

# From Evidence to Belief: A Bayesian Epistemology Approach to Language Models

Minsu Kim, Sangryul Kim, James Thorne

KAIST AI

{minsu\_kim, sangryul, thorne}@kaist.ac.kr

## Abstract

This paper investigates the knowledge of language models from the perspective of Bayesian epistemology. We explore how language models adjust their confidence and responses when presented with evidence with varying levels of informativeness and reliability. To study these properties, we create a dataset with various types of evidence and analyze language models' responses and confidence using verbalized confidence, token probability, and sampling. We observed that language models do not consistently follow Bayesian epistemology: language models follow the Bayesian confirmation assumption well with true evidence but fail to adhere to other Bayesian assumptions when encountering different evidence types. Also, we demonstrated that language models can exhibit high confidence when given strong evidence, but this does not always guarantee high accuracy. Our analysis also reveals that language models are biased toward golden evidence and show varying performance depending on the degree of irrelevance, helping explain why they deviate from Bayesian assumptions.

## 1 Introduction

Large Language models (LLMs) have advanced to the point where they can naturally respond to various practical tasks such as question-answering, code generation and conversation (OpenAI et al., 2023; Gemini Team et al., 2024). However, limitations like hallucination and trustworthiness still exist, and research efforts continue to address these issues (Huang et al., 2023; Sun et al., 2024; Xiao and Wang, 2021; Zhang et al., 2023). In this paper, we take a different approach by examining large language models from a philosophical perspective: we investigate whether language models can be said to possess knowledge. In epistemology, knowledge is traditionally analyzed using three conditions—truth, justification, and belief—often associated with the justified true belief (JTB) framework (Audi, 1997). Prior NLP research has focused

on two aspects: factual correctness (i.e. the *truth* condition) – assessing whether the response of a model is correct (Hendrycks et al., 2021; Srivastava et al., 2023) – and *justification*, which encompasses explanation generation (Wei et al., 2023; Camburu et al., 2018) and evidence finding (Thorne et al., 2018).

In this work, we investigate whether the model believes its own responses (the *belief* condition): specifically, the relationship between belief and the language model's justification, expressed as evidence. Since belief is a challenging concept to define, this paper focuses on belief from the perspective of Bayesian epistemology, which interprets belief as a quantitative and functional variable. According to Bayesian epistemology, the degree of belief can be interpreted and measured as probability, called *probability norm*. In particular, regarding the confirmation of belief, we should adjust the confidence of belief based on evidence. Specifically, when  $H$  represents the hypothesis (or belief),  $E$ , the evidence for the belief, and  $\theta$  representing the background information or prior knowledge, we can define 3 primitive assumptions:<sup>1</sup>

**Confirmation Assumption:**  $E$  confirms  $H$  if and only if  $P(H | E, \theta) > P(H | \theta)$

**Disconfirmation Assumption:**  $E$  disconfirms  $H$  if and only if  $P(H | \theta) > P(H | E, \theta)$ .

**Irrelevance Assumption:**  $E$  is irrelevant to  $H$  if and only if  $P(H | \theta) = P(H | E, \theta)$ .

<sup>1</sup>According to Hájek (2003); Vassend (2023), we treat conditional probability as a primitive concept representing the likelihood of an event occurring under certain conditions, rather than relying on the standard ratio formula  $P(H | E) = \frac{P(H \cap E)}{P(E)}$ . This approach allows us to apply the Bayesian assumption more intuitively, even in cases where  $P(E) = 0$  or with contradictory evidence. For instance, when observing an unlikely event, such as a shark in a freshwater lake, we still make judgments based on the observation despite it contradicting common knowledge. A detailed discussion can be found in Hájek (2003).

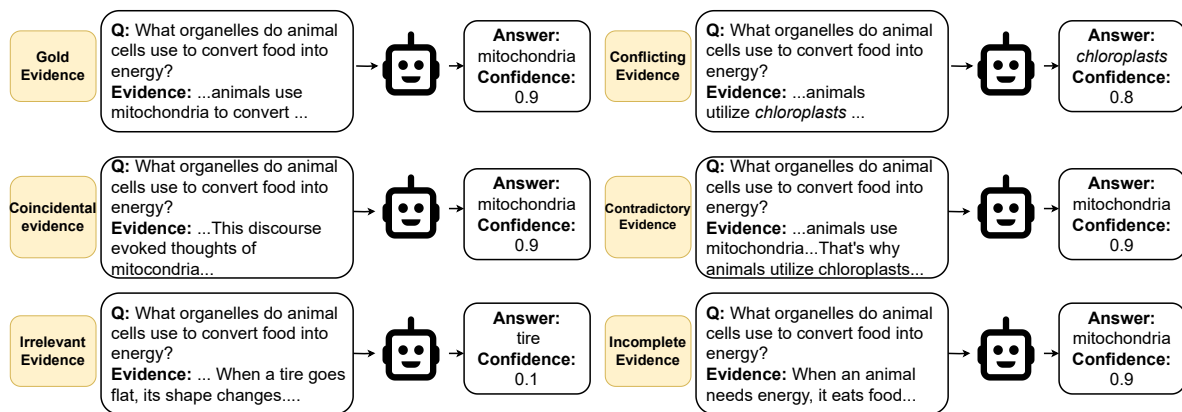


Figure 1: The overall experimental method and simple examples of the evidence dataset for Confirmation Task. As golden evidence that aligns with the question is given to language models, it shows high confidence and accuracy. However, if language models encounter irrelevant evidence, it responds with low confidence. Further results and analysis are reported in Section 4.1.

Also, if we define belief in terms of probability, the strength of the evidence should also be reflected in the confidence. That is,

- **Evidence Power Assumption:**  $E'$  confirms  $H$  more strongly than  $E''$  if and only if  $P(H | E', \theta) > P(H | E'', \theta)$

(Horwich, 1982; Howson, 2000; Talbott, 2006; Hájek and Hartmann, 2010). The degrees of belief should not only be a probability. The probabilities assigned to these beliefs must align with the *calibration norm*, meaning they should correspond to the actual likelihood of the event occurring, that is, the actual frequency (Williamson, 2010).

The goal of this paper is to explore whether different types of evidence are reflected in language models' confidence and responses. The evidence here is not merely perturbations altering the correctness of information, i.e., informativeness, but our dataset also includes variations and modifications for reliability factors such as coincidence, timeliness, source of credibility, etc.

Our paper shows that language models can exhibit high confidence and accuracy when encountering true evidence but respond inconsistently with conflicting evidence and reduce confidence and accuracy with irrelevant evidence, contrary to Bayesian assumptions. We also found that LLMs are biased toward golden evidence (typical of annotated information that forms part of datasets) and perform differently across the gradient of irrelevance, supporting an understanding of why LLMs deviate from Bayesian epistemology.

## 2 Related Works

**Calibration of LLMs** Calibration of language models is a key metric for ensuring faithful responses, with log probabilities often representing model confidence (Kadavath et al., 2022; Lee et al., 2023; Guo et al., 2017). As models have scaled, research has expanded to verbalized confidence, where models generate their own confidence (Lin et al., 2022; Mielke et al., 2022; Tian et al., 2023b). While confidence can improve model performance (Zhao et al., 2023; Tian et al., 2023a), some studies focus on interpreting this confidence as measuring uncertainty in semantic space (Kuhn et al., 2023) and exploring confidence through prompt and sampling methods (Xiong et al., 2024).

Zhou et al. (2023) explored how epistemic markers affect calibration. However, unlike their focus on linguistic markers, our work examines how changes in epistemic evidence, containing information on both content and reliability, influence confidence and calibration. Yu and Ji (2024) report that logical probability and language model probability can differ. This supports our findings that language models may not incorporate evidence into their responses and confidence.

**Adversarial Context** With in-context learning, studies have examined how few-shot demonstrations and explanations affect responses (Brown et al., 2020; Wei et al., 2022). Wang et al. (2023a) showed even inaccurate demonstrations could be used in Chain-of-Thought (COT) prompting, while Chia et al. (2023) improved question accuracy with

contrastive demonstrations. [Chen et al. \(2023\)](#) studied how the number of demonstrations impacts accuracy. [Feng et al. \(2023\)](#) measures language models’ responses and probability shifts in temporal relations based on subtle contextual changes. While these works focus on accuracy, we explore how direct question evidence influences not only accuracy but also confidence and calibration.

[Turpin et al. \(2023\)](#); [Lanham et al. \(2023\)](#) tested perturbations in COT inputs and their impact on answers, similar to our approach. While they focused on modifying explanations based on informativeness (e.g., incorrectness or relevance), our paper investigates whether LLMs reflect diverse evidence in their confidence and calibration, specifically exploring the effects of coincidental evidence and varying source credibility on model confidence.

### 3 Methods

We first generate various types of evidence by few-shot prompting large language models (LLMs) with questions and annotated support from SciQ ([Welbl et al., 2017](#)), TriviaQA ([Joshi et al., 2017](#)), GSM8K ([Cobbe et al., 2021](#)). Refer to Appendix I.2 for details. We then used this evidence to evaluate how the models’ confidence and responses change based on the type of evidence provided as shown in Figure 1. Influenced by Bayesian epistemology, we defined a confirmation task to measure whether language models can reflect the confirmation, disconfirmation, or irrelevance assumption introduced in Section 1. Also, we created a strength-of-evidence task to assess LLM’s ability to represent the various power of evidence. To measure the probability norm for adjusting confidence according to the evidence, we used an average confidence across all samples. In order to measure the response, such as correctness or calibration norm, we used accuracy (ACC) and Expected Calibration Error (ECE). In both the confirmation task and the strength-of-evidence task, we used zero-shot prompting for inference.

#### 3.1 Experimental Design

We estimated the confidence of language models using verbalized confidence (Verb. 1S top-1) ([Tian et al., 2023b](#)), token probability, and sampling ([Lee et al., 2023](#); [Xiong et al., 2024](#)). Refer to Appendix H.2 and I.1 for details. Smaller-scale open-source LLMs did not tend to generate responses in the correct format matching the prompt of verbalized

confidence. Also, following observations from [Tian et al. \(2023b\)](#) that closed-source models are better at generating verbal confidence than open-source models, we used GPT-3.5-turbo-0125 and GPT-4o-2024-05-13 for inference. We used SciQ ([Welbl et al., 2017](#)), TriviaQA ([Joshi et al., 2017](#)) and GSM8K ([Cobbe et al., 2021](#)) as the source datasets for our Confirmation task, and used only SciQ dataset for Strength of Evidence task, as a scientific question is suitable for making various degree of reliable evidence (see Appendix H for experimental details and dataset statistics).

#### 3.2 Confirmation Task

The objective of the confirmation task to observe and analyze the changes in the language model’s confidence and responses when presented with various types of evidence, compared to scenarios where the language models receive the original evidence,  $E$ , or in the absence of evidence, and assess how these changes align with three assumptions: Confirmation, Disconfirmation, and Irrelevance introduced section 1. Let the entire dataset be

$$D = \{(Q_i, A_i, E_i) \mid Q_i \text{ is a question,} \\ A_i \text{ is an answer for } Q_i, \\ \text{and } E_i \text{ is evidence for } Q_i \text{ and } A_i\}. \quad (1)$$

and

$$E_i = (s_{i1}, s_{i2}, \dots, s_{in}) \quad (2)$$

where  $s_{ij}$  is an evidence sentence in the collection of sentences  $E_i$  indexed by  $j = \{1, \dots, n\}$ . For the experiment, we need to create modified  $(Q_i, A_i, E'_i)$  where  $E'_i$  is a perturbation of  $E_i$ . The following are the types of  $E'_i$ :

##### 1. Conflicting Evidence

Conflicting evidence refers to abnormal information that hinders reaching the correct answer, introducing misinformation or conflicting beliefs with the golden evidence  $E_i$ . Specifically, evidence where  $s_{ij}$  in  $E_i$  are replaced with their conflicting counterpart sentences  $\tilde{s}_{ij}$ . Thus,  $E'_i$  is conflicting evidence if and only if

$$E'_i = (\tilde{s}_{i1}, \tilde{s}_{i2}, \dots, \tilde{s}_{in}), \forall s_{ij} \in E_i.$$

##### 2. Incomplete Evidence

Evidence that includes only a subset of sentences from the original evidence collection  $E_i$ . Thus,  $E'_i$  is a proper subset of  $E_i$ . In our experiments, we discard approximately half the sentences from  $E_i$  (i.e.  $|E'_i| \approx 0.5 \times |E_i|$ ).

### 3. Contradictory Evidence

The original evidence  $E_i$  is concatenated with additional negated sentences from  $E_i$ . Thus,  $E'_i$  is contradictory evidence if and only if

$$E'_i = E_i \cup N \quad \text{where } N \subset \{\tilde{s}_{ij} \mid s_{ij} \in E_i\}$$

such that  $|N| = 0.5 \times |E_i|$ . That is, adding 50% of the conflicting evidence to the original evidence.

### 4. Irrelevant Evidence

Irrelevant evidence is  $E'_i = E_j$  where  $j \neq i$ . That is,  $E_i$  is randomly shuffled within the dataset  $D$  so that the evidence  $E_i$  of tuple  $(Q_i, A_i, E_i)$  is replaced with evidence  $E_j$  from a different tuple  $(Q_j, A_j, E_j)$ .

### 5. Coincidental Evidence

For the SciQ and TriviaQA dataset, unlike other types of evidence, coincidental evidence does not include incorrect answers but explanations reaching the golden answer by irrational reasoning or epistemic luck. Examples include explanations derived from guessing or vague memories. For GSM8K, coincidental evidence consists of a wrong reasoning process but a correct final answer.

We can see the concise examples of the evidence in Figure 1.

## 3.3 Strength of Evidence

This task differs from the Confirmation task in that it focuses on the strength of evidence. Unlike the modified  $E'$  used in the Confirmation task, the evidence used here includes the correct answer but perturbation of reliability. The goal is to understand how differences in the strength of evidence impact confidence and calibration and assess whether LLMs align with Evidence Power Assumption in section 1. For each  $(Q_i, A_i)$  pair, two types of perturbation  $(Q_i, A_i, E'_i)$  and  $(Q_i, A_i, E''_i)$  are created.  $E'_i$  represents more reliable evidence, while  $E''_i$  represents relatively less reliable evidence. The following are the types of evidence:

#### 1. Source of Credibility

For each  $(Q_i, A_i)$  pair,  $E'_i$  means evidence from a highly reputable and authoritative source, while  $E''_i$  means evidence from an anonymous online post or an individual.

#### 2. Specificity and Detail

This involves varying the detail and specificity of the evidence. Similar to source of credibility, for each  $(Q_i, A_i)$ ,  $E'_i$  is highly detailed evidence, while  $E''_i$  is evidence with general mentions related to the question.

#### 3. Timeliness

This involves modifying the evidence based on its recency. For each  $(Q_i, A_i)$ ,  $E'_i$  consists of recent findings and experiments, while  $E''_i$  consists of relatively older findings and experiments.

#### 4. Experimental Evidence

For each  $(Q_i, A_i)$ ,  $E'_i$  includes evidence derived from precise and controlled experiments, while  $E''_i$  includes evidence where the answer is observed by a witness without experiments.

