CoT-UQ: Improving Response-wise Uncertainty Quantification in LLMs with Chain-of-Thought

Anonymous ACL submission

Abstract

Large language models (LLMs) excel in many tasks but struggle to accurately quantify uncertainty in their generated responses. This 004 limitation makes it challenging to detect misinformation and ensure reliable decision-making. Existing uncertainty quantification (UQ) methods for LLMs are primarily prompt-wise rather than response-wise, often requiring multiple response samples, which leads to inefficiency. Moreover, LLMs have been shown to be overconfident, particularly when using reasoning steps to derive their answers. In this work, we introduce a novel approach to quantify responsewise uncertainty by integrating LLMs' inherent reasoning capabilities through Chain-of-Thought (CoT) into the UQ process. Our CoT-UQ framework captures critical information 017 during inference by extracting keywords from each reasoning step and assessing their importance to the final answer. The uncertainty scores of keywords are then aggregated based on their significance to produce a final uncertainty estimate. We conduct extensive experiments based on LLaMA Family with model sizes varying from 8B to 13B across logical and mathematical reasoning tasks. Experimental results demonstrate that CoT-UQ significantly outperforms existing UQ methods, achieving an average improvement of 5.9% AUROC compared to current UQ methods.

1 Introduction

Large language models (LLMs) have demonstrated groundbreaking capabilities across a variety of applications (Ouyang et al., 2022; Chowdhery et al., 2023; OpenAI, 2024b). Particularly, prompting techniques like Chain-of-Thought (CoT) (Wei et al., 2022) have significantly enhanced LLMs reasoning capabilities, ranging from multi-round conversation (Long, 2023; Chen et al., 2023), logical reasoning (Creswell et al., 2022; Duan et al., 2024b) and mathematical reasoning (Yao et al., 2024; Shao



Figure 1: Comparison of existing UQ strategies with ours. Directly estimating the uncertainty of a generated incorrect answer leads to overconfidence, which is exacerbated by using CoT to derive the answer. We tackle this challenge by integrating CoT into the UQ process with keywords extraction and importance scores.

et al., 2024). However, LLMs often unpredictably hallucinate (Manakul et al., 2023), i.e., making plausible but incorrect statements (Ji et al., 2023), limiting their deployment in safety-critical applications (Clusmann et al., 2023).

To improve the reliability of LLMs, uncertainty quantification (UQ) has emerged as a key strategy for determining when humans can trust LLMgenerated outputs. However, existing UQ methods for LLMs are primarily prompt-wise (Malinin and Gales, 2021; Kuhn et al., 2023; Ling et al., 2024). That is, uncertainty is calculated at the prompt level rather than for each individual response. These methods require multiple response samples per prompt, leading to additional computational costs and inefficiency. Besides, some studies (Kadavath et al., 2022; Miao et al., 2023) propose leveraging an LLM's own ability to evaluate the uncertainty of its responses without relying on external knowledge. However, these approaches suffer from overconfidence issues, particularly when reasoning steps,

such as Chain-of-Thought (CoT), are used before deriving the final answer (Fu et al., 2025). Overconfidence has been attributed to the model's inherent bias toward trusting its own outputs (Mielke et al., 2022; Lin et al., 2022).

063

064

065

073

077

078

084

091

097

100

101

102

103

105

107

108

109

110

111

112

113

To enable response-wise reliable UQ, we propose leveraging reasoning steps not only for deriving final answers but also for UQ itself. Our motivation comes from an intuitive insight: providing access to the reasoning path allows the model to incorporate additional context for confidence calibration, leading to a more informed assessment of the final answer. As illustrated in Figure 1, while LLMs tend to be overconfident in their generated answers, and CoT can further amplify this issue, we believe incorporating key reasoning information into the UQ process can effectively mitigate inflated confidence scores, resulting in better-calibrated uncertainty estimates. This naturally leads to a critical research question: How to utilize the LLM reasoning path to estimate the uncertainty of its generations?

To answer this question, we propose a new framework, namely, Chain-of-Thought enhanced Uncertainty Quantification (CoT-UQ). At a high level, CoT-UQ follows the principle of *one response* \rightarrow *one uncertainty score*, integrating inference steps from CoT into the UQ process to mitigate overconfidence. In detail, as illustrated in Figure 2, CoT-UQ leverages the LLM's own reasoning process to extract keywords from each inference step and assess their importance in determining the final answer. By incorporating this critical information, CoT-UQ achieves better-calibrated uncertainty estimation, either by aggregating token probabilities of extracted keywords at intermediate steps or by integrating the reasoning path into the self-evaluation process.

We conducted extensive experiments to verify the effectiveness of our proposed framework. Under extensive evaluations, our CoT-UQ achieves superior performance compared with different baselines, which reveals that LLMs have the potential to use their own reasoning to better express the trustworthiness of their generations. We also conduct a range of ablation studies of the proposed framework and provide detailed further discussions from different perspectives. Our contributions can be summarized as the following:

 Conceptually, we introduce a novel perspective to quantify the response-wise uncertainty for LLMs by considering LLM's internal knowledge from the reasoning path. Technically, we propose a new UQ framework, namely, Chain-of-Thought enhanced
Uncertainty Quantification (CoT-UQ), which integrates the inference knowledge into the
UQ process through extracting keywords and evaluating corresponding importance scores.

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

• Empirically, we conduct extensive experiments on the LLaMA family across five datasets in two tasks and show that CoT-UQ achieves an average AUROC improvement of approximately 5.9% compared to baselines, which verifies the effectiveness of our method.

2 Related Works

2.1 Uncertainty Quantification in LLMs

Prior efforts to quantify uncertainty and confidence in LLMs can be categorized into four main approaches. The first approach is to derive calibrated confidence by examining agreement across multiple sampled responses (Malinin and Gales, 2021; Kuhn et al., 2023; Manakul et al., 2023; Tian et al., 2023a). However, as Qiu and Miikkulainen (2024) recently pointed out, these methods primarily quantify prompt-wise rather than response-wise uncertainty. While Qiu and Miikkulainen (2024) provides a method for response-wise uncertainty, it still relies on generating multiple response samples, making it computationally inefficient. The second approach is to leverage LLM's own ability to evaluate the confidence of its responses, often through self-probing techniques (Kadavath et al., 2022; Tian et al., 2023b; Xiong et al., 2023). The third approach is to aggregate token probabilities of its generated response, which includes adopting traditional UQ methods (Xiao et al., 2022; Ye et al., 2024) and assigning importance weights to tokens (Duan et al., 2024a; Bakman et al., 2024). However, the above two approaches still suffer from overconfidence due to the model's inherent bias to trust its own outputs (Mielke et al., 2022; Lin et al., 2022). The fourth approach is to fine-tune the original LLM to calibrate its confidence (Lin et al., 2022; Kapoor et al., 2024). However, the modelspecific tuning has limited their applications to new scenarios. In contrast to these four approaches, the proposed CoT-UQ is a response-wise uncertainty quantification method that does not require additional response sampling or model-specific tuning. Instead, it leverages the LLM's inherent knowledge



Figure 2: Illustration of our framework, CoT-UQ. Given a question and an incorrect response generated by an LLM, the *top* of the figure shows two common UQ strategies, which suffer from overconfidence issues. The *bottom* shows the four-step process of CoT-UQ: performing the reasoning process, extracting step-wise keywords, scoring the importance of keywords relative to the final answer, and leveraging reasoning information to enhance common UQ strategies. CoT-UQ leads to a better-calibrated response-wise uncertainty estimate.

and reasoning process to calibrate uncertainty/confidence scores, making it readily generalizable to new tasks and models.

