

SoftCoT: Soft Chain-of-Thought for Efficient Reasoning with LLMs

Yige Xu^{1,2,*}, Xu Guo^{1,*†}, Zhiwei Zeng^{1†}, Chunyan Miao^{1,2}

¹Joint NTU-UBC Research Centre of Excellence in Active Living for the Elderly

²College of Computing and Data Science

Nanyang Technological University, Singapore

{yige002,xu008}@e.ntu.edu.sg, {zhiwei.zeng,ascymiao}@ntu.edu.sg

Abstract

Chain-of-Thought (CoT) reasoning enables Large Language Models (LLMs) to solve complex reasoning tasks by generating intermediate reasoning steps. However, most existing approaches focus on hard token decoding, which constrains reasoning within the discrete vocabulary space and may not always be optimal. While recent efforts explore continuous-space reasoning, they often require full-model fine-tuning and suffer from catastrophic forgetting, limiting their applicability to state-of-the-art LLMs that already perform well in zero-shot settings with a proper instruction. To address this challenge, we propose a novel approach for continuous-space reasoning that does not require modifying the LLM. Specifically, we employ a lightweight fixed assistant model to speculatively generate instance-specific soft thought tokens as the initial chain of thoughts, which are then mapped into the LLM’s representation space via a trainable projection module. Experimental results on five reasoning benchmarks demonstrate that our method enhances LLM reasoning performance through supervised, parameter-efficient fine-tuning. Source code is available at <https://github.com/xuyige/SoftCoT>.

1 Introduction

In recent years, Large Language Models (LLMs) have become a cornerstone in Natural Language Processing (NLP), exhibiting advanced natural language understanding and generation (Brown et al., 2020; Du et al., 2022; Chowdhery et al., 2023; OpenAI, 2023; Touvron et al., 2023; Dubey et al., 2024; Yang et al., 2024). Scaling model sizes has not only improved instruction-following (Kojima et al., 2022) but also triggered emergent reasoning abilities, as first evidenced by chain-of-thought (CoT) prompting (Wei et al., 2022). CoT

prompts LLMs to generate intermediate reasoning steps before providing the final answer, which not only enhances interpretability but also improves a range of reasoning-intensive tasks (Zhang et al., 2023; Sprague et al., 2024). It has inspired many advanced prompting frameworks, marking a paradigm shift from scaling training-time compute (Kojima et al., 2022) to scaling inference-time compute (Wang et al., 2023; Yao et al., 2023) to further boost LLM performance.

Nevertheless, CoT’s effectiveness depends on the quality of intermediate thoughts, as the autoregressive generation process can propagate errors. To mitigate this challenge, methods like self-consistency (Wang et al., 2023) generate multiple reasoning paths, while Tree-of-Thought (Yao et al., 2023) and Graph-of-Thought (Besta et al., 2024) frameworks organize these paths to select higher-quality steps. Despite these improvements, such methods are computationally inefficient due to the need for extensive thought sampling.

To enhance CoT efficiency, recent research explores skipping the decoding of hard tokens at intermediate steps. Methods like Compressed CoT (Cheng and Durme, 2024) and Coconut (Hao et al., 2024) conduct reasoning in a continuous space by using latent representations instead of discrete token sequences. Their results have shown to be superior to long-sequence discrete reasoning chains using only a short length of continuous representation. Yet, these methods require full-model fine-tuning, which incurs substantial computational costs, risks catastrophic forgetting, and limits their transferability across tasks.

We empirically observed that fine-tuning LLaMA3.1-8B (Dubey et al., 2024) for continuous-space reasoning using a language modeling objective (as employed by Coconut and CCoT) results in performance degradation compared to zero-shot CoT (Tables 2 and 3). Drawing on a widely accepted definition of catastrophic forgetting (Kala-

*The first two authors contributed equally.

†Corresponding authors.

jdzievski, 2024; Lobo et al., 2024), defined as *the degradation of previously learned capabilities after fine-tuning on new data*, we conjecture that this drop in reasoning performance is attributable to **catastrophic forgetting**. This phenomenon appears particularly pronounced in already capable instruction-tuned models such as LLaMA-3.1-8B-Instruct and Qwen2.5-7B-Instruct, which exhibit strong zero-shot CoT reasoning abilities. Thus, the methodologies of Coconut, which is based on GPT-2 (Radford et al., 2019), may not be directly applicable to more recent models such as LLaMA3.1 and Qwen2.5 series. Therefore, it is crucial to explore alternative methodologies that mitigate catastrophic forgetting while effectively leveraging continuous reasoning techniques in large-scale, instruction-tuned models, which is the main research goal of this work. To the best of our knowledge, we are the first to systematically identify and address the forgetting issue.

To mitigate catastrophic forgetting, a straightforward approach is to freeze the backbone LLM and instead optimize an external model for reasoning. Inspired by prompt tuning (Lester et al., 2021) and speculative decoding (Leviathan et al., 2023), we propose to utilize an auxiliary small assistant model to generate a sequence of “thought” tokens conditioned on a task instruction followed by a specific instance (Li et al., 2023; Shao et al., 2023). These tokens serve as instance-specific prompts that adapt to different problems to boost LLM’s reasoning. Such an auxiliary prompting mechanism allows the LLM to achieve better generalization while preserving its pre-trained knowledge.

To exploit continuous-space reasoning, we use the last-layer hidden states from the small assistant model as the “soft” thought tokens, rather than the discrete tokens obtained after vocabulary mapping. Staying in the latent space avoids information loss inherent in autoregressive decoding. However, a representational gap between the assistant model and the LLM may hinder effective knowledge transfer. To bridge this gap, we train a projection module to map the soft thought tokens generated by the assistant model to the LLM’s representation space. Training the projection module for each task can be seen as *soft prompt tuning* for the LLM. The overall **Soft** thoughts for **CoT** (SoftCoT) reasoning framework is illustrated in Figure 1.

We evaluate SoftCoT on five reasoning benchmarks and two state-of-the-art LLM architectures.

The five benchmarks include mathematical reasoning, commonsense reasoning, and symbolic reasoning. For further exploration, we create a hard version of the ASDiv dataset (Miao et al., 2020), which requires stronger mathematical reasoning ability. The new dataset is named “ASDiv-Aug” in this paper. Experimental results show that SoftCoT consistently improves accuracy on both public datasets and our augmented ASDiv-Aug dataset, demonstrating the effectiveness of SoftCoT in enhancing LLM’s reasoning performance. Moreover, SoftCoT effectively mitigates catastrophic forgetting seen in previous methods based on full-model fine-tuning, and our results highlight the effectiveness of using an assistant model to generate soft thoughts. SoftCoT thus represents the first lightweight yet powerful approach that leverages the benefits of soft thought tokens while preserving LLM’s prior knowledge.

2 Related Works

Early research on chain-of-thought (CoT) reasoning can be traced back to Wei et al. (2022), who first introduced a prompting strategy that guides LLMs through decomposed intermediate reasoning steps using few-shot exemplars. Concurrently, Kojima et al. (2022) demonstrated that LLMs are capable of zero-shot CoT reasoning by simply appending the phrase “Let’s think step by step” to the prompt template. This discovery underscored the latent reasoning abilities of LLMs, even in the absence of explicit demonstrations.