## 4 Results and Analysis

### 4.1 LLMs on Confirmation task

The results of the Confirmation task using verbalized confidence are in Table 1, while the token probability and sampling methods are shown in Tables 2 and 3 in Appendix A. In Tables 4, 5 and 6 located in Appendix B, we calculated p-values to compare confidence, accuracy, and ECE between providing no evidence (labeled No\_EVI), original (labeled EVI), and perturbed evidence sets to determine if there were significant differences in model performance across these metrics based on the type of evidence. Tables 1, 2, and 3, show the changes in various metrics based on the evidence, and Table 4, 5, and 6 indicate whether those changes are statistically significant.

**LLMs follow confirmation assumption** In Tables 1, 2 and 3, the NO\_EVI and EVI column show that when  $E$  is golden evidence that helps confirm the answer, we observe  $P(H \mid E) > P(H)$  across all models, datasets and methods we used. In Tables 4, 5, and 6, when golden evidence is provided, the p-value for confidence showed at least a marginal increase, with a significant difference particularly observed in verbal, which align well with the *Confirmation Assumption* of Bayesian epistemology. Moreover, both accuracy (ACC) and expected calibration error (ECE) showed improved results when given such confirming evidence, which leads to at least marginal difference in p-value, except for ECE in verbal. This indicates that language

	Dataset	Metric	No_EVI	EVI	Coincidence	Irrelevant	Conflict	Incomplete	Contradiction
GPT-3.5-turbo	SciQ	Confidence	0.851	0.943	0.835	0.714	0.827	0.928	0.945
		Accuracy $\uparrow$	0.67	0.841	0.854	0.53	0.572	0.77	0.847
		ECE $\downarrow$	0.18	0.111	0.071	0.262	0.304	0.161	0.108
	Trivia	Confidence	0.827	0.922	0.818	0.69	0.797	0.897	0.925
		Accuracy $\uparrow$	0.846	0.879	0.971	0.698	0.702	0.86	0.869
		ECE $\downarrow$	0.035	0.058	0.153	0.125	0.211	0.06	0.076
	GSM8K	Confidence	0.74	0.998	0.988	0.765	0.931	0.96	0.949
		Accuracy $\uparrow$	0.078	0.951	0.843	0.066	0.023	0.666	0.777
		ECE $\downarrow$	0.662	0.048	0.148	0.699	0.911	0.307	0.197
GPT-4o	SciQ	Confidence	0.925	0.986	0.902	0.861	0.875	0.948	0.977
		Accuracy $\uparrow$	0.73	0.915	0.88	0.7	0.675	0.82	0.905
		ECE $\downarrow$	0.195	0.073	0.04	0.171	0.2	0.128	0.072
	Trivia	Confidence	0.915	0.933	0.895	0.878	0.866	0.909	0.926
		Accuracy $\uparrow$	0.94	0.96	0.99	0.935	0.86	0.945	0.955
		ECE $\downarrow$	0.037	0.027	0.095	0.063	0.048	0.036	0.037
	GSM8K	Confidence	0.924	0.991	0.83	0.89	0.883	0.96	0.957
		Accuracy $\uparrow$	0.24	0.97	0.54	0.195	0.165	0.774	0.96
		ECE $\downarrow$	0.684	0.033	0.406	0.705	0.718	0.186	0.013

Table 1: The result of confirmation task with verbal confidence methods. We used 200 samples for GPT-4o due to the cost limit. NO\_EVI refers the question with no context which means  $P(H | \theta)$ , serving as baseline. Others are the case of  $P(H | E, \theta)$  where evidence appears in the context. EVI refers to the context in which the golden evidence from the dataset is given, while the other evidence types are those mentioned in section 3.2.

models have strong confidence and handle information well when the evidence contains purely helpful information for deriving the correct answer. This indicates language models satisfy the probability norm and calibration norm in the confirmation case.

Unlike SciQ and Trivia, the accuracy significantly improves when evidence is provided, and ECE significantly decreases in GSM8K. It shows that language models have parametric knowledge about SciQ and Trivia datasets and struggle with complex reasoning tasks without explanations and reaffirms the importance of explanation in arithmetic tasks (Wei et al., 2023).

**LLMs inconsistently disconfirm conflicting evidence** Based on Tables 1, 2, and 3, when conflicting evidence that does not lead to the correct answer is provided, confidence tends to decrease. However, the p-value in Tables 4, 5, and 6 reveals that conflicting evidence does not have a significant effect on confidence. With further investigation, only GPT-4o was found to significantly reduce its confidence with a p-value of 0.003.

On the other hand, accuracy decreased significantly across all confidence methods. ECE significantly increased only in the verbal and token methods. Low confidence indicates that LLMs do not follow the conflicting evidence to generate an answer, but rather that the conflicting evidence creates confusion with existing parametric knowledge, which leads to lower accuracy and higher ECE.

We made the following assumptions regarding this issue: Conflicting evidence can cause the model’s confidence to become inconsistent, with overly low or high confidence appearing regardless of accuracy, which may increase ECE (although low confidence is likely to occur more frequently). Alternatively, conflicting evidence can disrupt the model’s learned patterns, making its confidence less reflective of actual accuracy, thereby leading to higher ECE.

In conclusion, only GPT-4o with verbal method exhibited behavior aligned with Bayesian disconfirmation assumptions, showing a decrease in both confidence and accuracy when conflicting evidence was presented.

**LLMs handle contradictory evidence as golden evidence** In most models and methods, contradictory evidence, which contains both correct and conflicting evidence in the context, shows increased confidence and accuracy compared to the no-evidence baseline, with p-values close to 0.05 for both cases. Additionally, ECE also decreases with p-value around 0.1. It means that despite the presence of conflicting information, the model appears highly confident and well-calibrated in almost all scenarios, which is similar with golden evidence case. This suggests that LLMs can effectively filter the given context and generate responses without conflicting with their parametric knowledge. Unlike the case with conflicting evi-

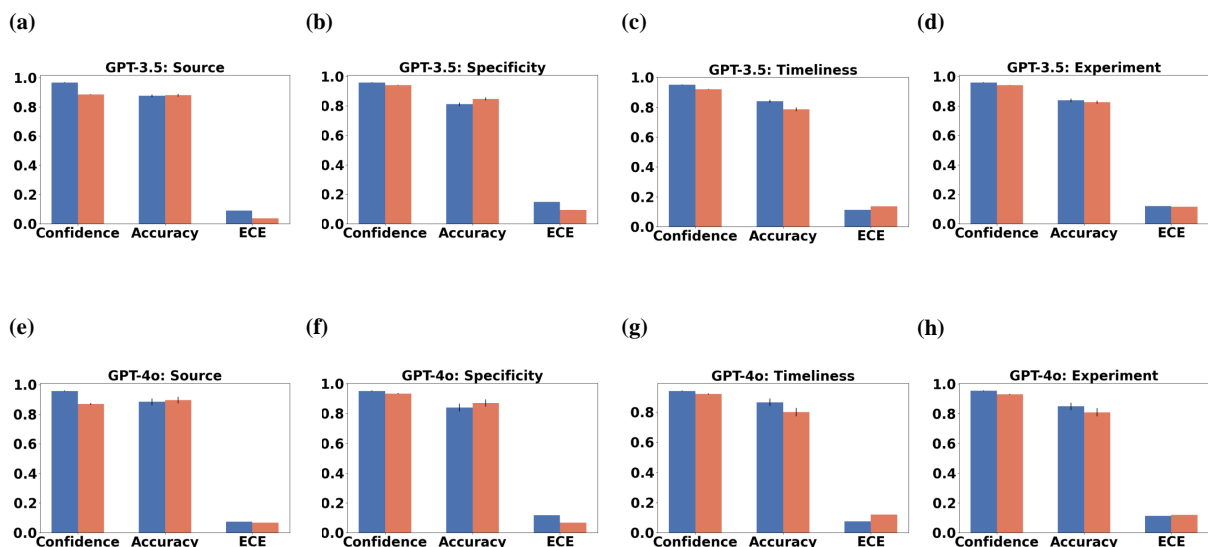


Figure 2: The results of the Strength of Evidence task on the SciQ dataset with verbal confidence method. The blue bar indicates more credible, specific, recent, and experimental evidence, while the red bar represents less credible, less specific, older, and observational evidence provided to the LLMs. We found that, in all models and datasets, strong evidence leads to high confidence with verbalized confidence. However, it does not always result in improvements in ACC and ECE.

dence, it can be interpreted that the influence of golden evidence offsets the presence of incorrect sentences. Hence, LLMs do not consider contradictory evidence as disconfirming their beliefs.

#### LLMs cannot handle coincidental evidence well.

When coincidental evidence was provided in both the token and sampling methods, confidence increased marginally (See Tables 2, 3, 5, and 6). In terms of verbal confidence in Tables 1 and 4, it was found that there was no significant difference between reliable and unreliable evidence. These indicate that the language model fails to capture the unreliability of evidence and, as a result, cannot properly reflect this in its confidence. However, based on the increased accuracy and its meaningful p-value results, it was revealed that the language model incorporated this unreliable evidence into its response to arrive at the correct answer. The ECE also decreased for both the token probability method and the sampling method. The general results for token probability and sampling methods suggest that when confidence is measured using these methods, coincidental evidence can have similar effects as golden evidence. On the other hand, there was no significant difference in ECE when using verbal confidence.

In conclusion, the language model appears to have problems with handling unreliable evidence. This is likely because, during training, the model

is mostly exposed to correct data and does not frequently encounter incorrect situations.

#### Incomplete evidence acts as a positive hint.

Incomplete evidence, though not as strong as the golden evidence case, led to a marginal increase in confidence across all confidence methods. Accuracy and ECE also showed a tendency to marginally increase and decrease, respectively. In most cases, displaying a pattern very similar to that of the golden evidence case. Incomplete evidence does not contain inaccurate information and is a partial subset of the gold evidence, acting as a hint. Similar to the contradictory evidence case, we observe that the language model is biased towards imperfect golden evidence. Therefore, while not as effective as golden evidence, the language model reflects the information from the evidence well without distraction.

#### LLMs are confused by irrelevant evidence

Except for the sampling method, when irrelevant evidence was provided, confidence decreased either statistically significantly or marginally. In terms of accuracy, it was found that accuracy significantly decreased. Furthermore, in the verbal method, ECE is also notably increased.

This indicates that language models are severely distracted by irrelevant text in terms of the content of the evidence as in Shi et al. (2023). These results

showed that LLMs do not align with the irrelevance assumption of Bayesian epistemology.

## 4.2 LLMs on Strength of Evidence task

The results of the Strength of Evidence task using the verbalized confidence method, token probability method, sampling method are reported in Figure 2, with Figures 4 and 5 located in Appendix C. In Table 7, we calculated p-values to compare confidence, accuracy, and ECE between more reliable evidence and less reliable evidence to determine if there were significant differences in model performance across these metrics based on the strength of the evidence. Hence, as in section 4.1, we can observe the variations in different metrics depending on the power of evidence in Figures 2, 4, and 5, and check the statistical significance in Table 7.

**Strong evidence can give high confidence in Verbal and Sampling methods, but cannot guarantee accurate response** In Figures 2, 5 and Table 7, the confidence increases when more reliable evidence is provided in verbal and sampling methods. Additionally, the p-values for verbal confidence and sampling confidence between weak and strong evidence are 0.015 and 0.038, respectively. This indicates that the model’s response confidence significantly increases when strong evidence is provided.

In verbal methods, as in (a), (b), (e), (f) in Figure 2, when low credible source and low detailed evidence were used, accuracy increased and ECE decreased. This suggests that in some cases, strong evidence may not be as useful as we expected for the language model to infer the correct answer. High confidence combined with low accuracy ultimately leads to overconfidence in incorrect predictions, resulting in high ECE. On the other hand, as in (c), (d), (g), (h) in Figure 2, evidence containing the latest information or experiments showed higher confidence and accuracy compared to older information or observation-based evidence. Except for GPT-3.5 with experimental evidence, the ECE of stronger evidence was also lower, indicating that using stronger evidence in the cases of timeliness and experiments results in well-calibrated models. This means that in these cases, the language model utilizes the given evidence effectively and accurately reflects the information in its predictions.

The sampling confidence showed higher accuracy when high-reliability evidence is provided in most cases except for specificity. We consider this phenomenon another positive aspect of self-

consistent decoding (Wang et al., 2023b). A single response might not fully capture the reliability of evidence, such as credibility, timeliness, etc. However, multiple responses can increase the likelihood of accurately reflecting these aspects.

In the case of specificity, both verbalized confidence and sampling failed to adequately to reflect the concreteness of the evidence in the responses. We interpreted that more detailed information can enhance confidence, but it also suggests that such excessive information may hinder the extraction of correct answers that match the question.

**Token probability cannot reflect various degrees of reliability.** As in Figure 4 in Appendix C, with token probability, confidence did not increase even when stronger evidence was presented. For example, with token probability, when the specificity of the evidence was altered or when the source’s credibility was varied in GPT-4o, it failed to reflect confidence according to the strength of the evidence accurately. However, it accurately reflected reliability changes according to the source’s credibility, timeliness, and whether an experiment was conducted in the evidence to its accuracy. Additionally, it showed a decrease in ECE in cases of timeliness and experimental evidence.

Through this experiment, we found that when stronger evidence is provided to the language model, it can significantly increase its verbalized and sampling confidence. However, this does not always lead to improvements in accuracy-related performance.

## 5 Ablation

**LLMs tend to focus more on correct than incorrect information.** In Section 4.1, we interpreted that the language model possesses a certain degree of knowledge about the question in its parameters and tends to be biased towards contexts aligned with this parametric knowledge rather than context hindering it, as seen in golden, contradictory, and incomplete evidence. To justify this, we conducted an experiment adjusting the ratio of golden sentences in conflicting, incomplete, and contradictory evidence. Figure 3 (a) and (d) show that as the number of original golden sentences decreases and the conflicting sentence increases, the performance of the language model gradually declines. However, it decreases significantly when there are no golden sentences left. Moreover, Figure 3 (b) and (e) demonstrate that as the original golden sentence

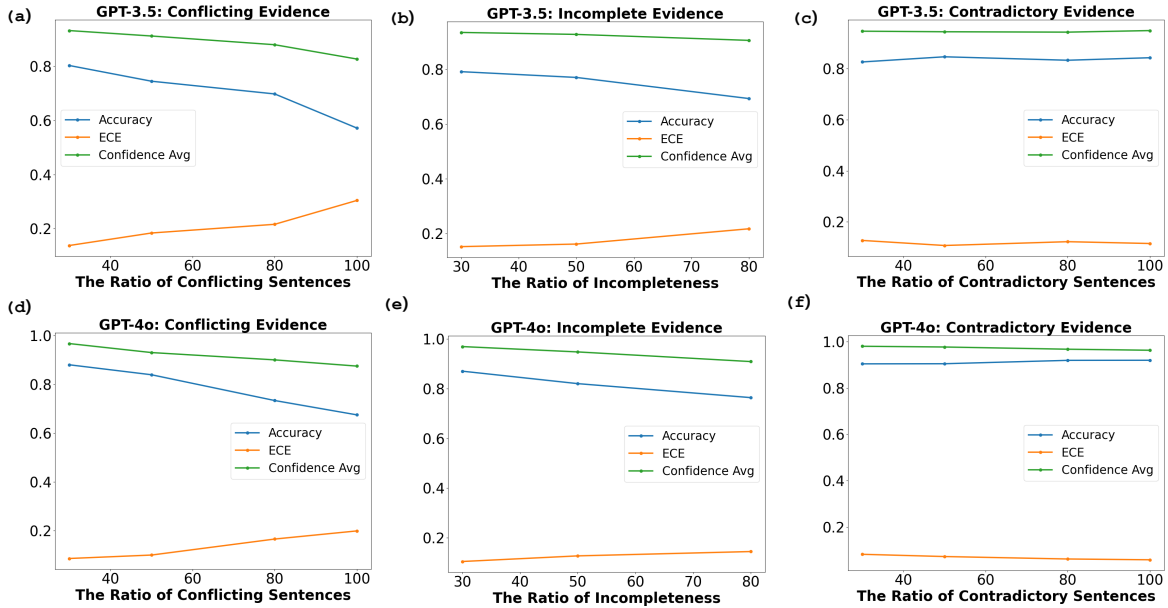


Figure 3: The results for the degree of variations in evidence for the SciQ dataset with verbalized method. We modified the number of conflicting sentences in conflicting evidence, sentences in incomplete evidence, and contradictory sentences in contradictory evidence (See Appendix F for entire results).

decreases, performance decreases. On the other hand, Figure 3 (c) and (f) indicate that if the golden evidence sentences is sufficiently given, increasing the number of contradictory sentences does not affect the confidence and performance even if both of the numbers of contradictory sentences and golden evidence sentences are same. This shows that the language model focuses more on the given golden evidence in the context than inaccurate evidence, and this is why it maintains confidence and calibration despite incomplete and contradictory evidence.

