

# Mind’s Mirror: Distilling Self-Evaluation Capability and Comprehensive Thinking from Large Language Models

Anonymous ACL submission

## Abstract

Large language models (LLMs) have achieved remarkable advancements in natural language processing. However, the sheer scale and computational demands of these models present formidable challenges when considering their practical deployment in resource-constrained contexts. While techniques such as chain-of-thought (CoT) distillation have displayed promise in distilling LLMs into small language models (SLMs), there is a risk that distilled SLMs may still inherit flawed reasoning and hallucinations from LLMs. To address these issues, we propose a twofold methodology: First, we introduce a novel method for distilling the self-evaluation capability from LLMs into SLMs, aiming to mitigate the adverse effects of flawed reasoning and hallucinations inherited from the LLM. Second, we advocate for a comprehensive distillation process that incorporates multiple distinct CoTs and self-evaluation outputs, to ensure a more thorough and robust knowledge transfer into SLMs. Experiments on three NLP benchmarks demonstrate that our method significantly improves the performance of distilled SLMs, offering a new perspective for developing more effective and efficient SLMs in resource-constrained environments. We will publicly release our code upon acceptance.

## 1 Introduction

With the gradual increase in the number of parameters, large language models (LLMs) have achieved significant successes in the field of natural language processing (Brown et al., 2020; Kaplan et al., 2020; Hoffmann et al., 2022; Chowdhery et al., 2023; OpenAI, 2023). However, LLMs’ tremendous model sizes and computational demands introduce challenges to their practical deployment, especially in resource-limited environments. To address these challenges, various studies have delved into the compression of LLMs into small language models

(SLMs) using knowledge distillation techniques and have led to significant reductions in computational complexity and inference costs (Jiang et al., 2020; Gu et al., 2023; Agarwal et al., 2023). This process involves traditional teacher-student learning methods and the more recent chain-of-thought (CoT) distillation method (Zhu et al., 2023). The CoT distillation methods use the CoT reasoning process of LLMs as supervision for training SLMs, rather than just labels. This allows SLMs to learn the reasoning process of LLMs, thereby improving the performance of SLMs.

While these CoT distillation methods have proven to be beneficial, they are not without their flaws, particularly:

1. Even during the CoT distillation process, the distilled SLMs remain vulnerable to the flawed supervision provided by LLMs, as observations suggest that chains of thought (CoTs) generated by LLMs may contain hallucinations (Zhang et al., 2023), accumulate errors (Shen et al., 2021), or lack robustness (Vaswani et al., 2017; Radford et al., 2019; Brown et al., 2020; Zhang et al., 2022). As shown in the example in Figure 1, “LLM Random CoT 2” incorrectly broadens the scope of the premise by arguing that “Being an animal welfare advocate means caring about all the animals that inhabit the planet.” In practice, it is not easy to exclude these flawed CoTs, since the ground truth of CoTs is not always easily obtainable (Zhang et al., 2023). Training SLMs with these flawed CoTs will result in SLMs inheriting these flaws and performance degradation (Alemohammad et al., 2023; Ho et al., 2023).
2. A single instance of CoT might not capture the diverse reasoning routes LLMs can explore, limiting the richness of the distilled knowledge of SLMs. Furthermore, relying solely

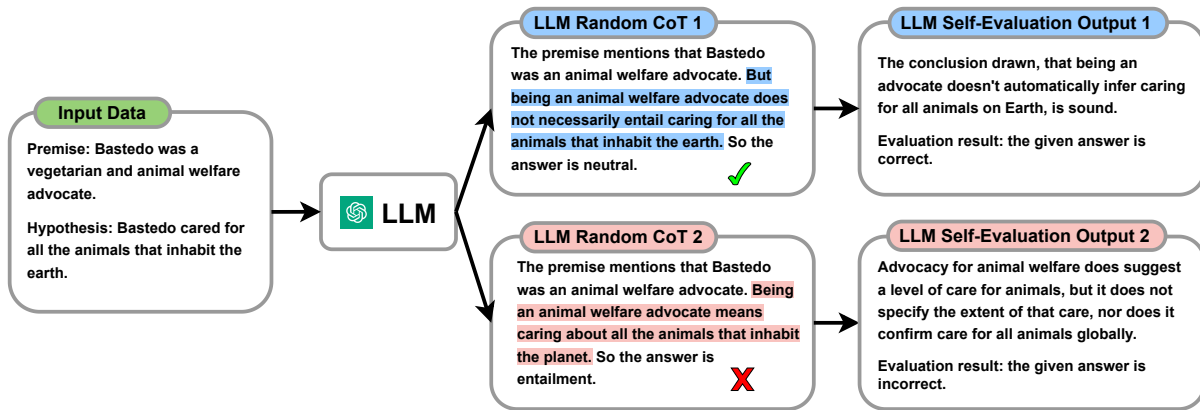


Figure 1: Examples of both the random CoT responses and their self-evaluation outputs generated by the LLM during natural language inference tasks. The human-like self-evaluation of the LLM enables the LLM to self-evaluate the correctness of its CoT reasoning processes, identifying which are correct (highlighted in blue) and which are incorrect (highlighted in red) in these randomly generated CoT reasoning.

on the CoT reasoning process as supervision for training SLMs is insufficient to distill the comprehensive capabilities of LLMs, such as the ability to check the correctness of answers.

To mitigate the impact of these flawed CoTs and allow SLMs to learn more comprehensive capabilities, we propose an innovative methodology that involves training SLMs to possess the self-evaluation capability. Humans often evaluate their reasoning processes to reduce errors in decision-making (Poole and Mackworth, 2010), and a similar self-evaluation capability has also been observed in LLMs (Kadavath et al., 2022; Shinn et al., 2023; Madaan et al., 2023; Paul et al., 2023), which recognizes and corrects the generated hallucinations, unreliable reasoning, and harmful content in a CoT (Pan et al., 2023). Figure 1 illustrates this with an example where incorrect reasoning in “LLM Random CoT 2” is identified and corrected in the self-evaluation. The advantage of self-evaluation is that it does not rely on external resources. However, it is constrained by the inherent capabilities of the model. To address this, we guide SLMs in distillation to learn the self-evaluation capability of LLMs. By learning the ability of LLMs to analyze right from wrong, SLMs can understand both what should and should not be generated, enhancing their predictive accuracy and reliability in various NLP tasks.

To facilitate comprehensive thinking and address the randomness and limitations of relying on a single CoT and a single self-evaluation, our second methodology insight involves distilling SLMs from diverse CoTs and multiple self-evaluation outputs generated by LLMs. This enables SLMs to inherit

a broader range of comprehensive thinking capabilities since multiple CoTs and self-evaluation collectively offer a more comprehensive perspective, derived from the varied state spaces of LLMs.

In summary, our contributions can be outlined as follows:

1. We distill the self-evaluation capability from LLMs into SLMs, primarily focusing on enhancing the accuracy of SLMs across various NLP tasks. This helps SLMs understand the potential reasons behind correct or incorrect reasoning and lays the foundation for mitigating errors (e.g., hallucinations) arising from flawed CoTs.
2. We distill a variety of CoTs and corresponding multiple self-evaluation outputs from LLMs into SLMs, leveraging extensive reasoning chains and self-evaluation outputs derived from the comprehensive state spaces of LLMs, thus enabling SLMs to encompass both enhanced reasoning and more comprehensive model capabilities.
3. Comprehensive experiments verified that our method significantly improves the performance and reliability of distilled SLMs, which enables SLMs to inherit the self-evaluation capability and comprehensive thinking of LLMs and outperforms previous CoT distillation methods.

