



# FIRST: Teach A Reliable Large Language Model Through Efficient Trustworthy Distillation

Anonymous ACL submission

## Abstract

Large language models (LLMs) have become increasingly prevalent in our daily lives, leading to an expectation for LLMs to be **trustworthy** — both accurate and well-calibrated (the prediction confidence should align with its ground truth correctness likelihood). Nowadays, fine-tuning has become the most popular method for adapting a model to practical usage by significantly increasing accuracy on downstream tasks. Despite the great accuracy it achieves, we found fine-tuning is still far away from satisfactory trustworthiness due to "tuning-induced mis-calibration". In this paper, we delve deeply into why and how mis-calibration exists in fine-tuned models, and how distillation can alleviate the issue. Then we further propose a brand new method named **EFFicient TRustworthy DiSTillation (FIRST)**, which utilizes a small portion of teacher's knowledge to obtain a reliable language model in a cost-efficient way. Specifically, we identify the "concentrated knowledge" phenomenon during distillation, which can significantly reduce the computational burden. Then we apply a "trustworthy maximization" process to optimize the utilization of this small portion of concentrated knowledge before transferring it to the student. Experimental results demonstrate the effectiveness of our method, where better accuracy (+2.3%) and less mis-calibration (-10%) are achieved on average across both in-domain and out-of-domain scenarios, indicating better trustworthiness.

## 1 Introduction

With the rapid development of large language models (LLMs), many powerful models have been deployed into our daily lives for practical usage to help us make decisions (Yao et al., 2023; Sha et al., 2023; Zhao et al., 2024). This makes it urgent for us to know to what extent we can trust the outputs of the models. Calibration is one of the most important indicators beyond accuracy which

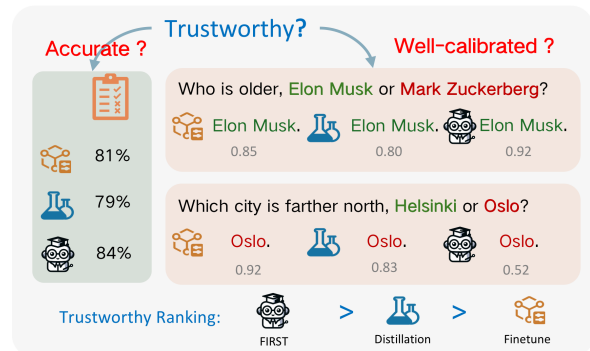


Figure 1: A trustworthy model should be both accurate (left) and well-calibrated (right). A well-calibrated model should produce high probabilities for the **correct answer** and low probabilities for the **wrong answer**.

provides a confidence measure to the model's predictions (Guo et al., 2017; Hsieh et al., 2023). In LLMs, confidence is exactly the probability for each generated token. Therefore, a well-calibrated model should align its prediction confidence with its ground-truth correctness likelihood. As an example, recent hallucination detection methods rely on model prediction confidence as a significant indicator of potential hallucination (Zhang et al., 2023; Varshney et al., 2023). If the model is incapable of giving accurate confidence levels, people may fail to detect hallucinations due to the model's over-confidence, or people may falsely identify hallucinations due to the model's under-confidence. Mis-calibration brings significant challenges for the deployment of LLMs in real-world applications.

Currently, there are two methods to obtain a language model for practical usage. First, fine-tuning, which fine-tunes pre-trained LLMs on specific datasets by matching each token entry with a target ground truth token. Although fine-tuning can consistently improve performance on downstream tasks (Dodge et al., 2020; Sun et al., 2020; Ziegler et al., 2020), we identify that the model obtained in this way exhibits a nature of "tuning-induced mis-calibration". Second, distillation-based methods

069	transfer knowledge (e.g., soft labels) from larger	calibration" phenomena, providing insights	119
070	LLMs to smaller models (Gu et al., 2023). Al-	into obtaining trustworthy models.	120
071	though distillation shows better calibration than		
072	fine-tuning as it matches each token entry with a	(ii) We propose <b>FIRST</b> , which maximizes the ef-	121
073	probability distribution instead of a hard label, we	fectiveness and trustworthiness of a relatively	122
074	find it is still biased because of the mis-calibration	small portion of knowledge transferred from	123
075	nature of teacher models. In addition, distilla-	the teacher by "trustworthy maximization" to	124
076	tion faces the challenge of determining the opti-	obtain a trustworthy student model.	125
077	mal amount of knowledge to transfer. Transferring		
078	all the teacher’s knowledge leads to high compu-	(iii) Extensive experiments demonstrate that mod-	126
079	tational costs while transferring too little knowl-	els obtained using <b>FIRST</b> consistently achieve	127
080	edge results in poor accuracy. Therefore, it is cru-	the highest level of trustworthiness across dif-	128
081	cial to balance between trustworthiness (accuracy	ferent settings.	129
082	and well-calibration) and efficiency for distillation-		
083	based methods.	<b>2 Related Work</b>	130
084	To address the challenge of obtaining a trustwor-	<b>2.1 Trustworthy Models</b>	131
085	thy model, we propose <b>eFFicient tRustworthy</b>	The current evaluation of LLMs predominantly fo-	132
086	<b>disTillation (FIRST)</b> , aiming to efficiently utilize	cuses on accuracy, overlooking whether the mod-	133
087	a relatively small amount of the teacher’s knowl-	els truly know the answer or are merely guess-	134
088	edge. Specifically, we first identify the "concent-	ing (i.e. trustworthy). Recent works (Sun et al.,	135
089	rated knowledge" phenomenon, which shows that	2024; Steyvers et al., 2024) have demonstrated that	136
090	in the context of LLMs, the probability distribution	accurate LLMs may not necessarily be "trustwor-	137
091	of generated tokens is not uniform but rather con-	thy" due to a significant calibration gap, so-called	138
092	centrated on a few high-probability tokens. Based	mis-calibration. This gap prevents us from trust-	139
093	on this finding, we propose to use the top-5 tokens	ing the output of the models, and it can further	140
094	as the knowledge to balance the trade-off between	cause LLMs to generate harmful content, especially	141
095	storage space and the amount of knowledge trans-	when subjected to adversarial attacks or jailbreak	142
096	ferred, achieving efficient distillation. Afterward,	prompts (Mo et al., 2024; Yao et al., 2024). Our	143
097	to eliminate the "tuning-induced mis-calibration"	work further reveals how mis-calibration exists in	144
098	of the teacher model, we applied a "trustworthy	different tuning methods and proposes a new trust-	145
099	maximization" to this portion of knowledge, en-	worthy evaluation metric that covers both accuracy	146
100	suring that it maximizes the enhancement of the	and calibration.	147
101	student model’s accuracy while also guaranteeing		
102	its well-calibration.	To achieve a well-calibrated LLM, recent work	148
103	We first validate our method in in-domain sce-	shows soft-label distillation shows better calibra-	149
104	narios, discovering that the models obtained by	tion ability (Gu et al., 2023). However, it still suf-	150
105	<b>FIRST</b> achieve excellent accuracy, even with the	fers from biased labels due to the mis-calibration	151
106	use of a relatively small amount of top-5 knowl-	nature of the fine-tuned teacher model. Our work	152
107	edge and the "trustworthy maximization" process	is an improvement on this line of work by applying	153
108	can significantly enhance these models’ calibra-	"concentrated knowledge" and "trustworthy max-	154
109	tion ability. Furthermore, we test our approach in	imization", leading to better accuracy, efficiency,	155
110	out-of-domain settings, demonstrating that models	and trustworthy.	156
111	obtained by <b>FIRST</b> still exhibit the best trustwor-	<b>2.2 Knowledge Distillation</b>	157
112	thiness and hold generalization ability. This indicates	Knowledge Distillation is a form of transfer learn-	158
113	that <b>FIRST</b> enables smaller models to genuinely	ing that facilitates the transfer of knowledge from	159
114	learn the capability of being trustworthy, rather	a larger teacher model to a smaller student model.	160
115	than being confined to in-domain scenarios.	The goal is to reduce the model size while main-	161
116	In summary, our key contributions include:	taining or even improving performance. Based on	162
117	(i) We discover that LLMs exhibit "concent-	whether we can access prediction probability, the	163
118	rated knowledge" and "tuning-induced mis-	existing distillation methods can be categorized	164
		into two types: Black-box Distillation and White-	165
		box Distillation.	166

