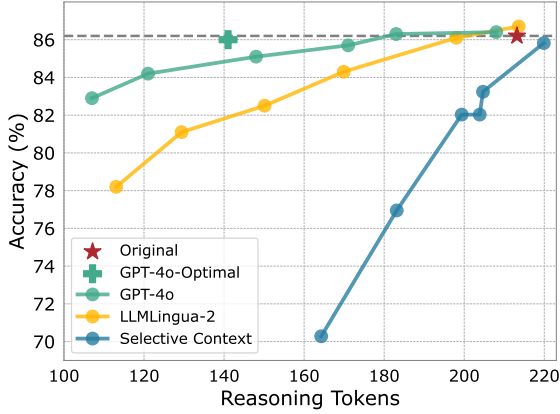


Figure 7: Performance comparison of TokenSkip using different token importance metrics, evaluated with LLaMA-3.1-8B-Instruct on GSM8K.

LLMLingua-2 (Pan et al., 2024) outperforms the variant with Selective Context (Li et al., 2023), which aligns with our demonstrations in Section 2.1. Additionally, incorporating GPT-4o for token importance measurement further enhances compression performance, suggesting that a more robust CoT compressor could improve TokenSkip even further. However, the API costs associated with GPT-4o make it impractical for processing large datasets. In contrast, LLMLingua-2, which includes a BERT-size model, offers a cost-effective and efficient alternative for training TokenSkip. Furthermore, GPT-4o-Optimal achieves a better balance between reasoning accuracy and CoT token reduction, emphasizing the potential of flexible compression ratios in CoT generation — an avenue we plan to explore in future work.

**Length Budget** As outlined in Section 4.1, we adjust the maximum length budget to  $\max\_len \times \gamma$  when evaluating TokenSkip on MATH-500, ensuring a fair comparison of compression ratios. However, this brute-force length truncation inevitably impacts the reasoning performance of LLMs, as LLMs are unable to complete the full generation. In this analysis, we explore whether LLMs can “think” more effectively using a compressed CoT format. Specifically, we evaluate TokenSkip under the same length budget as the original LLM (e.g., 1024 for MATH-500). The experimental results, shown in Figure 8, demonstrate a significant performance improvement of TokenSkip under this length budget, compared to those adjusted by compression ratios. Notably, with compression ratios of 0.7, 0.8, and 0.9, TokenSkip outperforms the original LLM, yielding an absolute performance in-

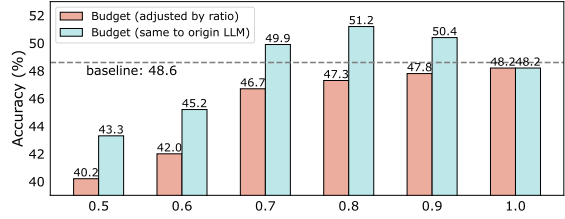


Figure 8: Performance comparison of TokenSkip with varying maximum length constraints, evaluated with LLaMA-3.1-8B-Instruct on the MATH-500 dataset.

crease of 1.3 to 2.6 points. These findings highlight TokenSkip’s potential to enhance the reasoning capabilities of LLMs within the same length budget.

**Case Study** Figure 9 presents several examples of TokenSkip, derived from the test sets of GSM8K and MATH-500. These examples clearly illustrate that TokenSkip allows LLMs to learn shortcuts between critical reasoning tokens, rather than generating shorter CoTs from scratch. For instance, in the first case, TokenSkip facilitates LLaMA-3.1-8B-Instruct to skip semantic connectors such as “of” and “the”, as well as expressions that contribute minimally to the reasoning, such as the first sentence. Notably, we observe that numeric values and mathematical equations are prioritized for retention in most cases. This finding aligns with recent research (Ma et al., 2024), which suggests that mathematical expressions may contribute more significantly to reasoning than CoT in natural language. Furthermore, we find that TokenSkip does not reduce the number of reasoning steps but instead trims redundant tokens within those steps.

