

# CritiCal: Can Critique Help LLM Uncertainty or Confidence Calibration?

Qing Zong, Jiayu Liu, Tianshi Zheng, Chunyang Li, Baixuan Xu, Haochen Shi,  
WeiQi Wang, Zhaowei Wang, Chunkit Chan, Yangqiu Song

Department of Computer Science and Engineering, HKUST

{qzong, yqsong}@cse.ust.hk

## Abstract

Accurate confidence calibration in Large Language Models (LLMs) is critical for safe use in high-stakes domains, where clear verbalized confidence enhances user trust. Traditional methods that mimic reference confidence expressions often fail to capture the reasoning needed for accurate confidence assessment. We propose natural language critiques as a solution, ideally suited for confidence calibration, as precise gold confidence labels are hard to obtain and often require multiple generations. This paper studies how natural language critiques can enhance verbalized confidence, addressing: (1) *What to critique*: uncertainty (question-focused) or confidence (answer-specific)? Analysis shows confidence suits multiple-choice tasks, while uncertainty excels in open-ended scenarios. (2) *How to critique*: self-critique or critique calibration training? We propose **Self-Critique**, enabling LLMs to critique and optimize their confidence beyond mere accuracy, and **CritiCal**, a novel **Critique Calibration** training method that leverages natural language critiques to improve confidence calibration, moving beyond direct numerical optimization. Experiments show that CritiCal significantly outperforms Self-Critique and other competitive baselines, **even surpassing its teacher model, GPT-4o**, in complex reasoning tasks. CritiCal also shows robust generalization in out-of-distribution settings, advancing LLM’s reliability. <sup>1</sup>

## 1 Introduction

Confidence calibration is crucial in ensuring the reliability and trustworthiness of LLMs in high-stakes applications (Vashurin et al., 2025; Xia et al., 2025). As LLMs increasingly interact with humans, verbalized confidence, such as "..., and my confidence is 80%," allows for clearer communication of response certainty, fostering trust and effective collaboration (Lin et al., 2022; Xiong et al., 2024).

<sup>1</sup><https://github.com/HKUST-KnowComp/CritiCal>

## Can Critique Help Calibration?

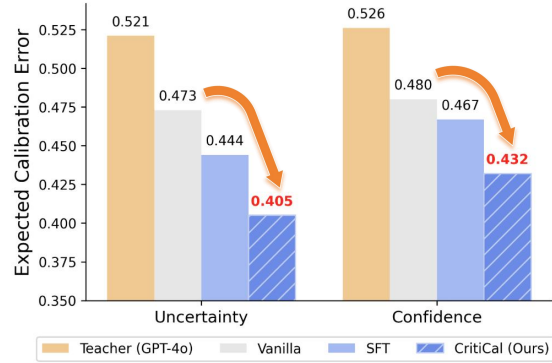


Figure 1: In-domain comparisons between CritiCal and other SFT methods by DeepSeek-R1-Distill-Qwen-7B on MATH-Perturb, showing CritiCal’s huge potential in improving LLM’s confidence calibration even with a teacher model having worse calibration performance.

Learning from critique (Wang et al., 2025b; Zhang et al., 2025) has proven highly effective in improving LLM’s accuracy. Natural language critiques clarify why answers are correct or incorrect, enabling more reasonable refinements rather than direct imitation of responses. This characteristic is ideal for confidence calibration, particularly verbalized confidence, as precise gold confidence labels are hard to obtain, but assessing whether confidence is too high or too low is straightforward based on reasoning and answer correctness. However, no related research has been conducted. This paper bridges this gap by investigating whether critique-based learning can improve uncertainty or confidence calibration, addressing two questions.

### What to critique: uncertainty or confidence?

Previous studies often treat uncertainty and confidence as antonyms, overlooking their distinction (Liu et al., 2025b): uncertainty pertains to the whole question, while confidence relates to the specific answer. Although Lin et al. (2024b) explored this difference using a consistency-based method, it was limited to the diversity of model

output and did not notice the difference between question types. We advance this by conducting a comprehensive study of LLMs’ direct outputs and their verbalized uncertainty and confidence. For brevity, we use "confidence" to broadly encompass both concepts, distinguishing them only when comparing their specific roles. *Extensive experiments show that verbalized confidence excels in multiple-choice questions, while uncertainty is better suited for open-ended tasks.*

**How to critique: self-critique or critique calibration training?** Unlike prior self-improvement methods (Huang et al., 2025c; Yang et al., 2025) that focus on accuracy, our **Self-Critique** approach targets confidence calibration. The model refines its confidence expression by analyzing the question, its reasoning steps, and final answer. But the results were unsatisfactory. Thus, we propose **CritiCal**, a supervised fine-tuning (SFT) **Critique Calibration** method that shifts from direct numerical optimization to critique-based learning (Wang et al., 2025b). During training, input consists of the question, model’s original response, and its confidence, while the output is a GPT-4o-generated critique of the confidence expression, based on the comparison of model’s reasoning process with a reference solution. Additionally, we explore replacing SFT with direct preference optimization (DPO) (Rafailov et al., 2023) for the training of CritiCal, using GPT-4o critiques as chosen responses and the model’s Self-Critique as rejected ones due to its suboptimal performance. Extensive experiments, both in-distribution and out-of-distribution, demonstrate that CritiCal significantly enhances confidence calibration for reasoning-intensive questions, *surpassing even GPT-4o’s calibration capabilities*, as is shown in Figure 1. *This suggests that a teacher model, with sufficient critique ability, can even enhance a student model’s confidence calibration beyond its own.* CritiCal also exhibits robust calibration improvements in out-of-distribution settings, with models trained on critique-suited data even outperforming those trained in-distribution, highlighting its exceptional transferability.

## 2 Related Works

### 2.1 Confidence Calibration

Confidence calibration methods for LLMs are divided into white-box and black-box approaches. White-box methods use internal model data, such as attention mechanisms (Lin et al., 2024a; Li et al.,

2023), hidden layers (Azaria and Mitchell, 2023), or token probabilities (Malinin and Gales, 2021; Zong et al., 2025) for precise confidence estimates. Conversely, black-box ones rely on model outputs without accessing internal structure. Consistency-based methods assess confidence by sampling multiple outputs and measuring their similarity, assuming consistent responses indicate higher certainty (Lin et al., 2024b; Huang et al., 2025a; Wang et al., 2024b; Su et al., 2024). Verbalization-based approaches train LLMs to explicitly express confidence through scores or epistemic markers (Li et al., 2025; Liu et al., 2025a; Zhang et al., 2024). SaySelf (Xu et al., 2024) uses a teacher model to generate reflective rationales and confidence scores by analyzing inconsistencies across numerous sampled reasoning chains. However, it focuses on imitating the reference reasoning and confidence expressions rather than learning from critiques of its own confidence and is computationally inefficient due to reliance on diverse outputs.

### 2.2 Critique Learning

Self-correction has recently emerged as a promising approach to enhance LLMs’ performance. Studies such as Madaan et al. (2023) and Welleck et al. (2023) utilize a model’s own feedback to refine outputs, though Huang et al. (2024), Zheng et al. (2025a) and Valmeekam et al. (2023) note limitations in its reliability for reasoning tasks. Alternatively, critique learning employs specialized models to provide feedback. Zhang et al. (2025) and Yang et al. (2024b) develop outcome-based reward models, while Wang et al. (2024a) and Lightman et al. (2024a) focus on process-based reward models to improve reasoning by evaluating intermediate steps. Damani et al. (2025) uses critique to train LLMs to reason about their uncertainty but is still limited to numerical critiques. Further work by Wang et al. (2025b) explicitly leverages natural language critiques as a training objective to encourage deeper understanding and reasoning. However, it focuses on using critique to improve LLM accuracy, whereas our work explores natural language critiques to improve confidence calibration.

