

CritiCal: Can Natural Language Critiques Help LLM’s Uncertainty or Confidence Calibration?

Anonymous ACL submission

Abstract

Confidence calibration is critical for safe use of Large Language Models (LLMs), where clear verbalized confidence enhances user trust. Traditional methods that mimic reference confidence expressions often fail to utilize the logic in model’s original reasoning chain. We propose natural language critique as a solution, which is ideal for confidence calibration, as precise gold confidence labels are hard to obtain and often require multiple generations but assessing whether the confidence is appropriate is easy by analyzing its internal logic and answer correctness. This paper studies how natural language critiques can enhance verbalized confidence: (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, and **CritiCal**, using natural language **Critiques** to train confidence **Calibration**. Experiments show that CritiCal significantly outperforms Self-Critique and other competitive baselines, **even surpassing its teacher, GPT-4o**, in complex reasoning tasks. CritiCal also shows robust generalization in out-of-distribution settings, proving its reliability.¹

1 Introduction

Confidence calibration is crucial in ensuring the trustworthiness of LLMs in high-stakes applications (Vashurin et al., 2025; Xia et al., 2025). As LLMs increasingly interact with humans, verbalized confidence (e.g., "I am 80% confident") clarifies response certainty, fostering trust and effective collaboration (Lin et al., 2022; Xiong et al., 2024).

Critique Fine-Tuning (CFT) (Wang et al., 2025) has proven highly effective in improving LLM’s accuracy. It makes LLMs learn from natural language

¹Data and code will be released upon acceptance.

Can Critique Help Calibration?

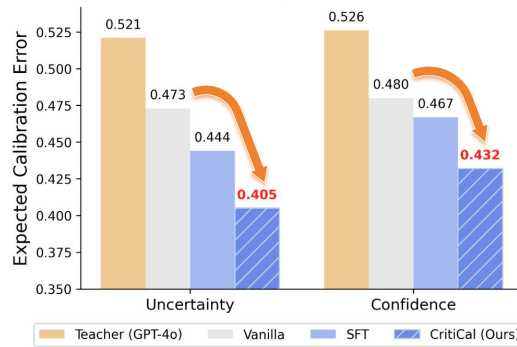


Figure 1: Comparisons between CritiCal and traditional SFT by DS-Qwen-7B on MATH-Perturb, showing CritiCal’s huge potential in improving LLM’s confidence calibration, since student can even outperform its teacher.

critiques which clarify why answers are correct or incorrect, enabling more reasonable and fine-grained refinements rather than direct imitation of gold responses. We find this characteristic ideal for confidence calibration, as precise gold confidence labels are hard to obtain, but assessing whether the confidence is appropriate is easy based on answer correctness and the logic shown in reasoning. So, this paper investigates whether natural language critiques can help uncertainty or confidence calibration, addressing the following two questions.

What to critique: uncertainty or confidence?

Previous studies often treat uncertainty and confidence as antonyms, overlooking their distinction (Liu et al., 2025b): While confidence evaluates the reliability of a specific response, uncertainty captures the inherent difficulty or ambiguity of the question itself. Crucially, uncertainty is invariant to the choice of answer for a given question. 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 confi-

dence. 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? In contrast to prior methods (Huang et al., 2025c; Yang et al., 2025) that focus on self-improving accuracy in the absence of external references, our **Self-Critique** approach also refines LLM’s confidence expressions by systematically analyzing the question, the reasoning chain, and the final answer. But we find that this zero-shot self-correction often yields unsatisfactory performance. Thus, we further introduce **CritiCal**, a supervised fine-tuning (SFT) method that uses natural language **Critiques** to train **Calibration**. During training, the model is prompted with the question and its own initial response; the target output is a GPT-4o-generated critique of the model’s confidence, based on the comparison between model’s reasoning chain and a reference solution. Additionally, we explore replacing SFT with direct preference optimization (DPO) (Rafailov et al., 2023) for the training of CritiCal, utilizing GPT-4o’s critiques as chosen responses and its own suboptimal Self-Critiques as rejected ones. Extensive experiments, both in-distribution and out-of-distribution, demonstrate that CritiCal significantly enhances confidence calibration for reasoning-intensive questions, *surpassing even its teacher, GPT-4o, in 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.* Finally, CritiCal exhibits robust transferability. In certain cases, models trained on critique-suited data even outperform those trained in-distribution, highlighting CritiCal’s generalizability to unseen domains.

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 information, 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 (Lin et al., 2024b; Huang et al., 2025a; Wang et al., 2024b; Su et al., 2024) assess confidence by sampling multiple outputs and measuring their similarity, assuming consistent responses indicate higher certainty. Verbalization-based approaches (Li et al., 2025; Liu et al., 2025a; Zhang et al., 2024) train LLMs to explicitly express confidence through scores or epistemic markers. 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 model’s own feedback to refine outputs, though Huang et al. (2024) and Valmeekam et al. (2023) note limitations in its reliability for reasoning tasks. Alternatively, critique learning uses specialized models to provide feedback. Zhang et al. (2025) and Yang et al. (2024b) develop outcome-based reward, while Wang et al. (2024a) and Lightman et al. (2024a) focus on process reward to improve reasoning by evaluating intermediate steps. Further work by Wang et al. (2025) 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. Damani et al. (2025) uses critique to train LLMs to reason about their uncertainty but is still limited to numerical critiques. Nair and Sinaga (2025) proves that learning calibrated confidence is impossible with only binary supervision. This motivates our exploration of natural language critiques in confidence calibration.

