

Learning to Maximize Mutual Information for Chain-of-Thought Distillation

Anonymous ACL submission

Abstract

Knowledge distillation, the technique of transferring knowledge from large, complex models to smaller ones, marks a pivotal step towards efficient AI deployment. Distilling Step-by-Step (DSS), a novel method utilizing chain-of-thought (CoT) distillation, has demonstrated promise by imbuing smaller models with the superior reasoning capabilities of their larger counterparts. In DSS, the distilled model acquires the ability to generate rationales and predict labels concurrently through a multi-task learning framework. However, DSS overlooks the intrinsic relationship between the two training tasks, leading to ineffective integration of CoT knowledge with the task of label prediction. To this end, we investigate the mutual relationship of the two tasks from Information Bottleneck perspective and formulate it as maximizing the mutual information of the representation features of the two tasks. We propose a variational approach to solve this optimization problem using a learning-based method. Our experimental results across four datasets demonstrate that our method outperforms the state-of-the-art DSS. Our findings offer insightful guidance for future research on language model distillation as well as applications involving CoT. Code and models will be released soon.

1 Introduction

The capabilities of larger language models (LLMs) tend to scale with their model size, leading to a substantial demand for memory and compute resources (Chowdhery et al., 2023; Wei et al., 2022a). Distilling knowledge from larger LLMs to smaller LLMs has been crucial for the efficient deployment of AI (Hinton et al., 2015; Phuong and Lampert, 2019). Chain-of-Thought (CoT) (Wei et al., 2022b) distillation represents a pivotal advance in the quest to endow smaller language models with the sophisticated reasoning capabilities of their larger counterparts. By distilling complex thought processes

into more compact models, this approach aims to democratize access to advanced natural language understanding and reasoning across a wider array of computational resources (Ma et al., 2023; Magister et al., 2023; Li et al., 2023).

Distilling Step-by-Step (DSS) (Hsieh et al., 2023) introduces a CoT distillation method that guides smaller models using rationales from LLMs within a multi-task learning (MTL) framework, training them for both label prediction and rationale generation tasks. While DSS brings out the benefits of reducing computational costs, it suffers from the same problem as the conventional MTL framework, that is the difficulty in effectively connecting the prediction and generation tasks. The intricacies inherent in training models within the MTL framework can undermine the effectiveness and reliability of the DSS process (Wang et al., 2023b). Despite the successful setup of an MTL framework in DSS, where the tasks of label prediction and rationale generation are intrinsically related, the current configuration may not optimally capture and maximize the mutual knowledge between these tasks. Furthermore, LLMs are prone to producing hallucinations and inconsistent rationales, which potentially mislead the student model toward incorrect answers and cause conflicts in MTL that destruct student model learning (Mueller et al., 2022).

To address this issue, we model the DSS using information bottleneck (IB) and investigate it from an information-theoretic viewpoint (Tishby and Zaslavsky, 2015). Subsequently, we formulate the DSS as an optimization problem to maximize mutual information (MI) of label prediction and rationale generation tasks. However, estimating MI from finite data is a known difficult problem (McAllester and Stratos, 2020; Belghazi et al., 2018). In this study, we introduce a variational method to estimate the MI. Recently, knowledge distillation has been used for solving the challenge

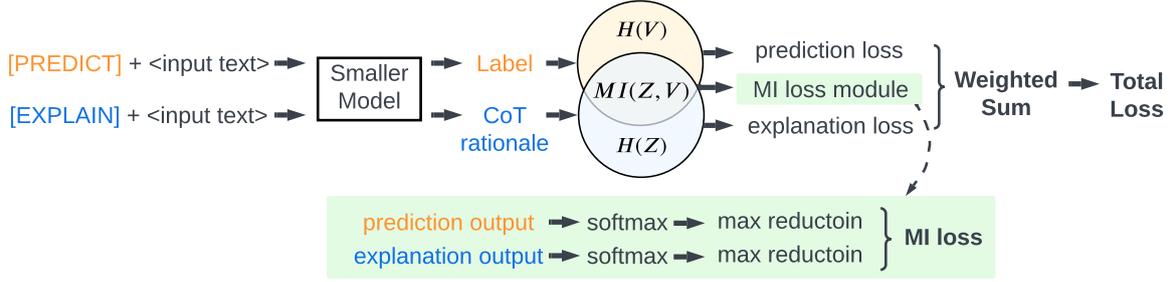


Figure 1: Overview of our approach: CoT distillation from an IB perspective and measurement of the intrinsic relationship between the two tasks by MI. The DSS is an MTL framework pipeline comprising label prediction and rationale generation tasks. H represents the entropy of representation features V and Z . Besides prediction loss and explanation losses used in conventional DSS, we design an auxiliary loss module to maximize MI between the two representation features. This process enhances CoT reasoning capacity through knowledge distillation.

of preserving task-specific training in MTL (Li and Bilen, 2020). Consequently, we propose a practical yet effective auxiliary loss to quantify the shared information between the prediction and the generation tasks, thereby enhancing the alignment between the two tasks and facilitating the knowledge transfer from CoT. We conduct comprehensive experiments with two smaller types of T5 models (Raffel et al., 2020), T5-base (220M) and T5-small (60M), on four popular datasets. Furthermore, we provide detailed analysis in Section 5. Our main contributions are summarized below:

- We reframe the MTL framework of DSS as a MI estimation challenge, aiming to maximize the MI between label prediction and rationale generation tasks. To achieve this, we introduce a variational approach grounded in the IB principle for effective MI estimation. To the best of our knowledge, we present the first work of improving CoT distillation from an IB perspective.
- Beyond establishing a theoretical foundation, we present a practical approach for MI estimation, incorporating a simple yet effective auxiliary loss to learning to maximize MI and enhance DSS.
- Our methodology demonstrably outperforms existing benchmarks across multiple datasets, evidencing the efficacy of our approach in enhancing the reasoning capabilities of distilled models.
- We conduct a systematic review of the relationship between predictive and explainable tasks under MTL training, presenting both qualitative and quantitative analysis results.