## 2 Related Work

**Chain-of-thought reasoning** Chain-of-thought (CoT) is a prompting method where a model generates intermediate reasoning steps to enhance its

problem-solving capabilities (Wei et al., 2022). The chain-of-thought with self-consistency (CoT-SC) (Wang et al., 2023b) builds upon CoT, sampling a set of diverse reasoning paths and selecting the most consistent answer as the final answer. This largely mitigates errors introduced by the inherent randomness of LLMs. The Tree of Thoughts (ToT) method (Yao et al., 2023) models problem-solving as a tree search process, enabling LLMs to explore different reasoning pathways and conduct self-evaluation to determine the solution taken at each step. Therefore, by leveraging the capability of LLMs to generate diverse reasoning paths and self-evaluation, ToT significantly enhances the performance of LLMs in solving tasks such as Game of 24, Creative Writing, and Mini Crosswords.

**Self-evaluation in LLMs** Many recent works have leveraged the self-evaluation capability of LLMs to enhance the reliability of their responses, such as Self-Refine (Madaan et al., 2023), Self-Check (Miao et al., 2023), SelfCheckGPT (Manakul et al., 2023), and Reflexion (Shinn et al., 2023). Concurrently, other studies have demonstrated the self-improvement potential of LLMs (Huang et al., 2022; Pan et al., 2023), as exemplified by RLAIIF (Lee et al., 2023). However, these methods are designed for LLMs and do not consider distilling the self-evaluation capability into SLMs.

**Knowledge distillation from LLMs** Knowledge distillation enhances the performance of smaller models by transferring knowledge from larger models (Hinton et al., 2015). This method has been widely adopted for the optimization and compression of models. Recent studies (Hsieh et al., 2023; Li et al., 2023; Ho et al., 2023; Wang et al., 2023a; Magister et al., 2023; Shridhar et al., 2023; Wang et al., 2023c; Chen et al., 2023; Fu et al., 2023) have been focusing on leveraging the CoT reasoning generated by LLMs to enhance the performance of SLMs. For instance, Hsieh et al. (2023) introduced a “Distilling step-by-step” method for extracting rationales from LLMs as additional supervision for training SLMs. Similarly, Li et al. (2023) proposed the Symbolic Chain-of-Thought Distillation (SCoTD) method, which trains SLMs to learn CoT reasoning. Additionally, Ho et al. (2023) presented “Fine-tune-CoT”, a method that generates reasoning samples from LLMs to fine-tune SLMs. However, these methods do not consider mitigating the impact of harmful content in CoTs generated by LLMs on smaller models, as well as distilling other

capabilities beyond CoTs. In contrast, our methodology incorporates the self-evaluation capability of LLMs into distillation, which can be utilized to mitigate the effects of flawed CoTs in a completely unsupervised manner and without relying on external resources, and further allows smaller models to learn the more comprehensive capabilities of LLMs. Furthermore, some related works utilize SLMs with up to several billion parameters and have not been able to validate their effectiveness on SLMs with as few as 220M parameters, so our approach exhibits lower resource requirements and broader applicability.

### 3 Distilling Self-Evaluation Capability and Comprehensive Thinking

We propose a new methodology for distilling the self-evaluation capability and comprehensive thinking of an LLM into an SLM. Our overall framework is illustrated in Figure 2, which operates in 4 steps: (1) Given an LLM and an unlabeled dataset, we utilize CoT prompts to generate diverse rationales and corresponding pseudo-labels from the LLM. (2) By devising self-evaluation prompts, we enable the LLM to evaluate the correctness of its generated CoTs, which also include both the rationales and labels in its self-evaluation outputs. (3) Leveraging the rationales and labels in the self-evaluation outputs generated by the LLM, we employ multi-task learning to train the SLM, enabling the SLM to distinguish right from wrong. (4) Utilizing the diverse rationales in CoTs and labels from either LLM-generated pseudo-labels or human-annotated labels, we employ multi-task learning to train the SLM’s reasoning capability.

#### 3.1 Obtaining diversity CoTs and self-evaluation outputs from the LLM

In our pipeline, an LLM functions as the teacher, while an SLM serves as the student. First, we let the LLM generate multiple different CoTs and self-evaluation outputs for a given task. We utilize few-shot CoT prompting to enhance the quality and standardize the formats of the CoTs generated by the LLM. This process is shown as step 1 and step 2 in Figure 2.

##### 3.1.1 Obtaining multiple CoTs

For an unlabeled dataset  $D$ , we devise a few-shot CoT prompt template  $p$  delineating how the task should be approached. We combine the concrete input data  $x_i$  with  $p$  and use this as input to prompt

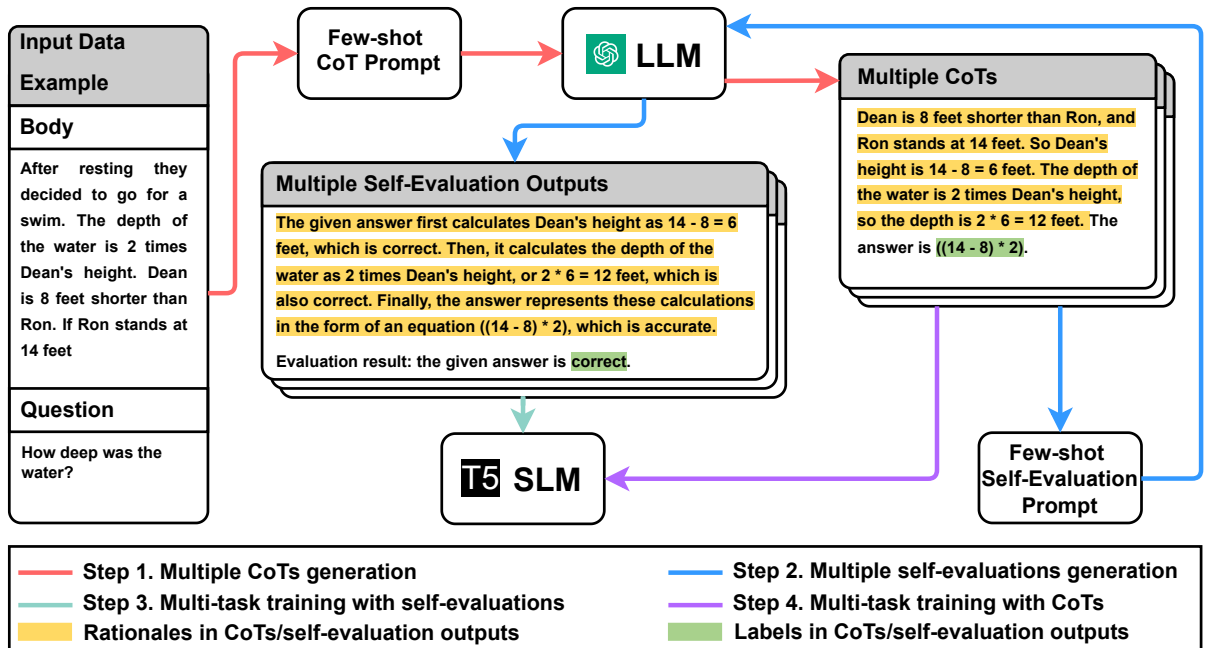


Figure 2: Detailed overview of our proposed methodology. **Step 1:** Obtain multiple CoTs from the LLM (Section 3.1.1). **Step 2:** Obtain multiple self-evaluation outputs from the LLM (Section 3.1.2). **Step 3:** Train the SLM with multiple self-evaluation outputs, enabling the SLM to distinguish right from wrong (Section 3.2.1). **Step 4:** Train the SLM with multiple CoTs to give the SLM comprehensive reasoning capabilities (Section 3.2.2).

the LLM. With examples from  $p$ , the LLM can simulate examples to generate the CoT response for  $x_i$  that contains a rationale  $r_i$  and a pseudo-label  $y_i$  (the yellow part and the green part of the “Multiple CoTs Outputs” in Figure 2). We let the LLM regenerate several times to get multiple different CoTs. Each CoT epitomizes a unique rationale underscoring the comprehensiveness of the state space generated by the LLM, hence broadening the explanation spectrum and laying a robust foundation for decision-making.