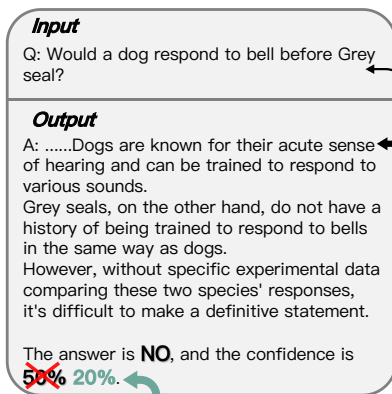
3 Method

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.

Question: Would a dog respond to bell before Grey seal?

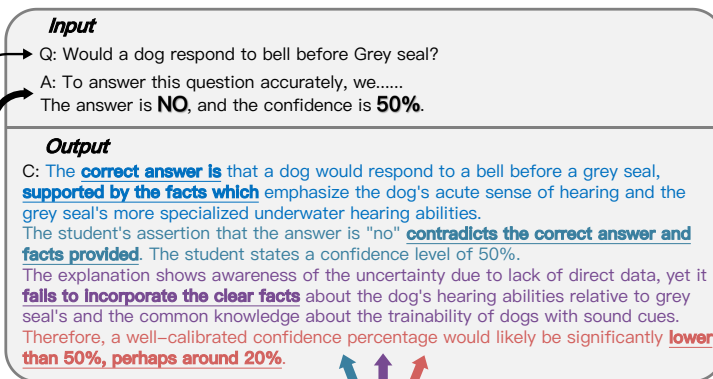
Reference Answer: Yes

Traditional Confidence Calibration



suggested confidence from correctness alignment or consistency in multiple sampling

CritiCal – Critique Confidence Calibration



summary of reference answer and facts
critique of answer correctness
critique of confidence
suggested confidence

from teacher model

original response
original question

Figure 2: Unlike traditional calibration methods that focus solely on the final answer, CritiCal introduces a teacher model to evaluate the student’s reasoning steps. By critiquing the logical path alongside facts and reference answers, it provides a more nuanced assessment of confidence than standard error-based approaches.

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 the reasoning chain. 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**, an SFT method that guides LLMs to refine their confidence expressions using critiques of their initial confidence expressions.

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 previous 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 and the student model’s original response, 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 each 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, for large reasoning models (LRMs), we instruct the teacher to structure their critiques with special "`</think>`" tokens, separating the explanation from the final judgment, to mitigate knowledge shift. 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; Compar-

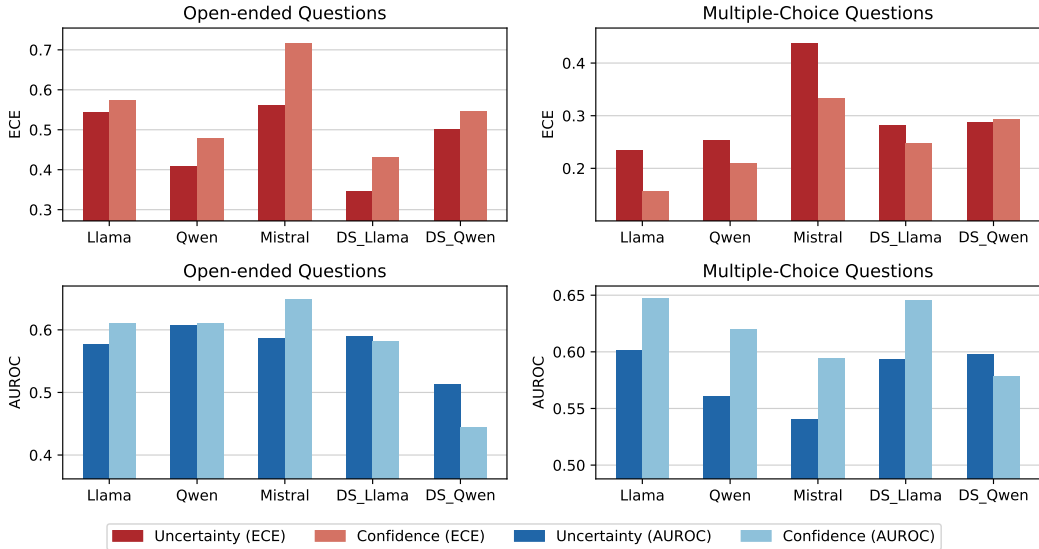


Figure 3: Mean ECE and AUROC values for each model across the same category of benchmarks. The dark bars are the result under uncertainty prompt, and the light ones are of confidence. Further analysis under the setting of multi-turn Self-Critique can be found in Appendix B.

isonQA (Zong et al., 2025) with multiple-choice, 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. (Appendix C)