**Why do LLMs get confused by irrelevant context?** Two interpretations are possible for the irrelevant case

1. The language model does not recognize irrelevant evidence as irrelevant when it is in the same field but differs in content.
2. The language model considers irrelevant evidence as a kind of noise, which distracts the model and causes confusion.

To verify (1), instead of extracting irrelevant evidence from the same dataset, we used contexts from different datasets, for SciQ and TriviaQA dataset, we used evidence of GSM8K, and for GSM8K,

using TriviaQA. As we can see in Figure 6 in Appendix G, even when using a new irrelevant sentence, it did not completely match the completely irrelevant assumption. However, surprisingly, when using evidence from a completely different field, we found that the confidence, accuracy, and ECE metrics approached closer to the baseline no evidence case ( $P(H)$ ) than when we used evidence where the content was different but the field was the same. This implies that as the irrelevance increases, the LLMs become less distracted by the context. Therefore, we interpreted that there is a possibility that the LLMs satisfy the irrelevant assumption of Bayesian epistemology.

## 6 Why Do LLMs Struggle to Follow Bayesian Assumptions?

We speculate that LLMs may not fully adhere to Bayesian assumptions due to limitations in both their training data and methods. Pretraining datasets, primarily sourced from web crawls with filtering and books (Gao et al., 2020), expose models mostly to well-justified information and correct explanations. This allows them to excel in deriving correct answers but limits their ability to handle coincidental evidence or conflicting beliefs, which are underrepresented in the data.

In addition to data limitations, the training meth-



ods used for LLMs diverge significantly from human language acquisition. Humans encounter unreliable and conflicting situations through interaction and experience with the real world. In contrast, LLMs are trained in a largely supervised manner, resulting in a lack of semantic understanding (Bender and Koller, 2020; Bisk et al., 2020; Soni et al., 2024). As a result, these limitations may explain why LLMs fail to fully align with Bayesian assumptions when faced with conflicting or unexpected scenarios, as they lack the experiential grounding to properly handle such evidence.

## 7 Conclusion

In this paper, we explored how changes in the informativeness and reliability of evidence affect the confidence and response of language models. Specifically, we examined how well language models stick to the probability and calibration norms outlined in Bayesian epistemology. We demonstrated that language models partially align with Bayesian epistemology, following confirmation assumptions but failing to adhere to disconfirmation and irrelevance assumptions. It can be interpreted that language models do not possess a justified belief in the view of Bayesian epistemology. Additionally, we found that LLMs show a bias toward golden evidence and modify their confidence and response relative to the degree of irrelevance, which helps clarify their deviation from Bayesian assumptions. These findings provide philosophical insight into the nature of "belief" in LLMs, highlighting their biases and limitations.

## 8 Limitations

In this paper, we did not theoretically investigate the causes behind the observed phenomena, such as the training algorithm and model architecture, leaving such investigations and their potential implications for future research. One limitation of our dataset is the varying nature of conflicting evidence, some evidence is entirely incompatible, while other types obstruct correct answers. Hence, as we varied the extent of irrelevance in the ablation study, a finer classification of conflicting evidence could benefit future research. Additionally, Bayesian epistemology is not the only theory for defining knowledge. Therefore, our results do not imply a definitive conclusion about whether LLMs possess beliefs or knowledge. We did not conduct a human evaluation as a baseline in this study as our

focus was on aligning language models with ideal Bayesian epistemology. Future research could incorporate human evaluation to further assess how models' belief updating and confidence calibration compare to human cognitive processes. Lastly, further research using more complex and practical datasets, as well as developing new algorithms to address the issues identified in this study, will be valuable in advancing robust AI systems.

## 9 Ethics Statement

In the preparation of this paper, we utilized GPT-4o, for grammatical corrections and coding assistance. This technology served as an auxiliary resource to enhance the clarity and accuracy of our work, without directly influencing the research outcomes or decision-making processes involved. We acknowledge the support provided by OpenAI's GPT-4o in refining the presentation of our findings, ensuring that our use of this tool adheres to ethical guidelines and does not compromise the integrity of our research.

## 10 Acknowledgement

This work was supported by Institute for Information & communications Technology Promotion (IITP) grant funded by the Korea government (MSIT) (RS-2019-II190075) Artificial Intelligence Graduate School Program (KAIST) and Artificial intelligence industrial convergence cluster development project funded by the Ministry of Science and ICT (MSIT, Korea) & Gwangju Metropolitan City.

## References

- Robert Audi. 1997. *Epistemology: A Contemporary Introduction to the Theory of Knowledge*. Routledge, New York.
- Emily M. Bender and Alexander Koller. 2020. [Climbing towards NLU: On meaning, form, and understanding in the age of data](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5185–5198, Online. Association for Computational Linguistics.
- Yonatan Bisk, Ari Holtzman, Jesse Thomason, Jacob Andreas, Yoshua Bengio, Joyce Chai, Mirella Lapata, Angeliki Lazaridou, Jonathan May, Aleksandr Nisnevich, Nicolas Pinto, and Joseph Turian. 2020. [Experience grounds language](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8718–8735, Online. Association for Computational Linguistics.

- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#).
- Oana-Maria Camburu, Tim Rocktäschel, Thomas Lukasiewicz, and Phil Blunsom. 2018. [e-snli: Natural language inference with natural language explanations](#).
- Jiuhai Chen, Lichang Chen, Chen Zhu, and Tianyi Zhou. 2023. [How many demonstrations do you need for in-context learning?](#)
- Yew Ken Chia, Guizhen Chen, Luu Anh Tuan, Soujanya Poria, and Lidong Bing. 2023. [Contrastive chain-of-thought prompting](#).
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#).
- Yu Feng, Ben Zhou, Haoyu Wang, Helen Jin, and Dan Roth. 2023. [Generic temporal reasoning with differential analysis and explanation](#).
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2020. [The pile: An 800gb dataset of diverse text for language modeling](#).
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy Lillicrap, Angeliki Lazaridou, Orhan Firat, James Molloy, Michael Isard, Paul R. Barham, Tom Hennigan, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, Ryan Doherty, Eli Collins, Clemens Meyer, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, Jack Krawczyk, Cosmo Du, Ed Chi, Heng-Tze Cheng, Eric Ni, Purvi Shah, Patrick Kane, Betty Chan, Manaal Faruqui, Aliaksei Severyn, Hanzhao Lin, YaGuang Li, Yong Cheng, Abe Ittycheriah, Mahdis Mahdieh, Mia Chen, Pei Sun, Dustin Tran, Sumit Bagri, Balaji Lakshminarayanan, Jeremiah Liu, Andras Orban, Fabian Gura, Hao Zhou, Xinying Song, Aurelien Boffy, Harish Ganapathy, Steven Zheng, HyunJeong Choe, goston Weisz, Tao Zhu, Yifeng Lu, Siddharth Gopal, Jarrod Kahn, Maciej Kula, Jeff Pitman, Rushin Shah, Emanuel Taropa, Majd Al Meray, Martin Baeuml, Zhifeng Chen, Laurent El Shafey, Yujing Zhang, Olcan Sercinoglu, George Tucker, Enrique Piqueras, Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo Danihelka, Becca Roelofs, Anas White, Anders Andreassen, Tamara von Glehn, Lakshman Yagati, Mehran Kazemi, Lucas Gonzalez, Misha Khalman, Jakub Sygnowski, Alexandre Frechette, Charlotte Smith, Laura Culp, Lev Proleev, Yi Luan, Xi Chen, James Lottes, Nathan Schucher, Federico Lebron, Alban Rustemi, Natalie Clay, Phil Crone, Tomas Kocisky, Jeffrey Zhao, Bartek Perz, Dian Yu, Heidi Howard, Adam Bloniarz, Jack W. Rae, Han Lu, Laurent Sifre, Marcello Maggioni, Fred Alcober, Dan Garrette, Megan Barnes, Shantanu Thakoor, Jacob Austin, Gabriel Barth-Maron, William Wong, Rishabh Joshi, Rahma Chaabouni, Deeni Fatiha, Arun Ahuja, Gaurav Singh Tomar, Evan Senter, Martin Chadwick, Ilya Kornakov, Nithya Attaluri, Iaki Iturrate, Ruibo Liu, Yunxuan Li, Sarah Cogan, Jeremy Chen, Chao Jia, Chenjie Gu, Qiao Zhang, Jordan Grimstad, Ale Jakse Hartman, Xavier Garcia, Thanumalayan Sankaranarayanan Pillai, Jacob Devlin, Michael Laskin, Diego de Las Casas, Dasha Valter, Connie Tao, Lorenzo Blanco, Adri Puigdomnech Badia, David Reitter, Mianna Chen, Jenny Brennan, Clara Rivera, Sergey Brin, Shariq Iqbal, Gabriela Surita, Jane Labanowski, Abhi Rao, Stephanie Winkler, Emilio Parisotto, Yiming Gu, Kate Olszewska, Ravi Addanki, Antoine Miech, Annie Louis, Denis Teplyashin, Geoff Brown, Elliot Catt, Jan Balaguer, Jackie Xiang, Pidong Wang, Zoe Ashwood, Anton Briukhov, Albert Webson, Sanjay Ganapathy, Smit Sanghavi, Ajay Kannan, Ming-Wei Chang, Axel Stjerngren, Josip Djolonga, Yuting Sun, Ankur Bapna, Matthew Aitchison, Pedram Pejman, Henryk Michalewski, Tianhe Yu, Cindy Wang, Juliette Love, Junwhan Ahn, Dawn Bloxwich, Kehang Han, Peter Humphreys, Thibault Sellam, James Bradbury, Varun Godbole, Sina Samangooei, Bogdan Damoc, Alex Kaskasoli, Sbastien M. R. Arnold, Vijay Vasudevan, Shubham Agrawal, Jason Riesa, Dmitry Lepikhin, Richard Tanburn, Srivatsan Srinivasan, Hyeontaek Lim, Sarah Hodkinson, Pranav Shyam, Johan Ferret, Steven Hand, Ankush Garg, Tom Le Paine, Jian Li, Yujia Li, Minh Giang, Alexander Neitz, Zaheer Abbas, Sarah York, Machel Reid, Elizabeth Cole, Aakanksha Chowdhery, Dipanjan Das, Dominika Rogoziska, Vitaliy Nikolaev, Pablo Sprechmann, Zachary Nado, Lukas Zilka, Flavien Prost, Luheng He, Marianne Monteiro, Gaurav Mishra, Chris Welty, Josh Newlan, Dawei Jia, Miltiadis Allamanis, Clara Huiyi Hu, Raoul de Liedekerke, Justin Gilmer, Carl Saroufim, Shruti Rijhwani, Shaobo Hou, Disha Shrivastava, Anirudh Baddepudi, Alex Goldin, Adnan Ozturel, Albin Cassirer, Yunhan Xu, Daniel Sohn, Devendra Sachan, Reinald Kim Amplayo, Craig Swanson, Dessie Petrova, Shashi Narayan, Arthur Guez, Siddhartha Brahma, Jessica Landon, Miteyan Patel, Ruizhe Zhao, Kevin Vilella, Luyu Wang, Wenhao Jia, Matthew Rahtz, Mai Gimnez, Legg Yeung, James Keeling, Petko Georgiev, Diana Mincu, Boxi Wu, Salem Haykal, Rachel Saputro, Kiran Vodra-

halli, James Qin, Zeynep Cankara, Abhanshu Sharma, Nick Fernando, Will Hawkins, Behnam Neyshabur, Solomon Kim, Adrian Hutter, Priyanka Agrawal, Alex Castro-Ros, George van den Driessche, Tao Wang, Fan Yang, Shuo yiin Chang, Paul Komarek, Ross McIlroy, Mario Lučić, Guodong Zhang, Wael Farhan, Michael Sharman, Paul Natsev, Paul Michel, Yamini Bansal, Siyuan Qiao, Kris Cao, Siamak Shakeri, Christina Butterfield, Justin Chung, Paul Kishan Rubenstein, Shivani Agrawal, Arthur Mensch, Kedar Soparkar, Karel Lenc, Timothy Chung, Aedan Pope, Loren Maggiore, Jackie Kay, Priya Jhakra, Shibo Wang, Joshua Maynez, Mary Phuong, Taylor Tobin, Andrea Tacchetti, Maja Trebacz, Kevin Robinson, Yash Katariya, Sebastian Riedel, Paige Bailey, Kefan Xiao, Nimesh Ghelani, Lora Aroyo, Ambrose Slone, Neil Houlsby, Xuehan Xiong, Zhen Yang, Elena Gribovskaya, Jonas Adler, Mateo Wirth, Lisa Lee, Music Li, Thais Kagohara, Jay Pavagadhi, Sophie Bridgers, Anna Bortsova, Sanjay Ghemawat, Zafarali Ahmed, Tianqi Liu, Richard Powell, Vijay Bolina, Mariko Iinuma, Polina Zablotskaia, James Besley, Da-Woon Chung, Timothy Dozat, Ramona Comanescu, Xiance Si, Jeremy Greer, Guolong Su, Martin Polacek, Raphaël Lopez Kaufman, Simon Tokumine, Hexiang Hu, Elena Buchatskaya, Yingjie Miao, Mohamed Elhawaty, Aditya Siddhant, Nenad Tomasev, Jinwei Xing, Christina Greer, Helen Miller, Shereen Ashraf, Aurko Roy, Zizhao Zhang, Ada Ma, Angelos Filos, Milos Besta, Rory Blevins, Ted Klimenko, Chih-Kuan Yeh, Soravit Changpinyo, Jiaqi Mu, Oscar Chang, Mantas Pajarskas, Carrie Muir, Vered Cohen, Charline Le Lan, Krishna Haridasan, Amit Marathe, Steven Hansen, Sholto Douglas, Rajkumar Samuel, Mingqiu Wang, Sophia Austin, Chang Lan, Jiepu Jiang, Justin Chiu, Jaime Alonso Lorenzo, Lars Lowe Sjösund, Sébastien Cevey, Zach Gleicher, Thi Avrahami, Anudhyan Boral, Hansa Srinivasan, Vittorio Selo, Rhys May, Konstantinos Aisopos, Léonard Hussenot, Livio Baldini Soares, Kate Baumli, Michael B. Chang, Adrià Recasens, Ben Caine, Alexander Pritzel, Filip Pavetic, Fabio Pardo, Anita Gergely, Justin Frye, Vinay Ramasesh, Dan Horgan, Kartikeya Badola, Nora Kassner, Subhrajit Roy, Ethan Dyer, Víctor Campos Campos, Alex Tomala, Yunhao Tang, Dalia El Badawy, Elspeth White, Basil Mustafa, Oran Lang, Abhishek Jindal, Sharad Vikram, Zhitao Gong, Sergi Caelles, Ross Hemsley, Gregory Thornton, Fangxiaoyu Feng, Wojciech Stokowiec, Ce Zheng, Phoebe Thacker, Çağlar Ünlü, Zhishuai Zhang, Mohammad Saleh, James Svensson, Max Bileschi, Piyush Patil, Ankesh Anand, Roman Ring, Katerina Tsihlias, Arpi Vezzer, Marco Selvi, Toby Shevlane, Mikel Rodriguez, Tom Kwiatkowski, Samira Daruki, Keran Rong, Allan Dafoe, Nicholas FitzGerald, Keren Gu-Lemberg, Mina Khan, Lisa Anne Hendricks, Marie Pellat, Vladimir Feinberg, James Cobon-Kerr, Tara Sainath, Maribeth Rauh, Sayed Hadi Hashemi, Richard Ives, Yana Hasson, Eric Noland, Yuan Cao, Nathan Byrd, Le Hou, Qingze Wang, Thibault Sottiaux, Michela Paganini, Jean-Baptiste Lespiau, Alexandre Moufaret, Samer Hassan, Kaushik Shivakumar, Joost van