## 5 Related Work

**Efficient CoT** While Chain-of-Thought (CoT) enhances task performance by simulating human-like reasoning patterns, its reasoning steps introduce significant computational overhead. As a result, researchers have sought methods to reduce this overhead while retaining the benefits of CoT. One intuitive approach is to simplify, skip (Marconato et al., 2024; Ding et al., 2024; Liu et al., 2024), or generate thinking steps in parallel (Ning et al., 2023) to improve efficiency. Another strategy involves compressing reasoning steps into continuous latent representations (Goyal et al., 2024; Deng et al., 2024; Hao et al., 2024; Cheng and Van Durme, 2024), allowing LLMs to reason without explicitly generating discrete word tokens. To minimize the generation of redundant natural lan-

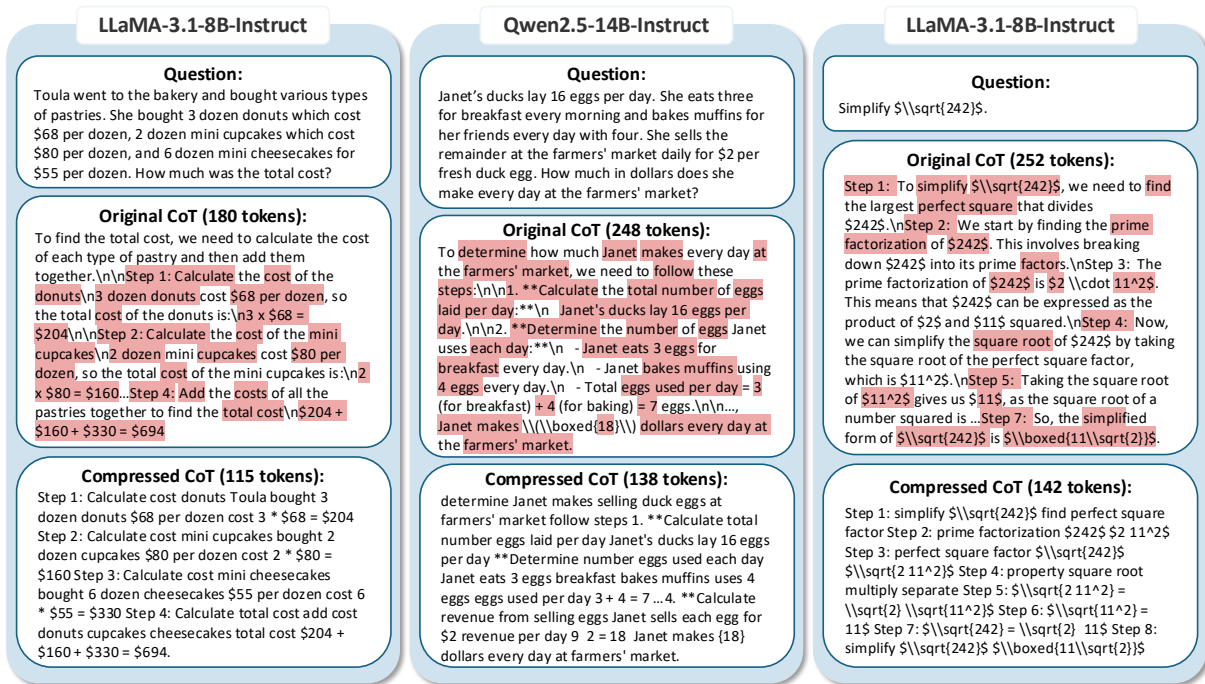


Figure 9: Three CoT compression examples from TokenSkip. For each sample, we list the question, original CoT outputs from corresponding LLMs, and the compressed CoT by TokenSkip. The tokens that appear in both the original CoT and the compressed CoT are highlighted in red.