## 3 Method

To investigate whether critique can enhance confidence calibration, we propose two methods: **Self-Critique**, a prompting-based approach, and **CritiCal**, a supervised fine-tuning (SFT) framework.

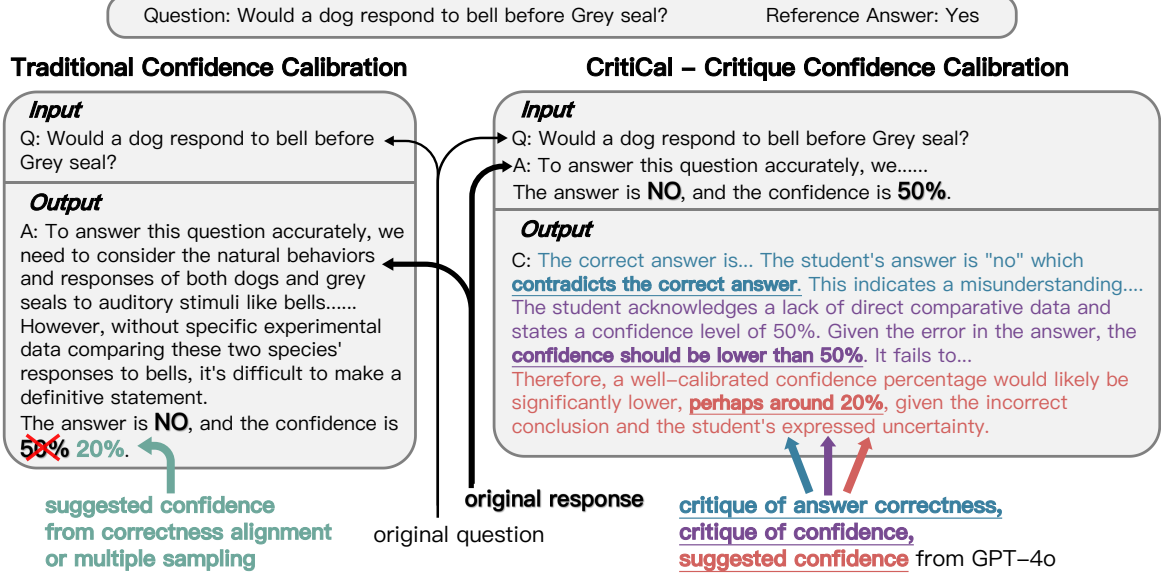


Figure 2: Comparisons between CritiCal and traditional confidence calibration methods.

### 3.1 Self-Critique

While prior studies on self-improvement (Huang et al., 2025c; Yang et al., 2025) focuses on refining reasoning processes and improving answer accuracy, **Self-Critique** targets confidence calibration, which aligns model’s verbalized confidence score with both the answer correctness and the uncertainty demonstrated in each reasoning steps. The model is prompted to reassess the question, its initial reasoning, and potential ambiguities or logical gaps, refining both the answer and confidence score to improve calibration. The detailed prompt is provided in Appendix A.

### 3.2 CritiCal

To further enable LLMs to express well-calibrated confidence aligned with their reasoning, we propose **CritiCal**, a SFT method that guides LLMs to refine their confidence expressions using critiques of their initial confidence scores.

As illustrated in Figure 2, CritiCal differs from traditional confidence calibration training methods (Zhang et al., 2024; Xu et al., 2024) in its input-output structure. In conventional methods, input is the original question, and output is the original model answer paired with a suggested confidence expression, which is derived from either the alignment of answer correctness or the generation probability of such an answer during multiple times of response generation. In contrast, CritiCal is a sampling-free approach that encourages LLMs to learn from their confidence estimation errors through critique-based training. Specifically, the

input consists of the question, the student model’s answer, and its associated confidence score, while the output is a critique from a teacher model (GPT-4o). This critique evaluates the calibration of the student’s confidence score, providing an explanation based on the clarity, strength, and correctness of the student’s reasoning compared to a reference solution.

In practice, we sample 2K questions from the training set and prompt the student model to generate answers along with confidence scores. These responses, paired with the questions and reference solutions from the benchmark, are provided to the teacher model to produce critiques assessing confidence calibration. The student model is then fine-tuned using the collected critique data. In particular, to mitigate knowledge shift in large reasoning models (LRMs), we instruct the teacher model to structure their critiques with special "</think>" tokens, separating the explanation from the final judgment. This structured critique format facilitates more effective learning. The detailed prompt for critique generation is provided in Appendix A.

## 4 Experiments

In this section, we answer the two questions: what to critique (§4.2) and how to critique (§4.3, §4.4).

### 4.1 Experimental Setup

**Datasets.** All the experiments involved a total of 7 datasets: TriviaQA (Joshi et al., 2017) with open-ended, single-hop factuality questions; ComparisonQA (Zong et al., 2025) with multiple-choice,

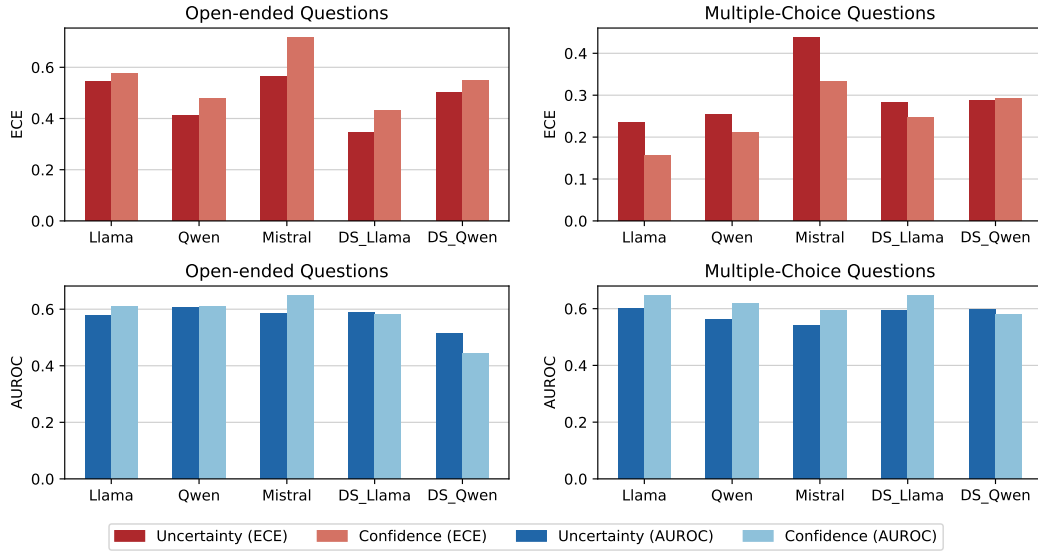


Figure 3: Mean ECE and AUROC values for each model across the same category of benchmarks. The dark bars are the result under uncertainty prompt, and the light ones are of confidence. Further analysis under the setting of multi-turn Self-Critique can be found in Appendix B.

single-hop factuality questions; StrategyQA (Geva et al., 2021) with yes/no, multi-hop factuality reasoning questions; HotpotQA (Yang et al., 2018) with open-ended, multi-hop factuality reasoning questions; MATH (Hendrycks et al., 2021) with open-ended, mathematical reasoning questions; MATH-500 (Lightman et al., 2024b) with harder ones selected from MATH test set; and MATH-Perturb (Huang et al., 2025b) with selected perturbed ones from MATH.