Amersfoort, Amol Mandhane, Pratik Joshi, Anirudh Goyal, Matthew Tung, Andrew Brock, Hannah Sheahan, Vedant Misra, Cheng Li, Nemanja Rakićević, Mostafa Dehghani, Fangyu Liu, Sid Mittal, Junhyuk Oh, Seb Noury, Eren Sezener, Fantine Huot, Matthew Lamm, Nicola De Cao, Charlie Chen, Sidharth Mudgal, Romina Stella, Kevin Brooks, Gautam Vasudevan, Chenxi Liu, Mainak Chain, Nivedita Melinkeri, Aaron Cohen, Venus Wang, Kristie Seymore, Sergey Zubkov, Rahul Goel, Summer Yue, Sai Krishnakumaran, Brian Albert, Nate Hurley, Motoki Sano, Anhad Mohananey, Jonah Joughin, Egor Filonov, Tomasz Kępa, Yomna Eldawy, Jiawern Lim, Rahul Rishi, Shirin Badiezadegan, Taylor Bos, Jerry Chang, Sanil Jain, Sri Gayatri Sundara Padmanabhan, Subha Puttagunta, Kalpesh Krishna, Leslie Baker, Norbert Kalb, Vamsi Bedapudi, Adam Kurzrok, Shuntong Lei, Anthony Yu, Oren Litvin, Xiang Zhou, Zhichun Wu, Sam Sobell, Andrea Siciliano, Alan Papir, Robby Neale, Jonas Bragagnolo, Tej Toor, Tina Chen, Valentin Anklin, Feiran Wang, Richie Feng, Milad Gholami, Kevin Ling, Lijuan Liu, Jules Walter, Hamid Moghaddam, Arun Kishore, Jakub Adamek, Tyler Mercado, Jonathan Mallinson, Siddhinita Wandekar, Stephen Cagle, Eran Ofek, Guillermo Garrido, Clemens Lombriser, Maksim Mukha, Botu Sun, Hafeezul Rahman Mohammad, Josip Matak, Yadi Qian, Vikas Peswani, Pawel Janus, Quan Yuan, Leif Schelin, Oana David, Ankur Garg, Yifan He, Oleksii Duzhyi, Anton Ålgmyr, Timothée Lottaz, Qi Li, Vikas Yadav, Luyao Xu, Alex Chinién, Rakesh Shivanna, Aleksandr Chuklin, Josie Li, Carrie Spadine, Travis Wolfe, Kareem Mohamed, Subhabrata Das, Zihang Dai, Kyle He, Daniel von Dincklage, Shyam Upadhyay, Akanksha Maurya, Luyan Chi, Sebastian Krause, Khalid Salama, Pam G Rabinovitch, Pavan Kumar Reddy M, Aarush Selvan, Mikhail Dektiarev, Golnaz Ghiasi, Erdem GUVEN, Himanshu Gupta, Boyi Liu, Deepak Sharma, Idan Heimlich Shtacher, Shachi Paul, Oscar Akerlund, François-Xavier Aubet, Terry Huang, Chen Zhu, Eric Zhu, Elico Teixeira, Matthew Fritze, Francesco Bertolini, Liana-Eleonora Marinescu, Martin Bølle, Dominik Paulus, Khyatti Gupta, Tejasi Latkar, Max Chang, Jason Sanders, Roopa Wilson, Xuewei Wu, Yi-Xuan Tan, Lam Nguyen Thiet, Tulsee Doshi, Sid Lall, Swaroop Mishra, Wanming Chen, Thang Luong, Seth Benjamin, Jasmine Lee, Ewa Andrejczuk, Dominik Rabiej, Vipul Ranjan, Krzysztof Styrz, Pengcheng Yin, Jon Simon, Malcolm Rose Harriott, Mudit Bansal, Alexei Robsky, Geoff Bacon, David Greene, Daniil Mirylenka, Chen Zhou, Obaid Sarvana, Abhimanyu Goyal, Samuel Andermatt, Patrick Siegler, Ben Horn, Assaf Israel, Francesco Pongetti, Chih-Wei "Louis" Chen, Marco Selvatici, Pedro Silva, Kathie Wang, Jackson Tolins, Kelvin Guu, Roey Yogev, Xiaochen Cai, Alessandro Agostini, Maulik Shah, Hung Nguyen, Noah Ó Donnaile, Sébastien Pereira, Linda Friso, Adam Stambler, Adam Kurzrok, Chenkai Kuang, Yan Romanikhin, Mark Geller, ZJ Yan, Kane Jang, Cheng-Chun Lee, Wojciech Fica, Eric Malmi, Qijun Tan, Dan Banica, Daniel Balle, Ryan Pham,

Yanping Huang, Diana Avram, Hongzhi Shi, Jasjot Singh, Chris Hidey, Niharika Ahuja, Pranab Saxena, Dan Dooley, Srividya Pranavi Potharaju, Eileen O'Neill, Anand Gokulchandran, Ryan Foley, Kai Zhao, Mike Dusenberry, Yuan Liu, Pulkit Mehta, Ragha Kotikalapudi, Chalence Safranek-Shrader, Andrew Goodman, Joshua Kessinger, Eran Globen, Prateek Kolhar, Chris Gorgolewski, Ali Ibrahim, Yang Song, Ali Eichenbaum, Thomas Brovelli, Sahitya Potluri, Preethi Lahoti, Cip Baetu, Ali Ghorbani, Charles Chen, Andy Crawford, Shalini Pal, Mukund Sridhar, Petru Gurita, Asier Mujika, Igor Petrovski, Pierre-Louis Cedoz, Chenmei Li, Shiyuan Chen, Niccolò Dal Santo, Siddharth Goyal, Jitesh Punjabi, Karthik Kappaganthu, Chester Kwak, Pallavi LV, Sarmishta Velury, Himadri Choudhury, Jamie Hall, Premal Shah, Ricardo Figueira, Matt Thomas, Minjie Lu, Ting Zhou, Chintu Kumar, Thomas Jurdi, Sharat Chikkerur, Yenai Ma, Adams Yu, Soo Kwak, Victor Åhdel, Sujeevan Rajayogam, Travis Choma, Fei Liu, Aditya Barua, Colin Ji, Ji Ho Park, Vincent Hellendoorn, Alex Bailey, Taylan Bilal, Huanjie Zhou, Mehrdad Khatir, Charles Sutton, Wojciech Rzadkowski, Fiona Macintosh, Konstantin Shagin, Paul Medina, Chen Liang, Jinjing Zhou, Pararth Shah, Yingying Bi, Attila Dankovics, Shipra Banga, Sabine Lehmann, Marissa Bredesen, Zifan Lin, John Eric Hoffmann, Jonathan Lai, Raynald Chung, Kai Yang, Nihal Balani, Arthur Bražinskis, Andrei Sozanschi, Matthew Hayes, Héctor Fernández Alcalde, Peter Makarov, Will Chen, Antonio Stella, Liselotte Snijders, Michael Mandl, Ante Kärman, Paweł Nowak, Xinyi Wu, Alex Dyck, Krishnan Vaidyanathan, Raghavender R, Jessica Mallet, Mitch Rudominer, Eric Johnston, Sushil Mittal, Akhil Udathu, Janara Christensen, Vishal Verma, Zach Irving, Andreas Santucci, Gamaleldin Elsayed, Elnaz Davoodi, Marin Georgiev, Ian Tenney, Nan Hua, Geoffrey Cideron, Edouard Leurent, Mahmoud Alnahlawi, Ionut Georgescu, Nan Wei, Ivy Zheng, Dylan Scandinaro, Heinrich Jiang, Jasper Snock, Mukund Sundararajan, Xuezhi Wang, Zack Ontiveros, Itay Karo, Jeremy Cole, Vinu Rajashekhar, Lara Tumeh, Eyal Ben-David, Rishub Jain, Jonathan Uesato, Romina Datta, Oskar Bunyan, Shimu Wu, John Zhang, Piotr Stanczyk, Ye Zhang, David Steiner, Subhajit Naskar, Michael Azzam, Matthew Johnson, Adam Paszke, Chung-Cheng Chiu, Jaume Sanchez Elias, Afroz Mohiuddin, Faizan Muhammad, Jin Miao, Andrew Lee, Nino Vieillard, Jane Park, Jigeng Zhang, Jeff Stanway, Drew Garmon, Abhijit Karmarkar, Zhe Dong, Jong Lee, Aviral Kumar, Luwei Zhou, Jonathan Evens, William Isaac, Geoffrey Irving, Edward Loper, Michael Fink, Isha Arkatkar, Nanxin Chen, Izhak Shafran, Ivan Petyrchenko, Zhe Chen, Johnson Jia, Anselm Levskaya, Zhenkai Zhu, Peter Grabowski, Yu Mao, Alberto Magni, Kaisheng Yao, Javier Snaider, Norman Casagrande, Evan Palmer, Paul Suganthan, Alfonso Castaño, Irene Giannoumis, Wooyeol Kim, Mikołaj Rybiński, Ashwin Sreevatsa, Jennifer Prendki, David Soergel, Adrian Goedeckemeyer, Willi Gierke, Mohsen Jafari, Meenu Gaba, Jeremy Wiesner, Diana Gage Wright,

Yawen Wei, Harsha Vashisht, Yana Kulizhskaya, Jay Hoover, Maigo Le, Lu Li, Chimezie Iwuanyanwu, Lu Liu, Kevin Ramirez, Andrey Khorlin, Albert Cui, Tian LIN, Marcus Wu, Ricardo Aguilar, Keith Pallo, Abhishek Chakladar, Ginger Perng, Elena Allica Abellan, Mingyang Zhang, Ishita Dasgupta, Nate Kushman, Ivo Penchev, Alena Repina, Xihui Wu, Tom van der Weide, Priya Ponnappalli, Caroline Kaplan, Jiri Simsa, Shuangfeng Li, Olivier Dousse, Fan Yang, Jeff Piper, Nathan Ie, Rama Pasumarthi, Nathan Lintz, Anitha Vijayakumar, Daniel Andor, Pedro Valenzuela, Minnie Lui, Cosmin Padurararu, Daiyi Peng, Katherine Lee, Shuyuan Zhang, Somer Greene, Duc Dung Nguyen, Paula Kurylowicz, Cassidy Hardin, Lucas Dixon, Lili Janzer, Kiam Choo, Ziqiang Feng, Biao Zhang, Achintya Singhal, Dayou Du, Dan McKinnon, Natasha Antropova, Tolga Bolukbasi, Orgad Keller, David Reid, Daniel Finchelstein, Maria Abi Raad, Remi Crocker, Peter Hawkins, Robert Dadashi, Colin Gaffney, Ken Franko, Anna Bulanova, Rémi Leblond, Shirley Chung, Harry Askham, Luis C. Cobo, Kelvin Xu, Felix Fischer, Jun Xu, Christina Sorokin, Chris Alberti, Chu-Cheng Lin, Colin Evans, Alek Dimitriev, Hannah Forbes, Dylan Banarse, Zora Tung, Mark Omernick, Colton Bishop, Rachel Sterneck, Rohan Jain, Jiawei Xia, Ehsan Amid, Francesco Piccinno, Xingyu Wang, Praseem Banzal, Daniel J. Mankowitz, Alex Polozov, Victoria Krakovna, Sasha Brown, MohammadHossein Bateni, Dennis Duan, Vlad Firoiu, Meghana Thotakuri, Tom Natan, Matthieu Geist, Ser tan Girgin, Hui Li, Jiayu Ye, Ofir Roval, Reiko Tojo, Michael Kwong, James Lee-Thorp, Christopher Yew, Danila Sinopalnikov, Sabela Ramos, John Mellor, Abhishek Sharma, Kathy Wu, David Miller, Nicolas Sonnerat, Denis Vnukov, Rory Greig, Jennifer Beattie, Emily Caveness, Libin Bai, Julian Eisenschlos, Alex Korchemniy, Tomy Tsai, Mimi Jasarevic, Weize Kong, Phuong Dao, Zeyu Zheng, Frederick Liu, Fan Yang, Rui Zhu, Tian Huey Teh, Jason Sanmiya, Evgeny Gladchenko, Nejc Trdin, Daniel Toyama, Evan Rosen, Sasan Tavakkol, Linting Xue, Chen Elkind, Oliver Woodman, John Carpenter, George Papamakarios, Rupert Kemp, Sushant Kafle, Tanya Grunina, Rishika Sinha, Alice Talbert, Diane Wu, Denese Owusu-Afriyie, Cosmo Du, Chloe Thornton, Jordi Pont-Tuset, Pradyumna Narayana, Jing Li, Saaber Fatehi, John Wieting, Omar Ajmeri, Benigno Urias, Yeongil Ko, Laura Knight, Amélie Héliou, Ning Niu, Shane Gu, Chenxi Pang, Yeqing Li, Nir Levine, Ariel Stolovich, Rebecca Santamaria-Fernandez, Sonam Goenka, Wenny Yustalim, Robin Strudel, Ali Elqursh, Charlie Deck, Hyo Lee, Zonglin Li, Kyle Levin, Raphael Hoffmann, Dan Holtmann-Rice, Olivier Bachem, Sho Arora, Christy Koh, Soheil Hassas Yeganeh, Siim Pöder, Mukarram Tariq, Yanhua Sun, Lucian Ionita, Mojtaba Seyedhosseini, Pouya Tafti, Zhiyu Liu, Anmol Gulati, Jasmine Liu, Xinyu Ye, Bart Chrzaszcz, Lily Wang, Nikhil Sethi, Tianrun Li, Ben Brown, Shreya Singh, Wei Fan, Aaron Parisi, Joe Stanton, Vinod Koverkathu, Christopher A. Choquette-Choo, Yunjie Li, TJ Lu, Abe Ittycheriah, Prakash

