# **Quantifying Uncertainty in Natural Language Explanations of Large Language Models**

#### Sree Harsha Tanneru

Chirag Agarwal

sreeharshatanneru@g.harvard.edu

chiragagarwall120gmail.com

#### Himabindu Lakkaraju

hlakkaraju@seas.harvard.edu

Harvard University

### **Abstract**

Large Language Models (LLMs) are increasingly used as powerful tools for several high-stakes natural language processing (NLP) applications. Recent prompting works claim to elicit intermediate reasoning steps and key tokens that serve as proxy explanations for LLM predictions. However, there is no certainty whether these explanations are reliable and reflect the LLM's behavior. In this work, we make one of the first attempts at quantifying the uncertainty in explanations of LLMs. To this end, we propose two novel metrics — Verbalized Uncertainty and *Probing Uncertainty* — to quantify the uncertainty of generated explanations. While verbalized uncertainty involves prompting the LLM to express its confidence in its explanations, probing uncertainty leverages sample and model perturbations as a means to quantify the uncertainty. Our empirical analysis of benchmark datasets reveals that verbalized uncertainty is not a reliable estimate of explanation confidence. Further, we show that the probing uncertainty estimates are correlated with the faithfulness of an explanation, with lower uncertainty corresponding to explanations with higher faithfulness. Our study provides insights into the challenges and opportunities of quantifying uncertainty in LLM explanations, contributing to the broader discussion of the trustworthiness of foundation models.

## 1 Introduction

Large Language Models (LLMs), such as GPT4 [18], Bard [16], Llama-2 [25], and Claude-2 [Anthropic], have garnered significant attention and are employed across a wide range of applications, including chat-bots, computational biology, creative work, and law [9] due to their impressive natural language understanding and generation capabilities. However, state-of-the-art LLMs are complex models with billions of parameters, where their inner working mechanisms are not fully understood yet, making them less trustworthy amongst relevant stakeholders. This lack of transparency causes hindrance to deploying LLMs in high-stakes decision-making applications, where the consequences of incorrect decisions are severe and could result in the generation of harmful content, misdiagnosis [35], and hallucinations [6, 31]. The lack of user trust demands the development of robust explanation techniques to gain insights into how these powerful LLMs work.

Previous works for explaining language models can be broadly categorized into perturbation-based methods [12, 13], gradient-based methods [10, 24]), attention-based methods [4, 28], example-based methods [8, 26, 29, 32], and Natural Language Explanations (NLEs) [30]. While most of the above methods require white-box access to models (e.g., model gradients and prediction logits), NLEs

37th R0-FoMo: Workshop on Robustness of Few-shot and Zero-shot Learning in Foundation Models at NeurIPS 2023(NeurIPS 2023).

generated by LLMs enable us to understand the behavior of these models even when the models are closed-source. For instance, Chain-of-Thought (CoT) [30] explanations, a popular class of NLEs generated by LLMs show the step-by-step reasoning process leading to the outputs generated by these models. While CoTs and other natural language explanations generated by LLMs often seem quite plausible and believable [27], recent works have demonstrated that these natural language explanations may not always faithfully capture the underlying behavior of these models [27]. However, there is little to no work that focuses on deciphering if and to what extent the generated NLEs faithfully capture the behavior of the underlying model. One way to address this problem is to quantify the uncertainty in the NLEs generated by LLMs. However, this critical direction remains unexplored.

Prior works on uncertainty estimation in the context of LLMs have only focused on providing uncertainty estimates (*i.e.*, confidence) corresponding to the responses (answers) generated by LLMs [34]. While uncertainty in LLM predictions has been studied using external calibrators [7], model fine-tuning [14], and non-logit-based approaches [34], there is little to no work on estimating the uncertainty of LLM explanations. Understanding the uncertainty in natural language explanations generated by LLMs is paramount to ensuring that these explanations are trustworthy and are not just plausible hallucinations.

**Present work.** In this work, we make one of the attempts at quantifying the uncertainty in natural language explanations generated by LLMs. In particular, we propose two novel approaches – *Verbalized uncertainty* and *Probing uncertainty* metrics – to quantify the confidence of NLEs generated by large language models and compare their reliability. While verbalized uncertainty metrics focus on prompting a language model to express its uncertainty in the generated explanations, probing uncertainty metrics leverage different kinds of input perturbations (e.g., replacing words with synonyms, paraphrasing inputs) and measure the consistency of the resulting explanations. Using our proposed metrics, we provide the first definition of uncertainty estimation of language model explanations. In addition, our work also demonstrates key connections between *uncertainty* and *faithfulness* of natural language explanations generated by LLMs.

We evaluate the effectiveness of our proposed metrics on three math word problems and two commonsense reasoning benchmark datasets and conduct experiments using different GPT variants. Our empirical results across these datasets and LLMs reveal the following key findings. 1) Verbalized uncertainty is not a reliable estimate of explanation confidence and LLMs often exhibit very high verbalized confidence in the explanations they generate. 2) Probing uncertainty is correlated with the predictive performance of the LLM, where correct answers from a model tend to generate more confident/less uncertain explanations. 3) A clear connection exists between the uncertainty and faithfulness of an explanation, where less uncertain explanations tend to be more faithful to the model predictions.

# 2 Preliminaries

**Notations.** Large language models typically have a single vocabulary  $\mathcal V$  that represents a set of unique "tokens" (words or sub-words). Let  $\mathcal M:Q\to A$  denote a language model mapping a sequence of n question tokens  $Q=(q_1,q_2,\ldots,q_n)$  to sequence of m answer tokens  $A=(a_1,a_2,\ldots,a_m)$ , where  $q_i$  and  $a_i$  are text tokens belonging to the model vocabulary  $\mathcal V$ . In addition to the original question Q, we design specific prompts  $Q_e$  to generate natural language explanation (NLE)  $A_e$  from the language model  $\mathcal M$ .