With the theoretical proofs and experiment results, we aim to provide stepping stones for future research on enhancing CoT distillation through an

effective MTL framework, guided by principles from information theory.

2 Related Work

We present an overview of previous work across three areas related to our study: knowledge distillation, multi-task learning and information bottleneck.

Knowledge Distillation (KD) While originally designed for training small models by leveraging the extensive knowledge of larger models (Hinton et al., 2015), KL has been extended to a variety of applications due to its effective transfer of knowledge between models and tasks (Chen et al., 2021; Wang and Yoon, 2021; Sanh et al., 2019; Jiao et al., 2020). An crucial yet open challenge is how to effectively transfer the knowledge. To address the issue, previous studies (Zhang et al., 2022b; Allen-Zhu and Li, 2023; Zhang et al., 2021) extract different features and design auxiliary loss functions to enhance KL. Our work focus on improving the model by acquiring mutual knowledge in addressing both label prediction and rationale generation tasks.

Multi-Task Learning (MTL) Through exploiting commonalities and differences of relevant tasks, MTL can enhance improve learning efficiency and prediction accuracy by learning multiple objectives from a share representation (Caruana, 1997; Zhang and Yang, 2021). In recent years, MLT has been broadly applied to NLP tasks (Worsham and Kalita, 2020; Zhang et al., 2023; Liu et al., 2019). However, some works figure out that multiple tasks trained simultaneously could conflict with each other and it is challenging to optimize the performance of all tasks simultaneously (Kendall et al., 2018; Lin et al., 2019). In recent years, KD has also

been applied within MTL frameworks, achieving state-of-the-art results in various applications (Li and Bilen, 2020; Xu et al., 2023; Yang et al., 2022). **Information Bottleneck (IB)** (Tishby and Zaslavsky, 2015; Slonim, 2002) provides a powerful statistical tool to learn representation to preserve complex intrinsic correlation structures over high dimensional data. As a general measure of the dependence between two random variables, MI is also widely used in deep learning to effectively represent the dependencies of features (Cover, 1999; Covert et al., 2023; Liu et al., 2009). MI estimation is known to be difficult, and recent progress has been made towards learning-based variational approaches (Tian et al., 2020; Covert et al., 2023; Bachman et al., 2019; Tschannen et al., 2019; Belghazi et al., 2018). Another challenge associated with the IB principle is the optimization process, which involves a trade-off between achieving a concise representation and maintaining strong predictive capabilities (Alemi et al., 2016; Wang et al., 2019). Consequently, optimizing IB becomes a complex task that heavily depends on the formulation of the problem and the provision of an effective optimization solution. Recent studies has applied IB to solve complex machine learning problems both in computer vision (Tian et al., 2021; Du et al., 2020; Wan et al., 2021) and NLP (Chen and Ji, 2020; Zhang et al., 2022a; Paranjape et al., 2020). In this paper, we formulate our CoT distillation problem with MTL training pipeline using IB and provide a learning-based solution to IB optimization for our CoT distillation problem, as detailed in Section 3.

3 Methodology

This section starts with preliminaries of IB. Following it, we formulate our CoT distillation idea within the IB and propose a learning approach to optimize MI.

3.1 Preliminaries

Under the DSS framework, a task prefix [PREDICT] and [EXPLAIN] will be pretended to the input text TEXT, corresponding to the label prediction and rationale generation task, respectively. In the label prediction task, given the input [PREDICT] + TEXT, and the predictive labels \mathbf{Y} , a representation feature \mathbf{V} , $\mathbf{V} \in \mathbb{R}^d$, is trained under \mathbf{Y} . In the rationale generation task, given the input [EXPLAIN] + TEXT and rationale label \mathbf{R} , a representation fea-

ture \mathbf{Z} , $\mathbf{Z} \in \mathbb{R}^d$, is trained under \mathbf{R} .

Considered the limited CoT capacity of smaller LLMs that are less than 10B parameters, our goal is to distill CoT knowledge of larger LLMs to Z which is also maximally informative to Y . In order to achieve this goal, on the basis of IB (Tishby and Zaslavsky, 2015; Zhang et al., 2022a; Wang et al., 2019), we can model the DSS as following:

$$I(Z; Y) = \int p(z, y) \log \frac{p(z, y)}{p(z)p(y)} dz dy. \quad (1)$$

where sampling observations $z \sim \mathbf{Z}$ and $v \sim \mathbf{V}$. $p(\cdot)$ is probability distribution.

To encourage CoT distillation to focus on the information represented in label \mathbf{Y} , we propose IB to enforce an upper bound I_c to the information flow from the representation features V to the representation features Z , achieved by maximizing the following objective:

$$\max I(Z; Y) \quad s.t. \quad I(Z; V) \leq I_c. \quad (2)$$

By using Lagrangian objective, IB allows Z to be maximally expressive about Y while being maximally compressive about input data by:

$$\mathcal{C}_{IB} = I(Z; V) - \beta I(Z; Y) \quad (3)$$

where β is the Lagrange multiplier. It is obvious that Eq. 3 is the trade-off optimization between high mutual information and high compression (Zhang et al., 2022a; Alemi et al., 2016). In our scenario, given a known small student model, the compression ratio is fixed. Therefore, we formulate the CoT distillation is an optimization problem as:

$$\max I(Z; V) \quad (4)$$

Due to symmetric property of MI, $I(Z; V) = I(V; Z)$. COT distillation can also enhance rationale generation task with the label knowledge. This is validated in Section 5.