162

163

164

2.2 Chain of Thought Reasoning in LLMs

To equip LLMs with capabilities to solve more com-166 plex and reasoning tasks, Wei et al. (2022) extended in-context learning by introducing the concept of 168 Chain of Thought (CoT) through a step-by-step reasoning process. Kojima et al. (2022) found that 170 simply adding a leading sentence "Let's think step 171 by step" to a cue allowed LLMs to perform zero-172 shot logical reasoning without any additional human prompts (Chu et al., 2023). Subsequently, CoT-SC 174 (Wang et al., 2022) introduces a self-consistency 175 strategy to replace the greedy decoding strategy. Feng et al. (2024) further reveals the underlying 177 mechanisms behind CoT through a theoretical perspective. Liu et al. (2024) refines CoT by capturing 179 relationships between entities to aid LLMs in understanding context. Although these studies highlight 181

the importance of CoT in enhancing LLMs' reasoning abilities in various situations, a recent study (Fu et al., 2025) observes that CoT can exacerbate the overconfident issues in LLMs when only measuring the final answer. To the best of our knowledge, CoT-UQ is the first approach to integrate reasoning knowledge into the UQ process for LLMs. 182

183

184

185

186

188

189

191

192

193

194

195

196

197

198

199

201

3 Preliminaries

In this section, we briefly introduce the preliminaries of response-wise uncertainty quantification (UQ) in LLMs, including problem settings and two popular UQ strategies, namely aggregated probabilities and self-evaluation. Further details on these two strategies can be found in Appendix A.2.

Problem Setups. Given an LLM M, an input prompt p, and the output sequence $\hat{y} = [y_1, y_2, ..., y_L]$, where L is the number of tokens generated by LLM, the task is to obtain a confidence score for users representing the probability that \hat{y} is correct. In the following, we illustrate the exist-

284

290

291

292

293

294

295

296

250

251

252

ing two paradigms for response-wise uncertaintyquantification and analyze their limitations.

Aggregated Probabilities (*AP*). Previous works 204 (Kadavath et al., 2022; Huang et al., 2023; Varsh-205 ney et al., 2023) based on aggregated probabilities 206 generally aggregate output token probabilities of the generated text tokens $\hat{y} = [y_1, y_2, ..., y_L]$ to measure the LLM's confidence for each response. For the type of aggregation techniques, we consider 210 several methods following Orgad et al. (2024), in-211 cluding the mean and the minimum of these values. 212 Formally, given an aggregation function $Aggr(\cdot)$, 213 the confidence score c can be abstracted as, 214

215

217

218

219

221

227

231

235

237

238

240

241

242

$$c = \underset{i=1}{\overset{N}{\text{Aggr}}} (\mathbb{P}(y_i | p, y_1, ..., y_{i-1})).$$
(1)

Self-Evaluation (*SE*). Self-evaluation strategies (Kadavath et al., 2022; Xiong et al., 2023) usually contain a two-stage process to elicit the confidence score from LLMs: 1) Using an input comprising of the question q combined with the prompt p to generate the text response \hat{y} . 2) Combining q and \hat{y} through a well-designed prompt p^t to instruct LLM to self-evaluate the correctness of \hat{y} . Among them, one representative baseline is P(True) (Kadavath et al., 2022). P(True) is straightforward yet effective by directly asking LLM whether the predicted \hat{y} is true or false to q via p^t and using the probability of "True" as confidence c, which is defined as follows,

 $c = \mathbb{P}(o = True), \text{ where } o = M(p^t(q, \hat{y})).$ (2)

Although previous UQ methods using AP and SE have shown promising results, they often suffer from overconfidence, particularly when using Chain-of-Thought reasoning for complex tasks. In this work, we explore how to leverage the intrinsic reasoning capabilities of LLMs to mitigate overconfidence. Specifically, we propose extracting step-wise keywords from the model's inference process and integrating this knowledge into AP and SE strategies to better assess the trustworthiness of LLM-generated outputs.

4 Methodology

243In this section, we introduce our new framework,244i.e., Chain-of-Thought enhanced Uncertainty Quan-245tification (CoT-UQ), as illustrated in Figure 2. Com-246pared to common UQ strategies implemented di-247rectly on the generated answer, CoT-UQ is a two-248stage paradigm containing four specific steps dur-249inference time. The first three steps focus on

refining the multi-step inference by extracting keywords and their corresponding importance scores to the final answer (Section 4.1). The fourth step illustrates how to integrate this crucial reasoning information into the two common UQ strategies, respectively (Section 4.2).

4.1 Stage 1: LLM Inference Refining

Step 1: Reasoning Extraction. We first instruct LLM to derive the reasoning for each response. Before inference, we add the step-wise Chain-of-Thought (CoT) prefix for prompting, i.e., "*Let's think step by step. Step 1:*". This ensures the model's inference results are structured into multiple reasoning steps, with each step explicitly starting with "*Step i:*". Upon completion of inference, we obtain a response \hat{y} for the question q, which contains a step-by-step reasoning $s_{1\sim k} = s_1, ..., s_k$ and a final answer a labeled with "*Final Answer:*".

Step 2: Keywords Extraction. After obtaining the step-by-step inference $s_{1\sim k}$ for each questionanswer pair (q, a), we choose to extract keywords from each step. Prior works generally consider the sum or average of all generated tokens (Slobodkin et al., 2023) to aggregate token-level uncertainty. However, these strategies potentially introduce redundant tokens, which can significantly compromise the accuracy of uncertainty scores.(Gupta et al., 2024). This motivates us to consider tokens from keywords, which better represent the most meaningful part of an inference step. Specifically, we request the LLM itself to complete the extraction, as (Ashok and Lipton, 2023; Orgad et al., 2024) have demonstrated LLMs' information extraction capability. Formally, we extract $n_i \ge 0$ keywords from each reasoning step $s_i \in s_{1 \sim k}$ (noted $n_i = 0$ means no effective keywords in a specific step, we explain this situation in Appendix B). The keywords set \mathcal{K} extracted from all steps can be formulated as,

$$\mathcal{K} = \bigcup_{i=1}^{k} \{w_j^i\}_{j=1}^{n_i}.$$
 (3)

Step 3: Importance Scoring. Relying on the selfevaluation (Ren et al., 2023) capability, we instruct the LLM to rate the importance of the keywords in deriving the final answer in a few-shot learning setup. We provide the context (question q, multistep reasoning $s_{1\sim k}$, and final answer a) combined with the keywords \mathcal{K} extracted in Step 2 to the LLM. Each keyword will be scored by the LLM, ranging

from 1 to 10, where 1 denotes the least critical and 297 10 is the most. For keywords that are more critical 298 in the inference time, i.e., require exact reasoning or imply vital elements to the final answer, we assign a higher score to this keyword. For instance, as shown in Figure 2, keywords that reveal the specific number of members in each band will get a higher importance score, even to 10, as they explicitly require reasoning and are crucially contributing to the final answer. After integrating with the corresponding importance indicator t, keywords set 307 308 \mathcal{K} can be updated as,

$$\mathcal{K} = \bigcup_{i=1}^{k} \{ (w_j^i, t_j^i) \}_{j=1}^{n_i}.$$
 (4)

In the first three steps, we deconstructed the redundant reasoning steps into keywords containing the most meaningful and critical information and evaluated their respective importance towards reaching the final answer, formalizing them as a keywords set containing dualist formulation. In the next step, we will use these keywords to help enhance existing uncertainty quantification strategies.