### 3.1.2 Obtaining multiple self-evaluation outputs

After forming multiple CoTs representing different thoughts, a self-evaluation phase is initiated to evaluate the correctness of the CoTs. This is essential to imitate the complete human thought process and correct mistakes in reasoning. Given an unlabeled dataset  $D$ , we devise a few-shot self-evaluation prompt template  $p_{eval}$ , which guides the LLM in evaluating each CoT’s correctness. For each CoT  $x_c$ , shown in “Multiple CoTs” in Figure 2, we add it to  $p_{eval}$  and use this as an input to prompt the LLM to generate the self-evaluation. With examples in  $p_{eval}$ , the LLM simulates examples to generate the self-evaluation output for  $x_c$  that also contains a rationale  $r_{eval_i}$  and a label  $y_{eval_i}$  (the yellow part and the green part of the “Multiple Self-Evaluation

Outputs” in Figure 2).

Similarly, to distill a more comprehensive self-evaluation capability of the LLM, we generate multiple different self-evaluation outputs for each CoT. Multiple self-evaluation outputs along with multiple CoTs represent a more comprehensive and complete thought process for the LLM.

### 3.2 Training the SLM with multiple self-evaluation outputs and diverse CoTs

After generating diverse CoTs and their corresponding self-evaluation outputs using the LLM, we begin to train the SLM. Our training methodology for SLMs first emphasizes distilling self-evaluation capability to lay the foundation for reducing the impact of errors in CoTs on SLMs, followed by incorporating comprehensive reasoning capability through diverse CoTs distillation. Hsieh et al. (2023) have demonstrated that multi-task learning can lead to better performance than simply treating rationale and label predictions as a single joint task, and can reduce computation overhead during inference since it allows the SLM to directly predict labels without generating rationales. Hence, we employ multi-task learning to train the SLM for self-evaluation capability and CoT reasoning capability. By appending different “task prefixes” at the beginning of the input, we can direct the SLM to generate either a label or a rationale (Raffel et al.,



2020). We train the SLM to generate a label when the prefix is “predict: ”, and to generate a rationale when the prefix is “explain: ”. This process is shown as step 3 and step 4 in Figure 2.

### 3.2.1 Distilling self-evaluation capability

Using the self-evaluation data generated by the LLM, we aim to distill this capability into the SLM. During this phase, the model is trained to predict the self-evaluation label  $y_{eval_i}$  as well as generate corresponding rationale  $r_{eval_i}$ . To guide the SLM in learning the self-evaluation outputs for each CoT, we employ a multi-task loss function:

$$L_{SE} = \frac{1}{N_{eval}} \sum_{c=1}^{N_{eval}} \left( \ell(f(x_c), y_{eval_c}) + \lambda \ell(f(x_c), r_{eval_c}) \right),$$

where  $f$  represents the SLM and  $\ell$  is the cross-entropy loss between the tokens predicted by the SLM and the target tokens.  $x_c$  is the CoT that needs to be evaluated.  $\lambda$  is a hyperparameter for weighing the rationale loss.  $y_{eval_c}$  indicates the self-evaluation label generated by the LLM,  $r_{eval_c}$  is the rationale in the  $c^{th}$  self-evaluation output, and  $N_{eval}$  is the total amount of self-evaluation outputs.

### 3.2.2 Distilling CoT reasoning capability

After successfully distilling self-evaluation capability, the focus shifts to leveraging diverse CoTs to train the comprehensive reasoning capability of SLMs. For each instance in the dataset, we also employ a multi-task loss function to guide the SLM in learning CoT reasoning by:

$$L_{CoT} = \frac{1}{N_{CoT}} \sum_{i=1}^{N_{CoT}} \left( \ell(f(x_i), \hat{y}_i) + \lambda \ell(f(x_i), r_{CoT_i}) \right),$$

where  $x_i$  indicates input data,  $\hat{y}_i$  indicates the pseudo-label  $y_i$  generated by the LLM or human-annotated label,  $r_{CoT_i}$  is the rationale in the  $i^{th}$  CoT, and  $N_{CoT}$  is the total amount of CoTs.

This two-pronged training regimen ensures that the SLM is not merely parroting the CoT reasoning but deeply understands introspective self-evaluation and nuanced reasoning, mirroring the powerful cognitive capabilities of the LLM.

## 4 Experiments

**Tasks and datasets** To evaluate our distillation method, we conduct comprehensive experiments

on three tasks: 1) Math Word Problems (MWP) task with the SVMAP dataset (Patel et al., 2021); 2) Commonsense Question Answering (CQA) task with the CQA dataset (Talmor et al., 2019; Rajani et al., 2019); 3) Natural Language Inference (NLI) task with the ANLI dataset (Nie et al., 2020). For dataset samples, we use either human-annotated labels from the dataset or LLM-generated pseudo-labels to explore the effect of human annotation availability on our method.

**Setup** In distillation, we utilize gpt-3.5-turbo as the LLM<sup>1</sup>. We utilize 5-shot CoT prompting to enhance the quality and standardize the formats of the responses generated by the LLM. We follow the CoT prompts from Wei et al. (2022) for the CQA dataset and devise similar prompts for other datasets and self-evaluation. To strike a balance between diversity and cost, in the main experiment, we obtain five CoTs for each training instance and five self-evaluation outputs of each CoT from the gpt-3.5-turbo model and choose the T5-Base model (220M) (Raffel et al., 2020) as the SLM. We provide more experimental details in Appendix A. We also explore the effect of the value of the hyperparameter  $\lambda$  on the results, which are presented in Appendix B. Therefore, we select  $\lambda = 0.5$  as the optimal hyperparameter for our main experiments. In all experiments, we report the mean results and standard deviations over 3 random runs.

### 4.1 Main results

Our results, presented in Table 1, show the advantages of our distillation method, which integrates multiple CoTs and self-evaluation capability into SLMs. The table shows consistent improvement across all tasks with our method over standard and CoT distillation baselines, whether using pseudo-labels or human-annotated labels. In particular, we observe significant leaps in model performance when simultaneously training with five CoTs and their corresponding self-evaluation outputs, reinforcing our hypothesis about the value of incorporating self-evaluation and comprehensive thinking during the distillation process. Moreover, the reduced standard deviation in performance metrics across multiple runs, especially in the “5 CoTs w/ self-evaluation” condition, suggests that our method provides a stable and reliable improvement over the baseline methods. This stability is cru-

<sup>1</sup>Most experiments were conducted in August 2023 using the gpt-3.5-turbo model provided by the OpenAI API.

Method	SVAMP		CQA		ANLI	
	P-labels	H-labels	P-labels	H-labels	P-labels	H-labels
Standard Distillation / Fine-tuning	49.2 ± 1.9	59.3 ± 1.2	58.7 ± 0.4	62.0 ± 0.4	37.7 ± 1.2	42.1 ± 5.0
1 CoT (i.e., CoT distillation)	51.7 ± 2.1	65.0 ± 1.1	59.7 ± 0.4	63.4 ± 0.2	39.8 ± 0.4	48.5 ± 1.2
1 CoT w/ Self-Evaluation	55.5 ± 0.4	67.8 ± 0.6	60.4 ± 0.2	63.7 ± 0.2	41.8 ± 0.4	49.2 ± 0.5
5 CoTs	54.8 ± 1.0	68.7 ± 0.2	61.2 ± 0.4	63.9 ± 0.2	41.7 ± 0.4	49.7 ± 0.8
5 CoTs w/ Self-Evaluation	<b>60.3 ± 0.6</b>	<b>72.7 ± 1.0</b>	<b>61.9 ± 0.3</b>	<b>65.0 ± 0.1</b>	<b>44.3 ± 0.2</b>	<b>50.8 ± 0.4</b>

Table 1: **Results of the main experiment.** We compare the accuracy (mean ± standard deviation, %) of different distillation methods on three different datasets (SVAMP, CQA, and ANLI) using 220M T5-Base models, utilizing pseudo-labels generated by the LLM or human-annotated labels. The P-labels in the table represent pseudo-labels, while the H-labels represent human-annotated labels. Across all datasets and label types, the method we propose consistently outperformed the baselines (standard distillation and CoT distillation), particularly when combining 5 CoTs and self-evaluation.

cial for real-world applications where consistent performance is necessary.