**Uncertainty.** Black-box LLMs do not provide access to parameter gradients or model logits, rendering traditional explainability techniques ineffective. To this end, most language models leverage NLEs, which are explanations generated from the language model to serve as proxy explanations and are a viable alternative. While NLEs are essentially a sequence of tokens sampled from the model that serve as explanations, there is an associated uncertainty for the generated explanations. Quantifying the uncertainty of these explanations is essential to estimate the reliability of generated NLEs. For the rest of the paper, we will use the term "confidence score" to refer to the uncertainty of an explanation, as determined by the language model.

**Explanation Methods.** We confine our study to two explanation methods — Token Importance and Chain of Thought (CoT) explanations. While token importance explanations [12, 33] aims to identify input tokens (refer to tokens t in an input text T for LLMs) that most contribute to a model's predictions, CoT explanations [30] focus on revealing the sequence of operations or reasoning steps  $S_i \in S$  the language model  $\mathcal{M}$  takes when processing the question Q and arriving at its predictions,

Read the question, and assign each word an importance score between 0 and 100 of how important it is for your answer. The output format is as follows:

Word: [Word 1 here], Importance: [Your importance score here]

. . .

Word: [Word N here], Importance: [Your importance score here]

**Final answer and overall confidence (0-100):** [Your answer as a number here], [Your confidence here]

Note: The importance scores of all words should add up to 100. The overall confidence score indicates the degree of certainty you have about your importance scores. For instance, if your confidence level is 80%, it means you are 80% certain that importance scores assigned are correct. Provide the answer in aforementioned format, and nothing else.

**Q:** Jake has 11 fewer peaches than Steven. If Jake has 17 peaches. How many peaches does Steven have?

**Answer:** 

Word: Jake, Importance: 20% Word: Steven, Importance: 20% Word: peaches, Importance: 60%

Final answer and overall confidence (0-100): 28, 100%

Figure 1: Template for generating token importance and its confidence. The prompt  $Q_e$  appended to the original question Q to elicit a token importance explanation TI. We ask the underlying LLM to verbally assign an importance score to each word in the question Q and then provide the final answer A with overall confidence.

where  $n_s = |S|$  denotes the total number of steps in a CoT explanation. For token importance explanation, we concatenate a prompt  $Q_e$  to the given question Q using the template: "Read the question and output the words important for your final answer...". Whereas, the prompt  $Q_e$  to generate CoT explanations uses the following template: "Read the question, give your answer by analyzing step by step, ...". Please refer to Figures 12-13 in appendix for more details.

We generate an answer from the LLM  $\mathcal{M}$  as follows:  $\mathcal{M}(Q) = A$ . We also generate an explanation  $A_e$  along with answer A using the aforementioned template question  $Q_e$  as:  $\mathcal{M}(Q_e + Q) = A + A_e$ .

# **3 Quantifying Uncertainty in Explanations**

Next, we describe our metrics which aim to estimate the uncertainty in token importance and CoT explanations generated by LLMs.

**Problem formulation (Uncertainty in Explanations).** Given a question-answer pair and prompt  $Q_e$  to generate natural language explanation  $A_e$  from the model  $\mathcal{M}:(Q,Q_e)\to(A,A_e)$ , we aim to develop an uncertainty function UNC:  $A_e\to[0,1]$ , which maps a generated explanation  $A_e$  to a scalar score that determines the uncertainty in the generated explanation, i.e.,

Uncertainty = 
$$UNC(A_e)$$
,

where  $\mathcal{M}(Q_e + Q) = A + A_e$ .

As mentioned before, we confine our study to two natural language explanation methods – Token Importance and CoT. We use  $\mathrm{TI}_q: \{w \mid w \in Q\}$  to denote a token importance explanation which is a subset of words in the question Q that are important for predicting the answer A and  $\mathrm{CoT}_q: \{(S_1,c_1)\to (S_2,c_2)\cdots\to (S_{n_s},c_{n_s})\}$  a CoT explanation for a prediction A from question Q. Here  $S_i=(s_1,s_2,\ldots s_{n_s})$  is a text sequence denoting the natural language reasoning and  $c_i\in[0,1]$  is the LLM's confidence of CoT step  $S_i$ .

## 3.1 Verbalized Uncertainty

A straightforward approach to elicit uncertainty of an explanation is to directly request the LLM  $\mathcal{M}$  to output a confidence score for the explanation ranging from 0% to 100%. By directly soliciting the

Read the question, give your answer by analyzing step by step, and assign a confidence level to each step and the final answer. The output format is as follows:

**Step 1**: [Your reasoning here], **Confidence**: [Your confidence here]% **Step 2**: [Your reasoning here], **Confidence**: [Your confidence here]%

•••

**Step N**: [Your reasoning here], **Confidence**: [Your confidence here]%

**Final answer and overall confidence (0-100):** [Your answer as a number here], [Your confidence here]%

Note: The confidence indicates the degree of certainty you have about your reasoning. For instance, if your confidence level is 80%, it means you are 80% certain that your reasoning is correct. Provide the answer in aforementioned format, and nothing else.

**Q:** Jake has 11 fewer peaches than Steven. If Jake has 17 peaches. How many peaches does Steven have?

#### **Answer:**

Step 1: Jake has 11 fewer peaches than Steven. Confidence: 100%

Step 2: Jake has 17 peaches. Confidence: 100%

Step 3: If Jake has 17 peaches, then Steven has 17 + 11 = 28 peaches. Confidence: 100%

Final answer and overall confidence (0-100): 28, 100%

Figure 2: Template for generating CoT explanation and its step-wise confidence. The prompt  $Q_e$  appended to the original question Q to elicit a CoT explanation. We ask the underlying LLM to verbally assign an importance score to each step of the CoT explanation and then provide the final answer A with overall confidence.

model's self-assessment of uncertainty, this approach seeks to extract explicit uncertainty information inherent in the model. We provide the template of the prompts for confidence elicitation for token importance and CoT explanations in Figures 1-2. For token importance, we ask the underlying LLM to verbally assign an importance score to each word in the question Q and then provide the final answer of the question with overall confidence (see Fig. 1). In contrast, for CoT explanations (see Fig. 2), we ask the LLM to assign verbalized confidence to each step in the CoT reasoning and the final answer.