3.2 Variational Bounds of MI

We rewrite MI $I(Z; V)$ of Equation 4 as:

$$I(Z; V) = \mathbb{E}_{p(z,v)} \left[\log \frac{p(v|z)}{p(v)} \right] \quad (5)$$

According to (Poole et al., 2019; Covert et al., 2023), a tractable variational upper bound can be built by introducing a variational approximation

245 $q(v)$ to the intractable marginal $p(v)$, thus:

$$\begin{aligned}
I(Z; V) &= \mathbb{E}_{p(z,v)} \left[\log \frac{p(v|z)q(v)}{p(v)q(v)} \right] \\
&= \mathbb{E}_{p(z,v)} \left[\log \frac{p(v|z)}{q(v)} \right] \\
&\quad - KL(p(v)||q(v))
\end{aligned}
\tag{6}$$

247 where $KL[\cdot||\cdot]$ denotes Kullback-Leibler diver-
248 gence. The bound is tight when $q(v) = p(v)$. Then
249 $KL(p(v)||q(v)) = KL(p(v)||p(v))$, and it is a
250 constant. Therefore, we can obtain the following
251 inequality:

$$I(Z; V) \leq \mathbb{E}_{p(z,v)} \left[\log \frac{p(z|v)}{p(v)} \right] \tag{7}$$

252 Then we can write the MI as the following:

$$\begin{aligned}
\mathbb{E}_{p(z,v)} \left[\log \frac{p(z|v)}{p(v)} \right] &= \sum p(z, v) \log \frac{p(v|z)}{p(v)} \\
&= \sum p(z|v)p(v) \log p(v|z) \\
&\quad - \sum p(v)p(z|v) \log p(v)
\end{aligned}
\tag{8}$$

255 Assuming that $p(v)$ is uniform distribution,
256 $\sum p(v)p(z|v) \log p(v)$ is a constant. Then maxi-
257 mizing the $I(Z; V)$ can defined as:

$$\begin{aligned}
\max I(Z; V) &\propto \max \sum p(z|v) \log p(v|z) \\
&= \max \left(- \sum p(z|v) \log \frac{1}{p(v|z)} \right) \\
&= \min \left(\sum p(z|v) \log \frac{1}{p(v|z)} \right) \\
&= \min \left(\sum CE(z|v, v|z) \right)
\end{aligned}
\tag{9}$$

259 here CE represents cross entropy. Therefore,
260 based on this equation, we propose a new MI loss,
261 which minimizes cross entropy to learn to maxi-
262 mize MI for CoT Distillation.

263 3.3 Training Loss

264 The training loss is given by

$$\mathcal{L}_{total} = \alpha_1 \mathcal{L}_{prediction} + \alpha_2 \mathcal{L}_{generation} + \alpha_3 \mathcal{L}_{CE} \tag{10}$$

265 where α_1 , α_2 and α_3 are regularization param-
266 eters, all of which are non-negative. $\mathcal{L}_{prediction}$
267 represents the loss of the label prediction task, and
268 $\mathcal{L}_{generation}$ represents the loss of the rationale gen-
269 eration task. Both are general cross-entropy loss as
270 defined in (Hsieh et al., 2023).
271

272 According to the last line of Equation 9, we
273 define the our MI loss as

$$\mathcal{L}_{CE} = l(f(\mathbf{Z}), f(\mathbf{V})) \tag{11}$$

274 l represents our proposed MI loss module, and l
275 denotes cross-entropy loss. As shown in Figure 1,
276 the MI loss module module consists of softmax and
277 max reduction layers. The softmax layer aims to
278 generate an output distribution over representation
279 learning space, while the max reduction aims to
280 reduce the dimension of representation features of
281 the label prediction task $\mathbb{R}^{m \times d} \rightarrow \mathbb{R}^{1 \times d}$ and the di-
282 mension of representation features of the rationale
283 generation task $\mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{1 \times d}$.
284

285 4 Experiments

286 4.1 Experimental Setting

287 **Datasets.** We conducted extensive the experi-
288 ments on 4 widely-used benchmark datasets across
289 3 different NLP tasks: e-SNLI (Camburu et al.,
290 2018) and ANLI (Nie et al., 2020) for natural lan-
291 guage inference; CQA (Talmor et al., 2018) for
292 commonsense question answering; SVAMP (Pa-
293 tel et al., 2021) for arithmetic math word problems.
294 We use rationale generated by PaLM 540B (Chowd-
295 hery et al., 2023) collected and open-sourced
296 by (Hsieh et al., 2023)¹.

297 **Setup.** Based on CoT properties and comparative
298 experimental study in (Hsieh et al., 2023), our work
299 adopted T5-base (220 million) and T5-small (60
300 million) to the student models. α_1 and α_2 are set
301 as 0.5 and α_3 is set as 0.1. We trained our models
302 on one A100 GPU with 80G memory. For T5 base
303 model, the training time on full-size four dataset
304 was about 14.4 hours. For T5 small model, the
305 training times was from 8.6 hours.

306 **Baselines.** We compare our work with the state-
307 of-the-art DSS (Hsieh et al., 2023) by running its
308 open-sourced code 1. We also compare with two
309 most common methods in learning task-specific
310 models: (1) Finetuning, which is the standard fine-
311 tuning with the prevailing pretrain-then-finetune
312 paradigm that finetunes a model with ground-truth
313 labels via standard label supervision (Howard and
314 Ruder, 2018). (2) Single-task, which the small
315 model is distilled to predict labels with the teacher
316 model’s predicted label. We run DSS, finetuning,

¹Data and DSS code are from <https://github.com/google-research/distilling-step-by-step>.

and single-task with the same settings as (Hsieh et al., 2023).