310

311

312

313

314

317

318

319

320

321

323

337

341

4.2 Stage 2 / Step 4: Reasoning Enhanced **Uncertainty Quantification Strategy**

Given the extracted reasoning path $s_{1\sim k}$ and keywords set \mathcal{K} defined in Section 4.1, we aim to elicit the confidence c and mitigate overconfidence issues generally exist in common UQ strategies, aggregated probabilities (AP) and self-evaluation (SE). We propose an integration method for each strategy to utilize the information provided by $s_{1 \sim k}$ and \mathcal{K} .

Reasoning Enhanced Aggregated Probabilities (AP) Strategy. Compared to directly combin-329 ing the token probabilities from output tokens, we choose to aggregate those from the extracted key-330 words to integrate the inference knowledge into the uncertainty quantification process. Since the keywords are generally short in token length, and for the sake of comparison, we use the same aggregation techniques $Aggr(\cdot)$ introduced in Section 335 3 to aggregate the token probabilities of a single 336 keyword into its prediction probability. Formally, given a single keyword w (text) with its corresponding token sequence $\hat{w} = [w_1, w_2, ..., w_l]$ of length l, the probability of the keyword w can be formulated as,

$$p(\hat{w}) = \underset{m=1}{\overset{l}{\operatorname{Aggr}}} (\mathbb{P}(w_m \mid p, w_1, \dots, w_{m-1})).$$

To consider their contributions to the final confidence, we propose to average the probabilities of keywords weighted by their importance scores. This ensures that more significant keywords have a greater influence on the final confidence estimation. The procedure can be formalized as follows:

$$c = \frac{\sum_{i=1}^{k} \sum_{j=1}^{n_i} t_j^i \cdot p(\hat{w_j^i})}{\sum_{i=1}^{k} \sum_{j=1}^{n_i} t_j^i}.$$
 (6)

Reasoning Enhanced Self-Evaluation (SE) Strat-We provide *four* approaches to instruct the egy. LLM to consider the reasoning information during the self-evaluation of uncertainty. The first two methods, namely, ALLSteps and ALLKeywords, directly add the extracted reasoning steps $s_{1\sim k}$ or keywords set \mathcal{K} to the self-evaluation process. To highlight the role of relevant importance of extracted keywords, we further introduce KEYStep and KEYKeywords strategies. KEYStep proposes to consider the most important step, where the importance is calculated from the average of the importance scores from each reasoning step. KEYStep can be abstracted as follows,

$$s^* = \arg\max_{1 \le i \le k} \left(\frac{1}{n_i} \sum_{j=1}^{n_i} t_j^i\right). \tag{7}$$

Meanwhile, the goal of KEYKeywords is to exclude redundant keywords and shortlist the most critical ones based on their importance. We formulate it as,

$$\mathcal{K}^* = \bigcup_{i=1}^k \left\{ (w_j^i, t_j^i) \mid t_j^i \ge \tau \right\}_{j=1}^{n_i}, \qquad (8)$$

where τ is a threshold to filter the sub-critical keywords during self-evaluation. We discuss the sensitivity to this hyper-parameter in Section 5.3.

We include the above information in the selfevaluation prompt p^t through an additional instruction, termed Considering <reasoning type> as additional information (See Appendix B.2 for the concrete realization). Formally, given a type of reasoning knowledge $z \in \{s_{1 \sim k}, s^*, \mathcal{K}, \mathcal{K}^*\}$, the self-evaluation process can be updated by the refined prompt p^r as,

$$c = \mathbb{P}(o = True), \text{ where } o = M(p^r(q, \hat{y}, z)).$$
(9)
381

(5)

364

343

344

345

347

348

349

350

351

352

353

354

355

356

357

358

359

361

362

363

365

366 367

368

369

371

372

373

374

375

376

377

378

379

Model	Strategy	Method	Logical	Reasoning	Mathematical Reasoning		
			HotpotQA	2WikiMHQA	GSM8K	SVAMP	ASDiv
	AP	Probas-mean	53.73	56.80	53.17	53.94	58.34
Llama 3.1-8B		w/ CoT-UQ	62.01	65.22	63.64	59.83	64.52
		Probas-min	58.34	56.81	54.95	54.79	58.69
		w/ CoT-UQ	64.37	70.02	63.09	60.49	64.84
		TOKEN <i>SAR</i>	53.57	56.92	54.46	55.01	58.71
		w/ CoT-UQ	61.07	65.38	65.10	62.11	66.91
	SE	P(True)	62.39	53.56	48.15	51.58	47.23
		w/ CoT-UQ	63.10	57.77	52.60	60.00	53.20
		Self-Probing	54.33	56.39	49.24	51.63	50.86
		w/ CoT-UQ	57.20	58.38	51.89	54.26	53.79
Llama 2-13B	AP	Probas-mean	56.27	51.54	53.96	54.48	57.73
		w/ CoT-UQ	66.56	63.29	58.54	57.37	59.44
		Probas-min	56.51	51.28	53.84	55.09	57.70
		w/ CoT-UQ	67.19	68.10	58.63	58.51	60.74
		TOKEN <i>SAR</i>	57.33	51.08	54.82	55.06	58.37
		w/ CoT-UQ	66.29	64.03	59.61	58.41	61.23
	SE	P(True)	51.13	47.52	46.06	46.36	48.02
		w/ CoT-UQ	57.10	53.52	52.59	56.87	56.10
		Self-Probing	60.63	59.81	52.72	47.27	52.35
		w/ CoT-UQ	64.03	62.63	55.14	50.53	57.58

Table 1: AUROC (\uparrow) comparison of Llama 3.1-8B and Llama 2-13B on various benchmarks for logical and mathematical reasoning tasks, where *AP* indicates *aggregated probabilities* and *SE* denotes *self-evaluation*.

5 Experiments

5.1 Experimental Setups

Datasets and Models. We consider the following reasoning scenarios: logical reasoning and mathematical reasoning, where existing UQ strategies suffer from overconfidence issues. Specifically, for logical reasoning, we use the HotpotQA dataset (Yang et al., 2018) and the 2WikiMultiHopQA dataset (Ho et al., 2020); for mathematical reasoning, we use the GSM8K dataset (Cobbe et al., 2021), the SVAMP dataset (Patel et al., 2021), and the ASDiv dataset (Miao et al., 2021). For models, we use Llama2-13B and Llama3.1-8B (Touvron et al., 2023). Details of dataset statistics and LLMs' hyper-parameters are provided in Appendix A.1.

Evaluation Metric. Following the common evaluation approach in Kuhn et al. (2023), we use the
Area Under the Receiver Operating Characteristic
curve (AUROC) (Davis and Goadrich, 2006) to
evaluate the performance of UQ methods, which
measures the likelihood that a positive sample will

receive a higher discriminating score than a negative sample (Fawcett, 2006). A higher AUROC score indicates better performance, while a score of 0.5 implies random guessing.

403

404

405

406

Baseline Methods. We compare the proposed 407 framework with various competitive uncertainty 408 quantification baseline methods. As mentioned 409 in Section 3, common response-wise UQ strate-410 gies include aggregated probabilities (AP) and self-411 evaluation (SE). For AP, inspired by (Kadavath 412 et al., 2022; Guerreiro et al., 2022), we first inves-413 tigate the most common aggregation techniques 414 Probas-mean and Probas-min. We also consider 415 the Toekn-level Shifting Attention to Relevance 416 (TOKENSAR) (Duan et al., 2024a) that evaluates 417 the relevance of each token in the final answer and 418 assigns higher weights to more relevant tokens. For 419 SE, we consider P(True) (Kadavath et al., 2022) and 420 Self-Probing (Xiong et al., 2023) that directly asks 421 LLM to evaluate the correctness of their generation 422 via prompting. Details on the implementations of 423 baseline methods are listed in Appendix A.2. 424



Figure 3: Comparison of different implementations of CoT-UQ for P(True). KEY*Keywords* is more effective for logical reasoning tasks (HotpotQA and 2WikiMultiHopQA), whereas ALL*Steps* works better for mathematical reasoning tasks (GSM8K and SVAMP).