## 3.2 Probing Uncertainty

Verbalized uncertainty elicits confidence in an explanation by directly requesting the underlying LLM to output a confidence score in a given range. In contrast, for estimating uncertainty using probing, we leverage the consistency of explanations as a measure to estimate the uncertainty in explanations generated by a language model  $\mathcal{M}$ . More specifically, let  $A_e$  denote the natural language explanation generated by the model  $\mathcal{M}$  for a given question Q and  $[A_{e_1}, A_{e_2}, \dots, A_{e_N}]$  be N explanations generated for N perturbation of the same question using its local neighborhood. Next, we describe two different perturbation strategies to generate N explanations for a given question and answer.

**Sample Probing.** Motivated by the local neighborhood approximation works in XAI [21, 22], we propose uncertainty metrics that leverage the consistency of a model in generating the explanation in a local neighborhood. Here, we presume that the local behavior of the underlying LLM is consistent for perturbed samples of the original question and gradually introduce perturbations in the questions by *paraphrasing* the original question Q. Given a question Q, we paraphrase the question into N different forms  $\{Q_1, Q_2, \ldots, Q_N\}$ , such that each paraphrased question  $Q_i$  is semantically equivalent to Q, and the true reasoning process remains the same, *i.e.*, given a question: "Jake has 11 fewer peaches than Steven. If Jake has 17 peaches. How many peaches does Steven have?", some of its local paraphrased counterparts used to calculate uncertainty in explanations are i)...What is the number of peaches Steven has? ii)...How many peaches is Steven in possession of? iii)...How many peaches does Steven possess? Next, we generate the explanations using the LLM by probing the model using the paraphrased questions  $Q_i$ . Mathematically,

$$\mathcal{M}(Q_e + Q_i) = A_i + A_{e_i} \; ; \; i = 1, 2, \dots, N$$
 (1)

where  $Q_i$  is a paraphrased form of question Q,  $Q_e$  is the prompt to generate explanations, and  $A_{e_i}$  is the corresponding generated explanation.

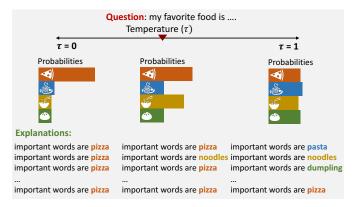


Figure 3: The impact of the temperature  $\tau$  on model stochasticity. We find that as  $\tau$  increases, the stochasticity in model responses increases.  $\tau$ =0 gives near-deterministic answers to a question, whereas  $\tau$ =1 gives a distribution of answers.

we use the "temperature" parameter  $\tau$  present in most LLMs that control the randomness in the generated answers by using the probability distribution of each generated token. A high value of temperature indicates an even distribution among all tokens, and a lower value of temperature indicates a sharper distribution (see Figure 3). As the temperature parameter increases, the language model becomes more creative and stochastic in the generated explanations. Intuitively, the temperature parameter affects the sampling process when generating answers from the model. For a given question Q (say "my favorite food is ..."), we sample N answers and their corresponding explanations,  $\{A_i, A_{e_i}\} \forall i \in 1, 2, \ldots, N$  from the language model M. Mathematically, we can denote this using:

$$\mathcal{M}(Q_e + Q) = A_i + A_{e_i} \; ; \; i \in \{1, 2, \dots, N\}$$
 (2)

where  $A_i$  is the  $i^{\text{th}}$  answer generated by the LLM for a given temperature  $\tau$  and  $A_{e_i}$  is its respective explanation.

#### 3.2.1 Token Importance Uncertainty

Using the above sample and model perturbation strategies, N perturbed natural language explanations  $A_{e_i}$  are generated for a given question Q, answer A, original explanation  $A_e$ . Next, we describe the metrics for estimating explanation confidence from these perturbed explanations.

We define the uncertainty in token importance explanations as the mean agreement between perturbed explanations and the original explanation. Two token importance explanations are said to agree with each other if they employ the same set of important words to arrive at a prediction. To quantify token importance uncertainty, we use token agreement and token rank metrics. While token agreement computes the fraction of important tokens that are common between two different explanations, token rank measures the fraction of important tokens that have the same position in their respective rank orders. The token rank (TR) metric is defined below:

$$TR(TI_{i}, TI_{j}, k) = \frac{1}{k} \Big( \bigcup_{s \in S} \{ s | s \in Tokens(TI_{i}, k) \land s \in$$

$$Tokens(TI_{j}, k) \land R(TI_{i}, s) = R(TI_{j}, s) \} \Big),$$
(3)

where  $TI_i$  and  $TI_j$  are any two given token important explanations,  $Tokens(TI_i, k)$  is the first k tokens in explanation  $TI_i$ , k denotes the topK tokens a user wants as explanations, and  $R(\cdot)$  function gives the rank of the word s in a token importance explanation TI. The uncertainty in token importance explanation is defined as the mean agreement between the perturbed explanations  $TI_{e_i}$  and the original explanation  $TI_{original}$ .

$$UNC_{TI} = \frac{1}{N} \sum_{i=1}^{N} TR(TI_{e_i}, TI_{original}, k),$$
(4)