500 guage information that has minimal impact on reason- 527  
501 ing, Hu et al. (2023) implements structured syntax 528  
502 and symbols, while Han et al. (2024) guides 529  
503 token consumption through dynamic token budget 530  
504 estimation. Similarly, Kang et al. (2024) prompts 531  
505 larger LLMs (i.e., GPT-4) to directly compress CoT, 532  
506 then fine-tunes LLMs to reason using these com- 533  
507 pressed CoTs. In contrast, this work focuses on 534  
508 pruning CoT tokens based on their semantic impor- 535  
509 tance. Additionally, TokenSkip leverages a small 536  
510 LM for token pruning, significantly reducing com- 537  
511 putational overhead. 538

512 **Prompt Compression** As LLMs advance in their 539  
513 zero-shot capabilities, the growing demand for 540  
514 complex instructions and long-context prompts 541  
515 has led to substantial computational and memory 542  
516 challenges in processing lengthy inputs. To ad- 543  
517 dress this bottleneck, researchers have explored 544  
518 various prompt compression techniques. One intu- 545  
519 itive approach involves using a lightweight LM 546  
520 to generate more concise, semantically similar 547  
521 prompts (Chuang et al., 2024). However, given that 548  
522 explicit natural language representations often con- 549  
523 tain redundant information, some researchers have 550  
524 turned to implicit continuous tokens to represent 551  
525 task prompts (Wingate et al., 2022; Mu et al., 2024) 552  
526 and long-context inputs (Chevalier et al., 2023; Ge

509 et al., 2024; Mohtashami and Jaggi, 2023). Other 527  
510 approaches focus on directly compressing input 528  
511 prompts by filtering and retaining high-informative 529  
512 tokens (Li et al., 2023; Jiang et al., 2023; Pan 530  
513 et al., 2024). For instance, Selective Context uses 531  
514 the perplexity of LLMs to measure token impor- 532  
515 tance and removes tokens deemed less important. 533  
516 LLMingua-2 (Pan et al., 2024) introduces a small 534  
517 bidirectional language model for token importance 535  
518 measurement and trains this LM with GPT-4 com- 536  
519 pression data, which serves as the token importance 537  
520 metric in this work. 538

## 539 6 Conclusion 540

541 This work introduces TokenSkip, a simple yet ef- 542  
543 fective approach for controllable Chain-of-Thought 543  
544 (CoT) compression. TokenSkip is built upon the 544  
545 semantic importance of CoT tokens — By selec- 545  
546 tively skipping less important tokens while pre- 546  
547 serving critical ones, TokenSkip enables LLMs 547  
548 to generate compressed CoTs with adjustable ra- 548  
549 tios, thereby striking an expected balance between 549  
550 reasoning efficiency and accuracy. Extensive ex- 550  
551 periments across various LLMs and tasks validate 551  
552 the effectiveness of TokenSkip. We hope our in- 552  
553 vestigations in *token skipping* will offer valuable 553  
554 insights for advancing efficient CoT research and 554  
555 inspire future studies in this area. 555



## 554 Limitations

555 Due to computational constraints, experiments with  
556 larger LLMs, such as Qwen2.5-32B-Instruct and  
557 Qwen2.5-72B-Instruct, were not conducted. We  
558 believe that TokenSkip could achieve a more fa-  
559 vorable trade-off between reasoning performance  
560 and CoT token usage on these models. Addition-  
561 ally, the token importance measurement used in  
562 our study, derived from the LLMLingua-2 com-  
563 pressor (Pan et al., 2024), was not specifically  
564 trained on mathematical data. This limitation may  
565 affect the compression effectiveness, as the model  
566 is not optimized for handling numerical tokens and  
567 mathematical expressions. Furthermore, experi-  
568 ments with long-CoT LLMs, such as QwQ-32B-  
569 Preview, were also excluded due to computational  
570 constraints. We plan to explore these aspects in  
571 future work, as we anticipate that TokenSkip’s  
572 potential can be further realized in these contexts.

## 573 Ethics Statement

574 The datasets used in our experiment are publicly  
575 released and labeled through interaction with hu-  
576 mans in English. In this process, user privacy is  
577 protected, and no personal information is contained  
578 in the dataset. The scientific artifacts that we used  
579 are available for research with permissive licenses.  
580 And the use of these artifacts in this paper is consis-  
581 tent with their intended use. Therefore, we believe  
582 that our research work meets the ethics of ACL.