Building upon these foundational works, the NLP community has extensively explored the potential of CoT reasoning. As summarized by Chu et al. (2024), recent advancements in CoT methods can be broadly categorized into three areas: (1) *Prompt Construction*, which aims to optimize prompts for improved CoT reasoning (Wei et al., 2022; Kojima et al., 2022; Zhang et al., 2023); (2) *Topological Variants*, which leverage structured representations such as trees and graphs to enhance CoT reasoning (Yao et al., 2023; Besta et al., 2024); and (3) *Enhancement Methods*, which introduce external strategies to further improve CoT reasoning, such as question decomposition (Zhou et al., 2023) and self-consistency decoding (Wang et al., 2023). Despite the effectiveness of these approaches, the majority of existing CoT methods rely on discrete token-by-token generation, which imposes inherent constraints and limits their expressiveness.

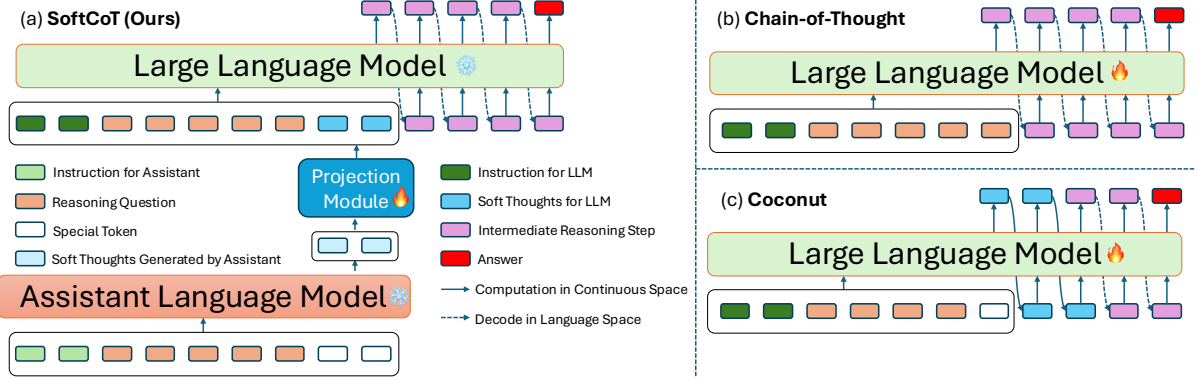


Figure 1: A comparison of SoftCoT, vanilla Chain-of-Thought, and Coconut.

To address the limitations of discrete language space, an effective approach is to leverage continuous representation space for reasoning. Coconut (Hao et al., 2024) introduces a Chain-of-Continuous-Thought, while CCoT (Cheng and Durme, 2024) employs Compressed Chain-of-Thought, generating content-rich and continuous contemplation tokens. Heima (Shen et al., 2025) further advances this idea by utilizing a single continuous vector to represent compressed reasoning tokens in multi-modal tasks. However, both Coconut and CCoT rely on a language modeling objective for supervised fine-tuning, which is infeasible for state-of-the-art LLMs due to the catastrophic forgetting problem. Moreover, Heima underperforms compared to its backbone model, LLaVA-CoT (Xu et al., 2024). These challenges underscore the need to develop methodologies that mitigate catastrophic forgetting in the application of continuous-space CoT reasoning.

3 Methodology

3.1 Problem Definition and Notations

Given a question $\mathcal{Q} = [q_1, q_2, \dots, q_{|\mathcal{Q}|}]$, the CoT reasoning will solve the problem on the following two steps: (1) Auto-regressively generate a list of rationale steps $\mathcal{R} = [r_1, r_2, \dots, r_{|\mathcal{R}|}]$ according to the question; (2) Auto-regressively generate the answer $\mathcal{A} = [a_1, a_2, \dots, a_{|\mathcal{A}|}]$ according to the question as well as the rationale steps. The generation process can be described as:

$$\begin{aligned} r_{i+1} &= \text{LLM}(\mathcal{Q}; \mathcal{R}_{\leq i}), \\ a_{j+1} &= \text{LLM}(\mathcal{Q}; \mathcal{R}; \mathcal{A}_{\leq j}), \end{aligned} \quad (1)$$

where $\text{LLM}(\cdot)$ indicates a large language model, $\mathcal{R}_{\leq i} = [r_1, \dots, r_i]$ indicates the previous gener-

ated i reasoning tokens, and $\mathcal{A}_{\leq j} = [a_1, \dots, a_j]$ indicates the previous generated j answer tokens.

The majority of recent works (Zhang et al., 2023; Zhou et al., 2023; Yao et al., 2023) focus on generating discrete hard tokens in \mathcal{R} , which is named as “**Hard-CoT**” in this paper. On contrast, some recent works (Hao et al., 2024; Cheng and Durme, 2024) focus on the continuous representations (*a.k.a* latent space) of \mathcal{R} , which is named as “**Soft-CoT**” in this paper.

In this paper, we manually define some rules (e.g., regular expression matching) to extract the framework answer $\hat{\mathcal{A}}$ from the answer $\bar{\mathcal{A}}$ generated by the LLM. Then we compute the accuracy of $\hat{\mathcal{A}}$ comparing with the ground-truth answer \mathcal{A} .

3.2 Overview of the SoftCoT Framework

SoftCoT is a novel framework designed to enhance reasoning capabilities in large language models (LLMs). Given an input question \mathcal{Q} , the framework produces both a sequence of reasoning steps \mathcal{R} and the final answer \mathcal{A} . SoftCoT consists of three key components: the soft thought token generation module, the projection module, and the CoT reasoning module. The overall architecture is illustrated in Figure 1(a).

The soft thought token generation module is inspired by prompt tuning techniques (Lester et al., 2021). In conventional prompt tuning, learnable prompts facilitate the retrieval of knowledge stored within the LLM (Xu et al., 2023). In SoftCoT, soft thought tokens are generated by an assistant language model, which is typically smaller than the backbone LLM (e.g., LLaMA-3.2-1B-Instruct as the assistant model and LLaMA-3.1-8B-Instruct as the backbone model).

A key challenge in this setup is that the assistant model can only generate discrete token sequences

as input to the backbone LLM, which imposes constraints and may not always yield optimal prompts. To address this limitation, we introduce continuous soft thought tokens that enable more expressive and flexible prompting. However, a representation gap exists between the assistant model and the backbone LLM, necessitating an intermediate transformation.

To bridge this gap, the projection module maps the soft thought tokens’ representations into a space more compatible with the backbone LLM. This ensures that the soft thought tokens effectively guide the reasoning process.

Finally, the CoT reasoning module leverages both the learned soft thought tokens and word embeddings to generate intermediate reasoning steps $\bar{\mathcal{R}}$ and the final answer $\bar{\mathcal{A}}$. The model is trained using a language modeling objective, optimizing the learnable parameters across the rationale steps and the answer spans.

3.3 Prompt Tuning for CoT Reasoning

Prompt tuning for CoT reasoning aims to optimize the structure and content of the prompt template to enhance the reasoning capabilities of a large language model (LLM). This process can be mathematically formulated as follows:

$$\begin{aligned} \hat{y} &= \text{LLM}(P_{\mathbf{p}}(x)), \\ \mathbf{p}^* &= \arg \min_{\mathbf{p}} \mathcal{L}(\hat{y}, y), \end{aligned} \quad (2)$$

where \hat{y} represents the predicted output, x denotes the input sequence, and $P_{\mathbf{p}}(x)$ is the input augmented with a prompt \mathbf{p} . The objective function $\mathcal{L}(\cdot)$ measures the discrepancy between the model’s prediction \hat{y} and the ground-truth label y . The primary goal of prompt tuning is to determine an optimal prompt configuration that effectively guides the LLM to perform CoT reasoning with improved accuracy and interpretability.