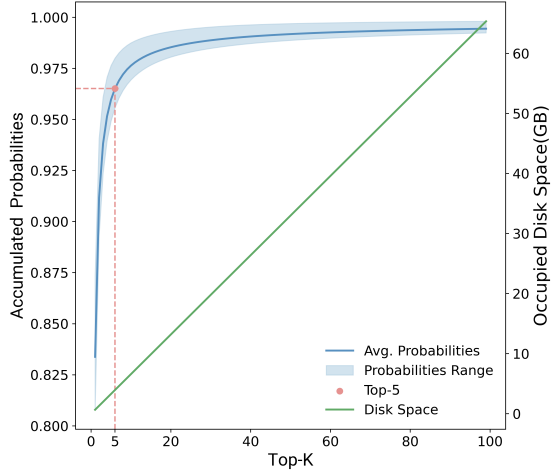


Figure 2: The blue line with range shows the averaged accumulated probability coverage for each token entry, from Top-1 to Top-100. **"Concentrated Knowledge"**: The red point represents accumulated probability for Top-5 tokens already exceed 95%. The green line describes the disk usage if use Top-K token distribution during distillation.

Black-box Distillation refers to distillation from models that we are unable to access the weight and prediction logits such as PaLM (Chowdhery et al., 2022). Recent studies have attempted to distill reasoning ability from GPT (Ho et al., 2023; Shridhar et al., 2023) or some emergent ability such as chain-of-thought (Hsieh et al., 2023; Li et al., 2023). However, these methods may still be categorized as the genre of data-augmentation-and-then-fine-tuning approaches.

White-box Distillation means the teacher models are either fully open-sourced such as Llama (Touvron et al., 2023a) or they can return partial probability distribution of the generated tokens, such as code-davinci-002. Instead of the hard token fine-tuning, white-box distillation typically uses more fine-grained signals by matching a distribution between teachers and students (Gu et al., 2023; Latif et al., 2023; Agarwal et al., 2024). Further, in the field of white-box distillation, there are two different ways: online distillation and offline distillation. Online distillation (Gu et al., 2023; Zhou et al., 2023) needs to keep both the teacher model and the student model on the GPU simultaneously during training. On the other hand, offline distillation typically involves obtaining knowledge from the teacher model beforehand. Our work is an extension of white-box offline distillation and focuses on how white-box offline distillation can be improved in terms of trustworthiness by re-calibrating the teacher distribution.

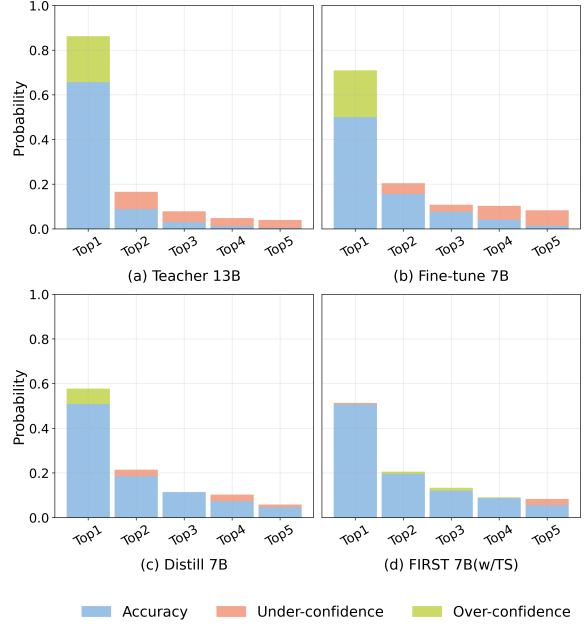


Figure 3: **"Tuning-Induced Mis-calibration"**: Position-wise prediction probabilities with corresponding actual accuracy of (a) fine-tuned teacher model and (b) fine-tuned small model, (c) distilled model and (d) model produced by FIRIST.