**Effect of label quality** A discernible pattern from the results is the gap in performance between models trained using LLM-generated pseudo-labels and human-annotated labels. Given the typically higher accuracy of human-annotated labels, which are considered the gold standard in supervised learning, this result is expected. However, regardless of the type of training labels used, our method exhibits consistent advantages, suggesting that the benefits of our distillation method are also robust to variations in label quality.

**Robustness across tasks** When considering performance on different tasks, our method’s superiority is consistently evident, although the degree of improvement varies. In tasks such as MWPs (SVAMP dataset) and NLI (ANLI dataset), where reasoning complexity and potential for hallucinatory content are higher, the benefits of our methodology are more pronounced. This suggests that the proposed method effectively mitigates flawed reasoning and hallucinations in complex reasoning scenarios. In tasks like CQA (CQA dataset), where the reasoning processes might be less convoluted, the increments in performance are smaller yet still notable. This showcases the adaptability of our method to different types of reasoning complexity within various NLP tasks.

**Effect of our method on model output** To investigate whether our method mitigates the flawed

reasoning and hallucinations of distilled SLMs, we conduct case studies on three datasets in the setting of using pseudo labels generated by LLMs. We compare the rationales and labels generated by the models trained using our method with those generated by the models trained using the CoT distillation method. The results indicate that our method effectively reduces flawed reasoning and hallucinations produced by distilled SLMs.

In the ANLI dataset case presented in Table 2, the task is to judge the relationship between the premise and hypothesis. The model trained by the baseline CoT distillation method incorrectly infers that the premise entails the hypothesis because superficially the geographic locations mentioned in the two statements match each other. This flawed reasoning likely results from a lack of critical evaluation of the information’s depth and relevance, a pitfall in models trained without a self-evaluation mechanism. Conversely, the model trained by our method identifies the lack of specific information about team members’ residences in the premise and correctly concludes that the premise is neutral to the hypothesis. This accurate judgment showcases our method’s strength in instilling a comprehensive and critical evaluation capability in the model, enabling it to discern the nuances and gaps in information that affect the reasoning. Case studies on other datasets are in Appendix C.

**Summary of main results** In conclusion, the experimental analysis demonstrates that our proposed distillation method, which emphasizes the distilla-

	ANLI Example	Ground Truth Label: neutral
Model Input	Premise: East Bengal Football Club is a professional football club, based in Kolkata, West Bengal, India. It currently competes in ... Hypothesis: All of the team members live in West Bengal.	
Model Output (CoT Distillation)	Rationale: The premise mentions that ... <b>The hypothesis states that all of the team members live in West Bengal. The hypothesis is consistent with the information provided in the premise.</b> Label: entailment	
Model Output (Our Method)	Rationale: The premise mentions that ... <b>However, it does not provide any information about the location of the team members. Therefore, the premise neither supports nor contradicts the hypothesis.</b> Label: neutral	

Table 2: **A case study on the ANLI dataset.** The model trained by the CoT distillation method incorrectly predicts the label as entailment due to the premise location matching the hypothesis statement (highlighted in red), while the model trained by our method correctly identifies the lack of information regarding the team members’ residences and correctly predicts the label as neutral (highlighted in blue). This indicates that our method effectively reduces flawed reasoning and hallucinations produced by distilled SLMs. The complete table is shown in Table 9.

460 tion of self-evaluation capability and the incorpora-  
461 tion of multiple CoTs, significantly improves the  
462 performance and reliability of SLMs. This affirms  
463 our hypothesis is essential for creating robust and  
464 efficient SLMs capable of high-quality reasoning  
465 in resource-constrained environments.

#### 466 4.2 Effect of the number of CoTs

467 Using the SVAMP dataset as an example, we fur-  
468 ther explore the effect of varying the number of  
469 CoTs on our method, where each CoT is accom-  
470 panied by five self-evaluation outputs. As shown  
471 in Figure 3, initially, as the number of CoTs in-  
472 creases from 1 to 5, there is a notable improve-  
473 ment in performance metrics across both pseudo-  
474 labels and human-annotated labels datasets. This  
475 trend underlines the benefit of exposing SLMs to a  
476 broader spectrum of reasoning processes and self-  
477 evaluation outputs, enhancing their ability to nav-  
478 igate complex reasoning landscapes and correct  
479 flawed reasoning. However, diminishing returns  
480 are observed when the number of CoTs exceeds  
481 five. In particular, when the number of CoTs ex-  
482 ceeded 7, performance degradation is observed  
483 using human-annotated labels training. It indi-  
484 cates that while multiple CoTs and self-evaluation  
485 outputs enrich the model’s reasoning capabilities,  
486 there is a threshold beyond which further complex-  
487 ity fails to enhance performance. This could be  
488 attributed to several factors: one possibility is that  
489 the integration of too many CoTs may introduce  
490 noise or conflicting reasoning patterns, thereby dis-

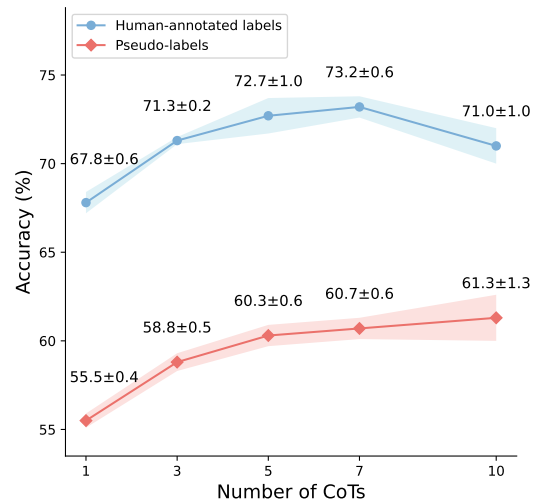


Figure 3: The experimental results of our method using the T5-Base model on the SVAMP dataset for different numbers of CoTs.

491 rupting the distilled SLM. Another factor could be  
492 the cognitive load on the SLM. Beyond a certain  
493 scope, the model may struggle to effectively learn  
494 from additional training data.

495 This observation underscores the importance of  
496 finding an optimal balance in the number of CoTs  
497 used for distillation. As the number of CoTs and  
498 self-evaluation outputs increases, there is a corre-  
499 sponding rise in data costs and training expenses.  
500 Therefore, we opted to use five CoTs in our main  
501 experiments, striking a balance between the cost  
502 and performance.

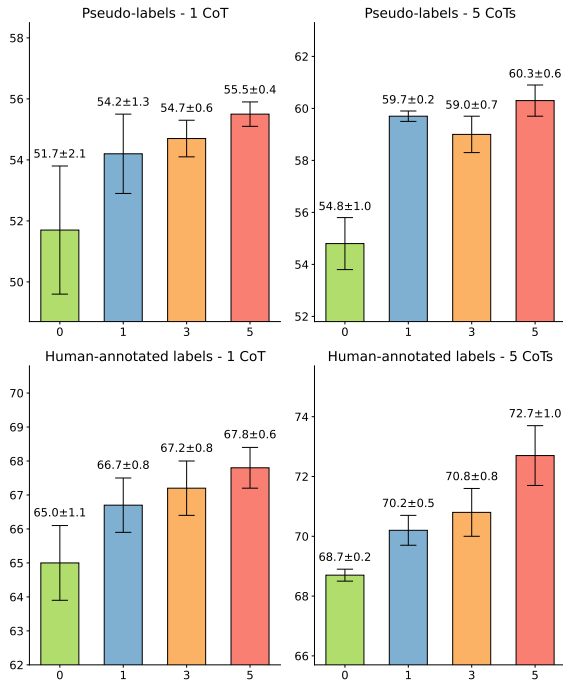


Figure 4: We present experimental results on the SVAMP dataset using the T5-Base model for different numbers of self-evaluation outputs for each CoT. Experiments were conducted separately for distilling a single CoT and five CoTs.