Evaluation Settings. Following the DSS work (Hsieh et al., 2023), we adopt the accuracy as the performance metrics on all four datasets. A high accuracy indicates that the generated results are better. Besides it, we also adopt Expected Calibration Errors (ECE) and Average Confidence Scores to evaluate T5 calibration. Both higher scores indicate better. We adopt GPT-4 to evaluate Quality of CoT examples and subjective analysis. Please refer to our codes for more details.

4.2 Results

Experiments of T5-base Model. We present our experimental results of the T5-base model in Table 1. In single-task training, the rationale and label are concatenated into a single sequence, which is treated as the target in training models (Hsieh et al., 2023). Our proposed method consistently achieves better performance than standard fine-tuning and single-task methods on all datasets. Compared to DSS, our method outperforms DSS on ANLI, CQA, and SVAMP, and achieves nearly the same accuracy on e-SNLI

Experiments of T5-small Model. The experimental results of T5-small model are shown in Table 2. The patterns of the results are similar to those of T5-base. Our proposed method consistently achieves better performance than standard finetuning across all dataset. Compared to DSS, our method outperforms DSS on ANLI, CQA and SVAMP, and is just 0.2% less accuracy on e-SNLI.

Distillation with LLM Labels. We conducted an experiment on e-SNLI and ANLI dataset with T5-base model to evaluate the effect of label quality. We distilled the student models using labels generated by 540B PaLM instead of the ground truth. The results are shown in Table 3. Comparing Table 1 and Table 3, we observe the label quality affects the distillation results in both methods. Even With the noisy LLM labels, our model still outperforms DSS on both datasets.

Distillation with smaller datasets. To evaluate the performance of our models on smaller datasets, we distilled T5-base and T5-small models on various sizes of four datasets and compared to DSS method. The results are shown in Figure. 2 and 3 respectively.

	e-SNLI	ANLI	CQA	SVAMP
Finetuning	88.38	43.58	62.19	62.63
Single-task	88.88	43.50	61.37	63.00
DSS	89.51	49.58	63.29	65.50
Ours	89.50	51.20	63.88	68.00

Table 1: CoT distillation results on T5-base model.

	e-SNLI	ANLI	CQA	SVAMP
Finetuning	82.90	42.00	43.16	45.00
DSS	83.43	42.90	43.24	48.00
Ours	83.23	43.70	43.90	52.50

Table 2: CoT distillation results on T5-small model.

Model	e-SNLI	ANLI
DSS	82.65	42.80
Ours	82.81	45.50

Table 3: Results on two dataset on T5-base model with LLM generated labels.

4.3 Ablation Study

Effectiveness of Difference Dimension Reduction Method In our proposed MI loss module, we employ maximum reduction to align the dimension of different features. Additionally, mean reduction serves as an alternative method for dimension reduction. We hypothesize that important features can represent better than average features. In Table 4, we present the results of two different layer of MI module. The results indicate the superiority of the MI module with maximum reduction.

Comparison with KL Divergence The KL divergence loss has been extensively utilized in KD tasks, serving as KL a metric for assessing the similarity between two data distributions (Hinton et al., 2015; Zhang et al., 2022b; Gou et al., 2021). While KL divergence has found widespread application in various KD scenarios, modeling DSS using IB framework proves to be more accurate than using similarity measures, as discussed in Section 3. To validate our hypothesis, we conduct experiments on T5-base model using all four datasets. As shown in Table 5, our proposed method consistently outperforms the KL divergence approach, demonstrating superior performance.

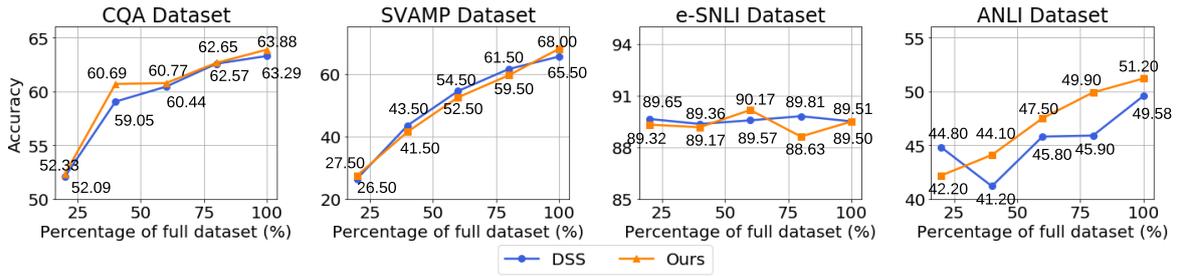


Figure 2: Comparison with DSS with varying sizes of training datasets on T5-base model.