### 3 Preliminaries

#### 3.1 Concentrated Knowledge

In the process of searching for a suitable trade-off between the amount of knowledge to transfer from the teacher model and efficiency, we begin by visualizing the probability distribution for each token entry. As illustrated in Figure 2, the blue line with range describes how averaged accumulated probabilities increase when we select more tokens (ranked from highest probability to lowest probability in one entry). The trend clearly shows a few top-position tokens take most of the probability information of a token entry. To be specific, the accumulated probabilities of Top-5 tokens can occupy over 95% probabilities while the remaining 49995 (i.e. a model with vocab. size of 50k) tokens have nearly 0 probability. We named this phenomenon "Concentrated Knowledge" as almost full knowledge of a token entry is stored in its top-k tokens where the remaining tokens have negligible information.

#### 3.2 Tuning-Induced Mis-calibration

In the context of LLMs, mis-calibration can be divided into two types: over-confidence and under-confidence. Over-confidence occurs when the predicted probability of a token is higher than its actual accuracy, while under-confidence takes place when

the predicted probability is lower than the actual accuracy.

During the fine-tuning process of LLMs, cross-entropy loss is commonly employed, which encourages the models to assign a probability of 1 to one token and 0 to all other tokens based on the ground-truth token. This training nature results in 1.) an over-estimation of the ground truth token’s probability and 2.) an under-estimation of all other token’s probability. As shown in Figure 3 (a) and (b), it is observed that both fine-tuned LLMs exhibit over-confidence in their top-1 token predictions, while demonstrating under-confidence in the subsequent tokens. This phenomenon, which we call "tuning-induced calibration", highlights the untrustworthy nature of fine-tuned models.

Since fine-tuned teacher models suffer from this tuning-induced mis-calibration, if the knowledge from the mis-calibrated teacher models is directly used in traditional distillation-based methods, the student models are very likely to inherit the same mis-calibration nature as depicted in Figure 3 (c). Motivated by the tuning-induced mis-calibration, our proposed method incorporates a "trustworthy maximization" procedure to re-calibrate the knowledge derived from the teacher models. This enables us to obtain a genuinely trustworthy student model.

### 3.3 Expected Calibration Error

To measure calibration in the context of LLMs, we adapt the expected calibration error (ECE) to the free-text generation task by treating the generation of a single token as a classification task. In this adaptation, we restrict the model to generate only one token from a set of candidate choices (e.g., A/B/C/D). For each token, we obtain the highest probability choice using  $\arg \max_{i \in C} P(i)$ , where  $C$  represents the set of candidates. The probability of the chosen token is taken as the predicted confidence, and we calculate the accuracy by comparing the predicted choice to the ground truth. Then we utilize a total  $M$  probability interval as bins and categorize each chosen token into  $m$ -th bin according to the predicted confidence. The ECE can be computed as follows:

$$ECE = \sum_{m=1}^M \frac{|B_m|}{n} |acc(B_m) - conf(B_m)| \quad (1)$$

Here,  $M$  is the number of bins.  $B_m$  represents the set of predictions in bin  $m$ ,  $|B_m|$  is the number

of prediction instances in bin  $m$ , and  $n$  is the total number of predictions.  $acc(B_m)$  is the average accuracy of predictions in bin  $m$ , and  $conf(B_m)$  is the average confidence of predictions in bin  $m$ . A lower ECE value indicates that the model’s predicted probabilities are more consistent with actual outcomes, meaning the model is better calibrated.

### 3.4 Trustworthy Score

In evaluating the trustworthiness of a model, it is essential to consider both high accuracy and effective calibration. Existing benchmarks primarily focus on accuracy, assuming that higher accuracy implies greater trustworthiness. However, our discovery of the widespread issue of "tuning-induced mis-calibration" has highlighted the inadequacy of relying solely on accuracy for a comprehensive evaluation of model reliability. To address this limitation, we propose Trust Score metric to quantify a model’s trustworthiness, which quantifies the trustworthiness of a model by considering two key aspects: its ability to provide accurate answers (measured by  $Acc$ ) and its capacity to align predicted confidences with actual accuracies (measured by  $ECE$ ). The Trust Score is defined as follows:

$$Trust = Acc - ECE \quad (2)$$

By incorporating the Trust Score, we achieve a more balanced evaluation of trustworthiness, taking into account both accuracy and calibration.

## 4 Efficient Trustworthy Distillation

In this section, we introduce eFFicient tRustworthy disTillation (**FIRST**), which can be divided into three parts. Firstly, we select Top-5 tokens as knowledge for transfer (Efficient Knowledge Selection) in Sec.4.1. Then, we adjust the knowledge for trustworthiness to ensure that the subsequent smaller models can maximize its utility (Knowledge Trustworthy Maximization) in Sec.4.2. Finally, we describe the learning process of the student model (Knowledge Matching) in Sec.4.3.

### 4.1 Efficient Knowledge Selection

Transferring knowledge directly from teachers to students can be computationally costly and storage-intensive. For example, if we consider a vocabulary size of 50,000 tokens, retrieving the complete probability distribution from a dataset of 100,000 samples, with an average length of 2,048, would require a staggering 120 TB of storage, which is impractical.

Based on the discovery of "concentrated knowledge" in teacher LLMs, we observe that the majority of knowledge is concentrated within a small portion of top-position tokens, as elaborated in Section §3.1. Therefore, considering that both computation and disk space increase linearly with the number of selected token entries, we argue that it is not necessary to use the complete probability distribution. Instead, by selecting a small amount of top-position tokens that contain majority of knowledge, we can strike the optimal balance between computational overhead and effectiveness. As depicted in figure 2, accumulated probability of Top-5 token entries occupy more than 95% probabilities while reducing storage from 120 TB to 1.2 GB.

## 4.2 Trustworthy Maximization

Once the top-5 tokens and their corresponding probabilities are collected from the teacher model, it is crucial to subject this knowledge to further processing to ensure proper calibration, as teacher models can also suffer from "tuning-induced Mis-calibration" due to fine-tuning (as we elaborate in Sec. §3.2). This additional calibration step ensures that the student model improves in both accuracy and trustworthiness.