Figure 4: Comparison of different implementations of CoT-UQ based on Self-Probing. KEY*Keywords* consistently achieves better performance on logical reasoning tasks HotpotQA and 2WikiMultiHopQA), while KEY*Step* performs better in mathematical reasoning (GSM8K and SVAMP).

5.2 Main Results

425

426

427

428

429

430

431

432

433

434

435

436

437

438

The overall comparison results are presented in Table 1. We evaluate various baseline methods by comparing their performance with and without CoT-UQ. As shown in Table 1, CoT-UQ consistently improves UQ performance across all tasks and datasets. This demonstrates that incorporating reasoning into uncertainty quantification enables LLMs to provide more calibrated assessments of the trustworthiness of their generated outputs. In general, CoT-UQ achieves greater improvements when applied to *AP* strategies compared to *SE* strategies, particularly for **Probas-min**, where it increases AUROC by up to **16.8%**.

Aggregated Probabilities (AP). When compar-439 ing CoT-UQ with three AP strategies, our framework 440 significantly outperforms them on the two logical 441 reasoning datasets, HotpotQA and 2WikiMulti-442 HopQA, achieving an average AUROC improve-443 444 ment of +10.3% across both models. Similarly, CoT-UQ improves performance on all three mathe-445 matical reasoning benchmarks. It is worth noting 446 that the recent work TOKENSAR gets limited im-447 provements or even worse performance compared 448

to Probas-mean on reasoning tasks, which may be caused by the limited length of the generated final answer in these tasks, especially in mathematical reasoning. CoT-UQ addresses this limitation by aggregating keyword token probabilities, mitigating the impact of response length on uncertainty estimation. 449

450

451

452

453

454

455

Self-Evaluation (SE). For the comparison with 456 SE strategies, we suggest different implementations 457 of CoT-UQ across different reasoning tasks and SE 458 strategies. Specifically, we propose KEYKeywords 459 for both two SE strategies in logical reasoning tasks, 460 ALLSteps for P(True) strategy and KEYStep for Self-461 Probing strategy in mathematical reasoning tasks. 462 These suggestions are based on the observation that 463 keywords extracted from mathematical reasoning 464 tend to be overly simplistic (e.g., single digits) and 465 lack informative content. As a result, step-level 466 strategies, which retain rich contextual information, 467 are more suitable. In contrast, keywords extracted 468 from logical reasoning generally retain meaningful 469 and logical information, while complete reasoning 470 steps may introduce redundant content that harms 471 the model's judgment. 472



Figure 5: Effect of Importance scoring in CoT-UQ. (a) Implement the aggregation strategy of Probas-min for comparison; (b) Implement Self-Probing with suggested realizations for comparison.

We report the results of recommended realizations in Table 1 and a detailed explanation and analysis for the above suggestions in Section 5.3. The results of these five datasets demonstrate an average of +4.4% improvement, highlighting the effectiveness of CoT-UQ compared to standard SE baselines. Notably, CoT-UQ applied to SE strategies shows a greater performance improvement in mathematical reasoning tasks (+5.3%) compared to logical reasoning tasks (+3.5%). This suggests that incorporating reasoning in UQ allows the model to identify critical errors in the thought process during self-evaluation, especially in mathematical problems where a single misstep can lead to incorrect conclusions. We provide detailed evidence and analysis for this observation in Appendix C.2.

473 474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

497

504

5.3 **Ablation and Future Discussions**

In this section, we provide a thorough understanding of our CoT-UQ. Additional results and discussions (e.g., sensitivity to hyper-parameters and case studies) can be found in Appendix C.

Effect of Different Implementation of CoT-UQ in SE Strategies. We observe an interesting 495 trend: the transition from step-level implemen-496 tations (ALLSteps and KEYStep) to keywords-level implementations (ALLKeywords and KEYKeywords) 498 shows opposite trends in logical and mathemat-499 ical reasoning tasks, as demonstrated in Figure 3 and Figure 4. Specifically, ALLKeywords and KEYKeywords perform significantly better than steplevel techniques on HotpotQA and 2WikiMHQA datasets. This suggests that step-level information usually contains redundant words in logical reason-505 ing tasks, and the keywords effectively filter out irrelevant information for uncertainty estimation. Conversely, ALLSteps consistently performs well when 508

applied to P(True) on GSM8K and SVAMP datasets, and KEYStep performs similarly when applied to Self-Probing. However, keywords-level methods show suboptimal performance, or even worse than the standard SE on mathematical datasets. This may be attributed to the necessity of including sufficient context when assessing mathematical answers, rather than relying on a few scattered keywords.

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

Effect of Importance Scoring. The effects of the importance scoring step are summarized in Figure 5, where we employ Probas-min and Self-Probing to represent the AP and SE strategies, respectively. For AP, w/o importance indicates directly calculating the mean of probabilities from keywords. For SE, we use the specific implementations as above suggested for different reasoning tasks. For instance, we adopt the KEYKeywords for HotpotQA, and the KEY*Step* for mathematical reasoning benchmarks. Here, w/o importance in SE demotes the corresponding realizations start with ALL, i.e, ALLKeywords and ALLSteps. The results highlight the necessity of evaluating the respective importance metric for each keyword in our method.

6 Conclusion

In this paper, we introduce Chain-of-Thought enhanced Uncertainty Quantification (CoT-UQ), a novel perspective for uncertainty quantification that leveraging LLM's internal knowledge through CoT to calibrate its confidence on each response. CoT-UQ consistently and significantly boosts the performance of current aggregated probabilities (AP) and self-evaluation (SE) strategies by using crucial information from the reasoning path. We have conducted extensive experiments to demonstrate the effectiveness of our framework and provided detailed discussions on it from various perspectives.

545

561

571

573

575

576

577

578

579

580

583

584 585

588

589

590

591

593

Ethics Statement

The datasets we used are sourced from the current public datasets. The prompts we used do not collect 547 or use personal information or information from 548 other individuals. Furthermore, they do not contain 549 any sensitive words or oppose any individual or 550 551 group. CoT-UQ has the potential to impact the credibility and reliability of LLMs, particularly in the context of reducing misinformation. LLMs 553 have the potential to generate highly plausible but 554 false information. Uncertainty quantification tech-555 niques can help distinguish between accurate and 556 misleading outputs. Successfully addressing this issue can help prevent the spread of misinformation and mitigate its potential societal impact.

Limitations 560

Our methods require access to token logits. Although commercial LLM providers widely support token logits, this still might restrict the potential application of our methods in black-box scenarios. In addition, our framework is limited to the 565 closed-ended question-answering domain, where a question has an objective ground-truth answer(s) so that we can justify the correctness of generated an-569 swer. Extensive analysis of CoT-UQ on open-ended question-answering tasks is beyond the scope of the 570 current study and is left as future work.