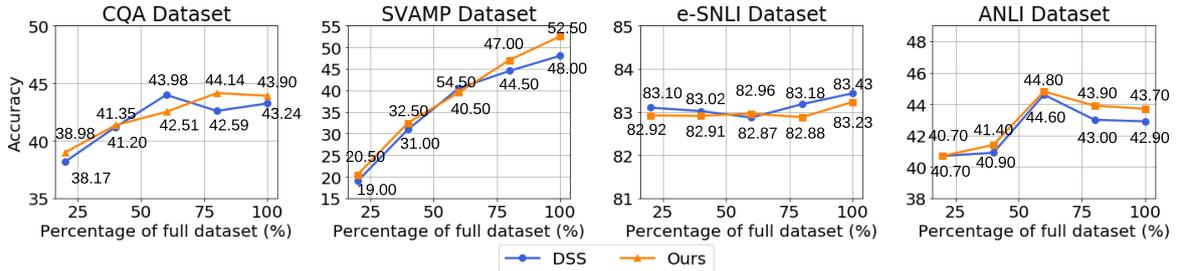


Figure 3: Comparison with DSS with varying sizes of training datasets on T5-small model.

	e-SNLI	ANLI	CQA	SVAMP
Mean	89.34	51.40	63.88	66.50
Max	89.50	51.20	63.88	68.00

Table 4: Results of Mean Reduction Vs Maximum Reduction on T5-based model.

	e-SNLI	ANLI	CQA	SVAMP
KL Divergence	89.42	42.00	62.49	67.00
Ours	89.50	51.2	63.88	68.00

Table 5: Results of KD loss VS our proposed cross entropy loss, on T5-base model.

5 Discussion

5.1 Analysis on T5 Calibration

Calibration measures the alignment between a model’s predicted accuracy and its confidence levels. Lee et al. (2022) introduced an innovative perspective on model distillation, positioning the teacher model not only as a source of knowledge but also as a tool for identifying mis-calibration during the training of the student model. This ability to maintain calibration and make reliable predictions is crucial for downstream applications and has been the focus of prior studies (Chen et al., 2023; Lee et al., 2022; Jiang et al., 2021). Here,

we apply the Expected Calibration Errors (ECE) and Average Confidence Scores to reflect the alignment between the model’s predicted probabilities and the actual outcomes, thereby gauging the reliability and certainty of its predictions. Despite the potential limitations inherent in these metrics, we still employ ECE in our experiments due to its simplicity and popularity, as in previous work on investigating the calibration quality of T5 (Chen et al., 2023; Lee et al., 2022).

We employ a 10-bin-based ECE metric and a softmax-based approach to compute average confidence scores from the test outputs across all four datasets. Given that e-SNLI and ANLI essentially represent the same task, we conduct an out-of-domain experiment by testing the model checkpoint trained on one dataset with the test set of the other. This analysis gives us insights into how well our model generalizes across similar tasks and the robustness of its predictions in out-of-domain scenarios and to assess the calibration quality of the model more comprehensively.

Table 6 presents the results of the distilled model calibration evaluation. Overall, both models report lower ECE and confidence scores on SVAMP and e-SNLI, indicating that these two tasks are more challenging and models are less certain about their prediction. Lower ECE values from our MI-based

Model	SVAMP		CQA		e-SNLI		ANLI		e-SNLI (Out)		ANLI (Out)	
	ECE	Conf.	ECE	Conf.	ECE	Conf.	ECE	Conf.	ECE	Conf.	ECE	Conf.
DSS	11.81	32.56	11.75	42.79	8.54	34.33	11.12	42.72	9.81	38.01	12.78	41.69
Ours	18.92	36.81	13.65	41.17	4.35	30.06	6.94	35.90	6.61	38.08	12.27	42.35

Table 6: Comparisons of our model and DSS on the expected calibration errors (ECE) and average confidence scores (Conf.).

distillation approach are presented for e-SNLI and ANLI, and their respective out-of-domain tests. Notably, our method achieves an ECE of 4.35 in e-SNLI, significantly lower than DSS’s 8.54. However, in SVAMP and CQA, our method records higher ECE, indicating potential areas for improvement in these domains. The trade-off in calibration accuracy in specific tasks like SVAMP and CQA compared to DSS suggests future directions for refining our approach.

Regarding average confidence scores (Conf.), our method generally maintains competitive confidence levels, with notable improvements in e-SNLI and ANLI. In e-SNLI, the confidence is lower (30.06) compared to DSS (34.33), which, combined with a lower ECE, suggests a more realistic confidence estimation. Conversely, in the out-of-domain scenarios for e-SNLI and ANLI, our method shows marginally higher confidence scores than DSS, which, coupled with the lower ECE, indicates robustness in out-of-domain generalization.

5.2 Analysis on CoT Output

5.2.1 Quality of CoT Examples by GPT-4 Evaluation

We evaluate the quality of CoT examples using GPT-4, as it achieves the state-of-the-art human alignment performance and is used for text generation evaluation in previous work (Liu et al., 2023; Hsu et al., 2023; Wang et al., 2023a). Inspired by (Wang et al., 2023a), we ask GPT-4 to evaluate the quality of the provided CoT examples based on their coherency and relevancy to the input questions and answers. We randomly sample 50 CoT examples from the four datasets and ask GPT-4 to score based on a scale from 1 to 5, where 1 indicates completely incoherent and irrelevant responses, and 5 represents highly coherent, relevant, and helpful responses. For each sample, we run the same sample for four times to obtain self-consistency to measure the reliability of the responses. Table 7 presents the prompt we use for GPT-4 evaluation, average scores and standard deviation on the scores ob-

tained over the four datasets. We report the scores on both provided CoT (“gold”) rationales and distilled model predicted rationales.

Prompt for GPT 4 Evaluation					
Given an input pair of a question and an answer of a <i>taskname</i> task, how good is the given Chain-of-thought example? From 1-5, where 1 is completely incoherent and irrelevant, 2 is somewhat incoherent and irrelevant, 3 is coherent, relevant but not helpful, 4 is somewhat helpful, and 5 is helpful and it explains the answer well.					
Average Scores and Standard Deviation					
Model	SVAMP	CQA	e-SNLI	ANLI	
Gold	4.63±1.05	3.95±1.16	2.42±1.23	3.82±1.26	
++	4.43±1.18	4.11±1.40	3.49±1.35	4.01±1.10	
DSS	2.50±1.42	3.60±1.61	3.24±1.27	3.48±1.40	
++	2.53±1.46	3.64±1.62	3.18±1.21	3.44±1.30	
Ours	2.30±1.54	3.70±1.45	3.03±1.47	3.42±1.37	
++	2.72±1.45	3.63±1.60	3.17±1.17	3.34±1.21	

Table 7: Prompt used and results of 50 randomly sampled CoT examples from the four datasets evaluated by GPT-4. We use ++ to denote the setting with *self-consistency* evaluation.