**Label Smoothing.** We first attempted to address tuning-induced mis-calibration" by applying a smoothing coefficient, denoted as  $\delta$ , to mitigate the teacher model's over-confidence in its top-1 token predictions while alleviating under-confidence in other predicted tokens as follows:

$$\begin{cases} P_T(i) := P_T(i) - \delta & \text{if } i = 1 \\ P_T(i) := P_T(i) + \frac{\delta}{4} & \text{if } 2 \leq i \leq 5 \end{cases} \quad (3)$$

Here,  $T$  denotes the teacher model,  $P_T(i)$  represents the probability of the  $i$ -th top token. While label smoothing can effectively mitigate over-confidence in top-1 token predictions, we have identified significant drawbacks associated with this approach. Firstly, directly applying label smoothing may compromise the preservation of token rankings, particularly between the top-1 and top-2 tokens. This can lead to a decline in model performance in certain cases. Secondly, label smoothing uses a constant probability, disregarding the varying levels of over-confidence or under-confidence in different token entries. Consequently, this can result in a transition from under-confidence to over-confidence among the top 2-5 tokens, making it

challenging to achieve a balanced calibration across all of them.

**Temperature Scaling.** Subsequently, we explore another approach using a temperature scaling technique to re-calibrate the probabilities:

$$P_T(i) = \frac{\exp(P_T(i)/c)}{\sum_j \exp(P_T(j)/c)} \quad (4)$$

This method offers several advantages. First, it allows for a more fine-grained adjustment of the probability distribution by controlling the temperature scaling parameter  $c$ , which can be optimized to achieve the lowest ECE values. Second, unlike label smoothing, temperature scaling can effectively balance the confidence levels of both top-1 and subsequent tokens, reducing both over-confidence and under-confidence issues. This results in a more consistent and reliable calibration across all tokens, thereby enhancing the overall trustworthiness of the knowledge. Additionally, we find that selecting the optimal  $c$  parameter on the validation set to maximize the knowledge significantly enhances the effectiveness of transferring trustworthy knowledge. The knowledge processed by using this  $c$  yields the best results for the student model (detailed in Sec. §5.5). Due to the low cost of selecting  $c$  on the validation set, we can tailor different  $c$  values for different tasks. This demonstrates "temperature scaling" excellent scalability and flexibility.

## 4.3 Knowledge Matching

After obtaining the re-calibrated probability data  $P_T$  that contains  $P_T(1), P_T(2), \dots, P_T(5)$ , we use the same training data to train the student model. Instead of utilizing language modeling loss on hard labels, the probabilities of the 5 tokens that correspond to the teacher's top-5 of the student model are retrieved as  $P_S$  which contains  $P_S(1), P_S(2), \dots, P_S(5)$ . Kullback-Leibler divergence is then used to measure the loss between the teacher model and the student model:

$$Loss(y_{1:N}) = \sum_{t=1}^N D_{KL}(P_T || P_S) \quad (5)$$

# 5 Experiment

## 5.1 Experimental Settings

Our experiments focus on both In-Domain and Out-of-Domain settings to ensure generalization abilities. In the **In-Domain setting**, we utilize CommonsenseQA (CSQA) (Talmor et al., 2019) and

	IN-DOMAIN						OUT-OF-DOMAIN					
	CSQA			BoolQ			CSQA			OBQA		
	<i>ECE</i> ↓	<i>Acc</i> ↑	<i>Trust</i> ↑	<i>ECE</i> ↓	<i>Acc</i> ↑	<i>Trust</i> ↑	<i>ECE</i> ↓	<i>Acc</i> ↑	<i>Trust</i> ↑	<i>ECE</i> ↓	<i>Acc</i> ↑	<i>Trust</i> ↑
LLAMA 1 : 33B → 7B												
Teacher <sub>33B</sub>	10.2	82.4	72.2	7.7	89.7	82	18.6	69.2	50.6	20.2	64.4	44.2
Fine-tune <sub>7B</sub>	11.8	79.9	68.1	6.5	82.5	76	12.5	48.2	35.7	21.9	43.4	21.5
Distill <sub>7B</sub>	9.4	78.9	69.5	<b>4.0</b>	85.3	81.3	5.3	43.1	37.8	18.1	39.8	21.7
Distill <sub>7B w/LS</sub>	9.1	78.1	69	19.0	85.3	66.3	5.2	43.9	38.7	19.0	37.6	18.6
<b>FIRST<sub>7B w/TS</sub></b>	<b>2.9</b>	<b>80.8</b>	<b>77.9</b>	<b>4.0</b>	<b>85.7</b>	<b>81.7</b>	<b>4.6</b>	<b>50.0</b>	<b>45.4</b>	<b>7.1</b>	<b>47.2</b>	<b>40.1</b>
FIRST to Fine-tune	↑8.9	↑0.9	↑9.8	↑2.5	↑3.2	↑5.7	↑7.9	↑1.8	↑8.7	↑14.8	↑3.8	↑18.6
LLAMA 2 : 13B → 7B												
Teacher <sub>13B</sub>	12.0	81.6	69.6	6.8	89.7	82.9	20.8	65.7	44.9	28.7	58.3	29.9
Fine-tune <sub>7B</sub>	14.0	76.8	62.8	8.4	87.5	79.1	21.2	50.0	28.8	30.1	45.6	15.5
Distill <sub>7B</sub>	10.9	80.0	69.1	4.0	85.3	81.3	7.7	50.9	43.2	12.5	46.6	34.1
Distill <sub>7B w/LS</sub>	10.3	<b>80.4</b>	70.1	3.9	87.5	83.6	7.5	51.1	43.6	16.2	47.6	31.4
<b>FIRST<sub>7B w/TS</sub></b>	<b>6.3</b>	80.3	<b>74</b>	<b>1.4</b>	<b>87.9</b>	<b>86.5</b>	<b>5.5</b>	<b>51.4</b>	<b>45.9</b>	<b>8.1</b>	<b>49.5</b>	<b>41.4</b>
FIRST to Fine-tune	↑7.7	↑3.5	↑11.2	↑7	↑0.4	↑7.4	↑15.7	↑1.4	↑17.1	↑22	↑3.9	↑25.9
OPENLLAMA : 13B → 7B												
Teacher <sub>13B</sub>	13.2	78.5	65.3	7.5	87.6	80.1	16.7	49.5	32.8	13.4	50.0	36.6
Fine-tune <sub>7B</sub>	10.5	75.0	64.5	3.6	81.5	77.9	21.6	28.3	6.7	16.1	30.4	14.3
Distill <sub>7B</sub>	9.2	75.2	66	6.2	83.8	77.6	9.7	27.7	18	13.7	29.8	16.1
Distill <sub>7B w/LS</sub>	9.6	74.5	65.9	3.3	83.3	80	4.1	29.2	25.1	14.2	29.8	15.6
<b>FIRST<sub>7B w/TS</sub></b>	<b>5.0</b>	<b>77.2</b>	<b>72.2</b>	<b>2.7</b>	<b>84.7</b>	<b>82</b>	<b>2.9</b>	<b>30.5</b>	<b>27.6</b>	<b>8.2</b>	<b>30.8</b>	<b>22.6</b>
FIRST to Fine-tune	↑5.5	↑2.2	↑7.7	↑0.9	↑3.2	↑4.1	↑18.7	↑2.2	↑20.9	↑7.9	↑0.4	↑8.3