References

- Dhananjay Ashok and Zachary C Lipton. 2023. Promptner: Prompting for named entity recognition. arXiv preprint arXiv:2305.15444.
- Yavuz Faruk Bakman, Duygu Nur Yaldiz, Baturalp Buyukates, Chenyang Tao, Dimitrios Dimitriadis, and Salman Avestimehr. 2024. Mars: Meaningaware response scoring for uncertainty estimation in generative llms. arXiv preprint arXiv:2402.11756.
- Zhipeng Chen, Kun Zhou, Beichen Zhang, Zheng Gong, Wayne Xin Zhao, and Ji-Rong Wen. 2023. Chatcot: Tool-augmented chain-of-thought reasoning on chat-based large language models. arXiv preprint arXiv:2305.14323.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. Journal of Machine Learning Research, 24(240):1-113.
 - Zheng Chu, Jingchang Chen, Qianglong Chen, Weijiang Yu, Tao He, Haotian Wang, Weihua Peng, Ming Liu,

Bing Qin, and Ting Liu. 2023. A survey of chain of thought reasoning: Advances, frontiers and future. arXiv preprint arXiv:2309.15402.

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

- Jan Clusmann, Fiona R Kolbinger, Hannah Sophie Muti, Zunamys I Carrero, Jan-Niklas Eckardt, Narmin Ghaffari Laleh, Chiara Maria Lavinia Löffler, Sophie-Caroline Schwarzkopf, Michaela Unger, Gregory P Veldhuizen, et al. 2023. The future landscape of large language models in medicine. Communications medicine, 3(1):141.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168.
- Antonia Creswell, Murray Shanahan, and Irina Higgins. 2022. Selection-inference: Exploiting large language models for interpretable logical reasoning. arXiv preprint arXiv:2205.09712.
- Jesse Davis and Mark Goadrich. 2006. The relationship between precision-recall and roc curves. In Proceedings of the 23rd international conference on Machine learning, pages 233-240.
- Jinhao Duan, Hao Cheng, Shiqi Wang, Alex Zavalny, Chenan Wang, Renjing Xu, Bhavya Kailkhura, and Kaidi Xu. 2024a. Shifting attention to relevance: Towards the predictive uncertainty quantification of free-form large language models. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 5050-5063.
- Jinhao Duan, Renming Zhang, James Diffenderfer, Bhavya Kailkhura, Lichao Sun, Elias Stengel-Eskin, Mohit Bansal, Tianlong Chen, and Kaidi Xu. 2024b. Gtbench: Uncovering the strategic reasoning limitations of llms via game-theoretic evaluations. arXiv preprint arXiv:2402.12348.
- Tom Fawcett. 2006. An introduction to roc analysis. Pattern recognition letters, 27(8):861–874.
- Guhao Feng, Bohang Zhang, Yuntian Gu, Haotian Ye, Di He, and Liwei Wang. 2024. Towards revealing the mystery behind chain of thought: a theoretical perspective. Advances in Neural Information Processing Systems, 36.
- Tairan Fu, Javier Conde, Gonzalo Martínez, María Grandury, and Pedro Reviriego. 2025. Multiple choice questions: Reasoning makes large language models (llms) more self-confident even when they are wrong. arXiv preprint arXiv:2501.09775.
- Nuno M Guerreiro, Elena Voita, and André FT Martins. 2022. Looking for a needle in a haystack: A comprehensive study of hallucinations in neural machine translation. arXiv preprint arXiv:2208.05309.

755

Neha Gupta, Harikrishna Narasimhan, Wittawat Jitkrittum, Ankit Singh Rawat, Aditya Krishna Menon, and Sanjiv Kumar. 2024. Language model cascades: Token-level uncertainty and beyond. *arXiv preprint arXiv:2404.10136*.

647

651

656

665

666

667

670

671

672

673

674

682

701

- Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing a multi-hop qa dataset for comprehensive evaluation of reasoning steps. *arXiv preprint arXiv:2011.01060*.
- Yuheng Huang, Jiayang Song, Zhijie Wang, Shengming Zhao, Huaming Chen, Felix Juefei-Xu, and Lei Ma. 2023. Look before you leap: An exploratory study of uncertainty measurement for large language models. *arXiv preprint arXiv:2307.10236*.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38.
- Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. 2022. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*.
- Sanyam Kapoor, Nate Gruver, Manley Roberts, Katherine Collins, Arka Pal, Umang Bhatt, Adrian Weller, Samuel Dooley, Micah Goldblum, and Andrew Gordon Wilson. 2024. Large language models must be taught to know what they don't know. *arXiv preprint arXiv:2406.08391*.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199– 22213.
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. *arXiv preprint arXiv:2302.09664*.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Teaching models to express their uncertainty in words. *arXiv preprint arXiv:2205.14334*.
- Chen Ling, Xujiang Zhao, Wei Cheng, Yanchi Liu, Yiyou Sun, Xuchao Zhang, Mika Oishi, Takao Osaki, Katsushi Matsuda, Jie Ji, et al. 2024. Uncertainty decomposition and quantification for in-context learning of large language models. *arXiv preprint arXiv:2402.10189*.
- Yanming Liu, Xinyue Peng, Tianyu Du, Jianwei Yin, Weihao Liu, and Xuhong Zhang. 2024. Era-cot: Improving chain-of-thought through entity relationship analysis. *arXiv preprint arXiv:2403.06932*.
- Jieyi Long. 2023. Large language model guided tree-ofthought. arXiv preprint arXiv:2305.08291.

- Andrey Malinin and Mark Gales. 2021. Uncertainty estimation in autoregressive structured prediction. In *International Conference on Learning Representations.*
- Potsawee Manakul, Adian Liusie, and Mark JF Gales. 2023. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models. *arXiv preprint arXiv:2303.08896*.
- Ning Miao, Yee Whye Teh, and Tom Rainforth. 2023. Selfcheck: Using llms to zero-shot check their own step-by-step reasoning. *arXiv preprint arXiv:2308.00436*.
- Shen-Yun Miao, Chao-Chun Liang, and Keh-Yih Su. 2021. A diverse corpus for evaluating and developing english math word problem solvers. *arXiv preprint arXiv:2106.15772*.
- Sabrina J Mielke, Arthur Szlam, Emily Dinan, and Y-Lan Boureau. 2022. Reducing conversational agents' overconfidence through linguistic calibration. *Transactions of the Association for Computational Linguistics*, 10:857–872.

OpenAI. 2024b. Learning to reason with llms.

- Hadas Orgad, Michael Toker, Zorik Gekhman, Roi Reichart, Idan Szpektor, Hadas Kotek, and Yonatan Belinkov. 2024. Llms know more than they show: On the intrinsic representation of llm hallucinations. *arXiv preprint arXiv:2410.02707*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are nlp models really able to solve simple math word problems? *arXiv preprint arXiv:2103.07191*.
- Xin Qiu and Risto Miikkulainen. 2024. Semantic density: Uncertainty quantification in semantic space for large language models. *arXiv preprint arXiv:2405.13845*.
- Jie Ren, Yao Zhao, Tu Vu, Peter J Liu, and Balaji Lakshminarayanan. 2023. Self-evaluation improves selective generation in large language models. In *Proceedings on*, pages 49–64. PMLR.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.
- Aviv Slobodkin, Omer Goldman, Avi Caciularu, Ido Dagan, and Shauli Ravfogel. 2023. The curious case of hallucinatory (un) answerability: Finding truths in the hidden states of over-confident large language models. In *Proceedings of the 2023 Conference on*

Empirical Methods in Natural Language Processing, pages 3607–3625.