#### 3.2.2 Chain of Thought Uncertainty

While the agreement between token importance explanations is intuitive, the agreement between the chain of thought explanations is non-trivial as each explanation has a sequence of steps  $S_i$  in natural language as output explanations. To check if the two steps in CoT explanations are equivalent, we propose using pre-trained sentence encoder models [20]. Let us consider two CoT explanations that generate  $N_a$  and  $N_b$  steps in their respective explanations, *i.e.*, (CoT<sub>a</sub> =  $(s_{a_1}, s_{a_2}, \ldots, s_{a_{N_a}})$  and CoT<sub>b</sub> =  $(s_{b_1}, s_{b_2}, \ldots, s_{b_{N_b}})$ . We define CoT agreement metric (CoTA) that measures the agreement between any two given CoT explanations as:

$$CoTA(CoT_{a}, CoT_{b}) = \frac{1}{N_{a} + N_{b}} \Big( \sum_{i=1}^{N_{a}} \max_{j \in 1, ..., N_{b}} E(s_{a_{i}}, s_{b_{j}}) + \sum_{i=1}^{N_{b}} \max_{i \in 1, ..., N_{a}} E(s_{a_{i}}, s_{b_{j}}) \Big),$$
(5)

The intuition behind the above metric is that for every step in the a CoT explanation, we check if there exists a step in other CoT explanation which agrees with it.  $E(\cdot,\cdot)$  denotes the entailment model that focuses on the task of textual entailment or natural language inference (NLI). The goal of NLI is to determine the logical relationship between two sentences, usually framed as "entailment", "contradiction", or "neutral". Formally, the entailment score between two explanation steps is defined as:

$$E(s_i, s_j) = \begin{cases} 1 & \text{if statements entail each other} \\ 0 & \text{if statements do not entail each other} \end{cases}$$

Finally, the uncertainty in the CoT explanation is calculated as the mean agreement of the perturbed chain of thought explanations with the original explanation.

$$UNC_{CoT} = \frac{1}{N} \sum_{i=1}^{N} CoTA(CoT_i, CoT_{original})$$
 (6)

To summarize, we introduce a metric for calculating the agreement between two CoT explanations (Eq. 5). In addition, we generate N perturbed explanations for a question, and calculate the mean agreement of perturbed explanations with the original explanation to estimate explanation uncertainty (Eq. 6).

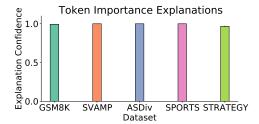
# 4 Experiments

Next, we validate the effectiveness of our proposed uncertainty metric which amounts to asking: What is the uncertainty in explanations generated by state-of-the-art LLMs with respect to different explanation methods? More specifically, we focus on the following research questions: RQ1) Does verbalized uncertainty estimation depict overconfidence in LLMs? RQ2) Is there a relation between uncertainty and faithfulness of an explanation? RQ3) How does explanation confidence vary for correct and incorrect answers? RQ4) Are changes in the metric parameters necessary for quantifying uncertainty in explanations?

#### 4.1 Datasets and Experimental Setup

We first describe the datasets and large language models used to study the uncertainty in explanations and then outline the experimental setup.

**Datasets.** We conduct experiments using three math word problem and two commonsense reasoning benchmark datasets. i) the **GSM8K** dataset that comprises several math word problems [3], ii) the **SVAMP** dataset contains math word problems with varying structures [19], iii) the **ASDiv** dataset consisting of diverse math word problems [17], iv) the **StrategyQA** [5] requires a language model to deduce a multi-step reasoning strategy to answer questions and v) the **Sports Understanding** dataset, which is a specialized evaluation set from the BIG-bench [23] that involves determining whether a sentence relating to sports is plausible or implausible.



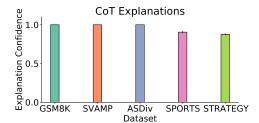


Figure 4: Verbalized explanation confidence of Token Importance and CoT explanations on three math word problems and two commonsense reasoning datasets. We observe that the verbalized explanation confidence is mostly high for explanations across all five datasets.

**Large language models.** We generate and evaluate the uncertainty in explanations by generating explanations using three large language models — InstructGPT, GPT-3.5, and GPT-4.

**Performance metrics.** Some recent works [2, 11, 15] have explored defining faithfulness for natural language explanations. i) Faithfulness of token importance explanations: We use the counterfactual test [2] for NLEs by intervening on input tokens and checking whether the explanation reflects these tokens. Specifically, we replace identified importance tokens in the explanation with synonyms and check whether the new explanation reflects these changes. Faithfulness is then quantified by the rank agreement (Eq. 3) between the new explanations and the expected explanation with intervened tokens. ii) Faithfulness of chain of thought explanations: Recent works that explored the topic of faithfulness in CoT explanations don't explicitly quantify the faithfulness of an individual explanation. Hence, we follow suit and follow Lanham et al. [11] to measure faithfulness at a dataset level. In our experiments, we use a strategy called "Early Answering" proposed by Lanham et al. [11] to measure the faithfulness of CoT explanations. It involves truncating the previously collected reasoning samples and prompting the model to answer the question with the partial CoT rather than the complete one, i.e., for a question Q and CoT  $[s_1, s_2, \dots s_n]$ , the model is prompted to answer with  $Q + s_1, Q + s_1 + s_2$ , until,  $Q + s_1 + s_2 \cdots + s_n$ . After collecting answers with each truncation of the CoT, we measure how often the model comes to the same conclusion as it did with the complete CoT. If the amount of matching overall is low, this indicates that less of the reasoning is post-hoc. If the reasoning is not post-hoc, there are fewer ways for it to be unfaithful than there are for reasoning which is post-hoc [11].

**Implementation details.** To run the paraphrase probing uncertainty, we formulate 10 semantically equivalent paraphrases of every question to measure uncertainty using sample probing. In the model probing uncertainty experiment, we sample five natural language explanations at a temperature of 1.0. To compute the rank agreement of token importance explanations, we use the top-3 words *i.e.*, k = 3. We run on a randomly sampled subset of 100 samples for each dataset. See the Appendix for more implementation details.