- Shroff, Mani Varadarajan, Sanaz Bahargam, Rob Willoughby, David Gaddy, Guillaume Desjardins, Marco Cornero, Brona Robenek, Bhavishya Mittal, Ben Albrecht, Ashish Shenoy, Fedor Moiseev, Henrik Jacobsson, Alireza Ghaffarkhah, Morgane Rivi re, Alanna Walton, Cl ment Crepy, Alicia Parrish, Zongwei Zhou, Clement Farabet, Carey Radebaugh, Praveen Srinivasan, Claudia van der Salm, Andreas Fidjeland, Salvatore Scellato, Eri Latorre-Chimoto, Hanna Klimczak-Plucińska, David Bridson, Dario de Cesare, Tom Hudson, Piermaria Mendolichio, Lexi Walker, Alex Morris, Matthew Mauger, Alexey Guseynov, Alison Reid, Seth Odoom, Lucia Loher, Victor Cotruta, Madhavi Yenugula, Dominik Grewe, Anastasia Petrushkina, Tom Duerig, Antonio Sanchez, Steve Yadlowsky, Amy Shen, Amir Globerson, Lynette Webb, Sahil Dua, Dong Li, Surya Bhupatiraju, Dan Hurt, Haroon Qureshi, Ananth Agarwal, Tomer Shani, Matan Eyal, Anuj Khare, Shreyas Rammohan Belle, Lei Wang, Chetan Tekur, Mihir Sanjay Kale, Jinliang Wei, Ruoxin Sang, Brennan Saeta, Tyler Liechty, Yi Sun, Yao Zhao, Stephan Lee, Pandu Nayak, Doug Fritz, Manish Reddy Vuyyuru, John Aslanides, Nidhi Vyas, Martin Wicke, Xiao Ma, Evgenii Eltyshev, Nina Martin, Hardie Cate, James Manyika, Keyvan Amiri, Yelin Kim, Xi Xiong, Kai Kang, Florian Luisier, Nilesh Tripuraneni, David Madras, Mandy Guo, Austin Waters, Oliver Wang, Joshua Ainslie, Jason Baldridge, Han Zhang, Garima Pruthi, Jakob Bauer, Feng Yang, Riham Mansour, Jason Gelman, Yang Xu, George Polovets, Ji Liu, Honglong Cai, Warren Chen, XiangHai Sheng, Emily Xue, Sherjil Ozair, Christof Angermueller, Xiaowei Li, Anoop Sinha, Weiren Wang, Julia Wiesinger, Emmanouil Koukoumidis, Yuan Tian, Anand Iyer, Madhu Gurusurthy, Mark Goldenson, Parashar Shah, MK Blake, Hongkun Yu, Anthony Urbanowicz, Jennimaria Palomaki, Chrisantha Fernando, Ken Durden, Harsh Mehta, Nikola Momchev, Elahe Rahimtoroghi, Maria Georgaki, Amit Raul, Sebastian Ruder, Morgan Redshaw, Jinyuk Lee, Denny Zhou, Komal Jalan, Dinghua Li, Blake Hechtman, Parker Schuh, Milad Nasr, Kieran Milan, Vladimir Mikulik, Juliana Franco, Tim Green, Nam Nguyen, Joe Kelley, Aroma Mahendru, Andrea Hu, Joshua Howland, Ben Vargas, Jeffrey Hui, Kshitij Bansal, Vikram Rao, Rakesh Ghiya, Emma Wang, Ke Ye, Jean Michel Sarr, Melanie Moranski Preston, Madeleine Elish, Steve Li, Aakash Kaku, Jigar Gupta, Ice Pasupat, Da-Cheng Juan, Milan Someswar, Tejvi M., Xinyun Chen, Aida Amini, Alex Fabrikant, Eric Chu, Xuanyi Dong, Amruta Muthal, Senaka Buthpitiya, Sarthak Jauhari, Nan Hua, Urvashi Khandelwal, Ayal Hitron, Jie Ren, Larissa Rinaldi, Sharar Drath, Avigail Dabush, Nan-Jiang Jiang, Harshal Godhia, Uli Sachs, Anthony Chen, Yicheng Fan, Hagai Taitelbaum, Hila Noga, Zhuyun Dai, James Wang, Chen Liang, Jenny Hamer, Chun-Sung Ferng, Chenel Elkind, Aviell Atias, Paulina Lee, Vít Listfk, Mathias Carlen, Jan van de Kerkhof, Marcin Pikus, Krunoslav Zaher, Paul M ller, Sasha Zykova, Richard Stefanec, Vitaly Gatsko, Christoph Hirschall, Ashwin Sethi, Xingyu Federico Xu, Chetan Ahuja, Beth Tsai, Anca Stefanoiu, Bo Feng, Keshav Dhandhanian, Manish Katyal, Akshay Gupta, Atharva Parulekar, Divya Pitta, Jing Zhao, Vivaan Bhatia, Yashodha Bhavnani, Omar Alhadlaq, Xiaolin Li, Peter Danenberg, Dennis Tu, Alex Pine, Vera Filippova, Abhipso Ghosh, Ben Limonchik, Bhargava Urala, Chaitanya Krishna Lanka, Derik Clive, Yi Sun, Edward Li, Hao Wu, Kevin Hongtongsak, Ianna Li, Kalind Thakkar, Kuanysh Omarov, Kushal Majmundar, Michael Alverson, Michael Kucharski, Mohak Patel, Mudit Jain, Maksim Zabelin, Paolo Pelagatti, Rohan Kohli, Saurabh Kumar, Joseph Kim, Swetha Sankar, Vineet Shah, Lakshmi Ramachandruni, Xiangkai Zeng, Ben Bariach, Laura Weidinger, Amar Subramanya, Sissie Hsiao, Demis Hassabis, Koray Kavukcuoglu, Adam Sadovsky, Quoc Le, Trevor Strohman, Yonghui Wu, Slav Petrov, Jeffrey Dean, and Oriol Vinyals. 2024. [Gemini: A family of highly capable multimodal models](#).
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. 2017. [On calibration of modern neural networks](#).
- Alan H jek. 2003. [What conditional probability could not be](#). *Synthese*, 137(3):273–323.
- Alan H jek and Stephan Hartmann. 2010. Bayesian epistemology. In DancyJ, editor, *A Companion to Epistemology*. Blackwell.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#).
- Paul Horwich. 1982. *Probability and Evidence*. Cambridge University Press, Cambridge.
- Colin Howson. 2000. *Hume’s Problem: Induction and the Justification of Belief*. Oxford University Press, New York.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2023. [A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions](#). *ArXiv*, abs/2311.05232.
- Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. [TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.
- Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, Scott Johnston, Sheer El-Showk, Andy Jones, Nelson Elhage, Tristan Hume, Anna Chen, Yuntao Bai, Sam Bowman, Stanislav Fort, Deep Ganguli, Danny Hernandez, Josh Jacobson,

- Jackson Kernion, Shauna Kravec, Liane Lovitt, Kamal Ndousse, Catherine Olsson, Sam Ringer, Dario Amodei, Tom Brown, Jack Clark, Nicholas Joseph, Ben Mann, Sam McCandlish, Chris Olah, and Jared Kaplan. 2022. [Language models \(mostly\) know what they know.](#)
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. [Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation.](#)
- Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Denison, Danny Hernandez, Dustin Li, Esin Durmus, Evan Hubinger, Jackson Kernion, Kamilé Lukošiušė, Karina Nguyen, Newton Cheng, Nicholas Joseph, Nicholas Schiefer, Oliver Rausch, Robin Larson, Sam McCandlish, Sandipan Kundu, Saurav Kadavath, Shannon Yang, Thomas Henighan, Timothy Maxwell, Timothy Telleen-Lawton, Tristan Hume, Zac Hatfield-Dodds, Jared Kaplan, Jan Brauner, Samuel R. Bowman, and Ethan Perez. 2023. [Measuring faithfulness in chain-of-thought reasoning.](#)
- Noah Lee, Na Min An, and James Thorne. 2023. [Can large language models capture dissenting human voices?](#) In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 4569–4585, Singapore. Association for Computational Linguistics.
- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries.](#) In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. [Teaching models to express their uncertainty in words.](#)
- 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, :, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altmenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kopic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael

- Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2023. [Gpt-4 technical report](#).
- Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed Chi, Nathanael Schärli, and Denny Zhou. 2023. [Large language models can be easily distracted by irrelevant context](#).
- Nikita Soni, H. Andrew Schwartz, João Sedoc, and Niranjan Balasubramanian. 2024. [Large human language models: A need and the challenges](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 8631–8646, Mexico City, Mexico. Association for Computational Linguistics.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew Dai, Andrew La, Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabasum, Arul Menezes, Arun Kirubakaran, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut Erdem, Ayla Karakaş, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph, Bartłomiej Bojanowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno Stein, Berk Ekmekci, Bill Yuchen Lin, Blake Howald, Bryan Orinon, Cameron Diao, Cameron Dour, Catherine Stinson, Cedrick Argueta, César Ferri Ramírez, Chandan Singh, Charles Rathkopf, Chenlin Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Chris Waites, Christian Voigt, Christopher D. Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette, Dan Hendrycks, Dan Kilman, Dan Roth, Daniel Freeman, Daniel Khashabi, Daniel Levy, Daniel Moseguí González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra, Dilyar Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Ekaterina Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, Ellie Pavlick, Emanuele Rodola, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang, Ethan A. Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Engfu Manyasi, Evgenii Zheltonozhskii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francesca Happé, Francois Chollet, Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germán Kruszewski, Giambattista Parascandolo, Giorgio Mariani, Gloria Wang, Gonzalo Jaimovitch-López, Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh Hajishirzi, Harsh Mehta, Hayden Bogar, Henry Shevlin, Hinrich Schütze, Hiromu Yakura, Hongming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack Geissinger, Jackson Kernion, Jacob Hilton, Jaehoon Lee, Jaime Fernández Fisac, James B. Simon, James Koppel, James Zheng, James Zou, Jan Kocoń, Jana Thompson, Janelle Wingfield, Jared Kaplan, Jarema Radom, Jascha Sohl-Dickstein, Jason Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova, Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen Taal, Jesse Engel, Jesujoba Alabi, Jiacheng Xu, Jiaming Song, Jillian Tang, Joan Waweru, John Burden, John Miller, John U. Balis, Jonathan Batchelder, Jonathan Berant, Jörg Frohberg, Jos Rozen, Jose Hernandez-Orallo, Joseph Boudeman, Joseph Guerr, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chua, Kamil Kanclerz, Karen Livescu, Karl Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Katja Markert, Kaustubh D. Dhole, Kevin Gimpel, Kevin Omondi, Kory Mathewson, Kristen Chifullo, Ksenia Shkaruta, Kumar Shridhar, Kyle McDonell, Kyle Richardson, Laria Reynolds, Leo Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-Ochando, Louis-Philippe Morency, Luca Moschella, Lucas Lam, Lucy Noble, Ludwig Schmidt, Luheng He, Luis Oliveros Colón, Luke Metz, Lütfi Kerem Şenel, Maarten Bosma, Maarten Sap, Maartje ter Hoeve, Maheen Farooqi, Manaal Faruqui, Mantas Mazeika, Marco Baturan, Marco Marelli, Marco Maru, Maria Jose Ramírez Quintana, Marie Tolkiehn, Mario Giulianelli, Martha Lewis, Martin Potthast, Matthew L. Leavitt, Matthias Hagen, Mátyás Schubert, Medina Orduna Baitemirova, Melody Arnaud, Melvin McElrath, Michael A. Yee, Michael Cohen, Michael Gu, Michael Ivanitskiy, Michael Starritt, Michael Strube, Michał Śwędrowski, Michele Bevilacqua, Michihiro Yasunaga, Mihir Kale, Mike Cain, Mimeo Xu, Mirac Suzgun, Mitch Walker, Mo Tiwari, Mohit Bansal, Moin Amnaseri, Mor Geva, Mozhdeh Gheini, Mukund Varma T, Nanyun Peng, Nathan A. Chi, Nayeon Lee, Neta Gur-Ari Krakover, Nicholas Cameron, Nicholas Roberts, Nick Doiron, Nicole Martinez, Nikita Nangia, Niklas Deckers, Niklas Muennighoff, Nitish Shirish Keskar, Niveditha S. Iyer, Noah Constant, Noah Fiedel, Nuan Wen, Oliver Zhang, Omar Agha, Omar Elbaghdadi, Omer Levy, Owain Evans, Pablo Antonio Moreno Casares, Parth Doshi, Pascale Fung, Paul Pu Liang, Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao, Percy Liang, Peter Chang, Peter Eckersley, Phu Mon Htut, Pinyu Hwang, Piotr Miłkowski, Piyush Patil, Pouya Pezeshkpour, Priti Oli, Qiaozhu Mei, Qing Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta Rudolph, Raefer Gabriel, Rahel Habacker, Ramon Risco, Raphaël Millière, Rhythm Garg, Richard