## 583 References

584 Liang Chen, Haozhe Zhao, Tianyu Liu, Shuai Bai, Jun-  
585 yang Lin, Chang Zhou, and Baobao Chang. 2024.  
586 [An image is worth 1/2 tokens after layer 2: Plug-and-  
587 play inference acceleration for large vision-language  
588 models.](#) In [Computer Vision - ECCV 2024 -  
589 18th European Conference, Milan, Italy, September  
590 29-October 4, 2024, Proceedings, Part LXXXI, vol-  
591 ume 15139 of Lecture Notes in Computer Science,](#)  
592 pages 19–35. Springer.

593 Jeffrey Cheng and Benjamin Van Durme. 2024. Com-  
594 pressed chain of thought: Efficient reasoning  
595 through dense representations. [arXiv preprint  
596 arXiv:2412.13171.](#)

597 Alexis Chevalier, Alexander Wettig, Anirudh Ajith,  
598 and Danqi Chen. 2023. [Adapting language mod-  
599 els to compress contexts.](#) In [Proceedings of the  
600 2023 Conference on Empirical Methods in Natural  
601 Language Processing,](#) pages 3829–3846, Singapore.  
602 Association for Computational Linguistics.

Yu-Neng Chuang, Tianwei Xing, Chia-Yuan Chang,  
Zirui Liu, Xun Chen, and Xia Hu. 2024. [Learn-  
ing to compress prompt in natural language for-  
mats.](#) In [Proceedings of the 2024 Conference of  
the North American Chapter of the Association  
for Computational Linguistics: Human Language  
Technologies \(Volume 1: Long Papers\),](#) pages 7756–  
7767, Mexico City, Mexico. Association for Compu-  
tational Linguistics. 603  
604  
605  
606  
607  
608  
609  
610  
611

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,  
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias  
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro  
Nakano, Christopher Hesse, and John Schulman.  
2021. [Training verifiers to solve math word prob-  
lems.](#) [CoRR](#), abs/2110.14168. 612  
613  
614  
615  
616  
617

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang,  
Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,  
Shirong Ma, Peiyi Wang, et al. 2025. [Deepseek-  
rl: Incentivizing reasoning capability in llms via  
reinforcement learning.](#) [Preprint](#), arXiv:2501.12948. 618  
619  
620  
621  
622

Yuntian Deng, Yejin Choi, and Stuart Shieber. 2024.  
From explicit cot to implicit cot: Learning to  
internalize cot step by step. [arXiv preprint  
arXiv:2405.14838.](#) 623  
624  
625  
626

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
Kristina Toutanova. 2019. [BERT: Pre-training of  
deep bidirectional transformers for language under-  
standing.](#) In [Proceedings of the 2019 Conference  
of the North American Chapter of the Association  
for Computational Linguistics: Human Language  
Technologies, Volume 1 \(Long and Short Papers\),](#)  
pages 4171–4186, Minneapolis, Minnesota. Asso-  
ciation for Computational Linguistics. 627  
628  
629  
630  
631  
632  
633  
634  
635

Mengru Ding, Hanmeng Liu, Zhizhang Fu, Jian Song,  
Wenbo Xie, and Yue Zhang. 2024. [Break the chain:  
Large language models can be shortcut reasoners.](#)  
[CoRR](#), abs/2406.06580. 636  
637  
638  
639

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,  
Abhishek Kadian, Ahmad Al-Dahle, et al. 2024. [The  
llama 3 herd of models.](#) [Preprint](#), arXiv:2407.21783. 640  
641  
642

Tao Ge, Hu Jing, Lei Wang, Xun Wang, Si-Qing Chen,  
and Furu Wei. 2024. [In-context autoencoder for  
context compression in a large language model.](#) In  
[The Twelfth International Conference on Learning  
Representations.](#) 643  
644  
645  
646  
647

Sachin Goyal, Ziwei Ji, Ankit Singh Rawat, Aditya Kr-  
ishna Menon, Sanjiv Kumar, and Vaishnavh Na-  
garajan. 2024. [Think before you speak: Train-  
ing language models with pause tokens.](#) In  
[The Twelfth International Conference on Learning  
Representations.](#) 648  
649  
650  
651  
652  
653