A straightforward yet effective approach to optimizing prompts involves leveraging an auxiliary assistant model to generate instance-specific prompts, which provide contextual hints or question summaries to facilitate reasoning (Li et al., 2023; Shao et al., 2023; Li et al., 2024). In this framework, the prompt \mathbf{p} can be decomposed into two components: (1) a fixed, task-specific prompt \mathbf{p}_{task} , which remains constant across all instances and encodes general problem-solving heuristics, and (2) a learnable, instance-specific prompt \mathbf{p}_{inst} , which dynam-

ically adapts to each input instance to provide tailored guidance.

Given the rapid advancements in LLMs, many LLMs are capable of solving complex reasoning tasks under zero-shot settings. Instead of fine-tuning the assistant model for each task, we adopt a more efficient strategy by employing a relatively small, frozen language model to generate \mathbf{p}_{inst} . This approach not only reduces computational costs but also ensures stability and generalization across different problem domains. By systematically integrating instance-specific prompting with fixed task-specific instructions, this method enhances the LLM’s reasoning process while maintaining adaptability to various contexts.

3.4 Soft Thought Tokens for CoT Reasoning

One of the advantages of Hard-CoT is that the generated discrete tokens can be tokenized by the LLMs, which does not require an external mapping module. However, there are two main limitations of Hard-CoT: (1) The decoded token space is discrete, which is constrained and sometimes not optimal; (2) The gradient cannot backpropagate to the assistant model since the decoding process cut off the gradient information. A convincing solution is to replace the hard tokens with soft thought tokens.

Generating Soft Thought Tokens with an Assistant Model

To generate instance-specific soft thoughts, we utilize an auxiliary assistant model that produces soft thoughts based on the given reasoning task. The input to the assistant model, denoted as $\mathbf{x}_{\text{assist}}$, consists of three main components:

$$\mathbf{x}_{\text{assist}} = \text{concat}[\mathcal{I}_{\text{assist}}, \mathcal{Q}, [\text{UNK}]_{1:N}], \quad (3)$$

where

- $\mathcal{I}_{\text{assist}}$ represents the instructional context provided to the assistant model, guiding it on how to generate relevant thoughts.
- \mathcal{Q} denotes the reasoning question that the primary LLM will solve, which has been defined in Section 3.1.
- N [UNK] tokens serve as placeholders for the soft thoughts.

Once the input sequence is constructed, the assistant model processes it, and the soft thought tokens

are obtained as follows:

$$\begin{aligned} \mathbf{h}^{\text{assist}} &= \text{Assistant}(\mathbf{x}_{\text{assist}}), \\ \mathbf{t}_{\text{assist}} &= \mathbf{h}^{\text{assist}}_{|I|+|Q|+1:|I|+|Q|+N}. \end{aligned} \quad (4)$$

Here $\mathbf{h}^{\text{assist}}$ denotes the final-layer hidden states of the assistant model, and $\mathbf{t}_{\text{assist}}$ corresponds to the segment of $\mathbf{h}^{\text{assist}}$ associated with the N [UNK] tokens. This extracted representation serves as the instance-specific soft thoughts, dynamically adapting to the input reasoning question.

Projection Module Since there exist both a representation gap and a dimensional gap between the assistant language model and the primary LLM, a direct utilization of $\mathbf{t}_{\text{assist}}$ may lead to suboptimal performance. The assistant model and the LLM often operate in different embedding spaces, with distinct hidden state distributions and dimensionalities. To bridge this gap, we introduce a projection module that maps the assistant-generated soft thoughts $\mathbf{t}_{\text{assist}}$ from the assistant model’s embedding space to the LLM’s embedding space:

$$\mathcal{T}_{\text{soft}} = \text{Linear}_{\theta}(\mathbf{t}_{\text{assist}}), \quad (5)$$

where $\text{Linear}_{\theta} : \mathbb{R}^{d_{\text{assist}}} \rightarrow \mathbb{R}^{d_{\text{LLM}}}$ is a trainable projection layer parameterized by θ . This layer ensures that the assistant-generated soft thoughts are transformed into a suitable format for the LLM, preserving relevant semantic information while adapting to the LLM’s feature space.

By incorporating this projection module, we effectively mitigate discrepancies between the assistant model and the LLM, enabling smooth integration of instance-specific thoughts into the CoT reasoning process. This design ensures that the learned thought tokens are both informative and compatible, thereby enhancing the overall reasoning performance of the LLM.

LLM Reasoning with Soft CoT With instance-specific soft thought tokens generated by the assistant model and mapped to the LLM’s embedding space, we proceed to the final step: applying these soft thoughts to aid LLMs in CoT reasoning.

The input to the LLM, denoted as \mathbf{x}_{LLM} , follows a structure similar to that of $\mathbf{x}_{\text{assist}}$:

$$\mathbf{x}_{\text{LLM}} = \text{concat}[\mathcal{I}_{\text{LLM}}, \mathcal{Q}, \mathcal{T}_{\text{soft}}], \quad (6)$$

where

- \mathcal{I}_{LLM} is the task-specific instruction, which is a fixed prompt shared across all instances of the same task. It provides general problem-solving heuristics and instructions relevant to the reasoning task.
- $\mathcal{T}_{\text{soft}}$ is the instance-specific soft thoughts computed by Eq (5). This component dynamically adapts soft thought tokens to each input question, enhancing contextual understanding.

With this structured input, the LLM generates step-by-step reasoning chains, following the principles of CoT reasoning. The reasoning process unfolds by systematically applying logical deductions or problem-solving heuristics, ultimately leading to the generation of the final answer:

$$\begin{aligned} \bar{\mathcal{R}} &= \text{LLM}(\mathbf{x}_{\text{LLM}}), \\ \bar{\mathcal{A}} &= \text{LLM}(\mathbf{x}_{\text{LLM}}, \bar{\mathcal{R}}), \\ \hat{\mathcal{A}} &= \mathcal{E}(\bar{\mathcal{A}}), \end{aligned} \quad (7)$$

where $\mathcal{E}(\cdot)$ is manual rules for answer extraction.

By integrating both fixed task-specific instructions and instance-specific soft thought tokens, our approach enables the LLM to systematically decompose complex reasoning tasks while leveraging auxiliary knowledge provided by the assistant model. The structured input ensures that the LLM benefits from both general domain knowledge and tailored instance-level guidance, ultimately improving its reasoning effectiveness.

Parameter-Efficient Training In this work, we focus on reasoning tasks that include annotated reasoning steps, which provide explicit intermediate reasoning trajectories leading to the final answer. To effectively train the model, we employ the standard language modeling objective (also known as next-token prediction) to supervise the generation of soft thoughts. During the training stage, the input sequence is structured as follows:

$$\mathbf{x}_{\text{train}} = \text{concat}[\mathcal{I}_{\text{assist}}, \mathcal{Q}, \mathcal{T}_{\text{soft}}, \mathcal{R}, \mathcal{A}]. \quad (8)$$

To effectively learn the soft thoughts, we apply the negative log-likelihood (NLL) loss over the reasoning steps and the answer span. Specifically, we mask the tokens before the intermediate reasoning steps to prevent the model from directly relying on them during loss computation. Instead, the model is trained to generate the reasoning steps \mathcal{R} and final answer \mathcal{A} in an autoregressive manner.