756

757

758

759

761

762

764

767

772

773

774

775

778

782 783

790

791

794

804

- Katherine Tian, Eric Mitchell, Huaxiu Yao, Christopher D Manning, and Chelsea Finn. 2023a. Finetuning language models for factuality. *arXiv preprint arXiv:2311.08401*.
- Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher D Manning. 2023b. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. *arXiv preprint arXiv:2305.14975*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971.
- Neeraj Varshney, Wenlin Yao, Hongming Zhang, Jianshu Chen, and Dong Yu. 2023. A stitch in time saves nine: Detecting and mitigating hallucinations of llms by validating low-confidence generation. *arXiv preprint arXiv:2307.03987*.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Yuxin Xiao, Paul Pu Liang, Umang Bhatt, Willie Neiswanger, Ruslan Salakhutdinov, and Louis-Philippe Morency. 2022. Uncertainty quantification with pre-trained language models: A large-scale empirical analysis. *arXiv preprint arXiv:2210.04714*.
- Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. 2023. Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms. *arXiv preprint arXiv:2306.13063*.
- 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. *arXiv preprint arXiv:1809.09600*.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2024. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36.

Fanghua Ye, Mingming Yang, Jianhui Pang, Longyue809Wang, Derek F Wong, Emine Yilmaz, Shum-
ing Shi, and Zhaopeng Tu. 2024. Benchmarking
llms via uncertainty quantification. arXiv preprint
arXiv:2401.12794.819

A Details about Considered Datasets and Baselines

A.1 Datasets.

814

815

816

817

818

819

820

822

823

824

829

830

831

832

835

837

839

841

842

847

848

852

853

Detailed introductions. We outline here all five datasets belonging to two reasoning domains that we investigate in our work. For each dataset that has been divided into training and test sets, we used samples from its test set unless otherwise instructed.

• HotpotQA (Yang et al., 2018): a dataset designed for diverse multi-hop question answering. Each entry includes Wikipedia documents that help answering the questions. We use the setting without context, where questions are asked directly without additional context.

• 2WikiMHQA (Ho et al., 2020): a dataset uses both structured and unstructured data. The dataset also introduces the evidence information containing a reasoning path for multi-hop questions. There are four types of questions: *comparison*, *inference*, *compositional*, and *bridge-comparison*. In this paper, we use the *inference* type of questions in all experiments.

• **GSM8K** (Cobbe et al., 2021): a dataset of 8.5K high quality linguistically diverse grade school math word problems created by human problem writers. These problems take between 2 and 8 steps to solve, and solutions primarily involve performing a sequence of elementary calculations using basic arithmetic operations to reach the final answer.

• **SVAMP**(Patel et al., 2021): a challenge dataset for elementary-level Math Word Problems (MWP). An MWP consists of a short Natural Language narrative that describes a state of the world and poses a question about some unknown quantities.

• ASDiv (Miao et al., 2021): a diverse (in terms of both language patterns and problem types) English math word problem (MWP) corpus for evaluating the capability of various MWP solvers. Each MWP is annotated with its problem type and grade level (for indicating the level of difficulty).

Dataset Statistics Table 2 provides detailed information about the data included in the experiment, with a minimum of 1000 samples and a total of 14563 samples taken.

Dataset	Num.	Length	Domain
HotpotQA	8447	23.2	Logical Reasoning
2WikiMHQA	1548	14.5	Logical Reasoning
GSM8K	1,319	58.9	Mathematical Reasoning
SVAMP	1000	39.4	Mathematical Reasoning
ASDiv	2249	38.2	Mathematical Reasoning

Table 2: Dataset statistics, where "Num." represents the number of sampled datasets, and "Length" is the number of average tokens in the sampled dataset.

A.2 Baselines.

Probas-mean & Probas-min: Based on the preliminaries introduced in Section 3, we exemplify the following formulation for common aggregated probabilities strategies to compute the Probas-mean baseline on the entire generated answer:

$$c = \frac{1}{N} \sum_{i=1}^{N} \mathbb{P}(y_i \mid p, y_1, \dots, y_{i-1})$$
(10) 867

861

862

863

864

865

866

868

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

Probas-min can be formalized as follows,

$$c = \min_{i \in \{1, \dots, N\}} \mathbb{P}(y_i \mid p, y_1, \dots, y_{i-1})$$
 (11) 80

TOKENSAR: Token-Level Shifting Attention to Relevance (**TOKENSAR**) is a component of the complete *SAR* method (Duan et al., 2024a) that corrects generative inequalities by reviewing the relevance of each token and emphasizing uncertainty quantification attention to those more relevant components. Formally, given a sentence s_j regarding prompt xwith the normalized relevance score for each token z_i termed as $\tilde{R}_T(z_i, s_j, x)$, the uncertainty proportions of relevant tokens are enlarged by re-weighting token entropy according to their respective \tilde{R}_T :

$$E_T(z_i, s_j, x) = -\log p(z_i \mid s_{
(12)$$

The token-level shifted (TOKENSAR) predictive entropy defined over s_j can be formulated as,

TOKEN
$$SAR(s_j, x) = \sum_{i}^{N_j} E_T(z_i, s_j, x).$$
 (13)

P(True):We follow Kadavath et al. (2022) and885prompt the LLM to judge whether its answer is correct. Our prompt followed the following template:886

P(True)

Question: [Question *q*] **A student submitted**: [LLM Answer]

- Is the student's answer:
- (A) True
- (B) False
- The student's answer is:

Self-Probing: We follow Xiong et al. (2023) and prompt the LLM with a question and its answer, then asked, *"How likely is the above answer to be correct"?* The procedure involves generating the answer in one chat session and obtaining its verbalized confidence in another independent chat session. Our prompt followed the following template:

Self-Probing

Question: [Question q] **Possible answer**: [LLM Answer]

Q: How likely is the above answer to be correct? Please first show your reasoning concisely and then answer with the following format:

Confidence: [the probability of answer [LLM Answer] to be correct, not the one you think correct, please only include the numerical number]%

B Implementation Details

B.1 LLM Hyperparameters.

For all LLMs, the max length of each generation is set to 128 tokens. The temperature of generation is respectively set to 1.0 for LlaMA3-8B and 1.2 for LLaMA2-13B, and other hyperparameters as default. Besides, the hyperparameters during the P(True) process are set following Kadavath et al. (2022), where the max token length is 1 and the temperature is followed as above setting.

B.2 Prompts.

900

901

903

904

905

906

907

908

909

910

911

912

913

914

915

916

Prompts in Stage 1 of CoT-UQ. As illustrated in Section 4.1 of the main text, Stage 1 of CoT-UQ contains three separate steps: Reasoning Extraction, Keywords Extraction, and Importance Scoring. To minimize the forward times of the LLM for the sake of computational efficiency, we merge Keyword Extraction and Importance Scoring into the same prompt under a one-shot setting. Specifically, the prompts for the first three steps are as follows,

Step 1: Reasoning Extraction

Please reason the following question step by step. Label each reasoning step as "Step i:", where "i" is the step number. You need to ensure that each step builds on the previous one and contributes meaningfully toward reaching the final answer. Once you finish all steps, put your final answer on a separate line after the reasoning steps, starting with "Final Answer:" (do not label it as a step).

Question: [Question q] **Response**: Let's think step by step.

917

Step 2 & Step 3: Keywords Extractions and Importance Scoring

You will be provided with a question and a multi-step response containing reasoning steps.

For each long reasoning step labeled "*Step i*:", extract the keywords, only the relevant tokens for that specific reasoning step.

You also need to evaluate the importance of each keyword to the final answer. Please evaluate the importance score following with the keyword by (/<importance score>/) on a scale of 1 to 10, where 1 is the least critical and 10 is the most critical.

If you find more than one keyword in a specific step, separate them with ";".