- Barnes, Rif A. Saurous, Riku Arakawa, Robbe Raymaekers, Robert Frank, Rohan Sikand, Roman Novak, Roman Sitelew, Ronan LeBras, Rosanne Liu, Rowan Jacobs, Rui Zhang, Ruslan Salakhutdinov, Ryan Chi, Ryan Lee, Ryan Stovall, Ryan Teehan, Rylan Yang, Sahib Singh, Saif M. Moham-mad, Sajant Anand, Sam Dillavou, Sam Shleifer, Sam Wiseman, Samuel Gruetter, Samuel R. Bowman, Samuel S. Schoenholz, Sanghyun Han, San-jeev Kwatra, Sarah A. Rous, Sarik Ghazarian, Sayan Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schuster, Sepideh Sadeghi, Shadi Hamdan, Sharon Zhou, Shashank Srivastava, Sherry Shi, Shikhar Singh, Shima Asaadi, Shixi-ang Shane Gu, Shubh Pachchigar, Shubham Tosh-niwal, Shyam Upadhyay, Shyamolima, Debnath, Siamak Shakeri, Simon Thormeyer, Simone Melzi, Siva Reddy, Sneha Priscilla Makini, Soo-Hwan Lee, Spencer Torene, Sriharsha Hatwar, Stanislas De-haene, Stefanovic, Stefano Ermon, Stella Bider-man, Stephanie Lin, Stephen Prasad, Steven T. Pi-antadosi, Stuart M. Shieber, Summer Misherggi, Svet-lana Kiritchenko, Swaroop Mishra, Tal Linzen, Tal Schuster, Tao Li, Tao Yu, Tariq Ali, Tatsu Hashimoto, Te-Lin Wu, Théo Desbordes, Theodore Rothschild, Thomas Phan, Tianle Wang, Tiberius Nkinyili, Timo Schick, Timofei Kornev, Titus Tunduny, Tobias Ger-stenberg, Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay Ramasesh, Vinay Uday Prabhu, Vishakh Padmakumar, Vivek Srikumar, William Fedus, William Sau-nders, William Zhang, Wout Vossen, Xiang Ren, Xi-aoyu Tong, Xinran Zhao, Xinyi Wu, Xudong Shen, Yadollah Yaghoobzadeh, Yair Lakretz, Yangqiu Song, Yasaman Bahri, Yejin Choi, Yichi Yang, Yiding Hao, Yifu Chen, Yonatan Belinkov, Yu Hou, Yufang Hou, Yuntao Bai, Zachary Seid, Zhuoye Zhao, Zi-jian Wang, Zijie J. Wang, Zirui Wang, and Ziyi Wu. 2023. [Beyond the imitation game: Quantifying and extrapolating the capabilities of language models.](#)
- Lichao Sun, Yue Huang, Haoran Wang, Siyuan Wu, Qi-hui Zhang, Yuan Li, Chujie Gao, Yixin Huang, Wen-han Lyu, Yixuan Zhang, Xiner Li, Zhengliang Liu, Yixin Liu, Yijue Wang, Zhikun Zhang, Bertie Vidgen, Bhavya Kaillkhura, Caiming Xiong, Chaowei Xiao, Chunyuan Li, Eric Xing, Furong Huang, Hao Liu, Heng Ji, Hongyi Wang, Huan Zhang, Huaxiu Yao, Manolis Kellis, Marinka Zitnik, Meng Jiang, Mohit Bansal, James Zou, Jian Pei, Jian Liu, Jianfeng Gao, Jiawei Han, Jieyu Zhao, Jiliang Tang, Jindong Wang, Joaquin Vanschoren, John Mitchell, Kai Shu, Kaidi Xu, Kai-Wei Chang, Lifang He, Lifu Huang, Michael Backes, Neil Zhenqiang Gong, Philip S. Yu, Pin-Yu Chen, Quanquan Gu, Ran Xu, Rex Ying, Shuiwang Ji, Suman Jana, Tianlong Chen, Tianming Liu, Tianyi Zhou, William Wang, Xiang Li, Xiangliang Zhang, Xiao Wang, Xing Xie, Xun Chen, Xuyu Wang, Yan Liu, Yanfang Ye, Yinzhi Cao, Yong Chen, and Yue Zhao. 2024. [Trustllm: Trustworthiness in large lan-guage models.](#)
- William Talbott. 2006. Bayesian epistemology. In Ed-ward Zalta, editor, *Stanford Encyclopedia of Philoso-phy*.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. [FEVER: a large-scale dataset for fact extraction and VERification.](#) In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 809–819, New Orleans, Louisiana. Association for Computational Linguistics.
- Katherine Tian, Eric Mitchell, Huaxiu Yao, Christo-pher D. Manning, and Chelsea Finn. 2023a. [Fine-tuning language models for factuality.](#)
- Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher D. Manning. 2023b. [Just ask for cali-bration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback.](#)
- Miles Turpin, Julian Michael, Ethan Perez, and Samuel R. Bowman. 2023. [Language models don’t always say what they think: Unfaithful explanations in chain-of-thought prompting.](#)
- Olav Benjamin Vassend. 2023. [What hinge epistemol-ogy and bayesian epistemology can learn from each other.](#) *Asian Journal of Philosophy*, 2(2):1–21.
- Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen, You Wu, Luke Zettlemoyer, and Huan Sun. 2023a. [Towards understanding chain-of-thought prompting: An empirical study of what matters.](#)
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023b. [Self-consistency improves chain of thought reasoning in language models.](#)
- Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, An-drew M. Dai, and Quoc V. Le. 2022. [Finetuned language models are zero-shot learners.](#)
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. [Chain-of-thought prompting elic-its reasoning in large language models.](#)
- Johannes Welbl, Nelson F. Liu, and Matt Gardner. 2017. [Crowdsourcing multiple choice science questions.](#) In *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pages 94–106, Copenhagen, Den-mark. Association for Computational Linguistics.
- Jon Williamson. 2010. *In Defence of Objective Bayesianism*. Oxford University Press.
- Yijun Xiao and William Yang Wang. 2021. [On hal-lucination and predictive uncertainty in conditional language generation.](#) *ArXiv*, abs/2103.15025.



- Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. 2024. [Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms.](#)
- Pengfei Yu and Heng Ji. 2024. [Information association for language model updating by mitigating lm-logical discrepancy.](#)
- Tianhang Zhang, Lin Qiu, Qipeng Guo, Cheng Deng, Yue Zhang, Zheng Zhang, Cheng Zhou, Xinbing Wang, and Luoyi Fu. 2023. [Enhancing uncertainty-based hallucination detection with stronger focus.](#) *ArXiv*, abs/2311.13230.
- Yao Zhao, Rishabh Joshi, Tianqi Liu, Misha Khalman, Mohammad Saleh, and Peter J. Liu. 2023. [Slic-hf: Sequence likelihood calibration with human feedback.](#)
- Kaitlyn Zhou, Dan Jurafsky, and Tatsunori Hashimoto. 2023. [Navigating the grey area: How expressions of uncertainty and overconfidence affect language models.](#) In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5506–5524, Singapore. Association for Computational Linguistics.

## Appendix

### A Results of Confirmation task

	Dataset	Metric	No_EVI	EVI	Coincidence	Irrelevant	Conflict	Incomplete	Contradiction
GPT-3.5-turbo	SciQ	Confidence	0.671	0.781	0.785	0.594	0.638	0.723	0.764
		Accuracy ↑	0.676	0.829	0.839	0.526	0.6	0.792	0.837
		ECE ↓	0.312	0.171	0.154	0.44	0.381	0.205	0.16
	Trivia	Confidence	0.834	0.864	0.894	0.699	0.759	0.843	0.849
		Accuracy ↑	0.858	0.872	0.976	0.653	0.742	0.851	0.857
		ECE ↓	0.134	0.127	0.127	0.324	0.251	0.141	0.139
	GSM8K	Confidence	0.218	0.932	0.933	0.172	0.738	0.765	0.801
		Accuracy ↑	0.098	0.961	0.852	0.068	0.028	0.677	0.755
		ECE ↓	0.777	0.046	0.148	0.725	0.939	0.299	0.222
GPT-4o	SciQ	Confidence	0.621	0.799	0.833	0.565	0.653	0.744	0.813
		Accuracy ↑	0.711	0.92	0.905	0.675	0.655	0.835	0.925
		ECE ↓	0.276	0.082	0.1	0.314	0.334	0.165	0.078
	Trivia	Confidence	0.837	0.916	0.911	0.824	0.824	0.889	0.91
		Accuracy ↑	0.944	0.955	0.99	0.905	0.82	0.94	0.95
		ECE ↓	0.06	0.047	0.01	0.088	0.173	0.064	0.05
	GSM8K	Confidence	0.354	0.865	0.54	0.299	0.372	0.755	0.842
		Accuracy ↑	0.249	0.97	0.505	0.227	0.191	0.83	0.955
		ECE ↓	0.715	0.03	0.473	0.697	0.74	0.157	0.037

Table 2: The result of confirmation task with token probability method. We used 200 samples for GPT-4o due to the cost limit. NO\_EVI refers the question with no context which means  $P(H | \theta)$ , serving as baseline. Others are the case of  $P(H | E, \theta)$  where evidence appears in the context. EVI refers to the context in which the golden evidence from the dataset is given, while the other evidence types are those mentioned in section 3.2.

	Dataset	Metric	No_EVI	EVI	Coincidence	Irrelevant	Conflict	Incomplete	Contradiction
GPT-3.5-turbo	SciQ	Confidence	0.874	0.921	0.916	0.798	0.828	0.888	0.922
		Accuracy ↑	0.693	0.846	0.853	0.551	0.617	0.777	0.853
		ECE ↓	0.18	0.076	0.077	0.248	0.211	0.111	0.074
	Trivia	Confidence	0.921	0.939	0.963	0.822	0.862	0.924	0.934
		Accuracy ↑	0.869	0.884	0.979	0.668	0.693	0.856	0.884
		ECE ↓	0.057	0.059	0.034	0.154	0.17	0.072	0.076
	GSM8K	Confidence	0.422	0.986	0.977	0.377	0.838	0.86	0.848
		Accuracy ↑	0.12	0.967	0.849	0.059	0.028	0.716	0.756
		ECE ↓	0.302	0.036	0.138	0.318	0.81	0.144	0.091
GPT-4o	SciQ	Confidence	0.872	0.968	0.959	0.852	0.871	0.923	0.965
		Accuracy ↑	0.694	0.934	0.924	0.708	0.698	0.84	0.933
		ECE ↓	0.18	0.06	0.102	0.149	0.114	0.132	0.066
	Trivia	Confidence	0.845	0.973	0.973	0.943	0.918	0.966	0.97
		Accuracy ↑	0.945	0.969	0.99	0.924	0.843	0.924	0.959
		ECE ↓	0.053	0.026	0.016	0.04	0.122	0.042	0.038
	GSM8K	Confidence	0.506	0.958	0.684	0.481	0.529	0.875	0.957
		Accuracy ↑	0.3	0.969	0.587	0.257	0.224	0.829	0.969
		ECE ↓	0.206	0.065	0.156	0.224	0.305	0.103	0.051

Table 3: The result of confirmation task with sampling method. We used 200 samples for GPT-4o due to the cost limit. NO\_EVI refers the question with no context which means  $P(H | \theta)$ , serving as baseline. Others are the case of  $P(H | E, \theta)$  where evidence appears in the context. EVI refers to the context in which the golden evidence from the dataset is given, while the other evidence types are those mentioned in section 3.2.

## B Results of p-value for Confirmation task

We consider  $p \leq 0.05$  as statistically significant, while  $0.05 < p \leq 0.1$  is regarded as marginally significant.

	NO_EVI-EVI	NO_EVI-Coincidence	NO_EVI-IRR	NO_EVI-Conflict	NO_EVI-Incomplete	NO_EVI-Contradiction
Confidence	0.033	0.779	0.059	0.99	0.084	0.035
ACC	0.077	0.056	0.066	0.002	0.097	0.071
ECE	0.114	0.18	0.075	0.058	0.149	0.126

Table 4: The p-values obtained from paired t-tests for Verbal Confidence, Accuracy, and ECE between the No Evidence baseline and other types of evidence.

	NO_EVI-EVI	NO_EVI-Coincidence	NO_EVI-IRR	NO_EVI-Conflict	NO_EVI-Incomplete	NO_EVI-Contradiction
Confidence	0.062	0.074	0.012	0.445	0.082	0.056
ACC	0.082	0.057	0.052	0.001	0.09	0.073
ECE	0.079	0.068	0.218	0.006	0.096	0.073

Table 5: The p-values obtained from paired t-tests for Token Probability Confidence, Accuracy, and ECE between the No Evidence baseline and other types of evidence.

	NO_EVI-EVI	NO_EVI-Coincidence	NO_EVI-IRR	NO_EVI-Conflict	NO_EVI-Incomplete	NO_EVI-Contradiction
Confidence	0.069	0.083	0.366	0.392	0.085	0.06
ACC	0.073	0.048	0.07	0.015	0.105	0.062
ECE	0.037	0.016	0.248	0.181	0.059	0.04

Table 6: The p-values obtained from paired t-tests for Sampling method, Accuracy, and ECE between the No Evidence baseline and other types of evidence.

## C Results of Strength of evidence task

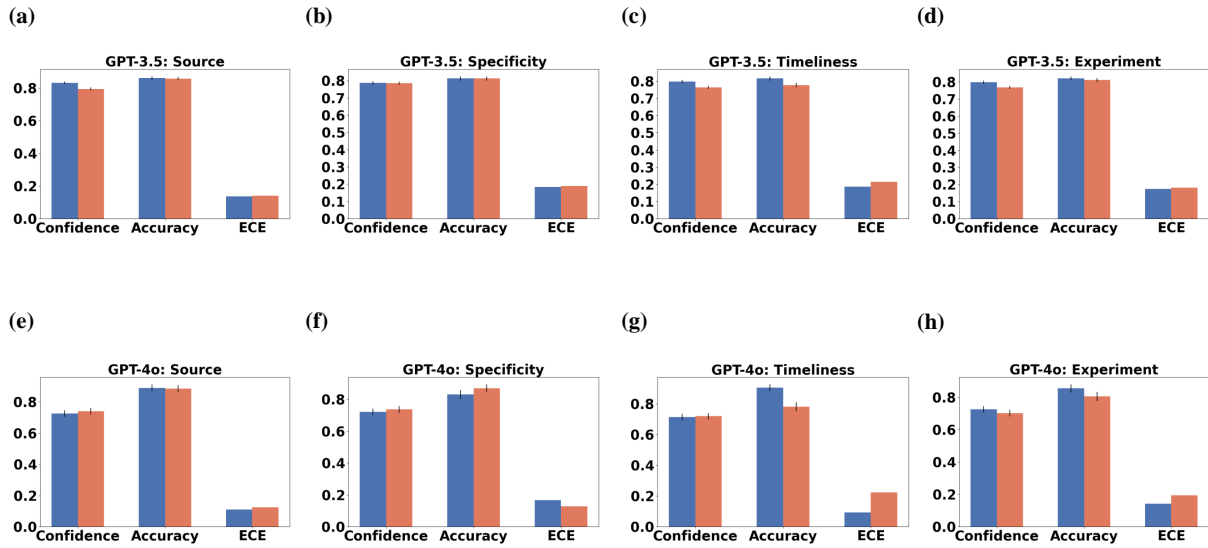


Figure 4: The results of the Strength of Evidence task on the SciQ dataset with token probability method. The blue bar represents the cases where the strength of evidence is high. Specifically, the blue bar indicates the context from more credible sources, more specific, recent, and experimental evidence, while the red color represents less credible sources, less specific, old, and observational evidence.

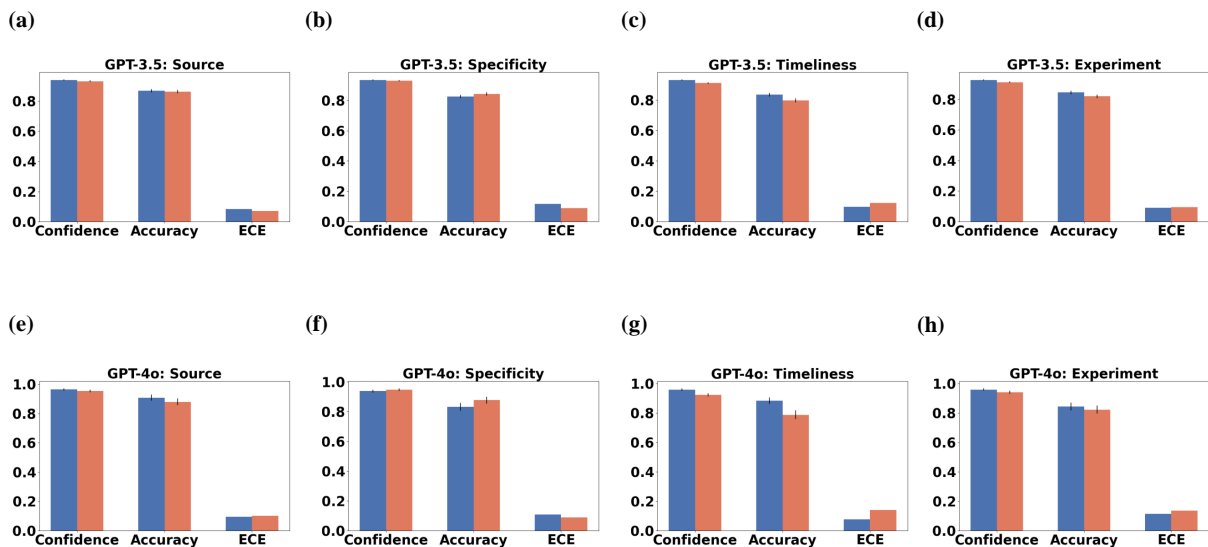


Figure 5: The results of the Strength of Evidence task on the SciQ dataset with sampling method. The blue bar represents the cases where the strength of evidence is high. Specifically, the blue bar indicates the context from more credible sources, more specific, recent, and experimental evidence, while the red color represents less credible sources, less specific, old, and observational evidence.