Models. Our test involves LLMs: Llama-3.1-8B-Instruct (Dubey et al., 2024), Qwen2.5-7B-Instruct (Yang et al., 2024a), Mistral-7B-Instruct-v0.3 (Jiang et al., 2023), LRMs: DeepSeek-R1-Distill-Llama-8B, DeepSeek-R1-Distill-Qwen-7B (DeepSeek-AI et al., 2025), and an API: GPT-4o (OpenAI, 2024), for their diverse architectures. **Metrics.** Following previous works (Xu et al., 2024; Huang et al., 2025c; Li et al., 2025), we use accuracy (via exact match for open-ended questions) for response correctness measurement, expected calibration error (ECE) with 10 bins 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

Confidence measures the reliability of a specific response, while uncertainty reflects the inherent difficulty or ambiguity of the question itself. Despite

their distinction, previous works often treat them casually in prompt (Liu et al., 2025b). Thus, we investigate the impact of their difference by explicitly distinguishing them within the prompt itself across various scenarios, as detailed in Appendix A. We evaluate 5 models across 6 benchmarks (TriviaQA, ComparisonQA, StrategyQA, HotpotQA, MATH-500, MATH-Perturb) grouped into open-ended and multiple-choice question types.

Figure 3 presents the mean ECE and AUROC for each model across benchmark categories. The results reveal that even minor definition distinctions in prompt can elicit statistically different behaviors from models based on question types. **Open-ended Questions:** Uncertainty consistently achieves a lower ECE across both LLMs and LRMs. This suggests that for expansive prediction spaces, uncertainty, which can capture the inherent ambiguity of the question, is more effective. **Multiple-Choice Questions:** The trend reverses, with confidence significantly outperforming uncertainty in both ECE and AUROC. In this constrained setting, models can leverage elimination strategies to provide more precise confidence estimates for specific options, even when the question remains ambiguous.

Thus, uncertainty and confidence should not be used interchangeably in prompts. Their definitions should be aligned with the task format to ensure optimal model calibration.

4.3 Self-Critique Analysis

This section investigates the impact of Self-Critique on confidence calibration and average confidence

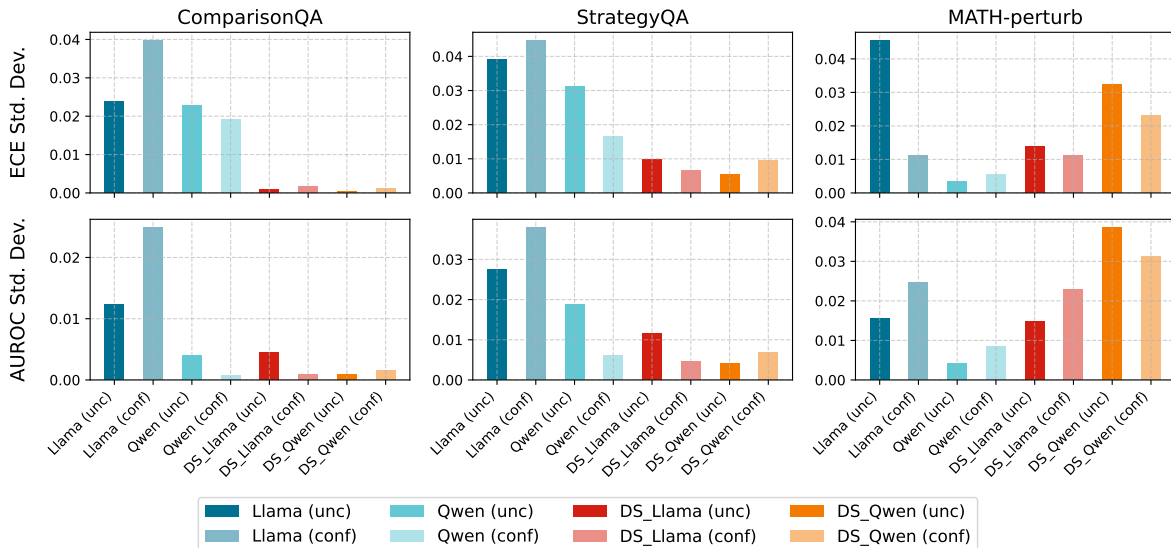


Figure 4: Standard deviation of multi-turn Self-Critique for ECE and AUROC across three benchmarks. Each bar represents the standard deviation of a model’s performance (uncertainty or confidence) across 6 iterations, where iteration 0 denotes the original response and iterations 1–5 indicate Self-Critique.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

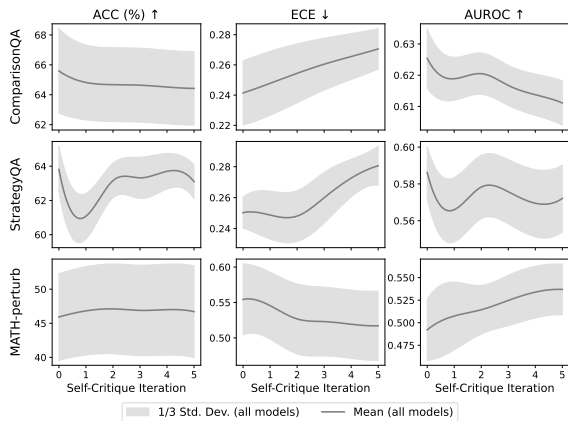


Figure 5: Multi-turn Self-Critique results on ComparisonQA, StrategyQA, and MATH-perturb benchmarks. Each plot shows the smoothed mean performance (solid line) and the corresponding 1/3 standard deviation range (shaded area) for ACC, ECE, and AUROC. Iteration 0 represents the original response without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

scores across multiple iterations.