Tingxu Han, Chunrong Fang, Shiyu Zhao, Shiqing  
Ma, Zhenyu Chen, and Zhenting Wang. 2024.  
[Token-budget-aware llm reasoning.](#) [arXiv preprint  
arXiv:2412.18547.](#) 654  
655  
656  
657

658	Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li,	Yucheng Li, Bo Dong, Frank Guerin, and Chenghua	713
659	Zhiting Hu, Jason Weston, and Yuandong Tian. 2024.	Lin. 2023. <a href="#">Compressing context to enhance inference efficiency of large language models</a> . In	714
660	Training large language models to reason in a continuous latent space. <a href="#">arXiv preprint arXiv:2412.06769</a> .	<a href="#">Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</a> , pages	715
661		6342–6353, Singapore. Association for Computational Linguistics.	716
662	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul	Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison	720
663	Arora, Steven Basart, Eric Tang, Dawn Song,	Edwards, Bowen Baker, Teddy Lee, Jan	721
664	and Jacob Steinhardt. 2021. <a href="#">Measuring mathematical problem solving with the MATH dataset</a> .	Leike, John Schulman, Ilya Sutskever, and Karl	722
665	In <a href="#">Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track</a>	Cobbe. 2024. <a href="#">Let’s verify step by step</a> . In	723
666	<a href="#">(Round 2)</a> .	<a href="#">The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May</a>	724
667		<a href="#">7-11, 2024</a> . OpenReview.net.	725
668			726
669	Le Hou, Richard Yuanzhe Pang, Tianyi Zhou, Yuexin	Zhenghao Lin, Zhibin Gou, Yeyun Gong, Xiao Liu,	727
670	Wu, Xinying Song, Xiaodan Song, and Denny Zhou.	yelong shen, Ruochen Xu, Chen Lin, Yujiu Yang,	728
671	2022. <a href="#">Token dropping for efficient BERT pretraining</a> .	Jian Jiao, Nan Duan, and Weizhu Chen. 2024. <a href="#">Not all tokens are what you need for pretraining</a> . In	729
672	In <a href="#">Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume</a>	<a href="#">The Thirty-eighth Annual Conference on Neural Information Processing Systems</a> .	730
673	<a href="#">1: Long Papers</a> ), pages 3774–3784, Dublin, Ireland.		731
674	Association for Computational Linguistics.		732
675			
676	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan	Tengxiao Liu, Qipeng Guo, Xiangkun Hu, Cheng Ji-	733
677	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and	ayang, Yue Zhang, Xipeng Qiu, and Zheng Zhang.	734
678	Weizhu Chen. 2022. <a href="#">Lora: Low-rank adaptation of large language models</a> . In <a href="#">The Tenth International</a>	2024. <a href="#">Can language models learn to skip steps?</a>	735
679	<a href="#">Conference on Learning Representations, ICLR</a>	In <a href="#">The Thirty-eighth Annual Conference on Neural</a>	736
680	<a href="#">2022, Virtual Event, April 25-29, 2022</a> . OpenRe-	<a href="#">Information Processing Systems</a> .	737
681	<a href="#">view.net</a> .		
682			
683	Hanxu Hu, Hongyuan Lu, Huajian Zhang, Yun-Ze Song,	Ilya Loshchilov and Frank Hutter. 2019. <a href="#">Decoupled</a>	738
684	Wai Lam, and Yue Zhang. 2023. Chain-of-symbol	<a href="#">weight decay regularization</a> . In <a href="#">7th International</a>	739
685	prompting elicits planning in large language models.	<a href="#">Conference on Learning Representations, ICLR</a>	740
686	<a href="#">arXiv preprint arXiv:2305.10276</a> .	<a href="#">2019, New Orleans, LA, USA, May 6-9, 2019</a> . Open-	741
		<a href="#">Review.net</a> .	742
687	Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing	Yiran Ma, Zui Chen, Tianqiao Liu, Mi Tian, Zhuo Liu,	743
688	Yang, and Lili Qiu. 2023. <a href="#">LLMLingua: Compressing</a>	Zitao Liu, and Weiqi Luo. 2024. <a href="#">What are step-level</a>	744
689	<a href="#">prompts for accelerated inference of large language</a>	<a href="#">reward models rewarding? counterintuitive findings</a>	745
690	<a href="#">models</a> . In <a href="#">Proceedings of the 2023 Conference on</a>	<a href="#">from mcts-boosted mathematical reasoning</a> . <a href="#">CoRR</a> ,	746
691	<a href="#">Empirical Methods in Natural Language Processing</a> ,	<a href="#">abs/2412.15904</a> .	747
692	pages 13358–13376, Singapore. Association for		
693	Computational Linguistics.	Emanuele Marconato, Stefano Teso, Antonio Vergari,	748
694	Mingyu Jin, Qinkai Yu, Dong Shu, Haiyan Zhao,	and Andrea Passerini. 2024. Not all neuro-symbolic	749
695	Wenyue Hua, Yanda Meng, Yongfeng Zhang, and	concepts are created equal: Analysis and mitiga-	750
696	Mengnan Du. 2024. <a href="#">The impact of reasoning step</a>	tion of reasoning shortcuts. <a href="#">Advances in Neural</a>	751
697	<a href="#">length on large language models</a> . In <a href="#">Findings of</a>	<a href="#">Information Processing Systems</a> , 36.	752
698	<a href="#">the Association for Computational Linguistics: ACL</a>	William Merrill and Ashish Sabharwal. 2024. <a href="#">The</a>	753
699	<a href="#">2024</a> , pages 1830–1842, Bangkok, Thailand. Associ-	<a href="#">expressive power of transformers with chain of</a>	754
700	ation for Computational Linguistics.	<a href="#">thought</a> . In <a href="#">The Twelfth International Conference</a>	755
701	Yu Kang, Xianghui Sun, Liangyu Chen, and Wei Zou.	<a href="#">on Learning Representations, ICLR 2024, Vienna,</a>	756
702	2024. C3ot: Generating shorter chain-of-thought	<a href="#">Austria, May 7-11, 2024</a> . OpenReview.net.	757
703	without compromising effectiveness. <a href="#">arXiv preprint</a>	Amirkeivan Mohtashami and Martin Jaggi. 2023.	758
704	<a href="#">arXiv:2412.11664</a> .	<a href="#">Random-access infinite context length for trans-</a>	759
		<a href="#">formers</a> . In <a href="#">Thirty-seventh Conference on Neural</a>	760
		<a href="#">Information Processing Systems</a> .	761
705	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yu-	Jesse Mu, Xiang Li, and Noah Goodman. 2024. Learn-	762
706	taka Matsuo, and Yusuke Iwasawa. 2022. <a href="#">Large</a>	ing to compress prompts with gist tokens. <a href="#">Advances</a>	763
707	<a href="#">language models are zero-shot reasoners</a> . In	<a href="#">in Neural Information Processing Systems</a> , 36.	764
708	<a href="#">Advances in Neural Information Processing Systems</a>	Xuefei Ning, Zinan Lin, Zixuan Zhou, Zifu Wang,	765
709	<a href="#">35: Annual Conference on Neural Information</a>	Huazhong Yang, and Yu Wang. 2023. Skeleton-of-	766
710	<a href="#">Processing Systems 2022, NeurIPS 2022, New</a>	<a href="#">thought: Large language models can do parallel de-</a>	767
711	<a href="#">Orleans, LA, USA, November 28 - December 9,</a>	<a href="#">coding</a> . <a href="#">Proceedings ENLSP-III</a> .	768
712	<a href="#">2022</a> .		