Table 1: Smaller models obtained by our method **FIRST** consistently achieves high accuracy *Acc* across various scenarios while maintaining a low expected calibration error *ECE* (see Eq. 1). The higher trust scores *Trust* (see Eq. 2), the more trustworthy models are. Note that in the out-of-domain setting, we only obtain smaller models by fine-tuning or distilling on Alpaca, with CSQA and OBQA being unseen in this context, validating the generalizability of our approach. ↑ represents the larger the better while the ↓ means the smaller the better. **Bold** represents the best.

BoolQ (Clark et al., 2019) for both training and testing. In the **Out-of-Domain setting**, we fine-tune and distill smaller models on a commonly used instruction-following dataset, Alpaca (Taori et al., 2023), while, testing the models’ performance over unseen task CommonsenseQA (CSQA) and OpenBook QA (OBQA) (Mihaylov et al., 2018). This approach allows us to assess the generalization abilities of the smaller models on unseen tasks, simulating real-world scenarios where these models need to perform on unfamiliar tasks.

To ensure the practicality of our approach, we select three widely used model families for our experiments: Llama-1 (Touvron et al., 2023a), Llama-2 (Touvron et al., 2023b), and OpenLlama (Geng and Liu, 2023). In our experiments, we test four types of smaller models obtained through different methods:

1) **Fine-tune<sub>7B</sub>**: Obtained by using fine-tuning with hard labels.

2) **Distill<sub>7B</sub>**: Obtained by distillation methods without "knowledge trustworthy maximization". For a fair comparison with our approach, we also use the top-5 tokens as knowledge in the latter comparison.

3) **FIRST<sub>7B w/TS</sub>**: Obtained by our proposed

method, primarily using temperature scaling (TS, see Eq. 4) within the trustworthy maximization phase.

4) **Distill<sub>7B w/LS</sub>**: We also explore the use of label smoothing (LS, see Eq. 3) to show why we ultimately adopt TS over LS in "knowledge trustworthy maximization". In the latter experiments, we pick up the popular smoothing coefficient 0.1 follow previous works (Müller et al., 2020). Additionally, we also provide the performance of **Teacher** models. For further implementation details, please refer to the Appendix.

## 5.2 Experiment Results

Based on the results shown in Table 1, we draw the following conclusions:

- **Fine-tuning lead to catastrophic mis-calibration**: We observed that although fine-tuned smaller models achieve relatively high accuracy in both in-domain and out-of-domain settings, their ECE values are notably high, resulting in overall low trust scores and lower reliability. This mis-calibration phenomenon is particularly pronounced in out-of-domain scenarios. For instance, we observe that the ECE of the model fine-tuned on OpenLlama 7B in the out-of-domain CSQA task reaches 21.6%, while its accuracy

is only 28.3%, indicating that smaller models obtained through fine-tuning tend to be unreliable on tasks they have not been trained on. In real-world scenarios, when smaller models are privately deployed, they will inevitably encounter tasks they have not been trained for. In such cases, there would be a mismatch between their confidence and true likelihood. They might confidently provide incorrect answers and even continuously emphasize their incorrect responses, thereby misleading users. This clearly does not meet the criteria of a trustworthy model.

• **Distillation brings bad calibration as well:** Furthermore, distilled models without "Knowledge Trustworthy Maximization" show relatively bad calibration ability. For in-domain tasks, the distilled Llama-1 7B and Llama-2 7B have ECE values of 9.4% and 10.9% on CSQA, a mis-calibration level similar to fine-tuned models. And distilled model of OpenLlama shows even worse calibration than fine-tuned models on BoolQ. While for accuracy, it generally has an improvement over standard fine-tuning, but on some settings such as Llama-1 on CSQA, it also shows worse performance than fine-tuning. This suggests that direct distillation without further process the knowledge does not consistently lead to better calibration and performance.

• **Temperature Scaling outperforms Label Smoothing:** Here, we compare the results of different methods used in the "Knowledge Trustworthy Maximization" phase. It is evident that  $\text{FIRST}_{7B \text{ w/TS}}$  performs significantly better than  $\text{Distill}_{7B \text{ w/LS}}$ . In the In-domain setting of BoolQ, the ECE values of  $\text{FIRST}_{7B \text{ w/LS}}$  astonishingly reached 19.0%, significantly worse than  $\text{Distill}_{7B}$ , which does not apply any additional processing to the knowledge. This highlights that LS cannot deliver stable performance across all scenarios. In contrast,  $\text{FIRST}_{7B \text{ w/TS}}$  consistently achieves lower ECE in both in-domain and out-of-domain scenarios. Additionally, they attain better accuracy in most cases, resulting in the highest Trust scores.