Dataset	Task Type	Answer Type	# Train samples	# Evaluation samples
GSM8K	Mathematical	Number	7,473	1,319
ASDiv-Aug		Number	4,183	1,038
AQuA		Option	97,467	254
StrategyQA	Commonsense	Yes/No	1,832	458
DU	Symbolic	Option	-	250

Table 1: Summary statistics of five datasets we used. “-” indicates that there is no training samples available. “DU” indicates the Date Understanding (BIG.Bench.authors, 2023) dataset.

4 Experiments

4.1 Datasets

We conduct experiments on five reasoning datasets spanning three categories of reasoning: mathematical reasoning, commonsense reasoning, and symbolic reasoning. For mathematical reasoning, we utilize **GSM8K** (Cobbe et al., 2021), **ASDiv** (Miao et al., 2020), and **AQuA** (Ling et al., 2017). For commonsense reasoning, we employ **StrategyQA** (Geva et al., 2021). For symbolic reasoning, we use **Date Understanding** (BIG.Bench.authors, 2023) from the BIG-benchmark.

Given that LLaMA-3.1-8B-Instruct is a well-trained LLM, we augment the ASDiv dataset to ensure that the model encounters novel instances. Specifically, we replicate each instance five times and systematically extract and replace numerical values in the questions with randomly selected alternatives. For example, the original dataset is “eat 6 apples per day”, we will duplicate this instance multiple times and **replace 6 with 7, or 8, or 12, or 18 randomly**. This augmentation is designed to evaluate the model’s reasoning capability rather than its ability to recognize patterns from memorized data. The augmented dataset is named as “**ASDiv-Aug**” in the following part of this paper. We open-source our ASDiv-Aug to the community at <https://huggingface.co/datasets/xuyige/ASDiv-Aug>. All detail statistics of the datasets is shown in Table 1.

4.2 Baselines

Considering that zero-shot CoT on modern LLMs already sets a strong baseline, often outperforming fine-tuned models due to their rich pretraining, we will consider zero-shot baseline as a fair comparison. Meanwhile, we also consider baselines the fine-tunes the LLM:

Zero-Shot CoT We adopt the prompt templates from Sprague et al. (2024) to test zero-shot CoT per-

formance. This baseline serves to assess whether the model experiences performance degradation after supervised fine-tuning.

Zero-Shot CoT-Unk We directly append some [UNK] tokens to represent the un-tuned prompts for the LLM to perform CoT reasoning. This baseline evaluates the effectiveness of projection tuning for soft thought tokens.

Zero-Shot Assist-CoT The assistant model is prompted to generate a hard-token sequence under standard CoT prompting, truncated at 24 tokens. This sequence is then used as a prompt for the LLM to perform CoT reasoning. This baseline evaluates the effectiveness of soft thoughts by comparing them with hard-token prompts.

Coconut Hao et al. (2024) propose training LLMs to reason in a continuous latent space by iteratively feeding hidden states from the previous step as input embeddings to the next step. The continuous thought encodes rich information, allowing the model to explore more effective reasoning paths. We use their official code¹ to implement this baseline. To adapt Coconut to larger Llama3.1 and Qwen2.5 models, we apply LoRA fine-tuning.

LoRA Fine-Tuning We apply LoRA fine-tuning (Hu et al., 2022) ($r = 16$) with the language modeling objective as our baseline. This baseline examines the effectiveness of appending soft thoughts to LLMs compared to traditional parameter-efficient methods like LoRA.

Implementation details for baselines as well as SoftCoT is shown in Appendix A.

5 Results and Discussions

5.1 Comparison with Baselines

To evaluate SoftCoT, we compare its performance against the baselines introduced in Section 4.2. The results are summarized in Table 2:

¹<https://github.com/facebookresearch/coconut>

Model	GSM8K	ASDiv-Aug	AQuA	StrategyQA	DU	Avg.
	Mathematical			Commonsense	Symbolic	
<i>GPT-2</i>						
Coconut (Hao et al., 2024)	34.10* _{±1.50}	38.92 [†] _{±0.00}	22.83 [†] _{±0.00}	-	-	-
<i>LLaMA-3.1-8B-Instruct</i>						
Zero-Shot CoT	79.61 _{±0.81}	86.78 _{±0.63}	54.65 _{±2.43}	65.63 _{±3.31}	54.40 _{±2.40}	68.21
Zero-Shot CoT-Unk	79.95 _{±0.59}	86.90 _{±0.41}	55.28 _{±1.88}	66.16 _{±2.70}	54.16 _{±1.46}	68.49
Zero-Shot Assist-CoT	80.76 _{±1.53}	86.96 _{±0.46}	55.83 _{±2.98}	66.55 _{±3.99}	58.24 _{±3.56}	69.67
<i>LoRA Fine-Tuning</i>						
Coconut (Hao et al., 2024) [†]	75.66 _{±0.00}	86.67 _{±0.00}	52.36 _{±0.00}	-	-	-
SoftCoT (Ours)	76.12 _{±0.00}	86.80 _{±0.00}	53.15 _{±0.00}	-	-	-
	81.03_{±0.42}	87.19_{±0.40}	56.30_{±1.67}	69.04_{±1.23}	59.04_{±1.93}	70.52

Table 2: Model comparison with baselines. “DU” indicates the Date Understanding (BIG.Bench.authors, 2023) dataset. The first row is under the backbone of GPT-2 (Radford et al., 2019) as backbone. The following rows are under the backbone of LLaMA-3.1-8B-Instruct (Dubey et al., 2024). The last three rows are models trained via the language modeling objective. We run for 5 random seeds and report the average accuracy as well as the standard variance. “*” indicates that the accuracy is reported by Hao et al. (2024). “†” indicates the results that we modify and run the official code of Coconut. ± 0.00 indicates that we only run once for baseline results.

(1) Supervised LoRA Fine-Tuning performs worse than zero-shot CoT, which make Coconut not applicable to larger language models:

We modify and run the official implementation of Coconut, adapting it to LLaMA-3.1-8B-Instruct. Our findings indicate that Coconut exhibits performance degradation following supervised fine-tuning with the language modeling objective. Further experiments on LoRA Fine-Tuning show that it fails to match the performance of zero-shot CoT, which demonstrates that the updated LLM after supervised fine-tuning will have a degradation of previously learned capabilities. We conjecture this drop in reasoning performance is attributable to **catastrophic forgetting**, which aligns with findings from prior studies, including Kalajdziewski (2024) and Lobo et al. (2024), which have reported similar issues. Thus, to mitigate the catastrophic forgetting, we fixed the whole LLM in SoftCoT.

(2) Incorporating [UNK] tokens mitigates performance variance: We examined the effect of directly adding [UNK] tokens as thoughts $\mathcal{T}_{\text{soft}}$ in Eq. (6). The results demonstrate a slight improvement in overall performance and a reduction in variance. The [UNK] token, also known as the “pause token” (Goyal et al., 2024), appears to expand the model’s computational capacity, leading to more stable and consistent outputs.