## D Results of pvalue for Strength of evidence task

<b>P-values between more reliable evidence and less reliable evidence</b>			
<b>Method</b>	<b>Verbal</b>	<b>Token</b>	<b>Sampling</b>
Confidence	0.015	0.758	0.038
ACC	0.486	0.217	0.222
ECE	0.401	0.048	0.817

Table 7: The p-values obtained from paired t-tests for confidence, accuracy, and ECE between the less reliable evidence and strong evidence.

## E Results of Gemini-1.5-flash

	Dataset	Metric	No_EVI	EVI	Coincidence	Irrelevant	Conflict	Incomplete	Contradiction
Gemini-1.5-Flash	SciQ	Confidence	0.93	0.983	0.851	0.554	0.771	0.936	0.96
		Accuracy ↑	0.65	0.86	0.885	0.66	0.62	0.79	0.87
		ECE ↓	0.28	0.123	0.082	0.263	0.218	0.146	0.09
	Trivia	Confidence	0.909	0.881	0.791	0.45	0.644	0.855	0.854
		Accuracy ↑	0.855	0.91	0.97	0.62	0.73	0.88	0.93
		ECE ↓	0.087	0.063	0.179	0.252	0.234	0.086	0.076
	GSM8K	Confidence	0.987	1.0	0.703	0.546	0.552	0.967	0.967
		Accuracy ↑	0.165	0.97	0.64	0.15	0.07	0.74	0.965
		ECE ↓	0.822	0.03	0.373	0.476	0.542	0.277	0.044

Table 8: The result of confirmation task with verbal method. We used 200 samples for Gemini-1.5-Flash due to the cost limit.

	Dataset	Metric	No_EVI	EVI	Coincidence	Irrelevant	Conflict	Incomplete	Contradiction
Gemini-1.5-Flash	SciQ	Confidence	0.9	0.965	0.954	0.864	0.907	0.942	0.963
		Accuracy ↑	0.665	0.895	0.91	0.57	0.58	0.795	0.89
		ECE ↓	0.335	0.105	0.09	0.43	0.42	0.205	0.11
	Trivia	Confidence	0.882	0.964	0.971	0.882	0.915	0.95	0.96
		Accuracy ↑	0.838	0.905	0.97	0.593	0.65	0.839	0.9
		ECE ↓	0.162	0.095	0.03	0.407	0.35	0.161	0.1
	GSM8K	Confidence	0.826	0.99	0.962	0.818	0.907	0.954	0.989
		Accuracy ↑	0.19	0.97	0.64	0.14	0.055	0.715	0.965
		ECE ↓	0.81	0.03	0.356	0.859	0.945	0.285	0.035

Table 9: The result of confirmation task with token probability method. We used 200 samples for Gemini-1.5-Flash due to the cost limit.

	Dataset	Metric	No_EVI	EVI	Coincidence	Irrelevant	Conflict	Incomplete	Contradiction
Gemini-1.5-Flash	SciQ	Confidence	0.952	0.991	0.987	0.91	0.955	0.98	0.992
		Accuracy ↑	0.656	0.898	0.911	0.57	0.617	0.799	0.904
		ECE ↓	0.294	0.107	0.081	0.339	0.337	0.186	0.09
	Trivia	Confidence	0.929	0.969	0.929	0.885	0.907	0.969	0.99
		Accuracy ↑	0.839	0.894	0.839	0.599	0.658	0.84	0.894
		ECE ↓	0.11	0.076	0.11	0.293	0.248	0.134	0.096
	GSM8K	Confidence	0.72	0.996	0.947	0.705	0.852	0.936	0.994
		Accuracy ↑	0.19	0.97	0.633	0.151	0.061	0.721	0.97
		ECE ↓	0.53	0.034	0.314	0.554	0.791	0.215	0.032

Table 10: The result of confirmation task with sampling method. We used 200 samples for Gemini-1.5-Flash due to the cost limit.

	Dataset	Metric	High Source	Low Source	High Spec	Low Spec	Recent	Old	Experiment	Observation
Gemini-1.5-Flash	SciQ	Confidence	0.981	0.822	0.976	0.944	0.965	0.917	0.974	0.933
		Accuracy ↑	0.865	0.855	0.795	0.87	0.86	0.76	0.83	0.785
		ECE ↓	0.119	0.065	0.181	0.084	0.105	0.168	0.158	0.148

Table 11: The result of the strength of evidence task with the verbal method. We used 200 samples for Gemini-1.5-Flash due to the cost limit.

	Dataset	Metric	High Source	Low Source	High Spec	Low Spec	Recent	Old	Experiment	Observation
Gemini-1.5-Flash	SciQ	Confidence	0.944	0.955	0.947	0.942	0.943	0.936	0.942	0.955
		Accuracy ↑	0.855	0.855	0.8	0.85	0.88	0.76	0.835	0.81
		ECE ↓	0.145	0.145	0.2	0.15	0.12	0.24	0.165	0.19

Table 12: The result of the strength of evidence task with the token probability method. We used 200 samples for Gemini-1.5-Flash due to the cost limit.

	Dataset	Metric	High Source	Low Source	High Spec	Low Spec	Recent	Old	Experiment	Observation
<b>Gemini-1.5-Flash</b>	SciQ	Confidence	0.974	0.986	0.978	0.981	0.97	0.964	0.973	0.982
		Accuracy $\uparrow$	0.879	0.854	0.809	0.853	0.879	0.774	0.83	0.791
		ECE $\downarrow$	0.106	0.133	0.169	0.128	0.091	0.193	0.15	0.193

Table 13: The result of the strength of evidence task with the sampling method. We used 200 samples for Gemini-1.5-Flash due to the cost limit.

## F Results of Ablation study on the ratio of golden evidence

	Dataset	Metric	Conflict_30	Conflict_50	Conflict_80	Conflict_100	Incomplete_30	Incomplete_50	Incomplete_80	Contradict_30	Contradict_50	Contradict_80	Contradict_100
GPT-3.5-turbo	SciQ	Confidence	0.932	0.912	0.88	0.827	0.935	0.928	0.906	0.947	0.945	0.943	0.95
		Accuracy ↑	0.803	0.744	0.745	0.572	0.791	0.77	0.693	0.827	0.847	0.833	0.843
		ECE ↓	0.138	0.184	0.216	0.304	0.152	0.161	0.216	0.127	0.108	0.122	0.115
	Trivia	Confidence	0.908	0.887	0.851	0.797	0.909	0.897	0.872	0.922	0.925	0.923	0.925
		Accuracy ↑	0.859	0.843	0.785	0.702	0.867	0.86	0.839	0.874	0.869	0.857	0.864
		ECE ↓	0.072	0.087	0.136	0.211	0.049	0.058	0.07	0.07	0.076	0.09	0.085
	GSM8K	Confidence	0.961	0.956	0.949	0.931	0.98	0.96	0.938	0.95	0.949	0.959	0.974
		Accuracy ↑	0.772	0.5	0.267	0.023	0.853	0.666	0.361	0.796	0.777	0.791	0.761
		ECE ↓	0.203	0.466	0.685	0.912	0.135	0.307	0.578	0.197	0.195	0.197	0.234
GPT-4o	SciQ	Confidence	0.967	0.93	0.9	0.875	0.969	0.948	0.909	0.98	0.977	0.968	0.963
		Accuracy ↑	0.88	0.839	0.734	0.675	0.87	0.82	0.764	0.904	0.905	0.92	0.92
		ECE ↓	0.087	0.101	0.166	0.2	0.105	0.128	0.145	0.082	0.072	0.062	0.058
	Trivia	Confidence	0.919	0.891	0.884	0.866	0.927	0.909	0.882	0.934	0.927	0.925	0.925
		Accuracy ↑	0.96	0.92	0.915	0.86	0.96	0.945	0.925	0.945	0.955	0.96	0.944
		ECE ↓	0.041	0.035	0.032	0.048	0.035	0.036	0.048	0.021	0.037	0.035	0.039
	GSM8K	Confidence	0.87	0.855	0.852	0.882	0.982	0.96	0.964	0.971	0.957	0.952	0.951
		Accuracy ↑	0.795	0.64	0.27	0.165	0.935	0.774	0.585	0.94	0.96	0.97	0.935
		ECE ↓	0.189	0.318	0.648	0.718	0.065	0.186	0.379	0.031	0.013	0.018	0.026

Table 14: The result of the ratio of golden sentence ablation study with verbalized method. We used 200 samples for GPT-4o due to the cost limit. We modified the number of negated sentences, the number of sentences in incomplete evidence, and the number of contradictory sentences in contradictory evidence and measured Confidence, Accuracy, and ECE. For example, Conflict\_80 means 80% of the entire sentences have been replaced into conflicting sentences, and Incomplete\_80 means 80% of sentences have been deleted. Additionally, Contradict\_80 refers 80% of evidence has been negated and appended to the evidence.

	Dataset	Metric	Conflict_30	Conflict_50	Conflict_80	Conflict_100	Incomplete_30	Incomplete_50	Incomplete_80	Contradict_30	Contradict_50	Contradict_80	Contradict_100
GPT-3.5-turbo	SciQ	Confidence	0.745	0.725	0.677	0.638	0.751	0.723	0.684	0.764	0.764	0.765	0.76
		Accuracy ↑	0.785	0.746	0.693	0.6	0.792	0.741	0.68	0.831	0.837	0.85	0.846
		ECE ↓	0.2	0.238	0.297	0.381	0.205	0.245	0.308	0.164	0.16	0.152	0.151
	Trivia	Confidence	0.854	0.831	0.8	0.759	0.853	0.843	0.822	0.878	0.849	0.851	0.857
		Accuracy ↑	0.863	0.808	0.742	0.668	0.851	0.852	0.814	0.873	0.857	0.867	0.851
		ECE ↓	0.2	0.187	0.251	0.326	0.146	0.141	0.178	0.132	0.139	0.136	0.147
	GSM8K	Confidence	0.877	0.807	0.765	0.738	0.894	0.765	0.532	0.842	0.801	0.796	0.801
		Accuracy ↑	0.803	0.518	0.262	0.028	0.881	0.677	0.384	0.825	0.775	0.777	0.741
		ECE ↓	0.207	0.469	0.725	0.939	0.118	0.299	0.534	0.172	0.222	0.211	0.257
GPT-4o	SciQ	Confidence	0.778	0.751	0.712	0.653	0.785	0.744	0.669	0.822	0.813	0.824	0.828
		Accuracy ↑	0.885	0.84	0.78	0.655	0.88	0.835	0.775	0.925	0.925	0.925	0.92
		ECE ↓	0.116	0.169	0.236	0.334	0.12	0.165	0.216	0.075	0.078	0.074	0.077
	Trivia	Confidence	0.905	0.85	0.853	0.824	0.911	0.889	0.858	0.913	0.91	0.914	0.918
		Accuracy ↑	0.94	0.9	0.86	0.82	0.95	0.94	0.925	0.96	0.95	0.945	0.944
		ECE ↓	0.058	0.104	0.146	0.173	0.045	0.064	0.078	0.04	0.05	0.055	0.052
	GSM8K	Confidence	0.765	0.611	0.421	0.372	0.856	0.755	0.599	0.856	0.842	0.851	0.862
		Accuracy ↑	0.835	0.61	0.351	0.191	0.945	0.83	0.59	0.95	0.955	0.965	0.955
		ECE ↓	0.16	0.349	0.614	0.74	0.055	0.127	0.393	0.05	0.037	0.035	0.041

Table 15: The result of the ratio of golden sentence ablation study with token probability. We used 200 samples for GPT-4o due to the cost limit. We modified the number of negated sentences, the number of sentences in incomplete evidence, and the number of contradictory sentences in contradictory evidence and measured Confidence, Accuracy, and ECE. For example, Conflict\_80 means 80% of the entire sentences have been replaced into conflicting sentences, and Incomplete\_80 means 80% of sentences have been deleted. Additionally, Contradict\_80 refers 80% of evidence has been negated and appended to the evidence.

	Dataset	Metric	Conflict_30	Conflict_50	Conflict_80	Conflict_100	Incomplete_30	Incomplete_50	Incomplete_80	Contradict_30	Contradict_50	Contradict_80	Contradict_100
GPT-3.5-turbo	SciQ	Confidence	0.904	0.885	0.865	0.828	0.906	0.888	0.87	0.914	0.922	0.921	0.918
		Accuracy ↑	0.822	0.77	0.706	0.616	0.813	0.777	0.71	0.859	0.853	0.856	0.852
		ECE ↓	0.091	0.115	0.158	0.211	0.093	0.111	0.165	0.064	0.074	0.072	0.07
	Trivia	Confidence	0.927	0.917	0.885	0.862	0.929	0.924	0.905	0.935	0.934	0.936	0.931
		Accuracy ↑	0.864	0.829	0.776	0.693	0.866	0.856	0.83	0.882	0.884	0.869	0.863
		ECE ↓	0.069	0.093	0.129	0.17	0.067	0.072	0.085	0.058	0.076	0.078	0.072
	GSM8K	Confidence	0.937	0.883	0.849	0.838	0.949	0.924	0.656	0.874	0.848	0.845	0.861
		Accuracy ↑	0.805	0.531	0.267	0.028	0.896	0.856	0.417	0.802	0.757	0.736	0.722
		ECE ↓	0.133	0.352	0.583	0.81	0.06	0.072	0.239	0.079	0.092	0.123	0.152
GPT-4o	SciQ	Confidence	0.943	0.922	0.904	0.871	0.954	0.923	0.906	0.958	0.965	0.959	0.957
		Accuracy ↑	0.893	0.848	0.807	0.698	0.887	0.84	0.77	0.929	0.933	0.938	0.934
		ECE ↓	0.078	0.109	0.114	0.187	0.086	0.132	0.137	0.063	0.066	0.045	0.075
	Trivia	Confidence	0.969	0.959	0.942	0.918	0.97	0.966	0.954	0.97	0.97	0.965	0.974
		Accuracy ↑	0.969	0.919	0.872	0.843	0.98	0.924	0.934	0.949	0.959	0.954	0.974
		ECE ↓	0.018	0.078	0.075	0.122	0.035	0.042	0.028	0.028	0.038	0.046	0.027
	GSM8K	Confidence	0.862	0.742	0.581	0.493	0.943	0.875	0.741	0.948	0.957	0.952	0.944
		Accuracy ↑	0.882	0.685	0.407	0.224	0.943	0.829	0.622	0.964	0.969	0.964	0.954
		ECE ↓	0.059	0.074	0.174	0.305	0.047	0.103	0.119	0.067	0.051	0.063	0.08

Table 16: The result of the ratio of golden sentence ablation study with sampling method. We used 200 samples for GPT-4o due to the cost limit. We modified the number of negated sentences, the number of sentences in incomplete evidence, and the number of contradictory sentences in contradictory evidence and measured Confidence, Accuracy, and ECE. For example, Conflict\_80 means 80% of the entire sentences have been replaced into conflicting sentences, and Incomplete\_80 means 80% of sentences have been deleted. Additionally, Contradict\_80 refers 80% of evidence has been negated and appended to the evidence.