If a specific step does not contribute meaningfully to deriving the final answer (e.g., repeating information already provided in the question, introducing irrelevant assumptions or speculations), return "*Step i: NO ANSWER*" for that step. For example:

Q: [Question q]

A: [Multi-Step Response $s_{1 \sim k}$]

Keywords for Each Reasoning Step: [Extracted Keywords and Corresponding Importance \mathcal{K}]

The following is your task: **Q**: [Question q] **A**: [Multi-Step Response $s_{1\sim k}$] **Keywords for Each Reasoning Step**:

It is worth noting that if the number of extracted keywords $n_i = 0$ for the *i*-th step, we will label that

step as "Step i: NO ANSWER", indicating it does
not contribute meaningfully to deriving the final
answer. We provide the example introduced in the
prompt template for each dataset in Table 3.

Prompts in Stage 2 of CoT-UQ. As illustrated in Section 4.2 of the body of the paper, the Stage 2 of 926 CoT-UQ is responsible for integrating the extracted 927 information from reasoning path into the UQ process, where only the self-evaluation (SE) strategy need to re-prompt the LLM. For the variations on prompts, we present them based on the four strate-931 gies proposed in the SE phase, namely, ALLSteps, KEYStep, ALLKeywords, and KEYKeywords. The 933 following is the refined prompt template p_r used by 935 CoT-UQ in the SE strategy.

• ALLSteps:

P(True) w/ ALL*Steps*

Question: [Question q] **A student submitted**: [LLM Answer a]

The student explained the answer, which included a step-by-step reasoning: [Multi-Step Response $s_{1\sim k}$]

Considering these reasoning steps as additional information, is the student's answer:

- (A) True
- (B) False

The student's answer is:

Self-Probing w/ ALLSteps

Question: [Question q] **Possible answer**: [LLM Answer a] **A step-by-step reasoning to the possible answer**: [Multi-Step Response $s_{1\sim k}$]

Q: Considering these reasoning steps as additional information, how likely is the above answer to be correct? Please first show your reasoning concisely and then answer with the following format: **Confidence**: [the probability of answer [LLM Answer] to be correct, not the one you think correct, please only include the numerical number]%

P(True) w/ KEYStep

Question: [Question q] **A student submitted**: [LLM Answer a]

The student explained the answer, where the most critical step is: [Key Step s^*]

Considering this critical reasoning step as additional information, is the student's answer: (A) True

(B) False

The student's answer is:

Self-Probing w/ KEYStep

Question: [Question q] **Possible answer:** [LLM Answer a] **The most critical step in reasoning to the possible answer:** [Key Step s^*]

Q: Considering this critical reasoning step as additional information, how likely is the above answer to be correct? Please first show your reasoning concisely and then answer with the following format:

Confidence: [the probability of answer [LLM Answer] to be correct, not the one you think correct, please only include the numerical number]%

• ALL*Keywords* & KEY*Keywords*: Prompts for keywords share the same template, except for the specific keywords content.

P(True) w/ Keywords

Question: [Question q] A student submitted: [LLM Answer a] The student explained the answer, which included the following keywords: [Keywords set $\mathcal{K}/\mathcal{K}^*$] Considering these keywords as additional information, is the student's answer: (A) True (B) False The student's answer is:

938 939 941

942

943

944

Self-Probing w/ Keywords

Question: [Question q] Possible answer: [LLM Answer a] Keywords during reasoning to the possible answer: [Keywords set $\mathcal{K}/\mathcal{K}^*$]

Q: Considering these keywords as additional information, how likely is the above answer to be correct? Please first show your reasoning concisely and then answer with the following format: **Confidence**: [the probability of answer [LLM Answer] to be correct, not the one you think correct, please only include the numerical number]%

946

- 947 948
- 949
- 951 952

953

955

957

960

961

962

964

965

966

967

969

970

971

973

974

975

977

B.3 Computational Costs Analysis

CoT-UQ is more generation-efficient compared to previous methods based on response sampling. We have counted the time consumed for each step in the overall uncertainty quantification pipeline, which takes about 12 seconds per sample and a total of approximately 50 GPU hours to derive all reported results for 14653 samples. All the experiments are conducted on a server with an Intel(R) Xeon(R) Gold 5218R CPU and 8 NVIDIA A6000 GPUs.

C Additional Experimental Results and Further Discussion

C.1 Sensitivity to Keywords Filtering threshold in SE.

To study how the KEY*Keywords* is affected by the importance filtering threshold τ , we conducted experiments for τ ranging from 1 to 10 on logical reasoning tasks that we have suggested using KEYKeywords before. To ensure that the KEY*Keywords* set \mathcal{K}^* is not empty, we apply the following strategy: if the number of keywords with an importance score above the threshold τ is fewer than three, we select the top three keywords in descending order of importance as KEYKeywords set \mathcal{K}^* . Figure 6 presents the correlations between the performance of KEYKeywords and τ . It is shown that our KEYKeywords method is not particularly sensitive to τ , but performs favorably when τ takes an intermediate value, and the results consistently outperform the ALLKeywords and baseline P(True) methods in the logical reasoning task.



Figure 6: Analysis of the sensitivity to the importance filtering threshold τ on the HotpotQA benchmark. The experiments are based on LLaMA3-8B.

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

1000

1001

1003

1004

1008

C.2 Case Study

How can CoT benefit UQ in AP? We first provide a case study on the HotpotQA dataset to visualize the effect of CoT-UQ on AP strategies. Table 4 shows an example using Probas-mean. In the standard Probas-mean method, the model assigns a probability of **1.0** to the incorrect answer, leading to an uncalibrated and misleading over-confidence score. However, by incorporating keyword-level probability adjustments, our approach assigns more nuanced confidence scores to key reasoning components (e.g., "France: 0.066", "Spain: 0.073", "located further east: 0.444"). This recalibration mitigates overconfidence in incorrect predictions and ensures a more reliable confidence estimation with a confidence of **0.387**, demonstrating the effectiveness of our method in refining uncertainty quantification through reasoning-aware adjustments.

How can CoT benefit UQ in SE? In Section 5.2, we noted that CoT-UQ applied to Self-Evaluation (SE) strategies achieves a greater performance improvement in mathematical reasoning tasks. We also hypothesized that this improvement is due to the model's ability to identify critical errors in its original thought process. Here, we provide detailed evidence supporting this assumption. First, we start the analysis with the following question:

Why is access to the reasoning path beneficial for UQ in self-evaluation, especially for mathematical problems?

To answer this question, we investigate the Self-

Probing strategy, which self-evaluates the credibil-1009 ity of an LLM's generated answers. As demon-1010 strated in a case from GSM8K (Table 5), given a 1011 math question, the LLM initially generates a multi-1012 step response that leads to an incorrect final answer. 1013 When applying the standard Self-Probing strategy, 1014 which assesses only the correctness of the final an-1015 swer, the model exhibits overconfidence, assigning 1016 an 80% certainty to its response. 1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037 1038

1039

1040

1041

1042 1043

1044

1045

However, when Self-Probing incorporates the reasoning path, it successfully calibrates its confidence to 10%, as it identifies problematic areas in the original thought process (highlighted in blue in Table 5). In this case, the incorrect response misinterpreted the problem by confusing *monthly salary after promotion* with *total annual salary*. By accessing its reasoning, the model correctly identifies this mistake, clarifying the key misunderstanding.

This step-by-step breakdown helps pinpoint the exact logical misstep, making it easier to adjust confidence accordingly. Furthermore, access to the reasoning path allows the LLM to distinguish between *calculation errors* and *conceptual errors*. If the mistake were a simple arithmetic error, the model's confidence might remain relatively high, as the reasoning itself would still be sound. However, in this case, the mistake is conceptual, requiring significantly adjusting the confidence.

Thus, incorporating the reasoning path into uncertainty quantification leads to more precise confidence calibration, enabling the model to differentiate between minor computational mistakes and fundamental conceptual misunderstandings.