Model	SVAMP	CQA	e-SNLI	ANLI
DSS	0.12	0.66	0.05	0.26
	$p > 0.05$	$p < 0.05$	$p > 0.05$	$p > 0.05$
Ours	0.42	0.53	0.03	0.26
	$p < 0.05$	$p < 0.05$	$p > 0.05$	$p > 0.05$

Table 8: Pearson correlation between CoT quality and accuracy of label prediction on the 50 random samples on the test set. We highlight the correlation with statistical significance ($p < 0.05$).

The effectiveness of our MI-based distillation method is closely linked to the quality of CoT reasoning in the training data. When the CoT quality is high, as in SVAMP, a strong correlation is observed between the model’s label prediction accuracy and the quality of its generated CoT. However, this correlation weakens significantly when the CoT quality is low (e-SNLI), suggesting that the model struggles to align label prediction with coherent rationale generation under poor training conditions. Interestingly, with average-quality CoT

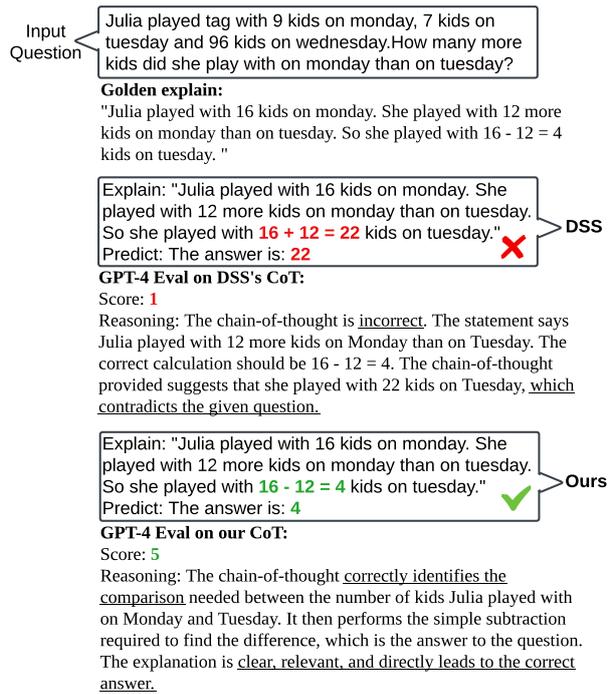


Figure 4: A case study of the output rationale on SVAMP.

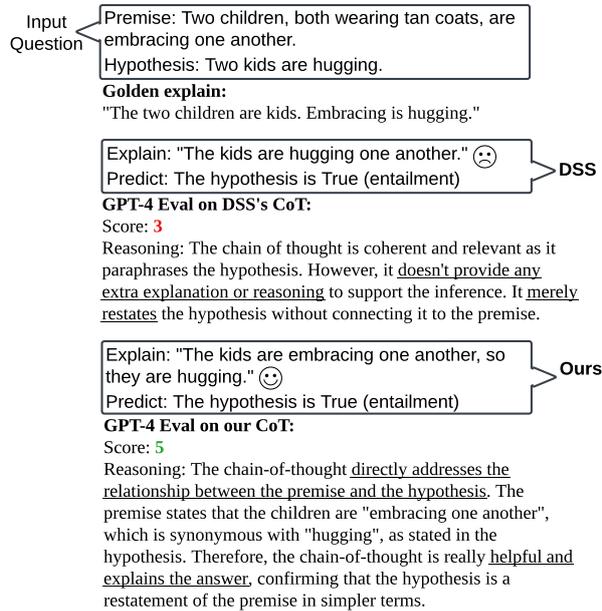


Figure 5: A case study of the output rationale on e-SNLI.

data (ANLI), the performance gap between our MI-based distillation and DSS is minimal, suggesting that the effectiveness of our approach is particularly reliant on the presence of high-quality reasoning in the training data.

5.2.2 Case Studies on the Output Rationale

We performed case studies on SVAMP and e-SNLI as illustrated in Figure 4 and 5. In the SVAMP example, the question asks the difference in the number of kids Julia played with from Monday to Tuesday, with specific numbers provided for Monday, Tuesday, and Wednesday. DSS generates an incorrect explanation, which contradicts the given question, resulting in to a wrong answer. Conversely, our method correctly identifies the comparison needed between the number of kids Julia played with on Monday and Tuesday, leading to the correct answer. Notably, our generated CoT reasoning is identical to the golden one, demonstrating that by precisely grasping the rationale, our approach effectively resolves the math problem. We also show the evaluation results (score and reasoning) from GPT-4, where our method gains a top score of 5 and DSS gains only a mere score of 1. This example showcases that the high-quality CoT generated by our method enhances problem-solving capabilities in math tasks like SVAMP.

Another example (Figure 5) is from e-SNLI, where the task is to identify whether the hypothesis is entailment, contradiction, or neutral, based on the given premise and hypothesis. Although both our method and DSS generate the correct label output, it is worth noting that, the CoT of our method points out the relationship between the premise and the hypothesis, while DSS only restates the hypothesis without providing any extra explanation or connecting the hypothesis to the premise. Our generated rationale also gains a higher score than DSS. A higher-quality rationale tends to facilitate more accurate label prediction, thereby enhancing overall task performance.