769	Maxwell I. Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, Charles Sutton, and Augustus Odena. 2021. <a href="#">Show your work: Scratchpads for intermediate computation with language models</a> . <a href="#">CoRR</a> , abs/2112.00114.	
776	OpenAI. 2023. <a href="#">GPT-4 technical report</a> . <a href="#">CoRR</a> , abs/2303.08774.	
778	OpenAI. 2024. <a href="#">Learning to reason with llms</a> .	
779	OpenAI et al. 2024. <a href="#">Openai o1 system card</a> . <a href="#">Preprint</a> , arXiv:2412.16720.	
781	Zhuoshi Pan, Qianhui Wu, Huiqiang Jiang, Menglin Xia, Xufang Luo, Jue Zhang, Qingwei Lin, Victor Rühle, Yuqing Yang, Chin-Yew Lin, H. Vicky Zhao, Lili Qiu, and Dongmei Zhang. 2024. <a href="#">LLMLingua-2: Data distillation for efficient and faithful task-agnostic prompt compression</a> . In <a href="#">Findings of the Association for Computational Linguistics: ACL 2024</a> , pages 963–981, Bangkok, Thailand. Association for Computational Linguistics.	
790	Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. <a href="#">Reflexion: language agents with verbal reinforcement learning</a> . In <a href="#">Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</a> .	
798	Charlie Victor Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2025. <a href="#">Scaling test-time compute optimally can be more effective than scaling LLM parameters</a> . In <a href="#">The Thirteenth International Conference on Learning Representations</a> .	
803	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. <a href="#">Self-consistency improves chain of thought reasoning in language models</a> . In <a href="#">The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</a> . <a href="#">OpenReview.net</a> .	
810	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. <a href="#">Chain-of-thought prompting elicits reasoning in large language models</a> . In <a href="#">Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022</a> .	
819	David Wingate, Mohammad Shoeybi, and Taylor Sorensen. 2022. <a href="#">Prompt compression and contrastive conditioning for controllability and toxicity reduction in language models</a> . In <a href="#">Findings of the Association for Computational Linguistics: EMNLP 2022</a> , pages 5621–5634, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	
	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. 2024. <a href="#">Qwen2 technical report</a> . <a href="#">CoRR</a> , abs/2407.10671.	826 827 828 829 830 831 832 833 834 835 836 837 838 839 840 841 842
	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. <a href="#">Tree of thoughts: Deliberate problem solving with large language models</a> . In <a href="#">Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</a> .	843 844 845 846 847 848 849 850
	Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song, Yuan-dong Tian, Christopher Re, Clark Barrett, Zhangyang Wang, and Beidi Chen. 2023. <a href="#">H2o: Heavy-hitter oracle for efficient generative inference of large language models</a> . In <a href="#">Thirty-seventh Conference on Neural Information Processing Systems</a> .	851 852 853 854 855 856 857
	Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, and Zheyang Luo. 2024. <a href="#">LlamaFactory: Unified efficient fine-tuning of 100+ language models</a> . In <a href="#">Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)</a> , pages 400–410, Bangkok, Thailand. Association for Computational Linguistics.	858 859 860 861 862 863 864 865
	Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V. Le, and Ed H. Chi. 2023. <a href="#">Least-to-most prompting enables complex reasoning in large language models</a> . In <a href="#">The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</a> . <a href="#">OpenReview.net</a> .	866 867 868 869 870 871 872 873