### 4.3 Effect of the number of self-evaluation outputs

Our experiments on the SVAMP dataset further delved into the effect of varying the number of self-evaluation outputs for each CoT. Experiments are conducted for a single CoT and five CoTs respectively, to investigate how the comprehensiveness of self-evaluation affects the performance of our methods. As shown in Figure 4, in both pseudo-labels and human-annotated labels settings, we can observe that: as the number of self-evaluation outputs per CoT increases, there is a notable enhancement in the accuracy of the distilled SLMs, although it may not strictly be monotonically increasing. This indicates that distilling more self-evaluation outputs enables SLMs to produce more accurate and reliable outputs. Notably, accuracy improves more with five CoTs than with a single CoT, underscoring the synergistic effect of combining multiple CoTs with corresponding self-evaluation. Overall, these findings emphasize the importance of incorporating self-evaluation in the distillation. The enhanced performance across different settings confirms the value of introspective self-evaluation in refining the reasoning and predictive capabilities of SLMs. Such introspective capabilities enable

models to refine internal representations, rectifying possible misconceptions or potential pitfalls in their reasoning.

### 4.4 Effect of model size

Method	SVAMP	
	P-labels	H-labels
Standard Distillation / Fine-tuning	60.2 ± 1.5	76.5 ± 1.2
1 CoT (i.e., CoT distillation)	66.2 ± 1.2	77.0 ± 1.2
1 CoT w/ Self-Evaluation	68.0 ± 1.1	79.0 ± 0.4
5 CoTs	66.5 ± 0.7	81.3 ± 0.8
5 CoTs w/ Self-Evaluation	<b>69.3 ± 0.6</b>	<b>83.7 ± 0.6</b>

Table 3: We compare the accuracy (mean ± standard deviation, %) of different distillation methods on the SVAMP dataset using T5-Large models (770M). The P-labels in the table represent pseudo-labels, while the H-labels represent human-annotated labels.

To further elucidate the impact of model size on the effectiveness of our proposed distillation methodology, we extend our experiments to include the larger T5-Large model on the SVAMP dataset. As presented in Table 3, an overarching observation is the T5-Large model consistently outperforms the T5-Base model across all methods and label types. Additionally, the comparison between the baselines (standard distillation and CoT distillation) and our proposed method (particularly 5 CoTs with self-evaluation) indicates that the benefits of our approach are scalable with the model size.

## 5 Conclusion

In this study, we have introduced an innovative method to effectively distill the more comprehensive capabilities from LLMs into SLMs, emphasizing both the transfer of self-evaluation capability and comprehensive thinking, to mitigate the shortcomings of previous CoT distillation methods. Comprehensive experiments demonstrate that our method outperforms prior distillation methods consistently in various NLP tasks. We hope that this study can promote the more effective and efficient utilization of SLMs, especially in resource-limited environments.



## 6 Limitations

Despite the promising results and advancements achieved in our study, certain limitations need acknowledgment and further investigation:

- Limited teacher and student models:** The experiments we conducted primarily utilized a single teacher model, GPT-3.5, and two student models, T5-Base and T5-Large. While these selections were influenced by their current popularity and efficacy, it is crucial to note that the landscape of LLMs and SLMs is rapidly evolving. As such, our distillation method may manifest differently when paired with other architectures or models. Future work will involve testing a wider range of models to confirm the universality of our method.
- Limited tasks:** Although we evaluated our methods on three different NLP tasks, NLP tasks are broad and complex. Therefore, future evaluations of our method’s performance on a wider range of tasks are needed to provide a more comprehensive evaluation of its strengths and potential weaknesses.
- Self-evaluation reliability:** One inherent limitation of the self-evaluation process is its reliance on the LLM’s own capacity for introspection. If the LLM’s self-evaluation mechanism is flawed or biased, it might adversely affect the distilled SLM. In future work, we will investigate the differences in self-evaluation capabilities among different LLMs, such as Llama 2 (Touvron et al., 2023), GPT-3.5, and GPT-4 (OpenAI, 2023), and how these differences affect the performance of distilled SLMs.

In conclusion, while we have made significant strides in advancing the distillation process from LLMs to SLMs, there exists a plethora of avenues for further refinement and exploration. Future endeavors should aim to address these limitations to ensure broader and more robust applicability.

## 7 Ethical Considerations

**Potential risks** While our approach is dedicated to reducing the flaws inherited by SLMs from LLMs, SLMs may still inherit harmful biases and discrimination from LLMs. Therefore, future work

will aim to further minimize the impact of harmful content from LLMs on SLMs.

**The use of closed source LLMs** Many related studies and open source models have already utilized data obtained from the GPT family of models provided by OpenAI. We also obtain CoTs and self-evaluation outputs from the gpt-3.5-turbo model. However, the purpose of this study is not to develop models that compete with general large language models like ChatGPT. Instead, it aims to enhance the effectiveness and efficiency of small language models in resource-constrained environments, promoting the democratization of NLP. We only use gpt-3.5-turbo as the LLM to validate the effectiveness of our method, and it is not required to use the gpt-3.5-turbo model in practical applications, so different LLMs can be employed according to the licenses.

**The use of AI assistants** We employed ChatGPT to assist us in polishing our paper and writing code.

## References

- Rishabh Agarwal, Nino Vieillard, Piotr Stanczyk, Sabela Ramos, Matthieu Geist, and Olivier Bachem. 2023. Gkd: Generalized knowledge distillation for auto-regressive sequence models. *arXiv preprint arXiv:2306.13649*.
- Sina Alemohammad, Josue Casco-Rodriguez, Lorenzo Luzi, Ahmed Imtiaz Humayun, Hossein Babaei, Daniel LeJeune, Ali Siahkoohi, and Richard G Baraniuk. 2023. Self-consuming generative models go mad. *arXiv preprint arXiv:2307.01850*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Hongzhan Chen, Siyue Wu, Xiaojun Quan, Rui Wang, Ming Yan, and Ji Zhang. 2023. **MCC-KD: Multi-CoT consistent knowledge distillation**. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 6805–6820, Singapore. Association for Computational Linguistics.
- 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.
- Yao Fu, Hao-Chun Peng, Litu Ou, Ashish Sabharwal, and Tushar Khot. 2023. Specializing smaller lan-