6 Conclusion

In this paper, we re-investigate the DSS framework from a information-theoretic perspective. We model it using Information Bottleneck and propose to strengthen it by maximizing the mutual information between rationale generation and label prediction tasks. The proposed learning-based method can automatically optimize the CoT distillation and bolster the reasoning ability of the distilled small models. Our qualitative and quantitative analysis demonstrate the rationale behind our method and shed light on language model distillation and CoT applications.

7 Limitation

Our comparative analysis primarily focuses on the Distilling Step-by-Step (DSS) framework, which serves as our main benchmark. This concentrated comparison, while valuable for a deep understanding of DSS’s nuances and our advancements over it, constitutes a limitation of our work. Specifically, our analysis does not extend to a broader range of knowledge distillation methods currently employed in the field. This focus may overlook the potential insights and contrasts that could emerge from evaluating our approach against a wider array of distillation techniques. Future research could benefit from a more expansive comparative study, incorporating diverse methodologies to fully contextualize our findings within the broader landscape of knowledge distillation practices. This broader comparison would not only validate the efficacy of our method in various settings but also illuminate areas for further refinement and innovation.

However, it is important to note that our contribution lies in providing an in-depth analysis from both theoretical and practical viewpoints to enhance the CoT distillation process. Our work delves into the intricacies of utilizing mutual information to improve distillation outcomes, offering significant advancements in understanding and applying CoT distillation techniques.

8 Ethical Issues

In this paper, we carefully considered the ethical implications in line with the ACL code of ethics. We evaluated the potential dual-use concerns, ensuring our research serves to benefit society and does not cause inadvertent harm. Our methodology and applications were thoroughly assessed for fairness, non-discrimination, and privacy, particularly in the context of data handling and model outputs. We also ensured our study did not expose any negative impact on individuals and groups. Moreover, we did not engage in academic dishonesty and adhered to high-quality processes and product standards in our professional work. We include this detailed discussion of these ethical considerations, affirming our commitment to responsible and beneficial computational linguistics research.

References

Alexander A Alemi, Ian Fischer, Joshua V Dillon, and Kevin Murphy. 2016. Deep variational information

bottleneck. In *International Conference on Learning Representations*.

Zeyuan Allen-Zhu and Yuanzhi Li. 2023. Towards understanding ensemble, knowledge distillation and self-distillation in deep learning. In *The Eleventh International Conference on Learning Representations*.

Philip Bachman, R Devon Hjelm, and William Buchwalter. 2019. Learning representations by maximizing mutual information across views. *Advances in neural information processing systems*, 32.

Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and Devon Hjelm. 2018. Mutual information neural estimation. In *International conference on machine learning*, pages 531–540. PMLR.

Oana-Maria Camburu, Tim Rocktäschel, Thomas Lukasiewicz, and Phil Blunsom. 2018. e-snli: Natural language inference with natural language explanations. *Advances in Neural Information Processing Systems*, 31.

Rich Caruana. 1997. Multitask learning. *Machine learning*, 28:41–75.

Hanjie Chen and Yangfeng Ji. 2020. Learning variational word masks to improve the interpretability of neural text classifiers. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4236–4251.

Pengguang Chen, Shu Liu, Hengshuang Zhao, and Jiaya Jia. 2021. Distilling knowledge via knowledge review. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5008–5017.

Yangyi Chen, Lifan Yuan, Ganqu Cui, Zhiyuan Liu, and Heng Ji. 2023. [A close look into the calibration of pre-trained language models](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1343–1367, Toronto, Canada. 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.

Thomas M Cover. 1999. *Elements of information theory*. John Wiley & Sons.

Ian Connick Covert, Wei Qiu, Mingyu Lu, Na Yoon Kim, Nathan J White, and Su-In Lee. 2023. Learning to maximize mutual information for dynamic feature selection. In *International Conference on Machine Learning*, pages 6424–6447. PMLR.