A complete procedure of CoT-UQ. To clarify each component of our approach, we present a case from the HotpotQA dataset in Table 6, illustrating the complete CoT-UQ process.

Table 3: Provided Examples in **Step 2 & Step 3** prompt template (Appendix B.2) for different datasets across the two reasoning domains.

Dataset	Provided Example				
HotpotQA	Q: Which band has more members, "We Are the Ocean" or "The Dream Academy"? A: Let's think step by step.				
	Step 1: The question is asking which band has more members.				
	Step 2: "We Are the Ocean" has 5 members.				
	Step 3: "The Dream Academy" has 3 members.				
	Step 4: 5 is greater than 3.				
	Step 5: Therefore, "We Are the Ocean" has more members.				
	Final Answer: We Are the Ocean				
	Keywords for Each Reasoning Step:				
	Step 1: NO ANSWER				
	Step 2: We Are the Ocean($\frac{5}{12}$; $\frac{5}{10}$)				
	Step 3: The Dream Academy(/5/); 3(/10/)				
	Step 4: greater(/7/)				
	Step 5: We Are the Ocean(/5/)				
2WikiMHQA	Q: Who is the paternal grandmother of Joseph Ferdinand Of Bavaria?				
	A: Let's think step by step.				
	Step 1: Joseph Ferdinand's father was Maximilian II Emanuel, Elector of Bavaria.				
	Step 2: Maximilian II Emanuel was the son of Ferdinand Maria, Elector of Bavaria, and his wife,				
	Henriette Adelaide of Savoy.				
	Step 3: As the mother of Maximilian II Emanuel, Henriette Adelaide of Savoy is the paternal				
	grandmother of Joseph Ferdinand of Bavaria.				
	Final Answer: Henriette Adelaide of Savoy				
	Keywords for Each Reasoning Step:				
	Step 1: father (/8/); Maximilian II Emanuel, Elector of Bavaria (/8/)				
	Step 2: son (/8/); Ferdinand Maria, Elector of Bavaria (/5/); Henriette Adelaide of Savoy (/9/)				
	Step 3: mother (/10/)				
GSM8K	Q: A robe takes 2 bolts of blue fiber and half that much white fiber. How many bolts in total				
SVAMP	does it take?				
ASDiv	A: Let's think step by step.				
	Step 1: Identify the amount of blue fiber needed. The robe requires 2 bolts of blue fiber.				
	Step 2: Determine the amount of white fiber needed. It is half the amount of blue fiber, which is				
	$2 \div 2 = 1$ bolt.				
	Step 3: Compute the total number of bolts. Add the bolts of blue fiber (2) and white fiber (1) to				
	get 3 bolts.				
	Final Answer: 3				
	Keywords for Each Reasoning Step:				
	Step 1: 2 bolts (/3/)				
	Step 2: 1 bolt (/10/)				
	Step 3: 3 bolts (/7/)				

Question: Which involved a country farther east, In Search of the Second Amendment or Without the King?

Response:

Step 1: In Search of the Second Amendment took place in France.

Step 2: Without the King took place in Spain.

Step 3: France is located further east than Spain.

Step 4: Therefore, In Search of the Second Amendment involved a country farther east.

Final Answer: In Search of the Second Amendment

Correct Answer: Without the King

Common Probas-mean: In Search of the Second Amendment 1.000 Keywords for Each Reasoning Step: *Step 1:* In Search of the Second Amendment(/5/); France(/10/) *Step 2:* Without the King(/5/); Spain(/10/) *Step 3:* located further east(/7/) *Step 4:* In Search of the Second Amendment(/5/)

Keywords Probabilities based on Probas-mean:

ney words i robublides bused on i robub inean.					
In Search of the Second Amendment	France	Without the King	Spain	located further east	
0.955	0.066	1.000	0.073	0.444	
In Search of the Second Amendment					
0.966					
Probas-mean with CoT-UQ: 0.387					

Table 4: A case study on HotpotQA demonstrating the effectiveness of CoT in AP strategy.

Question: Tim gets a promotion that offers him a 5% raise on his \$20000 a month salary. It also gives him a bonus worth half a month's salary. How much money will he make in a year?

Response: *Step 1*: Determine the amount of Tim's new salary per month. His monthly salary goes up by $5\% (0.05 \times 20000 = 1000)$ to 21000 per month.

Step 2: Compute the amount of Tim's bonus per month. The bonus is half a month's salary, which is 10500.

Step 3: Calculate the total amount Tim earns in a year. Add the monthly salary (21000) and bonus (10500) to get 31500 per year.

Final Answer: 31500

Correct Answer: 262500

Common Self-Probing: 100%.

Self-Probing with Reasoning Path: 10%.

You are confused between "Tim's new salary per month" and "Tim's total annual salary". This may seem like a trivial confusion, but reasoning is about critical reasoning. In the problem, there are two important numbers:

Tim's monthly salary = 20000

Tim's monthly salary after promotion = 21000

You are given 5% more than his monthly salary, which means he will get an additional 1000 a month. The problem asks about "Tim's total annual salary", not "Tim's monthly salary after promotion".

A possible answer is therefore $21000 \times 12 + 10500 = 262500$

Table 5: A case study on GSM8K demonstrating the effectiveness of CoT in SE strategy.

Question: Which of these board games was released earlier: Lord of the Rings or Finance? **Correct Answer**: Finance

Reasoning Extraction:

Please reason the following question step by step. Label each reasoning step as "*Step i*:", where "*i*" is the step number.

You need to ensure that each step builds on the previous one and contributes meaningfully toward reaching the final answer.

Once you finish all steps, put your final answer on a separate line after the reasoning steps, starting with "*Final Answer*:" (do not label it as a step)

Response: Step 1: Lord of the Rings was released in 2017.

Step 2: Finance was released in 2018.

Step 3: 2017 is earlier than 2018.

Step 4: Therefore, Lord of the Rings was released earlier.

Final Answer: Lord of the Rings

Keywords Extractions and Importance Scoring:

You will be provided with a question and a multi-step response containing reasoning steps.

For each long reasoning step labeled "*Step i:*", extract the keywords, only the relevant tokens for that specific reasoning step.

You also need to evaluate the importance of each keyword to the final answer. Please evaluate the importance score following with the keyword by (/<importance score>/) on a scale of 1 to 10, where 1 is the least critical and 10 is the most critical.

If you find more than one keyword in a specific step, separate them with ";".

If a specific step does not contribute meaningfully to deriving the final answer (e.g., repeating information already provided in the question, introducing irrelevant assumptions or speculations), return *"Step i: NO ANSWER"* for that step.

Question: Which of these board games was released earlier: Lord of the Rings or Finance?

Multi-Step Response: Step 1: Lord of the Rings was released in 2017.

Step 2: Finance was released in 2018.

Step 3: 2017 is earlier than 2018.

Step 4: Therefore, Lord of the Rings was released earlier.

Final Answer: Lord of the Rings

Keywords for Each Reasoning Step: Step 1: Lord of the Rings(/7/); 2017(/10/)

Step 2: Finance(/5/); 2018(/10/)

Step 3: earlier(/9/)

Step 4: Lord of the Rings(/8/)

Reasoning Enhanced Uncertainty Quantification Strategy (Exemplified by P(True) w/ KEY*Keywords***): Question**: Which of these board games was released earlier: Lord of the Rings or Finance?

A student submitted: Lord of the Rings

The student explained the answer, where the most critical step is: [2017, 2018, earlier]

Considering this critical reasoning step as additional information, is the student's answer:

(A) True

(B) False

The student's answer is:

Table 6: A case on HotpotQA demonstrates the complete procedure of CoT-UQ.