## 4.2 Results

Next, we discuss experimental results to answer questions (RQ1-RQ4) about uncertainty in explanations.

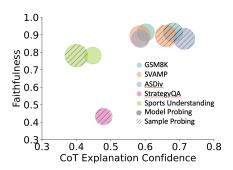
**RQ1)** Analyzing verbalized uncertainty. Verbalized confidence scores of both natural language explanation methods are almost always 100%. It raises questions about whether these uncertainty estimates are reliable. If the confidence in every explanation is the same, it is impossible to know when to trust the generated explanation and when not to. Our results in Figure 4 show that, on average, across both explanation methods and five datasets, the verbalized confidence is 94.46%. Our analysis of these methods uncovers that LLMs often exhibit a high degree of overconfidence when verbalizing their uncertainty in explanations. The verbalized uncertainty for commonsense reasoning datasets is lower than math word problem datasets but still very close to 100% with little standard deviation.

**RQ2**) Less uncertain explanations are more faithful. A model's explanation is said to be faithful if it reflects the true reasoning behind the prediction. For token importance explanations, we swap important words in explanations with synonyms and check if the corresponding replacements are reflected in the new explanation. In Fig. 7, we demonstrate that explanation confidence is correlated



Figure 5: Chain of thought explanation confidence distributions on three math word problems and two commonsense reasoning datasets using GPT-3.5. On average, across two probing strategies and five datasets, correct answers (in green) obtain higher explanation confidence than wrong answers (in red). See Table 1 in appendix for t-test statistics comparing explanation confidence scores of correct and incorrect answers to different datasets.

with faithfulness, and highly confident (certain) explanations are more faithful. In addition, we find a similar trend between the CoT explanation confidence and its faithfulness (see Fig. 6) and find that increased mean explanation confidence lead to an increase in the faithfulness of an explanation for most datasets. Our observations suggest that uncertainty estimation can be used as a test for the faithfulness of NLE, i.e., whether the explanation reflects the true reasoning process of the model.



0.8 Faithfulness 0.2 0.0 0.2 0.4 0.5 0.6 0.7 0.9 TI Explanation Confidence

Sports Understanding

Figure 6: Mean explanation confidence for CoT explanations generated using InstructGPT for five datasets. We find that the explanation confidence is positively correlated with faithfulness for four datasets, i.e., highly confident explanations tend to be more faithful. The circle size denotes the deviation in the confidence.

Figure 7: Mean explanation confidence for token importance explanations generated using InstructGPT for five datasets. We find that the explanation confidence is positively correlated with faithfulness, i.e., highly confident explanations tend to be more faithful.

**RQ3**) Correct answers have more certain explanations. Across five datasets and two probing uncertainty metrics, Fig. 5 shows that explanations of correct answers have higher explanation confidence compared to explanations of wrong answers. Our observation aligns with the general expectation that models tend to provide more reliable and confident explanations when they make correct predictions as opposed to incorrect ones.

**RQ4**) Ablation study. We conduct ablation on three key components of our proposed probing metrics i) the number of paraphrases we generate in sample probing, ii) the number of samples we generate at temperature  $\tau = 1$  in model probing, and iii) different LLMs (see Figs. 10-11 in appendix for more details and results). Results in Figure 8 show that the explanation confidence saturates as we increase the number of paraphrases of the original question Q and our chosen value of 10 is well justified. In addition, we observe that the explanation confidence using our proposed model probing technique shows similar behavior irrespective of the number of responses we generate using the LLM at  $\tau = 1$ (Figure 9). These findings explain our choice of hyperparameters in quantifying the uncertainty in explanations generated using different NLE techniques.

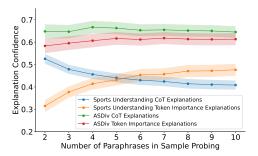


Figure 8: The effect of the number of paraphrased samples of the original Q on the mean explanation confidence of CoT and TI explanations generated from InstructGPT for Sports Understanding and ASDiv datasets. We observe that the confidence saturates as we increase the number of paraphrased samples.

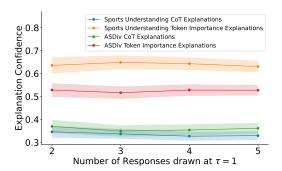


Figure 9: The effect of the number of responses drawn at  $\tau = 1$  on the mean explanation confidence of CoT and TI explanations generated from InstructGPT for Sports Understanding and ASDiv datasets. We observe that the confidence remains consistent irrespective of the number of responses generated using InstructGPT.

# 5 Conclusions

While improving the explainability of LLMs is crucial to establishing user trust, and better understanding the limitations and unintended biases present in LLMs, it is crucial to quantify the reliability of the generated explanations using uncertainty estimates. In this work, we present a novel way to estimate the uncertainty of natural language explanations (NLEs) using verbalized and probing techniques. Specifically, we propose uncertainty metrics to quantify the confidence of generated NLEs from LLMs and compare their reliability. We test the effectiveness of our metrics on math word problem and commonsense reasoning datasets and find that i) LLMs exhibit a high degree of overconfidence when verbalizing their uncertainty in explanations, ii) explanation confidence is positively correlated with explanation faithfulness, and iii) correct predictions tend to have more certain CoT explanations compared to incorrect predictions. Our work paves the way for several exciting future works in understanding the uncertainty of the natural language explanations generated by LLMs.

# References

[Anthropic] Anthropic. Model-card-claude-2.pdf. https://www-files.anthropic.com/production/images/Model-Card-Claude-2.pdf. (Accessed on 10/12/2023).