## Appendix

### A CoT Recovery

In this section, we provide the detailed prompt for our recovery experiments, which is illustrated in Figure 10. Besides, we present the CoT recovery result from GPT-4o (OpenAI, 2023) in Figure 11. The recovered results demonstrate that GPT-4o could understand the compressed CoT content and correctly restore the original CoT process.

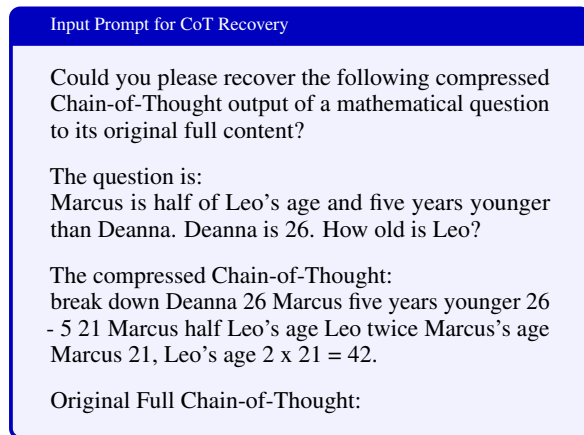


Figure 10: Input prompt for LLaMA-3.1-8B-Instruct designed to recover the compressed CoT from a GSM8K math problem.