**Models.** Our test involves LLMs: LLaMA (Dubey et al., 2024), Qwen (Yang et al., 2024a), Mistral (Jiang et al., 2023), LLMs: DeepSeek-Distill-Llama, DeepSeek-Distill-Qwen (DeepSeek-AI et al., 2025), and a proprietary API: GPT-4o (OpenAI, 2024), for their diverse architectures.

**Metrics.** We use accuracy (via exact match for open-ended questions) for response correctness measurement, expected calibration error (ECE) for confidence-accuracy alignment, and area under the receiver operating characteristic curve (AUROC) for confidence-based discrimination of correct and incorrect responses. For both accuracy and AUROC, the higher the better, but ECE is the opposite.

## 4.2 Uncertainty vs. Confidence

Uncertainty, which pertains to the question as a whole, and confidence, which relates to the specific answer generated, are distinct concepts in LLMs, yet often mixed up in previous works (Liu et al., 2025b). To address this, we investigate their differences and performance across various scenarios.

To ensure models distinguish between uncer-

tainty and confidence, we provide clear definitions in the prompts, as detailed in Appendix A. We evaluate five models across six benchmarks, grouped into open-ended and multiple-choice question types, with the MATH benchmark reserved exclusively for training.

Figure 3 presents the mean ECE and AUROC for each model across benchmark categories. The results reveal distinct performance patterns for uncertainty and confidence across question types. For open-ended questions, uncertainty consistently achieves lower ECE in both LLMs and LLMs. Conversely, for multiple-choice questions, confidence outperforms uncertainty in both ECE and AUROC.

This interesting discovery indicates that models exhibit better uncertainty calibration for open-ended questions, likely due to the expansive prediction space, where uncertainty captures the question’s inherent ambiguity. In contrast, confidence is better calibrated for multiple-choice questions, where the limited options allow models to leverage elimination strategies, enabling more precise confidence estimates for specific choices despite potential uncertainty about the question.

## 4.3 Self-Critique Analysis

This section investigates the impact of Self-Critique on confidence calibration and average confidence scores across multiple iterations.

### 4.3.1 Multi-Turn Self-Critique

To comprehensively evaluate Self-Critique performance, we conduct multi-turn Self-Critique exper-



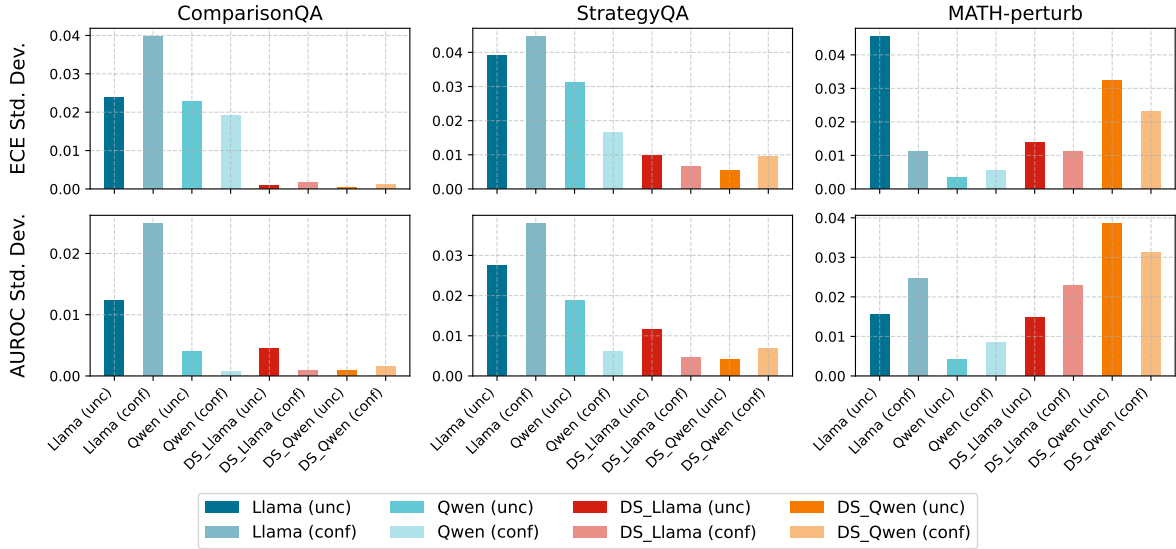


Figure 4: Standard deviation of multi-turn Self-Critique for ECE and AUROC across three benchmarks. Each bar represents the standard deviation of a model’s performance (uncertainty or confidence) across 6 iterations, where iteration 0 denotes the original response and iterations 1–5 indicate Self-Critique. Benchmarks are selected as representative of their task category due to question similarity under the same type. Full results are in Appendix B.

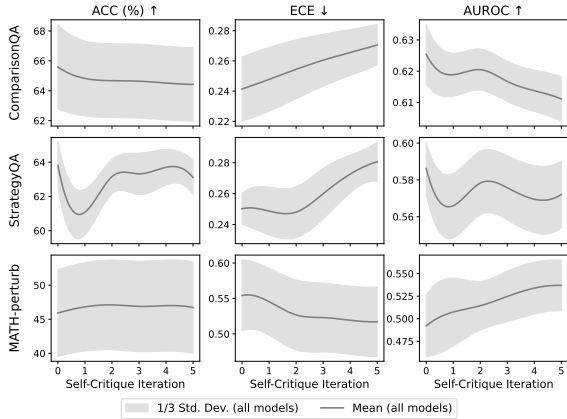


Figure 5: Multi-turn Self-Critique results on ComparisonQA, StrategyQA, and MATH-perturb benchmarks. Each plot shows the smoothed mean performance (solid line) and the corresponding 1/3 standard deviation range (shaded area) for ACC, ECE, and AUROC. Iteration 0 represents the original response without Self-Critique.

iments with 4 models, those have both reasoning and non-reasoning ones, across 6 benchmarks. In each iteration, the model receives the results of all previous iterations as context. Detailed prompt is provided in Appendix A.

Figure 4 illustrates the standard deviation of multi-turn Self-Critique for each model, focusing on ECE and AUROC to evaluate confidence calibration stability. Figure 5 presents the average performance of all models across three benchmarks, highlighting the impact of Self-Critique on different tasks. The specific variation curves of each model are shown in Figure 12 in Appendix B.

**Task Analysis.** The six benchmarks are catego-

rized into three tasks: one-hop factuality (ComparisonQA and TriviaQA), multi-hop factuality reasoning (StrategyQA and HotpotQA), and math reasoning (MATH500 and MATH-Perturb). As shown in Figure 5, the semi-transparent light gray area represents the average performance of all models with a one-third standard deviation. For accuracy, models exhibit greater stability on one-hop factuality and math reasoning tasks compared to multi-hop factuality reasoning. However, unlike prior self-improvement studies (Madaan et al., 2023; Welleck et al., 2023), Self-Critique shows no notable accuracy improvements, as it primarily targets confidence calibration rather than accuracy. For ECE and AUROC, Self-Critique exhibits relatively stable performance with slight improvements in math reasoning tasks. In contrast, factuality-related benchmarks experience negative impacts, with increased average ECE and decreased average AUROC. These findings suggest that Self-Critique has limited effectiveness, significantly worsening calibration for factuality-related tasks while only marginally enhancing it for math reasoning. Thus, prompting-based Self-Critique alone is inadequate for robust confidence calibration.