- [2] Atanasova, P., Camburu, O.-M., Lioma, C., Lukasiewicz, T., Simonsen, J. G., and Augenstein, I. (2023). Faithfulness tests for natural language explanations.
- [3] Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., et al. (2021). Training verifiers to solve math word problems. *arXiv*.
- [4] DeRose, J. F., Wang, J., and Berger, M. (2020). Attention flows: Analyzing and comparing attention mechanisms in language models.
- [5] Geva, M., Khashabi, D., Segal, E., Khot, T., Roth, D., and Berant, J. (2021). Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. *TACL*.
- [6] Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., Ishii, E., Bang, Y. J., Madotto, A., and Fung, P. (2023). Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38.
- [7] Jiang, Z., Araki, J., Ding, H., and Neubig, G. (2021). How can we know when language models know? on the calibration of language models for question answering. *Transactions of the Association for Computational Linguistics*, 9:962–977.
- [8] Jin, D., Jin, Z., Zhou, J. T., and Szolovits, P. (2020). Is bert really robust? a strong baseline for natural language attack on text classification and entailment.

- [9] Kaddour, J., Harris, J., Mozes, M., Bradley, H., Raileanu, R., and McHardy, R. (2023). Challenges and applications of large language models.
- [10] Kindermans, P.-J., Hooker, S., Adebayo, J., Alber, M., Schütt, K. T., Dähne, S., Erhan, D., and Kim, B. (2017). The (un)reliability of saliency methods.
- [11] Lanham, T., Chen, A., Radhakrishnan, A., Steiner, B., Denison, C., Hernandez, D., Li, D., Durmus, E., Hubinger, E., Kernion, J., Lukošiūtė, K., Nguyen, K., Cheng, N., Joseph, N., Schiefer, N., Rausch, O., Larson, R., McCandlish, S., Kundu, S., Kadavath, S., Yang, S., Henighan, T., Maxwell, T., Telleen-Lawton, T., Hume, T., Hatfield-Dodds, Z., Kaplan, J., Brauner, J., Bowman, S. R., and Perez, E. (2023). Measuring faithfulness in chain-of-thought reasoning.
- [12] Li, J., Chen, X., Hovy, E., and Jurafsky, D. (2016a). Visualizing and understanding neural models in nlp.
- [13] Li, J., Monroe, W., and Jurafsky, D. (2016b). Understanding neural networks through representation erasure. *CoRR*, abs/1612.08220.
- [14] Lin, S., Hilton, J., and Evans, O. (2022). Teaching models to express their uncertainty in words.
- [15] Lyu, Q., Havaldar, S., Stein, A., Zhang, L., Rao, D., Wong, E., Apidianaki, M., and Callison-Burch, C. (2023). Faithful chain-of-thought reasoning.
- [16] Manyika, J. (2023). An overview of bard: an early experiment with generative ai.
- [17] Miao, S.-Y., Liang, C.-C., and Su, K.-Y. (2021). A diverse corpus for evaluating and developing english math word problem solvers. *arXiv*.
- [18] OpenAI, R. (2023). Gpt-4 technical report. arXiv, pages 2303–08774.
- [19] Patel, A., Bhattamishra, S., and Goyal, N. (2021). Are nlp models really able to solve simple math word problems? *arXiv*.
- [20] Reimers, N. and Gurevych, I. (2019). Sentence-bert: Sentence embeddings using siamese bert-networks.
- [21] Ribeiro, M. T., Singh, S., and Guestrin, C. (2016). "Why should I trust you?" Explaining the predictions of any classifier. In *KDD*.
- [22] Smilkov, D., Thorat, N., Kim, B., Viégas, F., and Wattenberg, M. (2017). Smoothgrad: Removing noise by adding noise. *arXiv*.
- [23] Srivastava, A., Rastogi, A., Rao, A., Shoeb, A. A. M., Abid, A., Fisch, A., Brown, A. R., Santoro, A., Gupta, A., Garriga-Alonso, A., et al. (2022). Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv*.
- [24] Sundararajan, M., Taly, A., and Yan, Q. (2017). Axiomatic attribution for deep networks. In Precup, D. and Teh, Y. W., editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 3319–3328. PMLR.
- [25] Touvron, H. (2023). Llama 2: Open foundation and fine-tuned chat models.
- [26] Treviso, M., Ross, A., Guerreiro, N. M., and Martins, A. (2023). CREST: A joint framework for rationalization and counterfactual text generation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15109–15126, Toronto, Canada. Association for Computational Linguistics.
- [27] Turpin, M., Michael, J., Perez, E., and Bowman, S. R. (2023). Language models don't always say what they think: Unfaithful explanations in chain-of-thought prompting.
- [28] Vig, J. (2019). Visualizing attention in transformer-based language representation models.
- [29] Wang, B., Xu, C., Liu, X., Cheng, Y., and Li, B. (2022). Semattack: Natural textual attacks via different semantic spaces.

- [30] Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q., and Zhou, D. (2023). Chain-of-thought prompting elicits reasoning in large language models.
- [31] Weidinger, L., Mellor, J., Rauh, M., Griffin, C., Uesato, J., Huang, P.-S., Cheng, M., Glaese, M., Balle, B., Kasirzadeh, A., Kenton, Z., Brown, S., Hawkins, W., Stepleton, T., Biles, C., Birhane, A., Haas, J., Rimell, L., Hendricks, L. A., Isaac, W., Legassick, S., Irving, G., and Gabriel, I. (2021). Ethical and social risks of harm from language models.
- [32] Wu, T., Ribeiro, M. T., Heer, J., and Weld, D. (2021). Polyjuice: Generating counterfactuals for explaining, evaluating, and improving models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6707–6723, Online. Association for Computational Linguistics.
- [33] Wu, Z., Chen, Y., Kao, B., and Liu, Q. (2020). Perturbed masking: Parameter-free probing for analyzing and interpreting BERT. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4166–4176, Online. Association for Computational Linguistics.
- [34] Xiong, M., Hu, Z., Lu, X., Li, Y., Fu, J., He, J., and Hooi, B. (2023). Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms.
- [35] Zhang, J., Sun, K., Jagadeesh, A., Ghahfarokhi, M., Gupta, D., Gupta, A., Gupta, V., and Guo, Y. (2023). The potential and pitfalls of using a large language model such as chatgpt or gpt-4 as a clinical assistant.