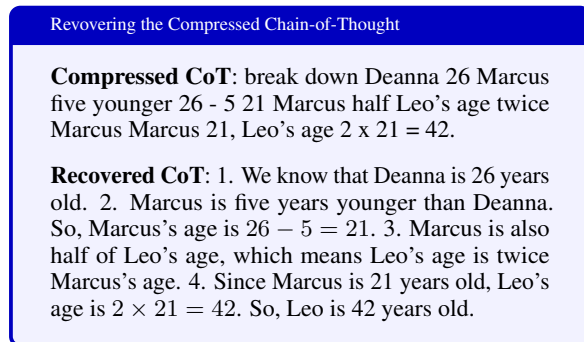


Figure 11: Recovering the compressed CoT for GSM8K math word problem using GPT-4o.

## B Experimental Details

### B.1 Implementation Details

We utilize LLMingua-2 (Pan et al., 2024) as the token importance metric to generate our compressed CoT training data. The compression ratio  $\gamma$  is randomly selected from  $\{0.5, 0.6, 0.7, 0.8, 0.9, 1.0\}$  for each training sample. We adopt LoRA (Hu et al., 2022) to train our models. The rank  $r$  is set to 8, and the scaling parameter  $\alpha$  is set to 16. We train the models for 3 epochs on both datasets. The

peak learning rate is set to  $5e-5$ , following a cosine decay schedule. We use AdamW (Loshchilov and Hutter, 2019) for optimization, with a warmup ratio of 0.1. We implement our training process using the LLaMA-Factory (Zheng et al., 2024) library. Inference for both our method and all baselines is performed using the Huggingface transformers package. During inference, the maximum number of tokens `max_len` is set to 512 for GSM8K and 1024 for MATH. All experiments are conducted using Pytorch 2.1.0 on  $2 \times$  NVIDIA GeForce RTX 3090 GPU (24GB) with CUDA 12.1, and an Intel(R) Xeon(R) Platinum 8370C CPU with 32 cores.

### B.2 Detailed Results with Qwen

We provide detailed experimental results of the Qwen2.5-Instruct series evaluated on GSM8K in Table 2. As the model scale increases, there is less performance degradation at higher compression ratios, indicating that larger LLMs are better at identifying shortcuts between critical reasoning tokens, enabling more efficient CoT generation.

Scale	Methods	Ratio	Accuracy	Tokens	ActRatio
3B	Original	-	83.7(0.0↓)	314.87	-
	TokenSkip	1.0	83.4(0.3↓)	318.79	1.00
		0.9	83.2(0.5↓)	262.99	0.83
		0.8	81.6(2.1↓)	250.71	0.79
		0.7	80.1(3.6↓)	233.03	0.73
		0.6	77.3(6.4↓)	199.55	0.63
		0.5	74.4(9.3↓)	170.55	0.54
7B	Original	-	91.4(0.0↓)	297.83	-
	TokenSkip	1.0	91.7(0.3↑)	295.78	1.00
		0.9	91.1(0.3↓)	254.77	0.86
		0.8	90.1(1.3↓)	237.27	0.80
		0.7	89.9(1.5↓)	216.73	0.73
		0.6	87.9(3.5↓)	178.07	0.60
		0.5	86.0(5.4↓)	151.44	0.51
14B	Original	-	93.1(0.0↓)	313.11	-
	TokenSkip	1.0	93.0(0.1↓)	314.55	1.00
		0.9	93.3(0.2↑)	269.22	0.86
		0.8	93.2(0.1↑)	247.24	0.79
		0.7	93.4(0.3↑)	218.62	0.70
		0.6	92.7(0.4↓)	180.68	0.57
		0.5	91.4(1.7↓)	156.85	0.50

Table 2: Experimental results on the Qwen2.5-Instruct series. We report accuracy, average CoT token count, and actual compression ratio (ActRatio) for comparison.