**Model Analysis.** In Figure 4 and Figure 12, LLMs are represented in cool colors, while LRMs are depicted in warm colors. For ECE and AUROC, LRMs demonstrate greater stability on factuality-related benchmarks, with significantly lower standard deviation than LLMs, whose calibration varies

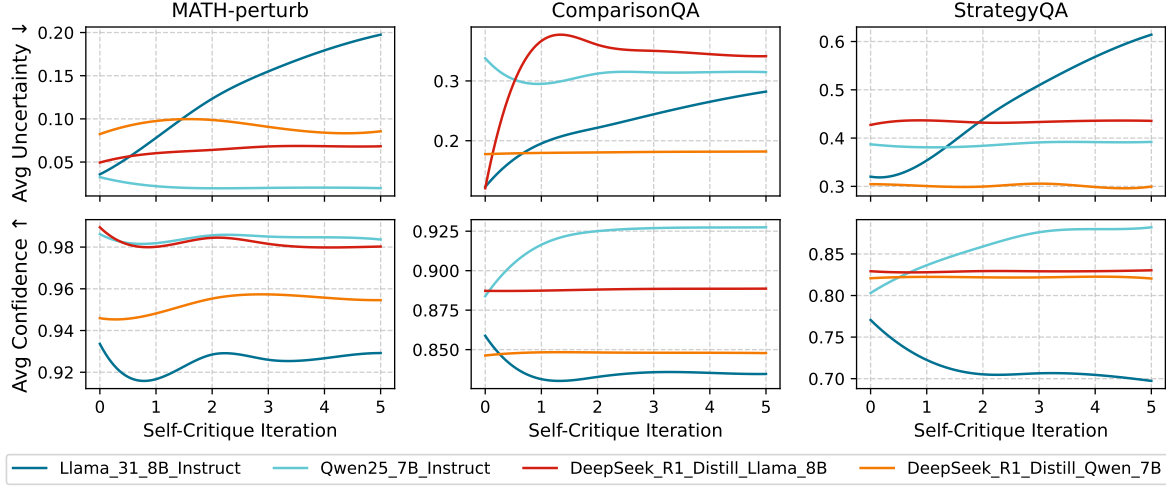


Figure 6: Curves of average uncertainty and confidence scores during multi-turn Self-Critique across 3 benchmarks.

widely. This stability in LRMs probably arises from their extended reasoning processes, enabling deeper reflection on initial responses and preventing erratic confidence shifts. Although LRMs show an increase in standard deviation in math-related questions, Figure 12 reveals that this stems from their progressively refined confidence calibration. Overall, LRMs exhibit more consistent and reliable confidence calibration compared to LLMs.

#### 4.3.2 Average Confidence Change during Self-Critique

Different from prior works (Huang et al., 2025c) finding models become more confident despite incorrect answers during self-improvement, Self-Critique focuses on refining confidence calibration, leading to more complex outcomes as it prioritizes confidence expression over answer correctness.

Results, shown in Figure 6, vary significantly by model. Llama consistently increases in uncertainty and decreases in confidence across all benchmarks, while Qwen shows the opposite trend, becoming more confident and less uncertain. This suggests that multi-turn Self-Critique amplifies these model-specific tendencies. The two DeepSeek distilled models generally maintain more consistent uncertainty and confidence scores compared to non-reasoning models, except for an increase in uncertainty of the distilled Llama on ComparisonQA. This indicates that extended reasoning processes enhance the robustness of confidence expressions.

### 4.4 CritiCal Analysis

We evaluate the performance of CritiCal in both in-distribution and out-of-distribution settings across multiple benchmarks.

We use the representative benchmark from each of the three tasks outlined in §4.3.1: one-hop factuality (ComparisonQA), multi-hop factuality reasoning (StrategyQA), and math reasoning (MATH-Perturb). For fair comparison, we randomly sample 2K questions from the training set to construct training data each time, using the method described in §3.2. For MATH-Perturb, where questions are perturbations of a subset of MATH, we sample from the original MATH training set to build training data and test only on perturbed questions from the original test set to prevent data leakage. Training is conducted using LlamaFactory (Zheng et al., 2024) with a batch size of 64 and other default hyperparameters, taking approximately half an hour per each dataset on a 45G single GPU.

#### 4.4.1 In Distribution

We first test the in-distribution performance of models fine-tuned with CritiCal, using Qwen and DeepSeek-Distill-Qwen as examples. Results are presented in Table 1.

For fair comparison, we include several sampling-free baselines: (1) **Vanilla**, a zero-shot prompt that directly asks model’s verbalized confidence (Xiong et al., 2024). (2) **Self-Critique**, the non-training method described in §3.1. (3) **SFT\_Hard**, a SFT approach using a suggested confidence score based on model’s original response (0% for incorrect answers, 100% for correct) for calibration (Zhang et al., 2024), with uncertainty as the inverse. (4) **SFT\_Soft**, a smoother SFT variant with confidence scores of 20% and 80%. (5) The performance of the teacher model, GPT-4o, is also included for reference.

Our key observations are as follows: (1) **Criti-**

Type	Method	Train	ComparisonQA			StrategyQA			MATH-Perturb		
			ACC (↑)	ECE (↓)	AUROC (↑)	ACC (↑)	ECE (↓)	AUROC (↑)	ACC (↑)	ECE (↓)	AUROC (↑)
GPT-4o											
Uncertainty	Vanilla	N	90.91	0.089	0.772	78.60	0.079	0.740	42.36	0.521	0.695
	Vanilla	N	91.97	0.036	0.787	79.48	0.103	0.716	44.54	0.526	0.683
Qwen-2.5-7B-Instruct (LLM)											
Uncertainty	Vanilla	N	69.65	<b>0.224</b>	0.615	64.63	0.283	0.507	39.57	0.587	0.525
	Self_Critique	N	68.24	0.268	0.605	<b>67.25</b>	0.308	0.464	40.00	0.583	0.542
	SFT_Hard	Y	69.49	0.229	0.616	65.07	0.288	0.537	36.52	0.605	0.554
	SFT_Soft	Y	69.68	0.228	0.615	64.19	0.245	0.564	38.70	0.593	0.558
	CritiCal	Y	<b>69.76</b>	<b>0.224</b>	<b>0.619</b>	<b>67.25</b>	<b>0.221</b>	<b>0.597</b>	<b>40.87</b>	<b>0.558</b>	<b>0.586</b>
Confidence	Vanilla	N	69.67	0.195	0.628	65.07	0.226	0.612	37.83	0.609	0.571
	Self_Critique	N	68.39	0.238	<b>0.630</b>	62.88	0.238	0.603	37.39	0.610	0.578
	SFT_Hard	Y	69.90	0.194	0.629	66.38	0.216	0.616	<b>41.30</b>	0.617	0.558
	SFT_Soft	Y	69.90	<b>0.193</b>	<b>0.630</b>	66.38	0.193	0.629	38.70	0.611	0.562
	CritiCal	Y	<b>69.97</b>	0.194	<b>0.630</b>	<b>69.00</b>	<b>0.179</b>	<b>0.644</b>	40.00	<b>0.588</b>	<b>0.593</b>
DeepSeek-R1-Distill-Qwen-7B (LRM)											
Uncertainty	Vanilla	N	52.18	0.331	0.586	58.52	0.247	<b>0.609</b>	62.54	0.473	0.380
	Self_Critique	N	52.12	0.330	<b>0.588</b>	59.83	0.242	0.604	64.87	0.491	0.383
	SFT_Hard	Y	<b>52.36</b>	<b>0.325</b>	0.578	59.83	0.281	0.558	65.65	0.446	0.413
	SFT_Soft	Y	52.51	<b>0.325</b>	0.580	62.45	0.272	0.516	66.09	0.444	0.437
	CritiCal	Y	52.30	0.326	0.579	<b>65.07</b>	<b>0.223</b>	0.572	<b>67.83</b>	<b>0.405</b>	<b>0.457</b>
Confidence	Vanilla	N	52.35	<b>0.326</b>	0.598	58.52	0.261	0.559	65.05	0.480	0.274
	Self_Critique	N	52.31	0.327	<b>0.602</b>	57.64	0.278	0.541	65.05	0.516	0.271
	SFT_Hard	Y	52.62	0.332	0.577	61.14	0.242	0.509	66.52	0.487	0.270
	SFT_Soft	Y	<b>52.66</b>	0.333	0.578	61.14	0.235	0.597	65.65	0.467	0.301
	CritiCal	Y	52.55	0.333	0.580	<b>66.81</b>	<b>0.176</b>	<b>0.630</b>	<b>69.13</b>	<b>0.432</b>	<b>0.328</b>