### 5.3 Reliability Analysis

**Reliability Diagrams.** To enhance our analysis and facilitate better comparisons, we employ reliability diagrams in addition to metric-based evaluations. As depicted in Figure 4, the reliability diagrams are divided into 10 bins based on the model's confidence. The bars represent the ex-

pected accuracy within each bin, and the colors indicate whether the model is under-confident (red) or over-confident (green) within each bin. A perfectly calibrated model would have a straight diagonal line from the bottom left to the top right of such a diagram, indicating that the confidence level is exactly consistent with expected accuracy.

The  $\text{Fine-tune}_{7B}$  model exhibits catastrophic mis-calibration, primarily characterized by over-confidence in its predictions. This means that the model tends to assign higher confidence levels to its predictions than what is justified by their actual accuracy. Although the  $\text{Teacher}_{33B}$  model also suffers from over-confidence, its overall high accuracy results in a much higher trust score. Additionally, the  $\text{Distill}_{7B}$  model demonstrates slightly improved calibration compared to the  $\text{Fine-tune}_{7B}$  model. Remarkably, our  $\text{FIRST}_{7B}$  model outperforms the other models, including the teacher model. It exhibits noticeably less under-confidence and over-confidence, as indicated by the smaller areas of the red and green bars, respectively, and its proximity to the perfect calibration line.

### 5.4 Analysis of Top-5 Selection.

Figure 2 illustrates the disk space usage and cumulative probability coverage for knowledge selection ranging from the top-1 to the top-100 tokens. The blue line represents the average accumulated probabilities, while the shaded area indicates the range of probabilities. The green line shows the corresponding disk space required. The reasons we finally adopted top-5 are as follows:

1. **Efficient Probability Coverage:** The figure demonstrates that selecting the top-5 tokens covers over 95% of the total probability. This high coverage ensures that the majority of relevant knowledge is captured, making the distillation process effective.
2. **Minimal Disk Space Usage:** The green line indicates the disk space required for storing the selected tokens. By selecting only the top-5 tokens, we significantly reduce the storage requirements compared to selecting more tokens. This efficiency is crucial for offline distillation, where disk space can be a limiting factor.
3. **Balancing Trade-offs:** The Top-5 selection strikes a balance between maximizing probability coverage and minimizing disk space

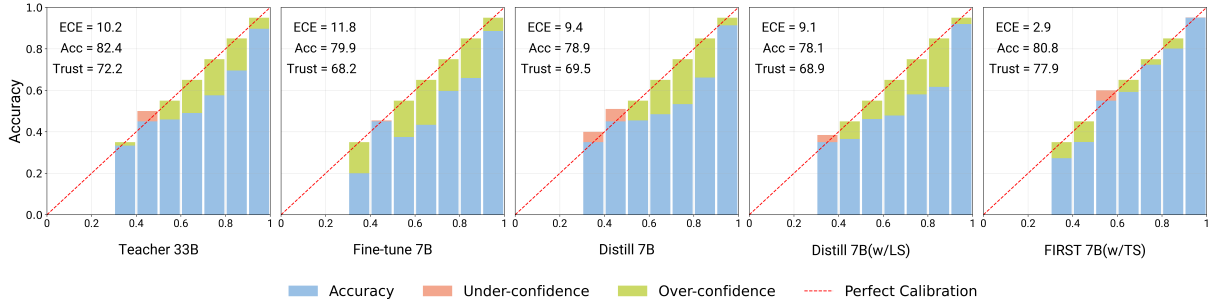


Figure 4: Reliability diagrams based on Llama-1 reveal the mis-calibration of various models on the CSQA dataset. In these diagrams, the X-axis is confidence divided into 10 bins, representing the model’s confidence levels for each question’s answer tokens. The Y-axis represents the accuracy within each bin. The red bar represents the degree to which the actual accuracy is higher than perfect calibration (under-confident), while the green bar means that the actual accuracy is lower than perfect calibration (over-confident).

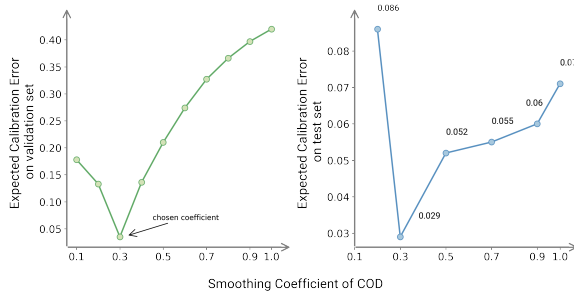


Figure 5: Left shows the comparison of different smoothing coefficients on the validation set, while the right part demonstrates its corresponding calibration effect on the test set.

usage. This balance ensures that the distilled knowledge is both comprehensive and storage-efficient, enabling practical implementation in various scenarios.

- Scalability:** Our method exhibits strong scalability. It is naturally extendable to distillation from models such as the GPT-3 series (text-davinci-003), which can only return top-5 token probabilities. This increases the range of LLMs that can be used as teacher models, allowing student models to be effectively trained even in semi-black box scenarios.

### 5.5 Temperature Scaling Parameter Analysis

As described in the section on Knowledge Trustworthy Maximization (Sec. §4.2), we employ a temperature scaling parameter to optimize the ECE (Expected Calibration Error) value on the validation set, as illustrated in the left part of Figure 5. We first divide the interval from 0 to 1 into steps of 0.1 and select the coefficient with the smallest ECE value. A larger coefficient results in all Top-5 tokens converging to the same probabilities, specifically 0.2. When the coefficient is set to 1, the probability of the top-1 token is dramatically com-

pressed, while the probabilities of the other tokens are enlarged accordingly. Conversely, a coefficient of 0.1 can even amplify the probabilities of over-confident tokens, leading to even worse calibration.

To further refine the search for the optimal smoothing coefficient, we narrow down the interval and use a smaller step size of 0.02. This allows us to pinpoint the best smoothing coefficient more precisely. Additionally, we compare the performance of FIRST using the selected optimal smoothing coefficient with other different smoothing coefficients as shown in the right part of Figure 5. FIRST with optimal smoothing coefficient do outperform those with other levels of smoothing coefficient with a large margin, indicating the effectiveness of selecting such optimal smoothing coefficient.