(3) Assistant model is effective to facilitate CoT reasoning: We utilize instruction to require the assistant model generate some hard prompts, which

can be regarded as the initial thoughts for CoT reasoning. Experiment results show that although it has a larger variance than CoT-Unk, it facilitates the LLM for more diverse CoT generation, which leads to the performance improvement from 68.49 to 69.67 in average.

(4) SoftCoT consistently benefits from the supervised fine-tuning: Overall, our proposed SoftCoT consistently outperforms baselines across all five reasoning datasets, involving the mathematical reasoning, the commonsense reasoning, and the symbolic reasoning. The experimental result highlights that our SoftCoT benefits from the supervised fine-tuning and mitigates the catastrophic forgetting problems in state-of-the-art LLMs.

5.2 Generalization to Other LLM Backbones

In addition to LLaMA-3.1, we evaluate SoftCoT on another state-of-the-art LLM family: Qwen2.5 (Yang et al., 2024). Specifically, we select Qwen2.5-7B-Instruct as the backbone LLM to assess the generalization capability of SoftCoT. As shown in Table 3, our analysis yields the following three key findings:

(1) SoftCoT is effective across different backbone models: Experimental results on Qwen2.5-7B-Instruct show that SoftCoT consistently improves performance across all reasoning datasets, underscoring its robustness. These findings suggest that SoftCoT serves as a generalizable framework that can be effectively adapted to diverse state-of-

Model	GSM8K	ASDiv-Aug	AQuA	StrategyQA	DU	Avg.
	Mathematical			Commonsense	Symbolic	
Zero-Shot CoT	83.70 \pm 0.78	87.19 \pm 0.28	64.53 \pm 3.27	49.65 \pm 3.18	66.40 \pm 2.26	70.29
Zero-Shot CoT-Unk	84.12 \pm 0.71	86.94 \pm 0.89	64.72 \pm 2.06	50.74 \pm 1.90	66.48 \pm 1.43	70.60
Zero-Shot Assist-CoT	84.85 \pm 0.35	88.63 \pm 1.05	64.96 \pm 2.83	52.71 \pm 2.65	67.04 \pm 2.84	71.64
LoRA Fine-Tuning	81.80 \pm 0.00	86.80 \pm 0.00	62.60 \pm 0.00	-	-	-
Coconut (Hao et al., 2024)	82.49 \pm 0.00	86.90 \pm 0.00	63.39 \pm 0.00	-	-	-
SoftCoT (Ours)	85.81\pm1.82	88.90\pm1.01	72.44\pm2.19	60.61\pm1.55	67.52\pm2.92	75.06

Table 3: Model performance using Qwen2.5-7B-Instruct. The first three rows are results without training, the last three rows are results trained by language modeling objective. Notably, only the parameters in the projection module is trained for our SoftCoT.

the-art LLM architectures.

(2) **SoftCoT enhances LLMs’ weaker areas while preserving their strengths:** Experiments on both LLaMA and Qwen LLMs reveal that SoftCoT yields the most significant improvements in commonsense reasoning tasks, where LLMs typically underperform compared to mathematical reasoning. This advantage may stem from SoftCoT’s ability to generate contextually relevant continuous thought processes, effectively activating the corresponding knowledge areas within the model. Furthermore, SoftCoT helps mitigate catastrophic forgetting in domains where LLMs already excel, such as mathematical reasoning, thereby preserving and reinforcing existing capabilities.

(3) **SoftCoT facilitates domain transfer:** Given that the Date Understanding dataset lacks training samples, we train the model on other similar datasets and apply zero-shot transfer to evaluate its generalization on Date Understanding. The results indicate that SoftCoT consistently enhances performance in zero-shot domain transfer scenarios, further demonstrating its adaptability.

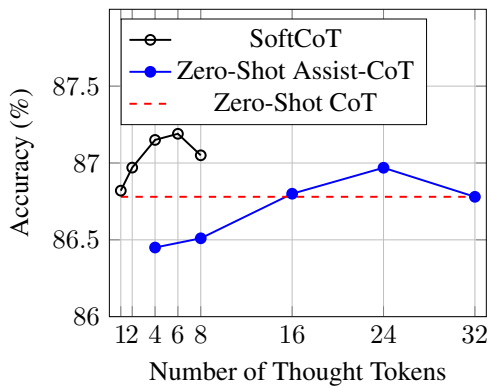


Figure 2: The impact of thought token numbers in ASDiv-Aug using LLaMA-3.1-8B-Instruct.

5.3 Model Analysis and More Studies

5.3.1 The Number of Thought Tokens

To better understand SoftCoT, we conduct experiments to examine the impact of the number of thought tokens. The results, presented in Figure 2, lead to the following key observations:

(1) **Soft thoughts reduce the required number of thought tokens:** We observe that SoftCoT achieves optimal performance with only six thought tokens, whereas Zero-Shot Assist-CoT requires 24 thought tokens to reach a similar level of effectiveness. This suggests that soft thoughts, which operate in a continuous space, exhibit a stronger representational capacity than hard thoughts expressed in the discrete language space. Our experiments indicate that the optimal number of hard thought tokens is approximately four times that of soft thought tokens, aligning with the 5x ratio reported by Cheng and Durme (2024).

(2) **SoftCoT mitigates the catastrophic forgetting problem:** Experimental results show that SoftCoT consistently outperforms Zero-Shot CoT across all tested numbers of soft thought tokens. In contrast, Zero-Shot Assist-CoT underperforms compared to Zero-Shot CoT when the number of thought tokens is insufficient. This is likely because the assistant model fails to generate a sufficiently informative set of thought tokens under these constraints, introducing noise and leading to confusion in the LLM’s reasoning process.

5.3.2 Varying the size of Assistant Model

In the zero-shot assist-CoT setting, the role of the assistant model is to extract or summarize key reasoning cues to guide the target reasoning LLM. While the assistant provides helpful context, the final answer is generated solely by the reasoning

Model	GSM8K		ASDiv-Aug		AQuA		StrategyQA		DU	
	$N = 1$	$N = 10$	$N = 1$	$N = 10$	$N = 1$	$N = 10$	$N = 1$	$N = 10$	$N = 1$	$N = 10$
Zero-Shot CoT	79.61 \pm 0.81	90.36 \pm 0.40	86.78 \pm 0.63	89.23 \pm 0.17	54.65 \pm 2.43	63.23 \pm 0.86	65.63 \pm 3.31	70.13 \pm 0.47	54.40 \pm 2.40	65.76 \pm 1.54
Zero-Shot Assist-CoT	80.76 \pm 1.53	90.43 \pm 0.69	86.96 \pm 0.46	89.48 \pm 0.36	55.83 \pm 2.98	63.62 \pm 0.99	66.55 \pm 3.99	70.48 \pm 0.68	58.24 \pm 3.56	65.84 \pm 1.93
SoftCoT (Ours)	81.03\pm0.42	90.63\pm0.39	87.19\pm0.40	89.75\pm0.29	56.30\pm1.67	65.51\pm0.72	69.04\pm1.23	71.14\pm0.10	59.04\pm1.93	67.36\pm1.12

Table 4: Self Consistency for SoftCoT on LLaMA-3.1-8B-Instruct. “ N ” indicates the number of reasoning chains.