753	Arkil Patel, Satwik Bhattamishra, and Navin Goyal.	Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen,	806
754	2021. Are nlp models really able to solve	You Wu, Luke Zettlemoyer, and Huan Sun. 2023a.	807
755	simple math word problems? <i>arXiv preprint</i>	Towards understanding chain-of-thought prompting:	808
756	<i>arXiv:2103.07191</i> .	An empirical study of what matters. In <i>Proceedings</i>	809
757	Mary Phuong and Christoph Lampert. 2019. Towards	of the 61st Annual Meeting of the Association for	810
758	understanding knowledge distillation. In <i>International</i>	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	811
759	<i>conference on machine learning</i> , pages 5142–	pages 2717–2739, Toronto, Canada. Association for	812
760	5151. PMLR.	Computational Linguistics.	813
761	Ben Poole, Sherjil Ozair, Aaron Van Den Oord, Alex	Lin Wang and Kuk-Jin Yoon. 2021. Knowledge distil-	814
762	Alemi, and George Tucker. 2019. On variational	lation and student-teacher learning for visual intelli-	815
763	bounds of mutual information. In <i>International Con-</i>	gence: A review and new outlooks. <i>IEEE transac-</i>	816
764	<i>ference on Machine Learning</i> , pages 5171–5180.	<i>tions on pattern analysis and machine intelligence</i> ,	817
765	PMLR.	44(6):3048–3068.	818
766	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	Peifeng Wang, Zhengyang Wang, Zheng Li, Yifan Gao,	819
767	Lee, Sharan Narang, Michael Matena, Yanqi	Bing Yin, and Xiang Ren. 2023b. SCOTT: Self-	820
768	Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the	consistent chain-of-thought distillation. In <i>Proceed-</i>	821
769	limits of transfer learning with a unified text-to-text	<i>ings of the 61st Annual Meeting of the Association for</i>	822
770	transformer. <i>Journal of Machine Learning Research</i> ,	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	823
771	21(140):1–67.	pages 5546–5558, Toronto, Canada. Association for	824
772	Victor Sanh, Lysandre Debut, Julien Chaumond, and	Computational Linguistics.	825
773	Thomas Wolf. 2019. Distilbert, a distilled version	Qi Wang, Claire Boudreau, Qixing Luo, Pang-Ning Tan,	826
774	of bert: smaller, faster, cheaper and lighter. <i>arXiv</i>	and Jiayu Zhou. 2019. Deep multi-view information	827
775	<i>preprint arXiv:1910.01108</i> .	bottleneck. In <i>Proceedings of the 2019 SIAM Inter-</i>	828
776	Noam Slonim. 2002. <i>The information bottleneck: The-</i>	<i>national Conference on Data Mining</i> , pages 37–45.	829
777	<i>ory and applications</i> . Ph.D. thesis, Hebrew Univer-	SIAM.	830
778	sity of Jerusalem Jerusalem, Israel.	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel,	831
779	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and	Barret Zoph, Sebastian Borgeaud, Dani Yogatama,	832
780	Jonathan Berant. 2018. Commonsenseqa: A question	Maarten Bosma, Denny Zhou, Donald Metzler, et al.	833
781	answering challenge targeting commonsense knowl-	2022a. Emergent abilities of large language models.	834
782	edge. <i>arXiv preprint arXiv:1811.00937</i> .	<i>arXiv preprint arXiv:2206.07682</i> .	835
783	Xudong Tian, Zhizhong Zhang, Shaohui Lin, Yanyun	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	836
784	Qu, Yuan Xie, and Lizhuang Ma. 2021. Farewell to	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	837
785	mutual information: Variational distillation for cross-	et al. 2022b. Chain-of-thought prompting elicits rea-	838
786	modal person re-identification. In <i>Proceedings of</i>	soning in large language models. <i>Advances in Neural</i>	839
787	<i>the IEEE/CVF Conference on Computer Vision and</i>	<i>Information Processing Systems</i> , 35:24824–24837.	840
788	<i>Pattern Recognition</i> , pages 1522–1531.	Joseph Worsham and Jugal Kalita. 2020. Multi-task	841
789	Yonglong Tian, Dilip Krishnan, and Phillip Isola. 2020.	learning for natural language processing in the 2020s:	842
790	Contrastive representation distillation. In <i>International</i>	where are we going? <i>Pattern Recognition Letters</i> ,	843
791	<i>Conference on Learning Representations</i> .	136:120–126.	844
792	Naftali Tishby and Noga Zaslavsky. 2015. Deep learn-	Yangyang Xu, Yibo Yang, and Lefei Zhang. 2023.	845
793	ing and the information bottleneck principle. In <i>2015</i>	Multi-task learning with knowledge distillation for	846
794	<i>iee information theory workshop (itw)</i> , pages 1–5.	dense prediction. In <i>Proceedings of the IEEE/CVF</i>	847
795	IEEE.	<i>International Conference on Computer Vision</i> , pages	848
796	Michael Tschannen, Josip Djolonga, Paul K Rubenstein,	21550–21559.	849
797	Sylvain Gelly, and Mario Lucic. 2019. On mutual in-	Chenxiao Yang, Junwei Pan, Xiaofeng Gao, Tingyu	850
798	formation maximization for representation learning.	Jiang, Dapeng Liu, and Guihai Chen. 2022. Cross-	851
799	In <i>International Conference on Learning Representa-</i>	task knowledge distillation in multi-task recommen-	852
800	<i>tions</i> .	dation. In <i>Proceedings of the AAAI Conference on</i>	853
801	Zhibin Wan, Changqing Zhang, Pengfei Zhu, and	<i>Artificial Intelligence</i> , volume 36, pages 4318–4326.	854
802	Qinghua Hu. 2021. Multi-view information-	Cenyuan Zhang, Xiang Zhou, Yixin Wan, Xiaoqing	855
803	bottleneck representation learning. In <i>Proceedings</i>	Zheng, Kai-Wei Chang, and Cho-Jui Hsieh. 2022a.	856
804	<i>of the AAAI conference on artificial intelligence</i> , vol-	Improving the adversarial robustness of nlp models	857
805	ume 35, pages 10085–10092.	by information bottleneck. In <i>Findings of the Asso-</i>	858
		<i>ciation for Computational Linguistics: ACL 2022</i> ,	859
		pages 3588–3598.	860

861 Linfeng Zhang, Chenglong Bao, and Kaisheng Ma.
862 2021. Self-distillation: Towards efficient and com-
863 pact neural networks. *IEEE Transactions on Pat-
864 tern Analysis and Machine Intelligence*, 44(8):4388–
865 4403.

866 Linfeng Zhang, Xin Chen, Junbo Zhang, Runpei Dong,
867 and Kaisheng Ma. 2022b. Contrastive deep supervi-
868 sion. In *European Conference on Computer Vision*,
869 pages 1–19. Springer.

870 Yu Zhang and Qiang Yang. 2021. A survey on multi-
871 task learning. *IEEE Transactions on Knowledge and
872 Data Engineering*, 34(12):5586–5609.

873 Zhihan Zhang, Wenhao Yu, Mengxia Yu, Zhichun Guo,
874 and Meng Jiang. 2023. A survey of multi-task learn-
875 ing in natural language processing: Regarding task
876 relatedness and training methods. In *Proceedings
877 of the 17th Conference of the European Chapter of
878 the Association for Computational Linguistics*, pages
879 943–956.