4.3.1 Multi-Turn Self-Critique

To comprehensively evaluate Self-Critique performance, we conduct multi-turn Self-Critique experiments with 4 models, those have both reasoning and non-reasoning ones, across the same 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 per-

formance 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 categorized 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 answer correctness. 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

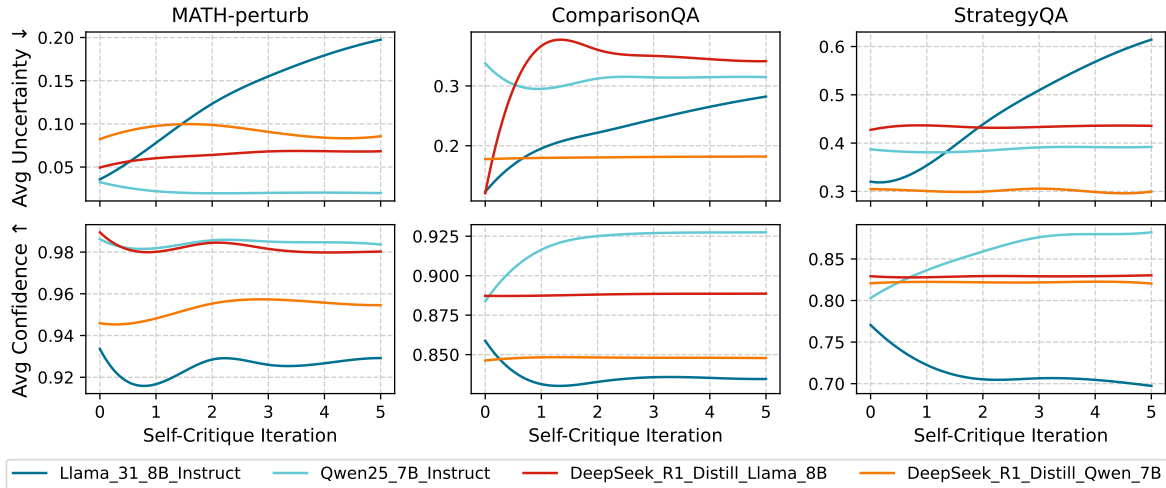


Figure 6: Curves of average uncertainty and confidence scores during multi-turn Self-Critique across 3 benchmarks.

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 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 select one representative benchmark from each of the 3 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

Model	Type	Method	Train	None Reasoning ComparisonQA			Require Reasoning					
							StrategyQA			MATH-Perturb		
				ACC	ECE	AUROC	ACC	ECE	AUROC	ACC	ECE	AUROC
(↑)	(↓)	(↑)	(↑)	(↓)	(↑)	(↑)	(↓)	(↑)				
GPT-4o	Unc	Vanilla	N	90.91	0.089	0.772	78.60	0.079	0.740	42.36	0.521	0.695
	Conf	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)	Unc	Vanilla	N	69.65	0.224	0.615	64.63	0.283	0.507	39.57	0.587	0.525
		Self_Critique	N	68.24	0.268	0.605	67.25	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	69.76	0.224	0.619	67.25	0.221	0.597	40.87	0.558	0.586
	Conf	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	0.630	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	41.30	0.617	0.558
		SFT_Soft	Y	69.90	0.193	0.630	66.38	0.193	0.629	38.70	0.611	0.562
		CritiCal	Y	69.97	0.194	0.630	69.00	0.179	0.644	40.00	0.588	0.593
DeepSeek-R1-Distill-Qwen-7B (LRM)	Unc	Vanilla	N	52.18	0.331	0.586	58.52	0.247	0.609	62.54	0.473	0.380
		Self_Critique	N	52.12	0.330	0.588	59.83	0.242	0.604	64.87	0.491	0.383
		SFT_Hard	Y	52.36	0.325	0.578	59.83	0.281	0.558	65.65	0.446	0.413
		SFT_Soft	Y	52.51	0.325	0.580	62.45	0.272	0.516	66.09	0.444	0.437
		CritiCal	Y	52.30	0.326	0.579	65.07	0.223	0.572	67.83	0.405	0.457
	Conf	Vanilla	N	52.35	0.326	0.598	58.52	0.261	0.559	65.05	0.480	0.274
		Self_Critique	N	52.31	0.327	0.602	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	52.66	0.333	0.578	61.14	0.235	0.597	65.65	0.467	0.301
		CritiCal	Y	52.55	0.333	0.580	66.81	0.176	0.630	69.13	0.432	0.328

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**. CritiCal fails for none-reasoning benchmarks (the gray columns), but shows effectiveness for those requiring reasoning.

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) **CritiCal 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 reduces the ECE of DeepSeek-Distill-Qwen, a model that already exhibits a lower ECE than GPT-4o itself. 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

We test OOD performance on StrategyQA and MATH-Perturb, with results presented in Table 2.

Compared with the baseline SFT methods, CritiCal still has the best OOD performance. Both baseline fine-tuning methods (SFT_Hard and SFT_Soft) exhibit degraded performance on OOD data, indicating limited generalization. In contrast, CritiCal even some times achieves improved calibration on OOD questions, with a lower ECE and a higher AUROC. The good generalization ability likely comes from the natural language critiques on 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.