Method	0.5B	1.5B	7B
Zero-Shot CoT	83.70	83.70	83.70
Zero-Shot Assist-CoT	84.78	84.85	84.90
SoftCoT	85.76	85.81	85.84

Table 5: Performance on GSM8K with different sizes of assistant model on Qwen2.5 series.

LLM. Empirically, we observe that the scale of the assistant model has limited impact on the accuracy of the final answer (see row “Zero-shot Assist-CoT” in Table 5).

A similar observation in the SoftCoT setting. Although the assistant model now produces continuous soft thought tokens instead of discrete hard tokens, its fundamental role—providing intermediate reasoning signals—remains unchanged. Our experiments similarly show that varying the assistant model’s scale does not significantly affect final task performance (see row “SoftCoT” in Table 5).

5.3.3 Model-Orthogonal Factors

Self-Consistency (Wang et al., 2023) is a widely adopted technique for enhancing Chain-of-Thought (CoT) reasoning by expanding the search space. One of the most straightforward implementations involves generating multiple CoT reasoning paths and determining the final answer through majority voting. This approach is effective in mitigating errors in reasoning steps by leveraging the diversity of generated thought processes.

To further assess the effectiveness of SoftCoT, we conduct experiments incorporating self-consistency. As shown in Table 4, SoftCoT consistently outperforms baseline models, demonstrating that its benefits are complementary to those of self-consistency rather than being redundant or conflicting. This suggests that SoftCoT introduces an independent improvement mechanism, which can be effectively combined with self-consistency for enhanced reasoning performance.

A key advantage of SoftCoT in this context is its ability to provide a more expressive and compact

representation of intermediate reasoning steps in continuous space. Unlike traditional CoT reasoning, where discrete thought tokens may introduce inconsistencies or redundant reasoning paths, SoftCoT enables more efficient reasoning trajectories with fewer tokens. This allows self-consistency methods to aggregate results from higher-quality reasoning paths, leading to a more robust and accurate final prediction.

6 Conclusion

In this paper, we introduce SoftCoT, a soft chain-of-thought prompting approach for efficient LLM reasoning. SoftCoT consists of three steps: (1) an assistant model generates soft thought tokens, (2) a projection module trained to map the soft thoughts to LLM’s representation space, and (3) the LLM applies soft thoughts for reasoning. To enhance efficiency, SoftCoT speculatively generates *all* the soft thought tokens in a single forward pass. To mitigate the catastrophic forgetting, SoftCoT freezes the backbone LLM and only tunes the projection module. Experiments on five datasets across three types of reason tasks demonstrate the effectiveness of our proposed SoftCoT. Experiments on multiple LLMs as well as orthogonal method such as self-consistency shows the robustness of SoftCoT, which can be adapted in widely scenarios.

Acknowledgements

This research is supported, in part, by the Joint NTU-WeBank Research Centre on Fintech (Award No. NWJ-2020-007), Nanyang Technological University, Singapore. This research is also supported, in part, by the National Research Foundation, Prime Minister’s Office, Singapore under its NRF Investigatorship Programme (NRFI Award No. NRF-NRFI05-2019-0002). Xu Guo thanks the Wallenberg-NTU Presidential Postdoctoral Fellowship. Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not reflect the views of National Research Foundation, Singapore.

Limitations

While SoftCoT represents a promising advancement in Chain-of-Thought (CoT) reasoning within a continuous space, several limitations must be acknowledged.

SoftCoT Cannot Fully Replace the Reasoning Path : Although SoftCoT employs soft thought tokens for reasoning, it does not entirely replace the reasoning path. The decoding stage functions as a search process, which is a crucial component of CoT reasoning. Soft thought tokens primarily serve to enrich the probability space for exploration rather than acting as the search mechanism itself.

Need for Further Empirical Evidence on Scalability : SoftCoT has been evaluated on LLaMA-3.1-8B-Instruct and Qwen2.5-7B-Instruct. However, larger backbone LLMs exist within the same model families. While its success on models with approximately 7–8 billion parameters suggests potential applicability to larger-scale models, its scalability to extremely large LLMs remains an open question and requires thorough empirical validation.

References

- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoeffler. 2024. [Graph of thoughts: Solving elaborate problems with large language models](#). In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2024, February 20-27, 2024, Vancouver, Canada*, pages 17682–17690. AAAI Press.
- BIG.Bench.authors. 2023. [Beyond the imitation game: Quantifying and extrapolating the capabilities of language models](#). *Trans. Mach. Learn. Res.*, 2023.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Jeffrey Cheng and Benjamin Van Durme. 2024. [Compressed chain of thought: Efficient reasoning through dense representations](#). *arXiv preprint arXiv:2412.13171*.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. [PaLM: Scaling language modeling with pathways](#). *Journal of Machine Learning Research*, 24(240):1–113.
- Zheng Chu, Jingchang Chen, Qianglong Chen, Weijiang Yu, Tao He, Haotian Wang, Weihua Peng, Ming Liu, Bing Qin, and Ting Liu. 2024. [Navigate through enigmatic labyrinth A survey of chain of thought reasoning: Advances, frontiers and future](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 1173–1203. Association for Computational Linguistics.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *arXiv preprint arXiv:2110.14168*.
- Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. [GLM: General language model pretraining with autoregressive blank infilling](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 320–335, Dublin, Ireland. Association for Computational Linguistics.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. [The llama 3 herd of models](#). *arXiv preprint arXiv:2407.21783*.
- Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. [Did aristotle use a laptop? A question answering benchmark with implicit reasoning strategies](#). *Trans. Assoc. Comput. Linguistics*, 9:346–361.
- Sachin Goyal, Ziwei Ji, Ankit Singh Rawat, Aditya Krishna Menon, Sanjiv Kumar, and Vaishnavh Nagarajan. 2024. [Think before you speak: Training language models with pause tokens](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li, Zhiting Hu, Jason Weston, and Yuandong Tian. 2024.