## 6 Conclusion

In conclusion, our proposed method, eFficient tRustworthy diSTillation (FIRST), effectively enhances both accuracy and calibration in large language models. By applying "trustworthy maximization", FIRST efficiently transfers the minimal yet most effective knowledge from teacher to student models. Experimental results show that FIRST consistently improves trustworthiness across various scenarios, demonstrating its potential to create reliable language models for practical applications.

## 7 Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. We address the critical issue of catastrophic mis-calibration in current training pipelines (supervised fine-tuning and knowledge distillation) and propose a pipeline to efficiently obtain a more trustworthy model. There are many potential societal consequences of our



621	work, none of which we feel must be specifically	Lewkowycz, Erica Moreira, Rewon Child, Oleksandr	672
622	highlighted here.	Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang,	673
623	<b>8 Limitations</b>	Brennan Saeta, Mark Diaz, Orhan Firat, Michele	674
624	It is shown that our efficient trustworthy distillation	Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas	675
625	(FIRST) demonstrates superior calibration ability	Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022.	676
626	and performance over direct distillation and stan-	<a href="#">Palm: Scaling language modeling with pathways.</a>	677
627	dard fine-tuning methods. However, despite these	Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom	678
628	exciting results, there are still some limitations to	Kwiatkowski, Michael Collins, and Kristina. Toutanova.	679
629	our current work, as well as potential opportunities	2019. <a href="#">Boolq: Exploring the surprising difficulty of</a>	680
630	for future research.	<a href="#">natural yes/no questions.</a> In <i>NAACL</i> .	681
631	<b>Extend to Large Teacher Model</b> : Due to the	Shizhe Diao, Rui Pan, Hanze Dong, Ka Shun Shum,	682
632	resource limitation, our largest teacher model is	Jipeng Zhang, Wei Xiong, and Tong Zhang. 2023. <a href="#">Lm-</a>	683
633	Llama 33B which is not very large but already	<a href="#">flow: An extensible toolkit for finetuning and inference</a>	684
634	achieving exciting results by distillation to a 7B	<a href="#">of large foundation models.</a>	685
635	student model. We expect that employing a large	Jesse Dodge, Gabriel Ilharco, Roy Schwartz, Ali	686
636	teacher model such as 70B can lead to better cali-	Farhadi, Hannaneh Hajishirzi, and Noah Smith. 2020.	687
637	bration ability and performance since a large model	<a href="#">Fine-tuning pretrained language models: Weight initial-</a>	688
638	learns a better distribution. However, we are unable	<a href="#">izations, data orders, and early stopping.</a>	689
639	to explore how very large teachers perform due to	Xinyang Geng and Hao Liu. 2023. <a href="#">Openllama: An open</a>	690
640	resource limitations.	<a href="#">reproduction of llama.</a>	691
641	<b>Top-K Chosen in Offline Distillation:</b> Another	Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. 2023.	692
642	limitation of this work is that it does not provide	<a href="#">Knowledge distillation of large language models.</a>	693
643	a rigorous study on how many token probabilities	Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Wein-	694
644	to choose for one entry is optimal for knowledge	berger. 2017. On calibration of modern neural networks.	695
645	distillation in large language models. Currently, we	In <i>Proceedings of the 34th International Conference</i>	696
646	consistently choose the top-5 token probability to	<i>on Machine Learning - Volume 70</i> , ICML'17, page	697
647	retrieve because of the reasons stated in §5.4. How-	1321–1330. JMLR.org.	698
648	ever, how much token probability to use is optimal	Namgyu Ho, Laura Schmid, and Se-Young Yun. 2023.	699
649	could be an important area for further exploration	<a href="#">Large language models are reasoning teachers.</a> In <i>Pro-</i>	700
650	and development.	<i>ceedings of the 61st Annual Meeting of the Association</i>	701
651	<b>References</b>	<i>for Computational Linguistics (Volume 1: Long Papers)</i> ,	702
652	Rishabh Agarwal, Nino Vieillard, Yongchao Zhou, Piotr	pages 14852–14882, Toronto, Canada. Association for	703
653	Stanczyk, Sabela Ramos, Matthieu Geist, and Olivier	Computational Linguistics.	704
654	Bachem. 2024. <a href="#">On-policy distillation of language mod-</a>	Cheng-Yu Hsieh, Chun-Liang Li, Chih-kuan Yeh,	705
655	<a href="#">els: Learning from self-generated mistakes.</a>	Hootan Nakhost, Yasuhisa Fujii, Alex Ratner, Ranjay	706
656	Aakanksha Chowdhery, Sharan Narang, Jacob De-	Krishna, Chen-Yu Lee, and Tomas Pfister. 2023. <a href="#">Distill-</a>	707
657	vlín, Maarten Bosma, Gaurav Mishra, Adam Roberts,	<a href="#">ing step-by-step! outperforming larger language models</a>	708
658	Paul Barham, Hyung Won Chung, Charles Sutton, Se-	<a href="#">with less training data and smaller model sizes.</a> In <i>Find-</i>	709
659	bastian Gehrmann, Parker Schuh, Kensen Shi, Sasha	<i>ings of the Association for Computational Linguistics:</i>	710
660	Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker	<i>ACL 2023</i> , pages 8003–8017, Toronto, Canada. Associ-	711
661	Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prab-	ation for Computational Linguistics.	712
662	hakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner	Ehsan Latif, Luyang Fang, Ping Ma, and Xiaoming	713
663	Pope, James Bradbury, Jacob Austin, Michael Isard,	Zhai. 2023. Knowledge distillation of llm for education.	714
664	Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm	<a href="#">arXiv preprint arXiv:2312.15842.</a>	715
665	Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk	Liunian Harold Li, Jack Hessel, Youngjae Yu, Xiang	716
666	Michalewski, Xavier Garcia, Vedant Misra, Kevin	Ren, Kai-Wei Chang, and Yejin Choi. 2023. <a href="#">Symbolic</a>	717
667	Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito,	<a href="#">chain-of-thought distillation: Small models can also</a>	718
668	David Luan, Hyeontaek Lim, Barret Zoph, Alexan-	<a href="#">“think” step-by-step.</a> <i>ArXiv</i> , abs/2306.14050.	719
669	der Spiridonov, Ryan Sepassi, David Dohan, Shivani	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish	720
670	Agrawal, Mark Omernick, Andrew M. Dai, Thanu-	Sabharwal. 2018. <a href="#">Can a suit of armor conduct electric-</a>	721
671	malayan Sankaranarayanan Pillai, Marie Pellat, Aitor	<a href="#">ity? a new dataset for open book question answering.</a> In	722
		<i>Proceedings of the 2018 Conference on Empirical Meth-</i>	723
		<i>ods in Natural Language Processing</i> , pages 2381–2391,	724
		Brussels, Belgium. Association for Computational Lin-	725
		guistics.	726