Model	Method	Uncertainty						Confidence					
		In-distribution			Out-of-distribution			In-distribution			Out-of-distribution		
		ACC (↑)	ECE (↓)	AUROC (↑)	ACC (↑)	ECE (↓)	AUROC (↑)	ACC (↑)	ECE (↓)	AUROC (↑)	ACC (↑)	ECE (↓)	AUROC (↑)
StrategyQA to MATH-Perturb													
Qwen-2.5-7B-Instruct (LLM)	SFT_Hard	36.52	0.605	0.554	37.83	0.603	0.531	41.30	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	40.87	0.558	0.586	39.57	0.574	0.595	40.00	0.588	0.593	42.17	0.571	0.593
DS-Distill-Qwen-7B (LRM)	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	67.83	0.405	0.457	67.83	0.375	0.465	69.13	0.432	0.328	69.13	0.434	0.350
MATH-Perturb to StrategyQA													
Qwen-2.5-7B-Instruct (LLM)	SFT_Hard	65.07	0.288	0.537	63.76	0.272	0.550	66.38	0.216	0.616	65.50	0.203	0.599
	SFT_Soft	64.19	0.245	0.564	64.63	0.288	0.531	66.38	0.193	0.629	65.94	0.216	0.591
	CritiCal	67.25	0.221	0.597	66.81	0.243	0.567	69.00	0.179	0.644	67.25	0.189	0.629
DS-Distill-Qwen-7B (LRM)	SFT_Hard	59.83	0.281	0.558	59.39	0.293	0.544	61.14	0.242	0.509	59.39	0.255	0.544
	SFT_Soft	62.45	0.272	0.516	60.26	0.281	0.535	61.14	0.235	0.597	61.14	0.234	0.568
	CritiCal	65.07	0.223	0.572	62.88	0.230	0.586	66.81	0.176	0.630	64.19	0.197	0.620

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
StrategyQA (Multi-hop)				
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
ComparisonQA (One-hop)				
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.

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 the same, DPO differs in its output, consisting of a chosen response, the same as SFT’s output, and a rejected response, which should have a similar structure to the chosen one. We use the model’s Self-Critique output as the rejected response 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 natural language critiques to enhance verbalized confidence calibration in LLMs, addressing two questions: (1) What to critique. Results 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, enabling LLMs to refine their own confidence, 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 natural language critiques to advance LLM reliability.

507 Limitations

508 While CritiCal demonstrates significant improve-
509 ments in confidence calibration for LLMs, several
510 limitations still exist that cannot be covered in this
511 single work.

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

522 Additionally, computational constraints restrict
523 our ability to evaluate all benchmarks in the com-
524 parison of training methods, SFT and DPO, where
525 only ComparisonQA and StrategyQA are tested.
526 Although these benchmarks are carefully chosen to
527 represent one-hop and multi-hop factuality reason-
528 ing tasks, this limitation may obscure potential vari-
529 ations in CritiCal’s effectiveness across other task
530 types. Future work could leverage greater computa-
531 tional resources to conduct a more comprehensive
532 analysis, incorporating additional benchmarks and
533 training configurations.

534 Ethics Statement

535 This paper utilizes several publicly available
536 datasets, including ComparisonQA, TriviaQA,
537 StrategyQA, HotpotQA, MATH, MATH-Perturb,
538 and MATH-500, which are accessible to the re-
539 search community under CC, Apache 2.0, MIT,
540 Apache 2.0, MIT, MIT, and MIT licenses, respec-
541 tively. The data is anonymized, ensuring our work
542 does not raise privacy concerns regarding specific
543 entities.

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

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

References

- 554 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. 555
- 556 Mehul Damani, Isha Puri, Stewart Slocum, Idan Shenfeld, Leshem Choshen, Yoon Kim, and Jacob Andreas. 2025. [Beyond binary rewards: Training lms to reason about their uncertainty](#). *CoRR*, abs/2507.16806. 557
- 558 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. 559
- 560 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. 561
- 562 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. 562
- 563 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*. 563
- 564 Chengsong Huang, Langlin Huang, Jixuan Leng, Jiacheng Liu, and Jiaxin Huang. 2025a. [Efficient test-time scaling via self-calibration](#). *CoRR*, abs/2503.00031. 564
- 565 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. 565
- 566 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. [Math-perturb: Benchmarking llms’ math reasoning abilities against hard perturbations](#). *CoRR*, abs/2502.06453. 566