- Training large language models to reason in a continuous latent space. *arXiv preprint arXiv:2412.06769*.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. [LoRA: Low-rank adaptation of large language models](#). In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.
- Damjan Kalajdzievski. 2024. [Scaling laws for forgetting when fine-tuning large language models](#). *arXiv preprint arXiv:2401.05605*.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. [Large language models are zero-shot reasoners](#). In *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*.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. [The power of scale for parameter-efficient prompt tuning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 3045–3059. Association for Computational Linguistics.
- Yaniv Leviathan, Matan Kalman, and Yossi Matias. 2023. [Fast inference from transformers via speculative decoding](#). In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 19274–19286. PMLR.
- Haochen Li, Jonathan Leung, and Zhiqi Shen. 2024. [Towards goal-oriented prompt engineering for large language models: A survey](#). *arXiv preprint arXiv:2401.14043*.
- Zekun Li, Baolin Peng, Pengcheng He, Michel Galley, Jianfeng Gao, and Xifeng Yan. 2023. [Guiding large language models via directional stimulus prompting](#). In *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*.
- Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017. [Program induction by rationale generation: Learning to solve and explain algebraic word problems](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, pages 158–167. Association for Computational Linguistics.
- Elita A. Lobo, Chirag Agarwal, and Himabindu Lakkaraju. 2024. [On the impact of fine-tuning on chain-of-thought reasoning](#). *arXiv preprint arXiv:2411.15382*.
- Ilya Loshchilov and Frank Hutter. 2019. [Decoupled weight decay regularization](#). In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.
- Shen-Yun Miao, Chao-Chun Liang, and Keh-Yih Su. 2020. [A diverse corpus for evaluating and developing english math word problem solvers](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 975–984. Association for Computational Linguistics.
- OpenAI. 2023. [GPT-4 technical report](#). *arXiv preprint arXiv:2303.08774*.
- Qwen Team. 2025. [Qwen3](#).
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie Huang, Nan Duan, and Weizhu Chen. 2023. [Synthetic prompting: Generating chain-of-thought demonstrations for large language models](#). In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 30706–30775. PMLR.
- Xuan Shen, Yizhou Wang, Xiangxi Shi, Yanzhi Wang, Pu Zhao, and Jiuxiang Gu. 2025. [Efficient reasoning with hidden thinking](#). *arXiv preprint arXiv:2501.19201*.
- Zayne Sprague, Fangcong Yin, Juan Diego Rodriguez, Dongwei Jiang, Manya Wadhwa, Prasann Singhal, Xinyu Zhao, Xi Ye, Kyle Mahowald, and Greg Durrett. 2024. [To cot or not to cot? chain-of-thought helps mainly on math and symbolic reasoning](#). *arXiv preprint arXiv:2409.12183*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. [Llama 2: Open foundation and fine-tuned chat models](#). *arXiv preprint arXiv:2307.09288*.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. [Self-consistency improves chain of thought reasoning in language models](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. [Chain-of-thought prompting elicits reasoning in large language models](#). In *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.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Guowei Xu, Peng Jin, Hao Li, Yibing Song, Lichao Sun, and Li Yuan. 2024. [Llava-cot: Let vision language models reason step-by-step](#). *arXiv preprint arXiv:2411.10440*.

Yige Xu, Zhiwei Zeng, and Zhiqi Shen. 2023. [Efficient cross-task prompt tuning for few-shot conversational emotion recognition](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 11654–11666, Singapore. Association for Computational Linguistics.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2024. [Qwen2.5 technical report](#). *arXiv preprint arXiv:2412.15115*.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. [Tree of thoughts: Deliberate problem solving with large language models](#). In *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*.

Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. 2023. [Automatic chain of thought prompting in large language models](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.

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. [Least-to-most prompting enables complex reasoning in large language models](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.

Appendix

A Implementation Details

We use the Huggingface Transformers framework (Wolf et al., 2020) for training. All models are trained on a single NVIDIA A100-80G GPU. The projection module is trained for 10 epochs and we use the checkpoint after the last epoch for evaluation. We use AdamW (Loshchilov and Hutter, 2019) as the optimizer. The learning rate is set as $1e-5$, the weight decay is set as 0.01. To fully utilize the GPU memory, we set the batch size as 8 for Qwen and 16 for LLaMA, which depends on the GPU memory usage. For Zero-Shot CoT-Unk, $L \in \{4, 8, 16, 24\}$ [UNK] tokens are appended and we report the best performance. For Zero-Shot Assist-CoT, we use the assistant model to generate 24 tokens. For SoftCoT, the number of soft thought tokens L is set as 4.

More details can be found in our source code at <https://github.com/xuyige/SoftCoT>. To better reproduce our work, we also release a reproduction that includes the training logs, the intermediate checkpoints, and the evaluation logs.

B Experiments on Qwen3-Series

In this section, we evaluate the performance of our proposed method on the recently released Qwen3-8B model (Qwen Team, 2025), a state-of-the-art large language model that continues the Qwen-series advancements. As presented in Table 6 and Table 7, our empirical results reveal several important findings:

(1) **SoftCoT remains consistently effective when applied to cutting-edge LLMs such as Qwen3-8B.** This demonstrates the generalizability and robustness of the SoftCoT framework across different LLM architectures. Despite architectural differences between Qwen3-8B and prior models such as LLaMA-3.1 or Qwen2.5, the underlying principles of SoftCoT, leveraging soft latent thought representations and their diversity, continue to provide performance gains across diverse reasoning tasks.

(2) **The key insights obtained from experiments on LLaMA-3.1 and Qwen2.5 seamlessly transfer to Qwen3-8B.** Specifically, we observe that SoftCoT helps alleviate the issue of catastrophic forgetting when scaling up reasoning chains, which is particularly critical when deploying CoT-based strategies in real-world inference scenarios. Addi-

Model	GSM8K		ASDiv-Aug		AQuA	
	$N = 1$	$N = 10$	$N = 1$	$N = 10$	$N = 1$	$N = 10$
Zero-Shot CoT	91.86 \pm 0.41	92.22 \pm 0.47	91.70 \pm 0.26	91.97 \pm 0.13	70.00 \pm 1.56	76.77 \pm 0.62
Zero-Shot Assist-CoT	91.90 \pm 0.50	92.68 \pm 0.17	91.64 \pm 0.16	91.91 \pm 0.28	70.16 \pm 2.35	76.77 \pm 0.79
Coconut (Hao et al., 2024)	87.95 \pm 0.00	90.37 \pm 0.00	89.40 \pm 0.00	90.37 \pm 0.00	68.50 \pm 0.00	76.38 \pm 0.00
SoftCoT (Ours)	92.48\pm0.29	93.19\pm0.32	91.83\pm0.19	92.14\pm0.15	75.04\pm2.68	80.63\pm1.90

Table 6: Model comparison with baselines on Qwen3-8B (Qwen Team, 2025) on mathematical reasoning datasets. “ N ” indicates the number of reasoning chains. For $N = 1$, we do not apply any self-consistency techniques. For $N = 10$, we use self-consistency with 10 return sequences to apply multiple reasoning chains. Majority vote is applied to obtain the final prediction answer.

Model	StrategyQA		DU	
	$N = 1$	$N = 10$	$N = 1$	$N = 10$
Zero-Shot CoT	69.87 \pm 0.35	70.96 \pm 0.15	80.32 \pm 2.32	84.56 \pm 0.61
Zero-Shot Assist-CoT	69.78 \pm 0.33	70.92 \pm 0.28	80.56 \pm 1.15	84.80 \pm 1.17
SoftCoT (Ours)	70.17\pm0.63	71.18\pm0.15	85.60\pm0.57	87.20\pm0.75

Table 7: Model comparison with baselines on Qwen3-8B (Qwen Team, 2025) on commonsense and symbolic reasoning datasets. “DU” indicates the Date Understanding (BIG.Bench.authors, 2023) dataset. “ N ” indicates the number of reasoning chains. For $N = 1$, we do not apply any self-consistency techniques. For $N = 10$, we use self-consistency with 10 return sequences to apply multiple reasoning chains. Majority vote is applied to obtain the final prediction answer. We run for 5 random seeds and report the average accuracy as well as the standard variance.

tionally, the orthogonality between SoftCoT and self-consistency, previously established in earlier models, holds true for Qwen3-8B. This suggests that the benefits of SoftCoT can be compounded with self-consistency sampling strategies, offering further improvements without redundancy or interference.

Overall, these findings support the claim that SoftCoT is a model-agnostic strategy that effectively adapts to new LLM releases, maintaining its strengths in enhancing reasoning performance and flexibility.

C Instruction Templates

In this section, we release the examples for GSM8K for reference. The instruction template for Zero-Shot CoT is from Sprague et al. (2024). We design the task-specific instruction template for the assistant model, which is listed under “**Input for Assistant Model**” in the examples of Zero-Shot Assist-CoT and SoftCoT.