# A Appendix

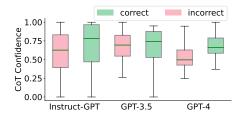


Figure 10: Comparison of chain of thought explanation uncertainty using sample probing across InstructGPT, GPT-3.5, and GPT-4 models on GSM8K dataset. We observe that the trend of correct answers having less uncertain explanations holds true across models.

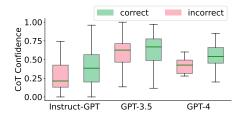


Figure 11: Comparision of chain of thought explanation uncertainty using model probing across InstructGPT, GPT-3.5, and GPT-4 models on GSM8K dataset. We observe that the trend of correct answers having less uncertain explanations holds true across models.

Read the question, and output the words important for your final answer, sorted in descending order of importance. The output format is as follows:

1. [Word 1 here]
2. [Word 2 here]
...
N. [Word N here]
Final Answer and Overall Confidence (0-100): [Your answer as a number here], [Your confidence here]%. Provide the answer in aforementioned format, and nothing else.

Figure 12: The prompt  $Q_e$  prepended to the question Q to elicit a token importance explanation TI along with an answer A.

## **Prompts**

The questions used to generate chain of thought and token importance explanations are described in 13 and 12 respectively. For sample probing and model probing uncertainty, we further tailor the prompt according to the dataset. Tailoring the question prompt helps in parsing answers and explanations from generated text. The prompts used are as follows GSM8K 15 14, ASDiv 17 16, SVAMP 19 18, StrategyQA 21 20, and Sports Understanding 23 22.

## Paraphrased Questions in Sample Probing

Semantically equivalent paraphrased questions are generated using INSTRUCTGPT using the following prompt - "Paraphrase the question into 10 different forms with the same meaning, and share them as a Python list of double quotes enclosed strings". An example is shown in 2.

Read the question, give your answer by analyzing step by step, and assign a confidence level to each step and the final answer. The output format is as follows:

**Step 1**: [Your reasoning here]

Step 2: [Your reasoning here]

Step 3:

•••

Step N: [Your reasoning here]

**Final Answer and Overall Confidence (0-100):** [Your answer as a number here], [Your confidence here]% Note: The confidence indicates the degree of certainty you have about your reasoning. For instance, if your confidence level is 80%, it means you are 80% certain that your reasoning is correct. Provide the answer in aforementioned format, and nothing else.

Figure 13: The prompt  $Q_e$  prepended to the question Q to elicit a chain of thought explanation CoT along with an answer A.

Table 1: T-Test Result Comparing Explanation Confidence Scores of Correct and Incorrect Answers using GPT-3.5 and InstructGPT models for Chain of Thought Explanations of GSM8K dataset.

Dataset	<b>Uncertainty Metric</b>	T-Statistic	P-Value
GSM8K	Sample Probing	-0.0977	0.9224
	Model Probing	0.7400	0.4611
SVAMP	Sample Probing	1.7913	0.0763
	Model Probing	1.2307	0.2214
ASDiv	Sample Probing	1.3031	0.1959
	Model Probing	1.7922	0.0765
StrategyQA	Sample Probing	-0.2752	0.7838
	Model Probing	-0.9779	0.3305
Sports Understanding	Sample Probing	1.3941	0.1665
	Model Probing	1.0851	0.2806

(i) GPT-3.5

Dataset	<b>Uncertainty Metric</b>	T-Statistic	P-Value
GSM8K	Sample Probing	1.5694	0.1198
	Model Probing	3.2404	0.0016
SVAMP	Sample Probing	2.6388	0.0097
	Model Probing	0.7660	0.4455
ASDiv	Sample Probing	3.7558	0.0003
	Model Probing	5.1783	0.0000
StrategyQA	Sample Probing	-0.1642	0.8699
	Model Probing	-0.1015	0.9194
Sports Understanding	Sample Probing	-0.8499	0.3975
	Model Probing	0.6971	0.4874

(ii) InstructGPT

```
Read the question, and output the words important for your final answer, sorted in descending order of importance. The output format is as follows:

1. [Word 1 here]
2. [Word 2 here]
...
N. [Word N here]
Final Answer and Overall Confidence (0-100): [Your answer as a number here], [Your confidence here]%. Provide the answer in aforementioned format, and nothing else.
```

Figure 14: **GSM8K** dataset. The prompt  $Q_e$  prepended to the paraphrased question Q to generate a token importance explanation TI along with an answer A in sample probing and model probing uncertainty experiments.

```
Read the question, give your answer by analyzing step by step, and assign a confidence level to each step and the final answer. The output format is as follows:

Step 1: [Your reasoning here], Confidence: [Your confidence here]%

Step 2: [Your reasoning here], Confidence: [Your confidence here]%

Step 3:
...

Step N: [Your reasoning here], Confidence: [Your confidence here]%

Final Answer and Overall Confidence (0-100): [Your answer as a number here], [Your confidence here]% Note: The confidence indicates the degree of certainty you have about your answer. For instance, if your confidence level is 80%, it means you are 80% certain that your answer is correct. Provide the answer in aforementioned format, and nothing else.
```

Figure 15: **GSM8K** dataset. The prompt  $Q_e$  prepended to the paraphrased question Q to elicit a chain of thought explanation CoT along with an answer A in sample probing and model probing uncertainty experiments.

```
Read the question, and output the words important for your final answer, sorted in descending order of importance. The output format is as follows:

1. [Word 1 here]
2. [Word 2 here]
...
N. [Word N here]
Final Answer and Overall Confidence (0-100): [Your answer as a number here], [Your confidence here]%. Provide the answer in aforementioned format, and nothing else.
```

Figure 16: **ASDiv** dataset. The prompt  $Q_e$  prepended to the paraphrased question Q to generate a token importance explanation TI along with an answer A in sample probing and model probing uncertainty experiments.