Table 1: Performance of various LLMs and LRMs on ComparisonQA, StrategyQA, and MATH-Perturb. The "Train" column indicates whether the method needs additional training, providing a fair comparison. The best performances among all methods are **bold-faced**.

**Cal excels in complex reasoning tasks.** Although CritiCal shows limited impact on ComparisonQA, it significantly improves calibration and accuracy on StrategyQA and MATH-Perturb, showing a huge decrease in ECE and increase in AUROC compared to all baselines, including Self-Critique. This improvement stems from the long structured reasoning processes elicited by multi-hop and math reasoning tasks, which provide robust cues for critiquing confidence calibration. **(2) CritiCal enables the student model to outperform even its teacher.** Notably, on MATH-Perturb, GPT-4o further reduces the ECE of DeepSeek-Distill-Qwen, a model whose ECE is already lower than its. This demonstrates that a teacher model, with sufficient critique capabilities, can continuously enhance a student model’s confidence calibration, highlighting CritiCal’s potential. **(3) Uncertainty and confidence distinctions persist in CritiCal.** Models trained with CritiCal maintain the pattern observed in §4.2: open-ended questions favor uncertainty, while multiple-choice questions favor confidence,

as evidenced by superior ECE and AUROC performance, indicating what to critique. **(4) Multi-hop factuality reasoning data is more suitable for critique than math reasoning.** CritiCal yields greater calibration improvements on StrategyQA than on MATH-Perturb, suggesting that factuality reasoning questions, with their explicit multi-hop reasoning steps, are more critique-suited.

#### 4.4.2 Out of Distribution

To evaluate CritiCal’s generalization, we focus on its performance in out-of-distribution (OOD) settings (Zheng et al., 2025b,c). Given CritiCal’s superior in-distribution performance on StrategyQA (as shown in Table 1), we train models on StrategyQA and test them on MATH-Perturb for OOD analysis. Results are presented in Table 2.

Both baseline fine-tuning methods (SFT\_Hard and SFT\_Soft) exhibit degraded performance on OOD data, indicating limited generalization. In contrast, CritiCal achieves improved calibration on OOD questions, with lower ECE and higher AUROC. This enhancement likely comes from Strat-

Method	Uncertainty						Confidence					
	In-distribution			Out-of-distribution			In-distribution			Out-of-distribution		
	ACC (↑)	ECE (↓)	AUROC (↑)	EM (↑)	ECE (↓)	AUROC (↑)	ACC (↑)	ECE (↓)	AUROC (↑)	EM (↑)	ECE (↓)	AUROC (↑)
<b>Qwen-2.5-7B-Instruct (LLM)</b>												
SFT_Hard	36.52	0.605	0.554	37.83	0.603	0.531	<b>41.30</b>	0.617	0.558	37.83	0.610	0.542
SFT_Soft	38.70	0.593	0.558	39.13	0.610	0.543	38.70	0.611	0.562	36.09	0.625	0.578
CritiCal	<b>40.87</b>	<b>0.558</b>	<b>0.586</b>	<b>39.57</b>	<b>0.574</b>	<b>0.595</b>	40.00	<b>0.588</b>	<b>0.593</b>	<b>42.17</b>	<b>0.571</b>	<b>0.593</b>
<b>DeepSeek-R1-Distill-Qwen-7B (LRM)</b>												
SFT_Hard	65.65	0.446	0.413	66.52	0.444	0.424	66.52	0.487	0.270	64.78	0.493	0.266
SFT_Soft	66.09	0.444	0.437	64.35	0.450	0.423	65.65	0.467	0.301	65.22	0.476	0.276
CritiCal	<b>67.83</b>	<b>0.405</b>	<b>0.457</b>	<b>67.83</b>	<b>0.375</b>	<b>0.465</b>	<b>69.13</b>	<b>0.432</b>	<b>0.328</b>	<b>69.13</b>	<b>0.434</b>	<b>0.350</b>

Table 2: Comparisons of CritiCal’s in-distribution and out-of-distribution performances. OOD Models are all trained on StrategyQA and tested on MATH-Perturb. The best performances among all methods are **bold-faced**.

Type	Method	ACC	ECE	AUROC
<b>StrategyQA (Multi-hop)</b>				
Uncertainty	CFT	67.25	0.221	0.597
	CPO	69.61	0.227	0.614
Confidence	CFT	69.00	0.179	0.644
	CPO	66.81	0.181	0.634
<b>ComparisonQA (One-hop)</b>				
Uncertainty	CFT	69.76	0.224	0.619
	CPO	69.61	0.227	0.614
Confidence	CFT	69.97	0.194	0.630
	CPO	69.94	0.192	0.630

Table 3: Comparisons of using SFT and DPO as the training method respectively for CritiCal.

egyQA’s critique-suited multi-hop reasoning data, which enables models to learn robust confidence calibration strategies based on their reasoning processes. These findings demonstrate CritiCal’s ability to foster reliable and generalizable confidence expressions across diverse tasks.

#### 4.4.3 Analysis of Training Method

We also explore another popular optimization method, DPO (Rafailov et al., 2023), for the training of CritiCal. While the input structure remains identical, DPO differs in its output, consisting of a chosen response, the same as SFT’s output, and a rejected response. For the rejected response, which should have a similar structure to the chosen one, we use the model’s Self-Critique output due to its suboptimal critique performance.

Since StrategyQA (multi-hop factuality reasoning) and MATH-Perturb (math reasoning) show similar performance trends in Table 1, we test only StrategyQA for multi-hop reasoning due to limited computing resources.

For clarity, we denote SFT-based CritiCal as CFT and DPO-based CritiCal as CPO, with results shown in Table 3. We can see that in both multi-hop and one-hop reasoning, the results of CFT and CPO differ very little compared to the huge improvement in Table 1. This suggests that CPO is also useful for reasoning-intensive tasks other than non-reasoning ones. Given DPO’s higher computational cost, SFT remains a sufficient and efficient training method for CritiCal.