612	Liangjie Huang, Dawei Li, Huan Liu, and Lu Cheng.	Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2024b.	669
613	2025c. Beyond accuracy: The role of calibration	Generating with confidence: Uncertainty quantifi-	670
614	in self-improving large language models.	cation for black-box large language models.	671
615	<i>CoRR</i> , abs/2504.02902.	<i>Trans. Mach. Learn. Res.</i> , 2024.	672
616	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch,	Jiayu Liu, Qing Zong, Weiqi Wang, and Yangqiu Song.	673
617	Chris Bamford, Devendra Singh Chaplot, Diego	2025a. Revisiting epistemic markers in confidence	674
618	de Las Casas, Florian Bressand, Gianna Lengyel,	estimation: Can markers accurately reflect large lan-	675
619	Guillaume Lample, Lucile Saulnier, L�el�io Renard	guage models’ uncertainty?	676
620	Lavaud, Marie-Anne Lachaux, Pierre Stock,	In <i>Proceedings of the</i>	677
621	Teven Le Scao, Thibaut Lavril, Thomas Wang, Timo-	<i>63rd Annual Meeting of the Association for Computa-</i>	678
622	th�ee Lacroix, and William El Sayed. 2023. Mistral	<i>tional Linguistics (Volume 2: Short Papers), ACL</i>	679
623	7b.	2025, Vienna, Austria, July 27 - August 1, 2025, pages	680
624	<i>CoRR</i> , abs/2310.06825.	206–221. Association for Computational Linguistics.	
625	Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke	Xiaou Liu, Tiejin Chen, Longchao Da, Chacha Chen,	681
626	Zettlemoyer. 2017. Triviaqa: A large scale distantly	Zhen Lin, and Hua Wei. 2025b. Uncertainty quantifi-	682
627	supervised challenge dataset for reading comprehen-	cation and confidence calibration in large language	683
628	sion.	models: A survey.	684
629	In <i>Proceedings of the 55th Annual Meeting of the</i>	<i>CoRR</i> , abs/2503.15850.	
630	<i>Association for Computational Linguistics, ACL</i>	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	685
631	<i>2017, Vancouver, Canada, July 30 - August 4, Volume</i>	Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon,	686
632	<i>1: Long Papers</i> , pages 1601–1611. Association for	Nouha Dziri, Shrimai Prabhumoye, Yiming Yang,	687
633	Computational Linguistics.	Shashank Gupta, Bodhisattwa Prasad Majumder,	688
634	Kenneth Li, Oam Patel, Fernanda B. Vi�egas, Hanspeter	Katherine Hermann, Sean Welleck, Amir Yazdan-	689
635	Pfister, and Martin Wattenberg. 2023. Inference-time	bakhsh, and Peter Clark. 2023. Self-refine: Itera-	690
636	intervention: Eliciting truthful answers from a lan-	tive refinement with self-feedback.	691
637	guage model.	In <i>Advances in Neural Information Processing Sys-</i>	692
638	In <i>Advances in Neural Information Processing Sys-</i>	<i>tems 36: Annual Conference on Neural Information</i>	693
639	<i>tems 36: Annual Conference on Neural Information</i>	<i>Processing Systems 2023, NeurIPS 2023, New Orleans,</i>	694
640	<i>Processing Systems 2023, NeurIPS 2023, New Orleans,</i>	<i>LA, USA, December 10 - 16, 2023.</i>	695
641	<i>LA, USA, December 10 - 16, 2023.</i>	Andrey Malinin and Mark J. F. Gales. 2021. Uncertainty	696
642	Yibo Li, Miao Xiong, Jiaying Wu, and Bryan Hooi.	estimation in autoregressive structured prediction.	697
643	2025. Conftuner: Training large language mod-	In <i>9th International Conference on Learning Representa-</i>	698
644	els to express their confidence verbally.	<i>tions, ICLR 2021, Virtual Event, Austria, May 3-7,</i>	699
645	<i>CoRR</i> , abs/2508.18847.	2021. OpenReview.net.	700
646	Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harri-	Arjun S. Nair and Kristina P. Sinaga. 2025. Disprov-	701
647	son Edwards, Bowen Baker, Teddy Lee, Jan Leike,	ing the feasibility of learned confidence calibration	702
648	John Schulman, Ilya Sutskever, and Karl Cobbe.	under binary supervision: An information-theoretic	703
649	2024a. Let’s verify step by step.	impossibility.	704
650	In <i>The Twelfth International Conference on Learning Representa-</i>	<i>CoRR</i> , abs/2509.14386.	
651	<i>tions, ICLR 2024, Vienna, Austria, May 7-11, 2024.</i>	OpenAI. 2024. Hello gpt-4o.	705
652	OpenReview.net.	<i>OpenAI.</i>	
653	Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harri-	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christo-	706
654	son Edwards, Bowen Baker, Teddy Lee, Jan Leike,	pher D. Manning, Stefano Ermon, and Chelsea Finn.	707
655	John Schulman, Ilya Sutskever, and Karl Cobbe.	2023. Direct preference optimization: Your language	708
656	2024b. Let’s verify step by step.	model is secretly a reward model.	709
657	In <i>The Twelfth International Conference on Learning Representa-</i>	In <i>Advances in Neural Information Processing Sys-</i>	710
658	<i>tions, ICLR 2024, Vienna, Austria, May 7-11, 2024.</i>	<i>tems 36: Annual Conference on Neural Information</i>	711
659	OpenReview.net.	<i>Processing Systems 2023, NeurIPS 2023, New Orleans,</i>	712
660	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022.	<i>LA, USA, December 10 - 16, 2023.</i>	713
661	Teaching models to express their uncertainty in	Jiayuan Su, Jing Luo, Hongwei Wang, and Lu Cheng.	714
662	words.	2024. API is enough: Conformal prediction for large	715
663	<i>Trans. Mach. Learn. Res.</i> , 2022.	language models without logit-access.	716
664	Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2024a.	In <i>Findings of the Association for Computational Linguistics:</i>	717
665	Contextualized sequence likelihood: Enhanced con-	<i>EMNLP 2024, Miami, Florida, USA, November 12-</i>	718
666	fidence scores for natural language generation.	<i>16, 2024, pages 979–995.</i>	719
667	In <i>Proceedings of the 2024 Conference on Empirical</i>	Association for Computational Linguistics.	720
668	<i>Methods in Natural Language Processing, EMNLP</i>	Karthik Valmeekam, Matthew Marquez, and Subbarao	721
	<i>2024, Miami, FL, USA, November 12-16, 2024,</i>	Kambhampati. 2023. Can large language models	722
	pages 10351–10368. Association for Computational Lin-	really improve by self-critiquing their own plans?	723
	guistics.	<i>CoRR</i> , abs/2310.08118.	724