Read the question, give your answer by analyzing step by step, and assign a confidence level to each step and the final answer. The output format is as follows:

Step 1: [Your reasoning here], Confidence: [Your confidence here]%

Step 2:

...

Step 3:

...

Step N:

...

Final Answer and Overall Confidence (0-100): [Your answer as a number here], [Your confidence here]% Note: The confidence indicates the degree of certainty you have about your answer. For instance, if your confidence level is 80%, it means you are 80% certain that your answer is correct.

Figure 17: **ASDiv** dataset. The prompt  $Q_e$  prepended to the paraphrased question Q to elicit a chain of thought explanation CoT along with an answer A in sample probing and model probing uncertainty experiments.

Provide the answer in aforementioned format, and nothing else.

Read the question, and output the words important for your final answer, sorted in descending order of importance. The output format is as follows:

1. [Word 1 here]
2. [Word 2 here]
...
N. [Word N here]
Final Answer and Overall Confidence (0-100): [Your answer as a number here], [Your confidence here]%. Provide the answer in aforementioned format, and nothing else.

Figure 18: **SVAMP** dataset. The prompt  $Q_e$  prepended to the paraphrased question Q to generate a token importance explanation TI along with an answer A in sample probing and model probing uncertainty experiments.

```
Read the question, give your answer by analyzing step by step, and assign a confidence level to each step and the final answer. The output format is as follows:

Step 1: [Your reasoning here], Confidence: [Your confidence here]%

Step 2:
...

Step 3:
...

Step N:
...

Final Answer and Overall Confidence (0-100): [Your answer as a number here], [Your confidence here]% Note: The confidence indicates the degree of certainty you have about your answer. For instance, if your confidence level is 80%, it means you are 80% certain that your answer is correct. Provide the answer in aforementioned format, and nothing else.
```

Figure 19: **SVAMP** dataset. The prompt  $Q_e$  prepended to the paraphrased question Q to elicit a chain of thought explanation CoT along with an answer A in sample probing and model probing uncertainty experiments.

Read the question, and output the words important for your final answer, sorted in descending order of importance. The output format is as follows:

1. [Word 1 here]
2. [Word 2 here]
...
N. [Word N here]
Final Answer and Overall Confidence (0-100): [Your answer Yes/No here], [Your confidence here]%. Provide the answer in aforementioned format, and nothing else.

Figure 20: **StrategyQA** dataset. The prompt  $Q_e$  prepended to the paraphrased question Q to generate a token importance explanation TI along with an answer A in sample probing and model probing uncertainty experiments.

Read the question, give your answer by analyzing step by step, and assign a confidence level to each step and the final answer. The output format is as follows:

Step 1: [Your reasoning here], Confidence: [Your confidence here]%

Step 2:
...

Step 3:
...

Step N:
...

Final Answer and Overall Confidence (0-100): [Your answer Yes/No here], [Your confidence here]% Note: The confidence indicates the degree of certainty you have about your answer. For instance, if your confidence level is 80%, it means you are 80% certain that your answer is correct. Provide the answer in aforementioned format, and nothing else.

Figure 21: **StrategyQA** dataset. The prompt  $Q_e$  prepended to the paraphrased question Q to elicit a chain of thought explanation CoT along with an answer A in sample probing and model probing uncertainty experiments.

Read the question, and output the words important for your final answer, sorted in descending order of importance. The output format is as follows:

1. [Word 1 here]
2. [Word 2 here]
...
N. [Word N here]
Final Answer and Overall Confidence (0-100): [Your answer plausible / implausible here], [Your confidence here]%. Provide the answer in aforementioned format, and nothing else.

Figure 22: **Sports Understanding** dataset. The prompt  $Q_e$  prepended to the paraphrased question Q to generate a token importance explanation TI along with an answer A in sample probing and model probing uncertainty experiments.

Read the question, give your answer by analyzing step by step, and assign a confidence level to each step and the final answer. The output format is as follows:

**Step 1**: [Your reasoning here], Confidence: [Your confidence here]%

Step 2: ... Step 3: ...

otep 5

Step N: ...

**Final Answer and Overall Confidence (0-100):** [Your answer plausible / implausible here], [Your confidence here]% Note: The confidence indicates the degree of certainty you have about your answer. For instance, if your confidence level is 80%, it means you are 80% certain that your answer is correct. Provide the answer in aforementioned format, and nothing else.

Figure 23: **Sports Understanding** dataset. The prompt  $Q_e$  prepended to the paraphrased question Q to elicit a chain of thought explanation CoT along with an answer A in sample probing and model probing uncertainty experiments.

Table 2: Paraphrased Samples for a question in GSM8K math word problem dataset. The original question is "How many signatures do the sisters need to collect to reach their goal?"

What is the number of signatures the sisters need to collect to reach their goal?		
How many signatures must the sisters acquire to reach their goal?		
What is the amount of signatures the sisters need to collect to reach their goal?		
How many signatures do the sisters have to collect to reach their goal?		
What is the total number of signatures the sisters need to collect to reach their		
goal?		
How many signatures do the sisters require to reach their goal?		
What is the quantity of signatures the sisters need to collect to reach their goal?		
How many signatures do the sisters need to gather to reach their goal?		
What is the sum of signatures the sisters need to collect to reach their goal?		
How many signatures do the sisters need to acquire to reach their goal?		