656	guage models towards multi-step reasoning. In <i>International Conference on Machine Learning</i> .	713
657		714
658	Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang.	715
659	2023. Knowledge distillation of large language mod-	716
660	els. <i>arXiv preprint arXiv:2306.08543</i> .	717
661	Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015.	718
662	Distilling the knowledge in a neural network. <i>arXiv</i>	719
663	<i>preprint arXiv:1503.02531</i> .	720
664	Namgyu Ho, Laura Schmid, and Se-Young Yun. 2023.	721
665	<a href="#">Large language models are reasoning teachers</a> . In	722
666	<i>Proceedings of the 61st Annual Meeting of the As-</i>	723
667	<i>sociation for Computational Linguistics (Volume 1:</i>	724
668	<i>Long Papers)</i> , pages 14852–14882, Toronto, Canada.	725
669	Association for Computational Linguistics.	
670	Jordan Hoffmann, Sebastian Borgeaud, Arthur Men-	726
671	sch, Elena Buchatskaya, Trevor Cai, Eliza Ruther-	727
672	ford, Diego de Las Casas, Lisa Anne Hendricks,	728
673	Johannes Welbl, Aidan Clark, et al. 2022. Train-	729
674	ing compute-optimal large language models. <i>arXiv</i>	730
675	<i>preprint arXiv:2203.15556</i> .	731
676	Cheng-Yu Hsieh, Chun-Liang Li, Chih-kuan Yeh,	732
677	Hootan Nakhost, Yasuhisa Fujii, Alex Ratner, Ranjay	733
678	Krishna, Chen-Yu Lee, and Tomas Pfister. 2023. <a href="#">Dis-</a>	734
679	<a href="#">tilling step-by-step! outperforming larger language</a>	735
680	<a href="#">models with less training data and smaller model</a>	736
681	<a href="#">sizes</a> . In <i>Findings of the Association for Computa-</i>	737
682	<i>tional Linguistics: ACL 2023</i> , pages 8003–8017,	738
683	Toronto, Canada. Association for Computational Lin-	739
684	guistics.	740
685	Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu,	741
686	Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2022.	742
687	Large language models can self-improve. <i>arXiv</i>	743
688	<i>preprint arXiv:2210.11610</i> .	744
689	Haoming Jiang, Pengcheng He, Weizhu Chen, Xi-	745
690	aodong Liu, Jianfeng Gao, and Tuo Zhao. 2020.	746
691	<a href="#">SMART: Robust and efficient fine-tuning for pre-</a>	747
692	<a href="#">trained natural language models through principled</a>	
693	<a href="#">regularized optimization</a> . In <i>Proceedings of the 58th</i>	
694	<i>Annual Meeting of the Association for Computational</i>	
695	<i>Linguistics</i> , pages 2177–2190, Online. Association	
696	for Computational Linguistics.	
697	Saurav Kadavath, Tom Conerly, Amanda Askell, Tom	750
698	Henighan, Dawn Drain, Ethan Perez, Nicholas	751
699	Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli	752
700	Tran-Johnson, et al. 2022. Language models	753
701	(mostly) know what they know. <i>arXiv preprint</i>	754
702	<i>arXiv:2207.05221</i> .	
703	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B	755
704	Brown, Benjamin Chess, Rewon Child, Scott Gray,	756
705	Alec Radford, Jeffrey Wu, and Dario Amodei. 2020.	757
706	Scaling laws for neural language models. <i>arXiv</i>	758
707	<i>preprint arXiv:2001.08361</i> .	759
708	Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie	760
709	Lu, Thomas Mesnard, Colton Bishop, Victor Car-	761
710	bune, and Abhinav Rastogi. 2023. Rlaif: Scaling	
711	reinforcement learning from human feedback with ai	
712	feedback. <i>arXiv preprint arXiv:2309.00267</i> .	
	Liunian Harold Li, Jack Hessel, Youngjae Yu, Xiang	713
	Ren, Kai-Wei Chang, and Yejin Choi. 2023. <a href="#">Sym-</a>	714
	<a href="#">bolic chain-of-thought distillation: Small models can</a>	715
	<a href="#">also “think” step-by-step</a> . In <i>Proceedings of the 61st</i>	716
	<i>Annual Meeting of the Association for Computational</i>	717
	<i>Linguistics (Volume 1: Long Papers)</i> , pages 2665–	718
	2679, Toronto, Canada. Association for Computa-	719
	tional Linguistics.	720
	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	721
	Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,	722
	Nouha Dziri, Shrimai Prabhumoye, Yiming Yang,	723
	et al. 2023. Self-refine: Iterative refinement with	724
	self-feedback. <i>arXiv preprint arXiv:2303.17651</i> .	725
	Lucie Charlotte Magister, Jonathan Mallinson, Jakub	726
	Adamek, Eric Malmi, and Aliaksei Severyn. 2023.	727
	<a href="#">Teaching small language models to reason</a> . In <i>Pro-</i>	728
	<i>ceedings of the 61st Annual Meeting of the Associa-</i>	729
	<i>tion for Computational Linguistics (Volume 2: Short</i>	730
	<i>Papers)</i> , pages 1773–1781, Toronto, Canada. Associ-	731
	ation for Computational Linguistics.	732
	Potsawee Manakul, Adian Liusie, and Mark JF Gales.	733
	2023. Selfcheckgpt: Zero-resource black-box hal-	734
	lucination detection for generative large language	735
	models. <i>arXiv preprint arXiv:2303.08896</i> .	736
	Ning Miao, Yee Whye Teh, and Tom Rainforth.	737
	2023. Selfcheck: Using llms to zero-shot check	738
	their own step-by-step reasoning. <i>arXiv preprint</i>	739
	<i>arXiv:2308.00436</i> .	740
	Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal,	741
	Jason Weston, and Douwe Kiela. 2020. <a href="#">Adversarial</a>	742
	<a href="#">NLI: A new benchmark for natural language under-</a>	743
	<a href="#">standing</a> . In <i>Proceedings of the 58th Annual Meet-</i>	744
	<i>ing of the Association for Computational Linguistics</i> ,	745
	pages 4885–4901, Online. Association for Computa-	746
	tional Linguistics.	747
	OpenAI. 2023. Gpt-4 technical report. <i>arXiv preprint</i>	748
	<i>arXiv:2303.08774</i> .	749
	Liangming Pan, Michael Saxon, Wenda Xu, Deepak	750
	Nathani, Xinyi Wang, and William Yang Wang. 2023.	751
	Automatically correcting large language models: Sur-	752
	veying the landscape of diverse self-correction strate-	753
	gies. <i>arXiv preprint arXiv:2308.03188</i> .	754
	Arkil Patel, Satwik Bhattamishra, and Navin Goyal.	755
	2021. <a href="#">Are NLP models really able to solve simple</a>	756
	<a href="#">math word problems?</a> In <i>Proceedings of the 2021</i>	757
	<i>Conference of the North American Chapter of the</i>	758
	<i>Association for Computational Linguistics: Human</i>	759
	<i>Language Technologies</i> , pages 2080–2094, Online.	760
	Association for Computational Linguistics.	761
	Debjit Paul, Mete Ismayilzada, Maxime Peyrard, Beat-	762
	riz Borges, Antoine Bosselut, Robert West, and	763
	Boi Faltings. 2023. Refiner: Reasoning feedback	764
	on intermediate representations. <i>arXiv preprint</i>	765
	<i>arXiv:2304.01904</i> .	766
	David L Poole and Alan K Mackworth. 2010. <i>Artificial</i>	767
	<i>Intelligence: foundations of computational agents</i> .	768
	Cambridge University Press, Cambridge, UK.	769