In the Zero-Shot Assist-CoT, we require the assistant model to generate 24 tokens to facilitate the reasoning process. In SoftCoT, we use two special token “<reserved_special_token_0>” and

“<reserved_special_token_1>” to represent the beginning as well as the ending of soft thought tokens, and “<end_of_text>” as the [UNK] token, which is the placeholder for soft thought tokens.

In the following examples, contents under the “**Answer**” section is the reasoning steps generated by the LLM, the contents under the “**Model Prediction**” section is number that represents the model prediction extracted by pre-defined manual rules, and the contents under the “**Evaluate**” section is the judgement about the model prediction is correct or not.

Example Input-Output for Zero-Shot CoT (GSM8K) with LLaMA-3.1-8B-Instruct

Input:

Solve the following math problem efficiently and clearly:

- For simple problems (2 steps or fewer): Provide a concise solution with minimal equation.
- For complex problems (3 steps or more):

Use this step-by-step format:

Step 1: [Brief calculations]

Step 2: [Brief calculations]

...

Regardless of the approach, always conclude with:

Therefore, the final answer is: $\boxed{\text{answer}}$. I hope it is correct.

Where [answer] is just the final number or expression that solves the problem.

Problem: John runs 60 miles a week. He runs 3 days a week. He runs 3 hours the first day and half as much the other two days he runs. How fast does he run?

Answer:

Step 1: Calculate the total hours John runs in a week

To find out how much he runs in total in a week, we multiply the number of hours he runs on the first day by 3 (for the 3 days in a week) and add the same multiplied by 0.5 (since he runs half as much on the other two days).

Step 2: Calculate the number of hours John runs in a week

John runs 3 hours on the first day, so he runs $(3 * 3) + (3 * 0.5 * 2)$ hours in a week.

Step 3: Calculate John's total hours and miles run in a week

$3 * 3 = 9$ hours for the first day, and $3 * 0.5 * 2 = 3$ hours for the other two days. His total hours are $9 + 3 = 12$ hours in a week.

Step 4: Calculate John's speed in miles per hour

John runs 60 miles in 12 hours. So, we divide $60 / 12$ to find his speed.

$60 / 12 = 5$ miles per hour.

Therefore, the final answer is: $\boxed{5}$.

Model Prediction:

5

Evaluate:

Incorrect, the ground-truth is 10 and the model prediction is 5.

Example Input-Output for Zero-Shot Assist-CoT (GSM8K) with LLaMA-3.1-8B-Instruct

Input for Assistant Model:

You are required to generate 24 tokens to help another language model to solve the following math reasoning task efficiently and clearly. Here are the requirements of your generated tokens:

- The tokens should include some useful information for the reasoning problem, for example, the numbers and the operations needed for calculation.
- Generate the tokens starts from the most important or the highest related tokens.
- ****Informative tokens are required****: (1) Do not need to generate a sentence or paragraph, (2) Do not need to generate the uninformative tokens such as serial number.
- The tokens should be useful for large language model to answer the question with the numbers.

...

Here is the problem: John runs 60 miles a week. He runs 3 days a week. He runs 3 hours the first day and half as much the other two days he runs. How fast does he run?

Input for Large Reasoning Model:

Solve the following math problem efficiently and clearly:

- For simple problems (2 steps or fewer): Provide a concise solution with minimal equation.
- For complex problems (3 steps or more):

Use this step-by-step format:

Step 1: [Brief calculations]

Step 2: [Brief calculations]

...

Regardless of the approach, always conclude with:

Therefore, the final answer is: $\boxed{\text{answer}}$. I hope it is correct.

Where [answer] is just the final number or expression that solves the problem.

Problem: John runs 60 miles a week. He runs 3 days a week. He runs 3 hours the first day and half as much the other two days he runs. How fast does he run?

There are some prompts generated by a weaker assistant model. Some prompts maybe useful while others maybe unuseful for your reasoning. If the prompts are correct, you can use it as reference. If the prompts are not correct, you can ignore them and focus back to solving the problem.

Here are prompts:

1. a
2. miles
3. week
4. days
5. 3-mile
6. run

Answer:

Step 1: Calculate the total hours John runs in a week

John runs 3 hours on the first day. On the other two days, he runs half as much, which is $3 / 2 = 1.5$ hours per day. So, the total hours John runs in a week is $3 + 1.5 + 1.5 = 6$ hours.

Step 2: Calculate John's speed

To find John's speed, divide the total distance he runs in a week (60 miles) by the total hours he runs in a week (6 hours). $\text{Speed} = \text{Distance} / \text{Time} = 60 \text{ miles} / 6 \text{ hours} = 10 \text{ miles per hour}$.

Therefore, the final answer is: $\boxed{10}$.

Model Prediction:

10

Evaluate:

Correct!

Example Input-Output for SoftCoT (GSM8K) with LLaMA-3.1-8B-Instruct

Input for Assistant Model:

<user>:

You are required to generate 4 tokens to help another language model to solve the following math reasoning task efficiently and clearly. Here are the requirements of your generated tokens:

- The tokens should include some useful information for the reasoning problem, for example, the numbers and the operations needed for calculation.
- Generate the tokens starts from the most important or the highest related tokens.
- ****Informative tokens are required****: (1) Do not need to generate a sentence or paragraph, (2) Do not need to generate the uninformative tokens such as serial number.
- The tokens should be useful for large language model to answer the question with the numbers.

...

Here is the problem: John runs 60 miles a week. He runs 3 days a week. He runs 3 hours the first day and half as much the other two days he runs. How fast does he run?

<assistant>:

Here are 4 tokens to help the language model solve this math reasoning task: <reserved_special_token_0> <end_of_text> <end_of_text> <end_of_text> <end_of_text> <reserved_special_token_1>

Input for Large Reasoning Model:

Solve the following math problem efficiently and clearly:

- For simple problems (2 steps or fewer): Provide a concise solution with minimal equation.
- For complex problems (3 steps or more):

Use this step-by-step format:

Step 1: [Brief calculations]

Step 2: [Brief calculations]

...

Regardless of the approach, always conclude with:

Therefore, the final answer is: $\boxed{\text{answer}}$. I hope it is correct.

Where [answer] is just the final number or expression that solves the problem.

Problem: John runs 60 miles a week. He runs 3 days a week. He runs 3 hours the first day and half as much the other two days he runs. How fast does he run?

There are some prompts generated by a weaker assistant model. Some prompts maybe useful while others maybe unuseful for your reasoning. If the prompts are correct, you can use it as reference. If the prompts are not correct, you can ignore them and focus back to solving the problem.

Here are prompts: <reserved_special_token_0> <end_of_text> <end_of_text> <end_of_text> <end_of_text> <reserved_special_token_1>

Answer:

Step 1: Calculate the total time John runs each week in hours. John runs 3 hours on the first day. For the other two days, he runs half as much, which is $3 / 2 = 1.5$ hours per day. So, his total running time each week is $3 + 2 * 1.5 = 6$ hours.

Step 2: Determine John's speed in miles per hour by dividing his total weekly distance by his total weekly running time. John runs 60 miles a week and his total running time each week is 6 hours. Thus, his speed is $60 / 6 = 10$ miles per hour.

Therefore, the final answer is: $\boxed{10}$.

Model Prediction:

10

Evaluate:

Correct!