## 5 Conclusions

This study investigates critique-based learning to enhance verbalized confidence calibration in LLMs, addressing two key questions: (1) What to critique. Our findings reveal that confidence expressions are better suited for multiple-choice tasks, while uncertainty is more effective for open-ended tasks, providing clear guidance for calibration strategies. (2) How to critique. We introduced Self-Critique, which enables LLMs to refine their own confidence assessments, and CritiCal, a novel critique calibration method that leverages natural language critiques from a teacher model to optimize calibration. Extensive experiments demonstrate that CritiCal significantly outperforms Self-Critique and other baselines, achieving superior calibration even beyond that of the teacher model, GPT-4o, in complex reasoning tasks. Moreover, CritiCal exhibits strong generalization, maintaining robust performance in both in-distribution and out-of-distribution settings, with notable transferability when trained on critique-suited multi-hop reasoning data. And compared to DPO, SFT is sufficient and efficient for CritiCal training. These findings underscore the potential of critique-based approaches to advance LLM reliability.



## Limitations

While CritiCal demonstrates significant improvements in confidence calibration for LLMs, several limitations still exist that cannot be covered in this single work.

The generalizability of CritiCal’s performance is potentially constrained by the specific benchmarks used in our experiments. Although we select diverse tasks (one-hop factuality, multi-hop factuality reasoning, and math reasoning), these benchmarks may not fully represent the broad range of real-world scenarios where LLMs are deployed, such as creative writing or multi-modality tasks. Further evaluation on a wider array of datasets could strengthen claims about CritiCal’s robustness.

Additionally, computational constraints restrict our ability to evaluate all benchmarks in the comparison of training methods, SFT and DPO, where only ComparisonQA and StrategyQA are tested. Although these benchmarks are carefully chosen to represent one-hop and multi-hop factuality reasoning tasks, this limitation may obscure potential variations in CritiCal’s effectiveness across other task types. Future work could leverage greater computational resources to conduct a more comprehensive analysis, incorporating additional benchmarks and training configurations.

## Ethics Statement

This paper utilizes several publicly available datasets, including ComparisonQA, TriviaQA, StrategyQA, HotpotQA, MATH, MATH-Perturb, and MATH-500, which are accessible to the research community under CC, Apache 2.0, MIT, Apache 2.0, MIT, MIT, and MIT licenses, respectively. The data is anonymized, ensuring our work does not raise privacy concerns regarding specific entities.

Our experiments involve the use of LLaMA, Qwen, Mistral, DeepSeek-Distill-Llama, DeepSeek-Distill-Qwen, and GPT-4o, so the same risks from LLMs research are also applicable to this work.

While CritiCal seeks to increase trust in AI by training on natural language critiques, there is a risk of users overly relying on its confidence estimates. These estimates may occasionally be inaccurate. Therefore, users are advised to treat these confidence expressions as a reference only.

## Acknowledgments

The authors of this paper were supported by the ITSP Platform Research Project (ITS/189/23FP) from ITC of Hong Kong, SAR, China, and the AoE (AoE/E-601/24-N), the RIF (R6021-20) and the GRF (16205322) from RGC of Hong Kong, SAR, China.

## References

- Amos Azaria and Tom M. Mitchell. 2023. [The internal state of an LLM knows when it’s lying](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 967–976. Association for Computational Linguistics.
- Mehul Damani, Isha Puri, Stewart Slocum, Idan Shenfelf, Leshem Choshen, Yoon Kim, and Jacob Andreas. 2025. [Beyond binary rewards: Training lms to reason about their uncertainty](#). *CoRR*, abs/2507.16806.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, and 81 others. 2025. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#). *CoRR*, abs/2501.12948.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, and 81 others. 2024. [The llama 3 herd of models](#). *CoRR*, abs/2407.21783.
- Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. [Did aristotle use a laptop? A question answering benchmark with implicit reasoning strategies](#). *Trans. Assoc. Comput. Linguistics*, 9:346–361.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. [Measuring mathematical problem solving with the MATH dataset](#). In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*.
- Chengsong Huang, Langlin Huang, Jixuan Leng, Jiacheng Liu, and Jiaxin Huang. 2025a. [Efficient test-time scaling via self-calibration](#). *CoRR*, abs/2503.00031.
- Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. 2024. [Large language](#)

- models cannot self-correct reasoning yet. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Kaixuan Huang, Jiacheng Guo, Zihao Li, Xiang Ji, Jiawei Ge, Wenzhe Li, Yingqing Guo, Tianle Cai, Hui Yuan, Runzhe Wang, Yue Wu, Ming Yin, Shange Tang, Yangsibo Huang, Chi Jin, Xinyun Chen, Chiyuan Zhang, and Mengdi Wang. 2025b. [Mathperturb: Benchmarking llms' math reasoning abilities against hard perturbations](#). *CoRR*, abs/2502.06453.
- Liangjie Huang, Dawei Li, Huan Liu, and Lu Cheng. 2025c. [Beyond accuracy: The role of calibration in self-improving large language models](#). *CoRR*, abs/2504.02902.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. 2023. [Mistral 7b](#). *CoRR*, abs/2310.06825.
- Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. 2017. [Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, pages 1601–1611. Association for Computational Linguistics.
- Kenneth Li, Oam Patel, Fernanda B. Vi  gas, Hanspeter Pfister, and Martin Wattenberg. 2023. [Inference-time intervention: Eliciting truthful answers from a language model](#). In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.
- Yibo Li, Miao Xiong, Jiaying Wu, and Bryan Hooi. 2025. [Conftuner: Training large language models to express their confidence verbally](#). *CoRR*, abs/2508.18847.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2024a. [Let's verify step by step](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2024b. [Let's verify step by step](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. [Teaching models to express their uncertainty in words](#). *Trans. Mach. Learn. Res.*, 2022.
- Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2024a. [Contextualized sequence likelihood: Enhanced confidence scores for natural language generation](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 10351–10368. Association for Computational Linguistics.
- Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2024b. [Generating with confidence: Uncertainty quantification for black-box large language models](#). *Trans. Mach. Learn. Res.*, 2024.
- Jiayu Liu, Junhao Tang, Hanwen Wang, Baixuan Xu, Haochen Shi, Weiqi Wang, and Yangqiu Song. 2024. [GProofT: A multi-dimension multi-round fact checking framework based on claim fact extraction](#). In *Proceedings of the Seventh Fact Extraction and Verification Workshop (FEVER)*, pages 118–129, Miami, Florida, USA. Association for Computational Linguistics.
- Jiayu Liu, Qing Zong, Weiqi Wang, and Yangqiu Song. 2025a. [Revisiting epistemic markers in confidence estimation: Can markers accurately reflect large language models' uncertainty?](#) In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pages 206–221. Association for Computational Linguistics.
- Xiaou Liu, Tiejun Chen, Longchao Da, Chacha Chen, Zhen Lin, and Hua Wei. 2025b. [Uncertainty quantification and confidence calibration in large language models: A survey](#). *CoRR*, abs/2503.15850.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. [Self-refine: Iterative refinement with self-feedback](#). In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.
- Andrey Malinin and Mark J. F. Gales. 2021. [Uncertainty estimation in autoregressive structured prediction](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- OpenAI. 2024. [Hello gpt-4o](#). OpenAI.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Ermon, and Chelsea Finn. 2023. [Direct preference optimization: Your language model is secretly a reward model](#). In *Advances in Neural Information Processing Systems 36: Annual*

*Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023.*

- Jiayuan Su, Jing Luo, Hongwei Wang, and Lu Cheng. 2024. [API is enough: Conformal prediction for large language models without logit-access](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024*, pages 979–995. Association for Computational Linguistics.
- Karthik Valmeekam, Matthew Marquez, and Subbarao Kambhampati. 2023. [Can large language models really improve by self-critiquing their own plans?](#) *CoRR*, abs/2310.08118.
- Roman Vashurin, Ekaterina Fadeeva, Artem Vazhentsev, Lyudmila Rvanova, Daniil Vasilev, Akim Tsvigun, Sergey Petrakov, Rui Xing, Abdelrahman Boda Sadallah, Kirill Grishchenkov, Alexander Panchenko, Timothy Baldwin, Preslav Nakov, Maxim Panov, and Artem Shelmanov. 2025. [Benchmarking uncertainty quantification methods for large language models with lm-polygraph](#). *Trans. Assoc. Comput. Linguistics*, 13:220–248.
- Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang Sui. 2024a. [Math-shepherd: Verify and reinforce llms step-by-step without human annotations](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2024, Bangkok, Thailand, August 11-16, 2024, pages 9426–9439. Association for Computational Linguistics.
- Rui Wang, Qihan Lin, Jiayu Liu, Qing Zong, Tianshi Zheng, Weiqi Wang, and Yangqiu Song. 2025a. [Prospect theory fails for llms: Revealing instability of decision-making under epistemic uncertainty](#). *CoRR*, abs/2508.08992.
- Yubo Wang, Xiang Yue, and Wenhui Chen. 2025b. [Critique fine-tuning: Learning to critique is more effective than learning to imitate](#). *CoRR*, abs/2501.17703.
- Zhiyuan Wang, Jinhao Duan, Lu Cheng, Yue Zhang, Qingni Wang, Xiaoshuang Shi, Kaidi Xu, Heng Tao Shen, and Xiaofeng Zhu. 2024b. [Conu: Conformal uncertainty in large language models with correctness coverage guarantees](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024*, pages 6886–6898. Association for Computational Linguistics.
- Sean Welleck, Ximing Lu, Peter West, Faeze Brahman, Tianxiao Shen, Daniel Khashabi, and Yejin Choi. 2023. [Generating sequences by learning to self-correct](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.
- Zhiqiu Xia, Jinxuan Xu, Yuqian Zhang, and Hang Liu. 2025. [A survey of uncertainty estimation methods on large language models](#). In *Findings of the Association for Computational Linguistics, ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pages 21381–21396. Association for Computational Linguistics.
- Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. 2024. [Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Tianyang Xu, Shujin Wu, Shizhe Diao, Xiaozhe Liu, Xingyao Wang, Yangyi Chen, and Jing Gao. 2024. [Sayself: Teaching llms to express confidence with self-reflective rationales](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 5985–5998. Association for Computational Linguistics.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jixai Yang, Jingren Zhou, Junyang Lin, Kai Dang, and 22 others. 2024a. [Qwen2.5 technical report](#). *CoRR*, abs/2412.15115.
- An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang Ren, and Zhenru Zhang. 2024b. [Qwen2.5-math technical report: Toward mathematical expert model via self-improvement](#). *CoRR*, abs/2409.12122.
- Zhe Yang, Yichang Zhang, Yudong Wang, Ziyao Xu, Junyang Lin, and Zhifang Sui. 2025. [Confidence v.s. critique: A decomposition of self-correction capability for llms](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2025, Vienna, Austria, July 27 - August 1, 2025, pages 3998–4014. Association for Computational Linguistics.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. [Hotpotqa: A dataset for diverse, explainable multi-hop question answering](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 2369–2380. Association for Computational Linguistics.
- Hanning Zhang, Shizhe Diao, Yong Lin, Yi R. Fung, Qing Lian, Xingyao Wang, Yangyi Chen, Heng Ji, and Tong Zhang. 2024. [R-tuning: Instructing large language models to say 'i don't know'](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, NAACL 2024, Mexico City, Mexico,

## Appendices

### A Prompt

Xiaoying Zhang, Hao Sun, Yipeng Zhang, Kaituo Feng, Chaochao Lu, Chao Yang, and Helen Meng. 2025. Critique-grpo: Advancing LLM reasoning with natural language and numerical feedback. *CoRR*, abs/2506.03106.

Tianshi Zheng, Yixiang Chen, Chengxi Li, Chunyang Li, Qing Zong, Haochen Shi, Baixuan Xu, Yangqiu Song, Ginny Y. Wong, and Simon See. 2025a. The curse of cot: On the limitations of chain-of-thought in in-context learning. *CoRR*, abs/2504.05081.

Tianshi Zheng, Zheyang Deng, Hong Ting Tsang, Weiqi Wang, Jiaxin Bai, Zihao Wang, and Yangqiu Song. 2025b. From automation to autonomy: A survey on large language models in scientific discovery. *CoRR*, abs/2505.13259.

Tianshi Zheng, Kelvin Kiu-Wai Tam, Newt Hue-Nam K. Nguyen, Baixuan Xu, Zhaowei Wang, Jiayang Cheng, Hong Ting Tsang, Weiqi Wang, Jiaxin Bai, Tianqing Fang, Yangqiu Song, Ginny Y. Wong, and Simon See. 2025c. Newtonbench: Benchmarking generalizable scientific law discovery in llm agents.

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, and Yongqiang Ma. 2024. Llamafactory: Unified efficient fine-tuning of 100+ language models. *CoRR*, abs/2403.13372.

Qing Zong, Zhaowei Wang, Tianshi Zheng, Xiyu Ren, and Yangqiu Song. 2025. Comparisonqa: Evaluating factuality robustness of llms through knowledge frequency control and uncertainty. In *Findings of the Association for Computational Linguistics, ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pages 4101–4117. Association for Computational Linguistics.

We design our prompt according to previous works (Xiong et al., 2024; Liu et al., 2024; Wang et al., 2025a). Figure 7, 8, 9, 10 illustrate the prompt we use for vanilla uncertainty inquiry, vanilla confidence inquiry, Self-Critique with confidence, and critique generation, respectively, with StrategyQA as an example.

#### Vanilla prompt using uncertainty

Answer the following yes/no question and provide your uncertainty score. Your response should end with 'The answer is [your\_answer], and the uncertainty is [uncertainty\_percentage]%' where [your\_answer] is yes or no, and the uncertainty percentage is a number between 0 and 100, indicating how uncertain you are about the question. If you are not sure, you should give a higher uncertainty percentage.

Question: [Question]

Figure 7: The vanilla prompt using uncertainty on StrategyQA. Placeholders [Question] will be replaced with the real one.

#### Vanilla prompt using confidence

Answer the following yes/no question and provide your confidence score. Your response should end with 'The answer is [your\_answer], and the confidence is [confidence\_percentage]%' where [your\_answer] is yes or no, and the confidence percentage is a number between 0 and 100, indicating how sure you are about your answer. If you are not sure, you should give a lower confidence percentage.

Question: [Question]

Figure 8: The vanilla prompt using confidence on StrategyQA. Placeholders [Question] will be replaced with the real one.



### Multi-turn Self-Critique prompt using confidence on StrategyQA

You previously answered the following yes/no question, and your responses have gone through one or more rounds of refinement. Below is the question, your initial response, and all subsequent refined responses. Now, reassess the question and your previous reasoning, including the initial and all refined responses. Consider any potential ambiguities, logical steps, or overlooked aspects that could improve the accuracy of your response and the calibration of your confidence score. Answer the question and provide a new confidence score.