770	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	828
771	Dario Amodei, Ilya Sutskever, et al. 2019. Language	pages 5546–5558, Toronto, Canada. Association for	829
772	models are unsupervised multitask learners. <i>OpenAI</i>	Computational Linguistics.	830
773	<i>blog</i> , 1(8):9.		
774	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V.	831
775	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	Le, Ed H. Chi, Sharan Narang, Aakanksha Chowd-	832
776	Wei Li, and Peter J Liu. 2020. Exploring the limits	hery, and Denny Zhou. 2023b. Self-consistency im-	833
777	of transfer learning with a unified text-to-text trans-	proves chain of thought reasoning in language mod-	834
778	former. <i>The Journal of Machine Learning Research</i> ,	els. In <i>Proceedings of the 11th International Confer-</i>	835
779	21(1):5485–5551.	<i>ence on Learning Representations, ICLR’2023</i> .	836
780	Nazneen Fatema Rajani, Bryan McCann, Caiming	Zhaoyang Wang, Shaohan Huang, Yuxuan Liu, Jiahai	837
781	Xiong, and Richard Socher. 2019. <a href="#">Explain your-</a>	Wang, Minghui Song, Zihan Zhang, Haizhen Huang,	838
782	<a href="#">self! leveraging language models for commonsense</a>	Furu Wei, Weiwei Deng, Feng Sun, and Qi Zhang.	839
783	<a href="#">reasoning</a> . In <i>Proceedings of the 57th Annual Meet-</i>	2023c. <a href="#">Democratizing reasoning ability: Tailored</a>	840
784	<i>ing of the Association for Computational Linguistics</i> ,	<a href="#">learning from large language model</a> . In <i>Proceed-</i>	841
785	pages 4932–4942, Florence, Italy. Association for	<i>ings of the 2023 Conference on Empirical Methods</i>	842
786	Computational Linguistics.	<i>in Natural Language Processing</i> , pages 1948–1966,	843
787	Jianhao Shen, Yichun Yin, Lin Li, Lifeng Shang, Xin	Singapore. Association for Computational Linguis-	844
788	Jiang, Ming Zhang, and Qun Liu. 2021. <a href="#">Generate &amp;</a>	tics.	845
789	<a href="#">rank: A multi-task framework for math word prob-</a>	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	846
790	<a href="#">lems</a> . In <i>Findings of the Association for Computa-</i>	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	847
791	<i>tional Linguistics: EMNLP 2021</i> , pages 2269–2279,	et al. 2022. Chain-of-thought prompting elicits rea-	848
792	Punta Cana, Dominican Republic. Association for	soning in large language models. <i>Advances in Neural</i>	849
793	Computational Linguistics.	<i>Information Processing Systems</i> , 35:24824–24837.	850
794	Noah Shinn, Beck Labash, and Ashwin Gopinath.	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,	851
795	2023. Reflexion: an autonomous agent with dyn-	Thomas L Griffiths, Yuan Cao, and Karthik	852
796	amic memory and self-reflection. <i>arXiv preprint</i>	Narasimhan. 2023. Tree of thoughts: Deliberate	853
797	<i>arXiv:2303.11366</i> .	problem solving with large language models. <i>arXiv</i>	854
798	Kumar Shridhar, Alessandro Stolfo, and Mrinmaya	<i>preprint arXiv:2305.10601</i> .	855
799	Sachan. 2023. <a href="#">Distilling reasoning capabilities into</a>	Muru Zhang, Ofir Press, William Merrill, Alisa	856
800	<a href="#">smaller language models</a> . In <i>Findings of the Asso-</i>	Liu, and Noah A Smith. 2023. How language	857
801	<i>ciation for Computational Linguistics: ACL 2023</i> ,	model hallucinations can snowball. <i>arXiv preprint</i>	858
802	pages 7059–7073, Toronto, Canada. Association for	<i>arXiv:2305.13534</i> .	859
803	Computational Linguistics.	Susan Zhang, Stephen Roller, Naman Goyal, Mikel	860
804	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and	Artetxe, Moya Chen, Shuohui Chen, Christopher De-	861
805	Jonathan Berant. 2019. <a href="#">CommonsenseQA: A ques-</a>	wan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022.	862
806	<a href="#">tion answering challenge targeting commonsense</a>	Opt: Open pre-trained transformer language models.	863
807	<a href="#">knowledge</a> . In <i>Proceedings of the 2019 Conference</i>	<i>arXiv preprint arXiv:2205.01068</i> .	864
808	<i>of the North American Chapter of the Association for</i>	Xunyu Zhu, Jian Li, Yong Liu, Can Ma, and Weip-	865
809	<i>Computational Linguistics: Human Language Tech-</i>	ing Wang. 2023. A survey on model compres-	866
810	<i>nologies, Volume 1 (Long and Short Papers)</i> , pages	sion for large language models. <i>arXiv preprint</i>	867
811	4149–4158, Minneapolis, Minnesota. Association for	<i>arXiv:2308.07633</i> .	868
812	Computational Linguistics.		
813	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-		
814	bert, Amjad Almahairi, Yasmine Babaei, Nikolay		
815	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti		
816	Bhosale, et al. 2023. Llama 2: Open founda-		
817	tion and fine-tuned chat models. <i>arXiv preprint</i>		
818	<i>arXiv:2307.09288</i> .		
819	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob		
820	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz		
821	Kaiser, and Illia Polosukhin. 2017. Attention is all		
822	you need. <i>Neural Information Processing Systems</i> ,		
823	30.		
824	Peifeng Wang, Zhengyang Wang, Zheng Li, Yifan Gao,		
825	Bing Yin, and Xiang Ren. 2023a. <a href="#">SCOTT: Self-</a>		
826	<a href="#">consistent chain-of-thought distillation</a> . In <i>Proceed-</i>		
827	<i>ings of the 61st Annual Meeting of the Association for</i>		

## A Experimental Details

**Datasets** The dataset statistics are shown in Table 4. Following Hsieh et al. (2023), for the SVAMP dataset, 20% of the original data is used as the test set. For the CQA dataset, the original validation set is used as the test set. Then, for both datasets, 10% of the data from the original training set is sampled to serve as the validation set. The ANLI dataset follows the original split. The language of all datasets is English. To the best of our knowledge, all datasets used have been widely employed in NLP research and do not contain any information that names or uniquely identifies individual people or offensive content.

Dataset	Train	Validation	Test
SVAMP	720	80	200
CQA	8,766	975	1,221
ANLI	16,946	1,000	1,000

Table 4: Dataset statistics.

**LLM performance** In Table 5, we report the accuracy of LLM (gpt-3.5-turbo) on three datasets in our experiments, including accuracy on the training set (i.e., the accuracy of pseudo-labels used for training SLMs) and accuracy on the test set.

Dataset	SVAMP	CQA	ANLI
Training Set	85.6	69.1	68.6
Test Set	84.3	72.4	55.1

Table 5: The accuracy (%) of LLM (gpt-3.5-turbo).

**Models & Training** The T5-Base<sup>2</sup> (220M) and T5-Large<sup>3</sup> (770M) models are initialized with pre-trained weights obtained from Hugging Face, and the hyperparameter settings for their training are shown in Table 6. We perform the main experiments on 4 A100 GPUs.

## B Effect of the hyperparameter $\lambda$

As shown in Figure 5, our experiments reveal trends regarding the effect of the hyperparameter  $\lambda$  on the accuracy of the small language models (SLMs) trained using both pseudo-labels and human-annotated labels.

<sup>2</sup>[https://huggingface.co/google/t5-v1\\_1-base](https://huggingface.co/google/t5-v1_1-base)

<sup>3</sup>[https://huggingface.co/google/t5-v1\\_1-large](https://huggingface.co/google/t5-v1_1-large)

Hyperparameter	T5-Base	T5-Large
Total Batch Size	64	32
Learning Rate	$5 \times 10^{-5}$	$5 \times 10^{-5}$
Max Input Length	1,024	1,024
Maximum Steps (for SVAMP)	4,000	9,000
Maximum Steps (for CQA & ANLI)	12,000	-

Table 6: Training hyperparameter settings.