725	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 . <i>Trans. Assoc. Comput. Linguistics</i> , 13:220–248.	
726		
727		
728		
729		
730		
731		
732		
733		
734	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 <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 9426–9439. Association for Computational Linguistics.	
735		
736		
737		
738		
739		
740		
741		
742		
743	Yubo Wang, Xiang Yue, and Wenhui Chen. 2025. Critique fine-tuning: Learning to critique is more effective than learning to imitate . <i>CoRR</i> , abs/2501.17703.	
744		
745		
746	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 <i>Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024</i> , pages 6886–6898. Association for Computational Linguistics.	
747		
748		
749		
750		
751		
752		
753		
754	Sean Welleck, Ximing Lu, Peter West, Faeze Brahman, Tianxiao Shen, Daniel Khachabi, and Yejin Choi. 2023. Generating sequences by learning to self-correct . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	
755		
756		
757		
758		
759		
760	Zhiqiu Xia, Jinxuan Xu, Yuqian Zhang, and Hang Liu. 2025. A survey of uncertainty estimation methods on large language models . In <i>Findings of the Association for Computational Linguistics, ACL 2025, Vienna, Austria, July 27 - August 1, 2025</i> , pages 21381–21396. Association for Computational Linguistics.	
761		
762		
763		
764		
765		
766		
767	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 <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	
768		
769		
770		
771		
772		
773	Tianyang Xu, Shujin Wu, Shizhe Diao, Xiaoze Liu, Xingyao Wang, Yangyi Chen, and Jing Gao. 2024. Sayself: Teaching llms to express confidence with self-reflective rationales . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024</i> , pages 5985–5998. Association for Computational Linguistics.	
774		
775		
776		
777		
778		
779		
780		
	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, Jixi Yang, Jingren Zhou, Junyang Lin, Kai Dang, and 22 others. 2024a. Qwen2.5 technical report . <i>CoRR</i> , abs/2412.15115.	781
		782
		783
		784
		785
		786
		787
	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 . <i>CoRR</i> , abs/2409.12122.	788
		789
		790
		791
		792
		793
		794
	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 <i>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</i> , pages 3998–4014. Association for Computational Linguistics.	795
		796
		797
		798
		799
		800
		801
		802
	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 <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018</i> , pages 2369–2380. Association for Computational Linguistics.	803
		804
		805
		806
		807
		808
		809
		810
		811
	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 <i>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, June 16-21, 2024</i> , pages 7113–7139. Association for Computational Linguistics.	812
		813
		814
		815
		816
		817
		818
		819
		820
		821
	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 . <i>CoRR</i> , abs/2506.03106.	822
		823
		824
		825
		826
	Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, and Yongqiang Ma. 2024. Llamafactory: Unified efficient fine-tuning of 100+ language models . <i>CoRR</i> , abs/2403.13372.	827
		828
		829
		830
	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 <i>Findings of the Association for Computational Linguistics, ACL 2025, Vienna, Austria, July 27 - August 1, 2025</i> , pages 4101–4117. Association for Computational Linguistics.	831
		832
		833
		834
		835
		836
		837
		838

Appendices

A Prompt

We design our prompt according to previous works (Xiong et al., 2024). 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.

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.

C Benchmark Statistics

For training, we randomly sampled 2K questions for each benchmark. For testing, we used the full sets of most benchmarks, except for MATH-Perturb. Since MATH-Perturb is constructed from the original MATH dataset, we excluded those derived from the MATH training set to avoid data leakage. Instead, we only evaluated on perturbations from the MATH test set. Note that the training regarding math reasoning relied on the original MATH training set. Table 4 details the question counts for each benchmark.

Benchmark	# Question
TriviaQA	11,313
ComparisonQA	56,696
StrategyQA	229
HotpotQA	7,405
MATH-500	500
MATH-Perturb	230

Table 4: Benchmark statistics: the number of questions we tested in each benchmark.

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]*

All Previously Refined responses: *[All_Previously_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.

Critique generation prompt for LLMs 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 is well-calibrated, considering the clarity and strength of the reasoning provided in the student's response and also 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 for LLMs on StrategyQA. Placeholders *[Question]*, *[Correct_Answer]*, *[Facts]*, and *[Student's_Response]* will be replaced with the real ones. For LRMs, we will add the sentence "Use "</think>" to separate your thinking process and your final response" at the end of above prompt.