727	Lingbo Mo, Boshi Wang, Muhao Chen, and Huan Sun. 2024. <a href="#">How trustworthy are open-source llms? an assessment under malicious demonstrations shows their vulnerabilities.</a>	784
728		785
729		786
730		787
731	Rafael Müller, Simon Kornblith, and Geoffrey Hinton. 2020. <a href="#">When does label smoothing help?</a>	788
732		789
733	Hao Sha, Yao Mu, Yuxuan Jiang, Li Chen, Chenfeng Xu, Ping Luo, Shengbo Eben Li, Masayoshi Tomizuka, Wei Zhan, and Mingyu Ding. 2023. <a href="#">Languagempc: Large language models as decision makers for autonomous driving.</a> <i>arXiv preprint arXiv:2310.03026</i> .	790
734		791
735		792
736		793
737		794
738	Kumar Shridhar, Alessandro Stolfo, and Mrinmaya Sachan. 2023. <a href="#">Distilling reasoning capabilities into smaller language models.</a> In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 7059–7073, Toronto, Canada. Association for Computational Linguistics.	795
739		796
740		797
741		798
742		799
743		800
744	KaShun Shum, Shizhe Diao, and Tong Zhang. 2023. <a href="#">Automatic prompt augmentation and selection with chain-of-thought from labeled data.</a>	801
745		802
746		803
747		804
748		805
749	Mark Steyvers, Heliodoro Tejada, Aakriti Kumar, Catarina Belem, Sheer Karny, Xinyue Hu, Lukas Mayer, and Padhraic Smyth. 2024. <a href="#">The calibration gap between model and human confidence in large language models.</a> <i>arXiv preprint arXiv:2401.13835</i> .	806
750		807
751		808
752		809
753		810
754	Chi Sun, Xipeng Qiu, Yige Xu, and Xuanjing Huang. 2020. <a href="#">How to fine-tune bert for text classification?</a>	811
755		812
756		813
757		814
758		815
759	Lichao Sun, Yue Huang, Haoran Wang, Siyuan Wu, Qihui Zhang, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, Xiner Li, et al. 2024. <a href="#">Trustllm: Trustworthiness in large language models.</a> <i>arXiv preprint arXiv:2401.05561</i> .	816
760		817
761		818
762		819
763		820
764		821
765		822
766		823
767		824
768	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. <a href="#">CommonsenseQA: A question answering challenge targeting commonsense knowledge.</a> In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.	825
769		826
770		827
771		828
772		829
773		830
774		831
775		832
776		833
777		834
778		835
779		836
780		837
781		838
782		839
783		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

	STANDARD FINE-TUNING	DIRECT DISTILLATION	FIRST
Question	Which city is farther north, Oslo or Helsinki?		
Correct Answer	Helsinki		
Generated Confidence	Oslo is farther north than Helsinki. 0.92 → over-confident	Oslo is farther north than Helsinki. 0.83 → over-confident	Oslo is farther north than Helsinki. 0.52
Question	Is Donald Trump a Neo-con American politician and businessman for the Republicans, with a long and varied career?		
Correct Answer	No		
Generated Confidence	Yes. 0.91 → over-confident	Yes. 0.85 → over-confident	Yes. 0.54

Table 2: A case study on how fine-tuned model and direct distilled model tend to over-confident on the wrong answer with high confidence. While FIRST though outputs a wrong answer, it produces low confidence to show its uncertainty.

## A Detailed Experimental Setting

839

### A.1 Implementation Details

840

We train our models on 8 GPU (RTX A6000 48G) using the Adam optimizer with beta set to be [0.9, 0.999] and epsilon fixed to be 1e-6 and cosine annealing scheduler with a warm-up ratio of 0.03. For fine-tuning, we utilize LMFlow (Diao et al., 2023) package to obtain a well fine-tuned model by a standard 3-epoch training and control the batch size to be 32 on each GPU and the learning rate for teacher models to be 2e-5. For question-answering tasks, we follow Shum et al. (2023)’s format and fine-tune the model in a zero-shot setting. For out-of-domain tasks, we directly follow Alpaca’s (Taori et al., 2023) setting to obtain the fine-tuned model. In both settings, we make use of the next token strategy for inferencing and answer generation. Finally, for distillation, the batch size is set to 32 on each GPU and we train our model for 3 epochs, the last checkpoint is used for evaluation since it has the best performance.

841

842

843

844

845

846

847

848

849

## B Additional Analysis

850

### B.1 Case Study

851

We further conduct a case study to see whether FIRST indeed helps mitigate mis-calibration in real-world question answering. As shown in Table 2, we ask the models of three different tuning methods on Alpaca to answer the question: which city is farther north, Oslo or Helsinki? The correct answer is Helsinki and the wrong answer is Oslo.

852

853

854

855

From the output confidence, we can see that standard fine-tuned models and direct distillation give high confidence in the wrong answer, which is far from satisfactory for trustworthy in real-world settings, especially when additional post-processing procedures were expected to be applied to filter wrong answers by identifying unconfident responses. In comparison, FIRST greatly mitigates this mis-calibration by producing a confidence of around 50% which indicates the model is not sure about the generated answer, allowing systems to filter those undesirable answers by a hard confidence threshold.

856

857

858

859

860

861