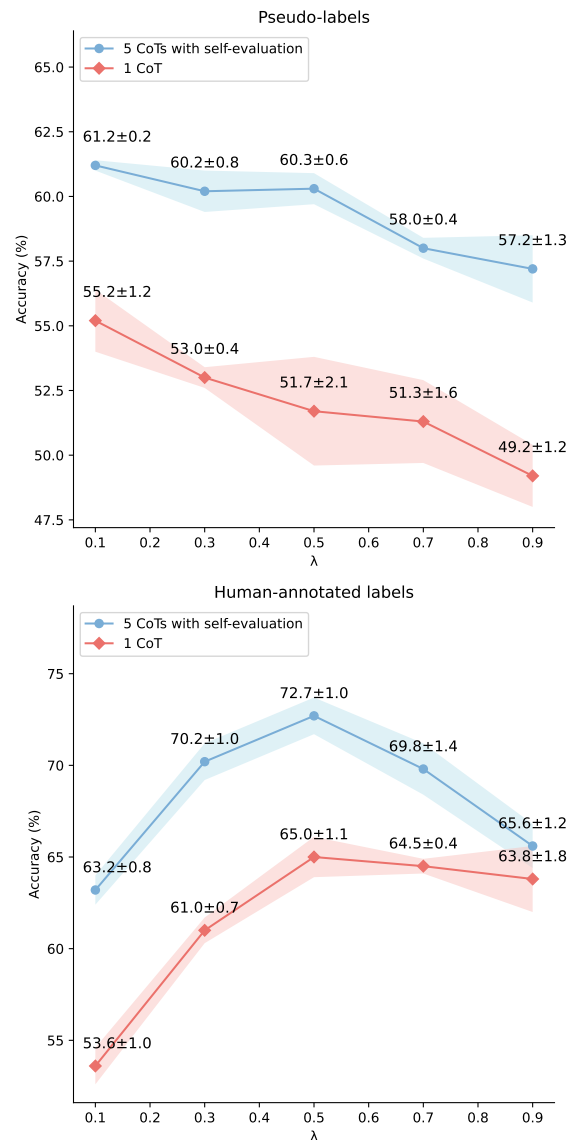


Figure 5: We present experimental results of distillation using the T5-Base model on the SVAMP dataset with different  $\lambda$  values for “1 CoT” and “5 CoTs with self-evaluation” respectively.

For pseudo-labels, the average accuracy with a single CoT is relatively low and slightly declines with an increase in  $\lambda$ . In the case of combining

900  
901  
902



self-evaluation with 5 CoTs, there is a significant improvement in accuracy. Moreover, although accuracy slightly decreases as  $\lambda$  increases, the decline is less when  $\lambda \leq 0.5$ . This trend indicates that the incorporation of 5 CoTs and self-evaluation significantly enhances accuracy in tasks based on pseudo-labels, although an excessive focus on rationale may not be conducive to improved accuracy.

Contrastingly, in the case of human-annotated labels, we observe a different trend. The accuracy initially increases with  $\lambda$ , peaking at  $\lambda = 0.5$ , and then begins to decline. This pattern underscores a critical observation: up to a certain point ( $\lambda \leq 0.5$ ), increasing the weight on rationale positively impacts the model’s ability to predict labels in human-annotated data. However, beyond this optimal point, further emphasis on rationale seems to divert too much focus from the primary task, leading to a decrease in label prediction accuracy.

Based on these observations, we select  $\lambda = 0.5$  as the optimal hyperparameter for our main experiments, as it shows good accuracy in both types of data (pseudo-labels and human-annotated labels) and strikes a balance between fostering a deep understanding of the reasoning process through rationale and maintaining high accuracy in label prediction.

### C Case Study

The detailed case studies presented in Tables 2, 7, and 8 provide insightful examples demonstrating the effectiveness of our methodology compared to the baseline CoT distillation method. These cases highlight the importance of incorporating both self-evaluation and comprehensive thinking in the distillation process, which significantly reduces reasoning errors and hallucinations in small language models (SLMs).

In the SVAMP example (Table 7), the model trained by the baseline CoT distillation method exhibits flawed reasoning in its calculation, erroneously summing the hours for learning Chinese and Spanish only, resulting in an incorrect total. This illustrates a common issue with CoT distillation, where the model may focus on a part of the problem, leading to incomplete reasoning. In stark contrast, the model trained by our method correctly identifies and sums the hours for all three languages, demonstrating a more comprehensive understanding and accurate reasoning process. This accurate reasoning underscores the effectiveness

of our method, which incorporates both multiple CoTs and self-evaluation capability. By exposing the model to diverse reasoning processes and enabling it to evaluate its reasoning, our method equips the model to consider all relevant information comprehensively and to avoid flawed reasoning paths.

Similarly, in the CQA example (Table 8), the model trained by the baseline CoT distillation method incorrectly concludes that the most logical result of dying is a change of color, showcasing a clear case of flawed reasoning or hallucination. This error is likely due to a superficial association between the concepts of dying and color change, without a deeper understanding of the context of organic material decay. The model trained by our method, on the other hand, correctly identifies “death and decay” as the logical result of dying in the context of organic material, reflecting a deeper and more accurate comprehension of the question’s essence. The comprehensive thinking instilled by our method, coupled with the ability to critically evaluate its reasoning, enables the model to select the most logical answer from the provided choices.

These case studies unequivocally demonstrate the effectiveness of our method in mitigating reasoning flaws and hallucinations often observed in SLMs. By distilling multiple CoTs and self-evaluation outputs from LLMs, we enable SLMs to engage in more comprehensive and critical reasoning.

	SVAMP Example	Ground Truth Label: $((6.0 + 7.0) + 3.0)$
Model Input	Every day Ryan spends 6 hours on learning English 7 hours on learning Chinese and 3 hours on learning Spanish. How many hours does he spend on learning English, Chinese, and Spanish in all?	
Model Output (CoT Distillation)	Rationale: Ryan spends 6 hours on learning English and 7 hours on learning Chinese. <b>To find out how many hours he spends on learning English, we need to add the number of hours spent on learning Chinese and Spanish together: <math>7 + 3 = 6</math>.</b> Label: $(7 + 3)$	
Model Output (Our Method)	Rationale: Ryan spends 6 hours on learning English, 7 hours on learning Chinese, and 3 hours on learning Spanish in total. <b>To find the total number of hours he spends on learning English, Chinese, and Spanish in all, we need to add the number of hours spent on learning English, Chinese, and Spanish together. So, the equation is <math>6 + 7 + 3</math>.</b> Label: $(6 + 7) + 3$	

Table 7: **A case study on the SVAMP dataset.** The model trained by the CoT distillation method incorrectly calculates the total time to learn the three languages (highlighted in red), while the model trained by our method correctly sums the time to learn the three languages (highlighted in blue).

	CQA Example	Ground Truth Label: death and decay
Model Input	Question: As a result of dying, what happens to organic material? Choices: [change of color, stop breathing, wake up, death and decay, getting cold]	
Model Output (CoT Distillation)	Rationale: The question is asking about the result of dying. Of the given choices, <b>the most logical result of dying is a change of color.</b> Label: change of color	
Model Output (Our Method)	Rationale: The question is asking about the result of dying and the process of organic material. Of the given choices, <b>the most logical result of dying is death and decay, as organic material undergoes the process of decay.</b> Label: death and decay	

Table 8: **A case study on the CQA dataset.** The model trained by the CoT distillation method incorrectly considers that the most logical result of dying is a change of color (highlighted in red), while the model trained by our method correctly identifies the most logical result as death and decay (highlighted in blue).

	ANLI Example	Ground Truth Label: neutral
Model Input	<p>Premise: East Bengal Football Club is a professional football club, based in Kolkata, West Bengal, India. It currently competes in the I-League, the top tier of Indian football. The club has won three National Football League (India) (later named as I league) titles, eight Federation Cups, and three Indian Super Cups, among others. The club is the current holder of the Calcutta Football League.</p> <p>Hypothesis: All of the team members live in West Bengal.</p>	
Model Output (CoT Distillation)	<p>Rationale: The premise mentions that East Bengal Football Club is based in Kolkata, West Bengal, India. <b>The hypothesis states that all of the team members live in West Bengal. The hypothesis is consistent with the information provided in the premise.</b></p> <p>Label: entailment</p>	
Model Output (Our Method)	<p>Rationale: The premise mentions that East Bengal Football Club is based in Kolkata, West Bengal, India. <b>However, it does not provide any information about the location of the team members. Therefore, the premise neither supports nor contradicts the hypothesis.</b></p> <p>Label: neutral</p>	

Table 9: **A case study on the ANLI dataset.** The full version of Table 2.