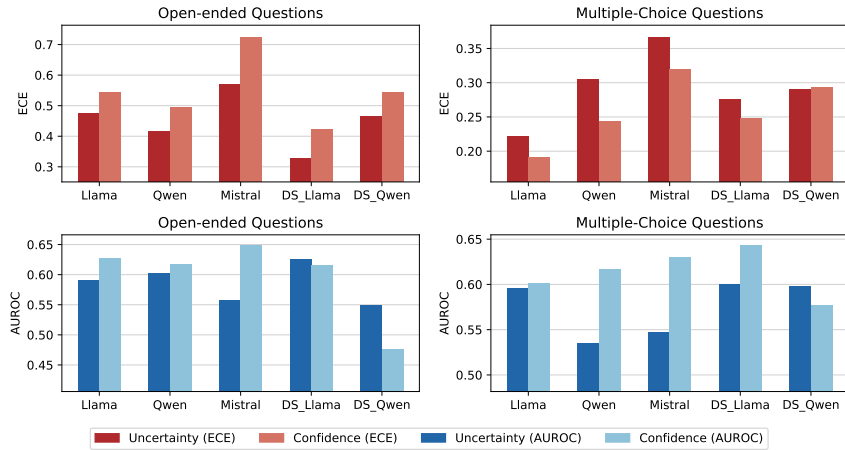


Figure 11: Mean ECE and AUROC values for each model across the same category of benchmarks, which are taken the average of across the 5 turns of Self-Critique. The dark bars are the result under uncertainty prompt, and the light ones are of confidence.

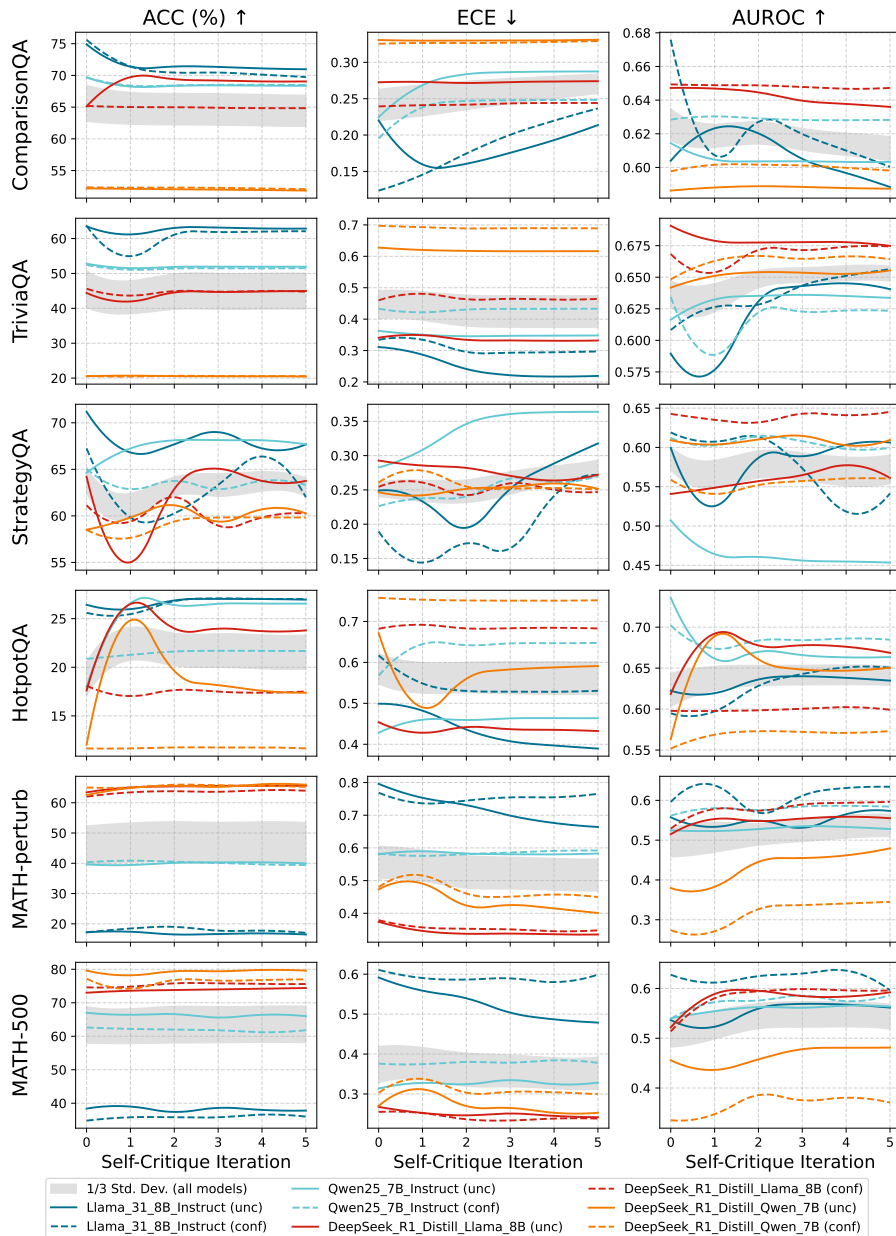


Figure 12: Multi-turn Self-Critique results on all the six benchmarks. The 0 iteration means the original response without Self-Critique. The semi-transparent light gray area represents the average performance of all models with a one-third standard deviation.