Question: *[Question]*

Initial response: *[Initial\_Responses]*

Refined responses: *[Refined\_Responses]*

Your response should end with 'The refined answer is [your\_answer], and the confidence is [confidence\_percentage]%' where [your\_answer] is yes or no, and the confidence percentage is a number between 0 and 100, indicating how sure you are about your refined answer. If you are not sure about your refined answer, you should give a lower confidence percentage.

Figure 9: The prompt for multi-turn Self-Critique using confidence on StrategyQA. Placeholders *[Question]*, *[Initial\_Responses]*, and *[Refined\_Responses]* will be replaced with the real ones.

## B Detailed Self-Critique Results

Figure 11 shows the difference between uncertainty and confidence after Self-Critique. The distinctions between these two concepts still remain evident after applying Self-Critique.

Figure 12 displays the performance trajectories of each model across all six benchmarks. Self-Critique demonstrates relatively stronger improvements on mathematical reasoning tasks but falls short on factuality-related tasks, highlighting its limitations and lack of robustness.

### Critique generation prompt on StrategyQA

Confidence indicates how how sure the student is about his answer. If he is not sure, he should give a lower confidence percentage. You are a teacher expert in confidence calibration. A student previously answered a question and provided his confidence score. Please evaluate the calibration of his confidence score for the question based on his response. If his response is incorrect, the confidence percentage should be low.

Question: *[Question]*

Correct Answer: *[Correct\_Answer]*

Facts: *[Facts]*

Student's Response: *[Student's\_Response]*

Using the facts and the correct answer as a reference, assess whether the confidence percentage in the student's response is well-calibrated, considering the clarity and strength of the reasoning provided and your own knowledge of the question. Is the confidence percentage appropriate, too high, or too low? Provide a brief explanation of your evaluation, focusing on how well his confidence aligns with the strength of his reasoning and the context of the question.

Figure 10: The prompt we use to generate confidence calibration critique on StrategyQA. Placeholders *[Question]*, *[Correct\_Answer]*, *[Facts]*, and *[Student's\_Response]* will be replaced with the real ones.

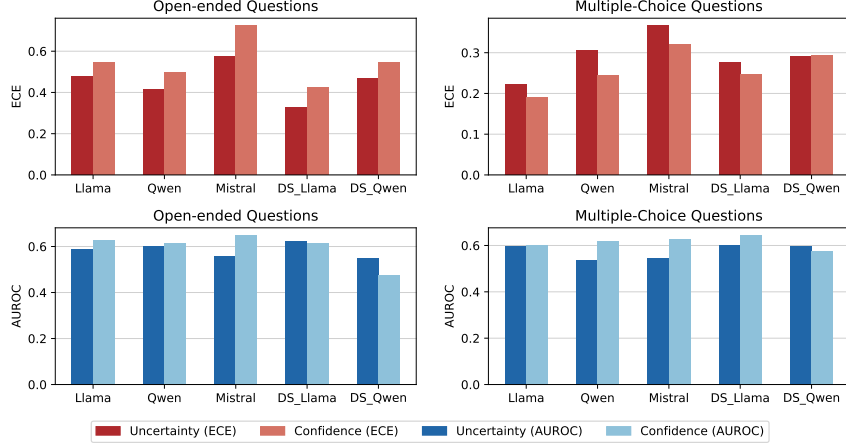


Figure 11: Mean ECE and AUROC values for each model across the same category of benchmarks, which are taken the average of across the 5 turns of Self-Critique. The dark bars are the result under uncertainty prompt, and the light ones are of confidence.

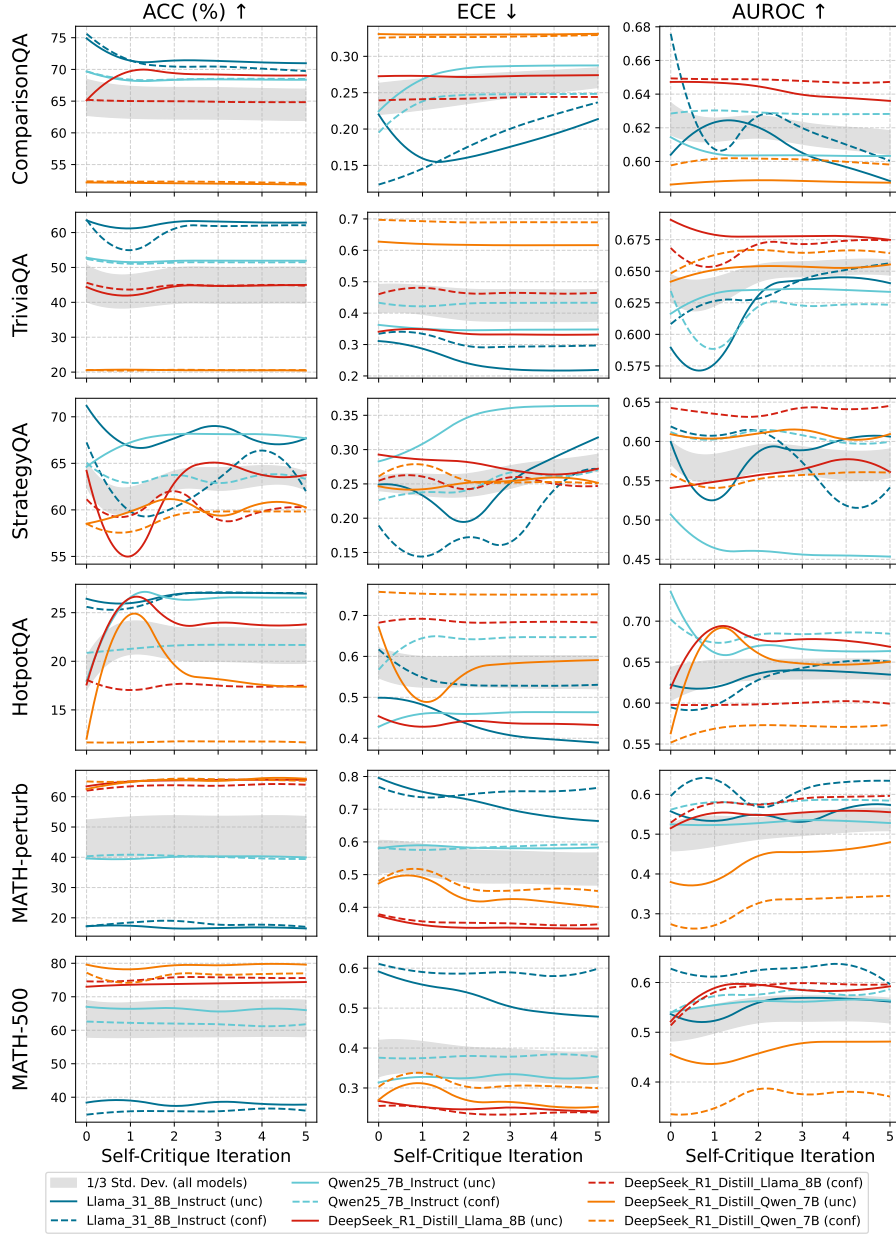


Figure 12: Multi-turn Self-Critique results on all the six benchmarks. The 0 iteration means the original response without Self-Critique. The semi-transparent light gray area represents the average performance of all models with a one-third standard deviation.