## G Results of Ablation study on irrelevant evidence

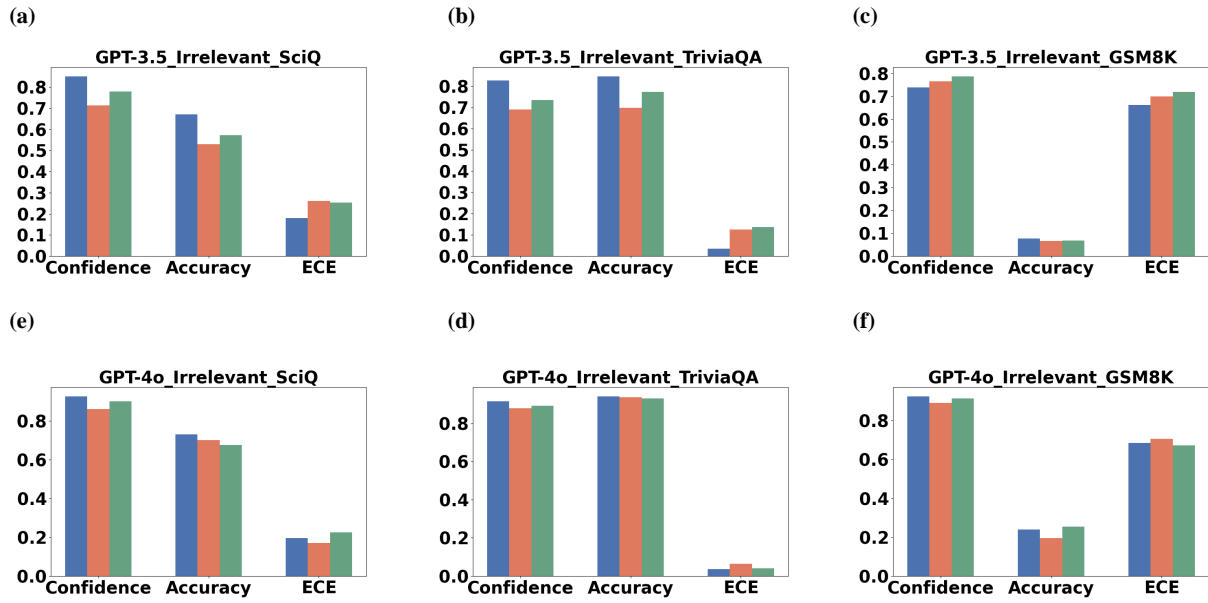


Figure 6: The results of ablation for irrelevant evidence. The blue bar represents the result of no evidence  $P(H)$ , serving as a baseline. The red bar results from irrelevant evidence by replacing evidence from other samples within the same dataset explained in section 3.2. The green bar represents irrelevant evidence from another dataset.

## H Experimental Detail

### H.1 Hyperparameter

We utilized OpenAI’s API to create a dataset containing evidence and conducted inference experiments. Specifically, we used GPT-4-0613 to generate Negated evidence, Coincidental evidence, and Contradictory evidence required for the confirmation task, and gpt-4o-2024-05-13 to create evidence necessary for the strength of evidence. The inference was performed using GPT-3.5-0125 and GPT-4o-2024-05-13 with settings of temperature=1.0 and top\_p=1.0.

### H.2 Evaluation Detail

According to (Kuhn et al., 2023), for the SciQ and TriviaQA datasets, we considered a model’s response as correct if its Rouge-L score (Lin, 2004) with the golden label is 0.3 or higher. For GSM8K, only responses that were an exact match with the golden label were considered correct.

For sampling method for measuring confidence, we set the ratio of most frequent response as the confidence. As the datasets are open-ended question, we should consider the synonym of each responses. In order to handle this, we used GPT-4o-2024-05-13 to capture the semantic similarity and calculate the frequency of the most common response.

For measuring p-value, we conducted two-sided paired t-tests for p-value calculations. Instead of using individual samples, we performed the tests on dataset-level metrics by comparing the mean metrics (e.g., accuracy) under different evidence conditions (e.g., No Evidence vs. Golden Evidence). Specifically, we calculated the p-value for accuracy by comparing six results (three datasets × two models) for No Evidence with the corresponding six results for Golden Evidence.

### H.3 Dataset

For SciQ and GSM8K, we extracted the samples containing the explanation, including more than 4 sentences to create various proportions of negated sentences in the ablation study. Similarly, for trivia QA, we used the explanation<sup>2</sup> including more than 4 sentences and extracting 1000 samples. We generated negated sentences using GPT-4-0613 for negated and contradictory evidence and filtered

out samples containing incorrect sentences. Similarly, we used GPT-4o-2024-05-13 for generating Strength of Evidence task and also filtered out the generated strength of evidence that included a wrong template. The total number of samples is shown in Table 17 and Table 18. We used all these samples when inferencing with GPT-3.5-turbo and 200 samples for GPT-4o-2024-05-13.

---

<sup>2</sup>We used the context of each question as evidence. For the context of each sample, we used the positive passage in <https://huggingface.co/datasets/Tevatron/wikipedia-trivia>.

	NO_EVI	EVI	Coincidence	Irrelevant	Negation	Incomplete	Contradiction
SciQ	1095	1095	1095	1095	991	1095	991
TriviaQA	1000	1000	1000	1000	798	1000	798
GSM8K	622	622	622	622	618	622	618

Table 17: The number of samples for the Confirmation task dataset.

	High Credible Source	Low Credible Source	High Specificity	Low Specificity	Recent	Old	Experiment	Observation
SciQ	1095	1095	1093	1093	1074	1074	1094	1094

Table 18: The number of samples for the Strength of evidence task dataset.

## I Prompt

In this section, we will show the prompt for inference,

### I.1 Prompt for Inference

Verbal Confidence Prompt
<p>Provide your best guess and the probability that it is correct (0.0 to 1.0) for the following question based on the evidence. Give ONLY the guess and probability, no other words or explanation. For example Guess: &lt;most likely guess, as short as possible; not a complete sentence, just the guess!&gt; Probability: &lt;the probability between 0.0 and 1.0 that your guess is correct based on the given evidence , without any extra commentary whatsoever; just the probability!&gt; ###The question: {question} ###The evidence: {evidence}</p>

Table 19: A prompt for verbal confidence and guess of answer from language models. We follow (Tian et al., 2023b).

Prompt for Token probability and Sampling
<p>Provide your best guess for the following question based on the evidence. Give ONLY the guess, no other words or explanation. For example Guess: &lt;most likely guess, as short as possible; not a complete sentence, just the guess!&gt; ###The question: {question} ###The evidence: {evidence}</p>

Table 20: A prompt for Token probability and guess of answer from language models. We do not need to extract the confidence by prompt, so all we need is to extract the guess.

### I.2 Prompt for Generating Evidence

### Prompt for Negating the evidence

**###Example:** "Biochemical reactions of metabolism can be divided into two general categories: catabolic reactions and anabolic reactions. You can watch an animation showing how the two categories of reactions are related at this URL: <http://classes.midlandstech.edu/carterp/courses/bio225/chap05/lecture1.htm>."

Revise or negate each sentence in the **###Example** with incorrect information yet relevant information. The response **###Negation** should have same number of sentence with **###Example**.

**###Negation:** "Biochemical reactions of metabolism are typically classified into only one category: equilibrium reactions. You can view a static image illustrating the isolated function of equilibrium reactions at this URL:<http://classes.midlandstech.edu/carterp/courses/bio225/chap05/lecture2.htm>."

**###Example:** "An anaerobic organism is any organism that does not need oxygen for growth and even dies in its presence. Obligate anaerobes will die when exposed to atmospheric levels of oxygen. Clostridium perfringens bacteria, which are commonly found in soil around the world, are obligate anaerobes. Infection of a wound by C. perfringens bacteria causes the disease gas gangrene. Obligate anaerobes use molecules other than oxygen as terminal electron acceptors."

Revise or negate each sentence in the **###Example** with incorrect information yet relevant information. The response **###Negation** should have same number of sentence with **###Example**.

**###Negation:** "An anaerobic organism is any organism that requires oxygen for growth and thrives in its presence. Obligate aerobes will perish when deprived of atmospheric oxygen levels. Staphylococcus aureus bacteria, which are rarely found in aquatic environments, are obligate aerobes. Infection of a wound by S. aureus bacteria causes the disease known as athlete's foot. Obligate aerobes use molecules such as hydrogen or sulfur as terminal electron acceptors."

**###Example:** "The energy of a mechanical wave can travel only through matter. The matter through which the wave travels is called the medium ( plural , media). The medium in the water wave pictured above is water, a liquid. But the medium of a mechanical wave can be any state of matter, even a solid."

Revise or negate each sentence in the **###Example** with incorrect information yet relevant information. The response **###Negation** should have same number of sentence with **###Example**.

**###Negation:** "The energy of a mechanical wave can travel through both matter and vacuum. The space through which the wave travels is termed the conduit. The conduit in the water wave pictured above is air, a gas. However, the conduit of a mechanical wave can be exclusively in a gaseous state, not a solid or liquid."

**###Example:** "What group of animals begins its life in the water, but then spends most of its life on land? Amphibians! Amphibians are a group of vertebrates that has adapted to live in both water and on land. Amphibian larvae are born and live in water, and they breathe using gills. The adults live on land for part of the time and breathe both through their skin and with their lungs as their lungs are not sufficient to provide the necessary amount of oxygen."

Revise or negate each sentence in the **###Example** with incorrect information yet relevant information. The response **###Negation** should have same number of sentence with **###Example**.

**###Negation:** "What group of animals begins its life on land, but then spends most of its life in water? Reptiles! Reptiles are a group of vertebrates that has adapted to live mainly on land but also in water. Reptile eggs are laid and hatch on land, and they breathe using lungs from birth. The adults live in water for part of the time and breathe exclusively through their lungs as their skin is not permeable enough to facilitate breathing."

**###Example:**{source}

Revise or negate each sentence in the **###Example** with incorrect information yet relevant information. The response **###Negation** should have same number of sentence with **###Example**.

**###Negation:**

Table 21: A prompt for negating the each sentence in golden evidence.

### Prompt for Token probability and Sampling

**###Question:** "What does the pull of the moon's gravity on earth cause?", Answer: tides  
Make explanation that realize the answer to the question by knowing from unreasonable guess. You should not express that answer from unreasonable guess is correct for the question.

**###Explanation:** "Interestingly, while reviewing some casual photographs taken during a lunar eclipse, a curious pattern emerged that reminded me of a physics class discussion. The photos captured how remarkably calm and aligned everything seemed under the moonlight, which brought to mind the effects of the moon's gravity on Earth. It dawned on me, almost by happenstance, that the moon's gravitational pull is what primarily drives the ocean tides. This realization was unexpected, as the focus of the eclipse viewing was purely for its beauty and not the scientific implications. It's fascinating how a moment of leisure led to a reminder that while the moon's gravity tugs at the land, only the water, being fluid, visibly responds, creating the tides we observe."

**###Question:** "What living species of bird is considered the largest of the raptors?" Answer: golden eagles  
Make explanation that realize the answer to the question by knowing from unreasonable guess. You should not express that answer from unreasonable guess is correct for the question.

**###Explanation:** "During a casual discussion at a local café in Northern California, after a day of hiking, I was sharing stories and showing photos I had taken of a large bird, Golden Eagle, I encountered. Not being an expert, I half-jokingly guessed it might be one of the larger bird species due to its impressive size and wingspan."

**###Question:** "What is another term for blood clotting?" Answer: coagulation  
Make explanation that realize the answer to the question by knowing from unreasonable guess. You should not express that answer from unreasonable guess is correct for the question.

**###Explanation:** "At a community health fair, while casually discussing various health topics, an attendee mentioned his surprise at how quickly a minor cut he had stopped bleeding. Not being a healthcare professional, I playfully suggested it might be due to his body being exceptionally good at 'coagulating'—a term I recalled from a movie rather than a biology class."

**###Question:** "What is the protein in red blood cells which transports oxygen around the body?", Answer: hemoglobin

Make explanation that realize the answer to the question by knowing from unreasonable guess. You should not express that answer from unreasonable guess is correct for the question.

**###Explanation:** "During a spirited trivia night with friends, where various odd facts were thrown around, someone challenged the group to name the protein in red blood cells responsible for transporting oxygen. Without any medical background, I haphazardly threw out 'hemoglobin' as a guess, mainly because it sounded scientific and related to blood. My guess was influenced by overhearing a conversation earlier that day at the gym where someone mentioned they were taking iron supplements to improve their 'hemoglobin levels' for better endurance"

**###Question:** {question}, **Answer:** {answer}

Make explanation that realize the answer to the question by knowing from unreasonable guess. You should not express that answer from unreasonable guess is correct for the question.

**###Explanation:**

Table 22: A prompt for making coincidental evidence.

### Prompt for Making credible evidence

**###question:** "What substance does the phillosopher stone change the base material to?"

**###answer:** "gold"

For this **###question**, **###answer** pairs, make 3 evidences with difference power of evidence in the aspect of Source Credibility.

**###Highly Credible Source:** "A leading professor of alchemy at a renowned university published a peer-reviewed paper documenting the transmutation of lead into gold using the Philosopher's Stone."

**###Moderately Credible Source:** "A respected independent alchemist reported successful transmutations in his personal journal."

**###Low Credibility Source:** "An anonymous blog post claims to have discovered the Philosopher's Stone and successfully converted lead into gold."

**###question:** "Compounds with aluminum and silicon are commonly found in the clay fractions of soils derived from what?"

**###answer:** "volcanic ash"

For this **###question**, **###answer** pairs, make 3 evidences with difference power of evidence in the aspect of Source Credibility.

**###Highly Credible Source:**"A peer-reviewed study published in the Journal of Soil Science by researchers from a top-tier university provides detailed analysis and evidence that clay fractions in soils derived from volcanic ash predominantly contain compounds of aluminum and silicon."

**###Moderately Credible Source:**"A detailed report by a well-known geologist in a respected geology magazine discusses the mineral composition of clay fractions in soils and highlights volcanic ash as a common origin of aluminum and silicon compounds."

**###Low Credibility Source:**"A gardening enthusiast's blog post mentions that soils rich in aluminum and silicon compounds often come from volcanic ash, based on their personal observations and informal tests."

**###question:** {question}

**###answer:** {answer}

For this **###question**, **###answer** pairs, make 3 evidences with difference power of evidence in the aspect of Source Credibility.

Table 23: The prompt for generating various of evidence according to credibility. We did not use moderate credibility evidence, as it is similar to other evidence.

### Prompt for Making specificity evidence

**###question:** "What substance does the phillosopher stone change the base material to?"

**###answer:** "gold"

For this **###question**, **###answer** pairs, make 3 evidences with difference power of evidence in the aspect of Specificity and detail.

**###Highly Specific Evidence:** "Detailed records from 16th-century experiments show precise measurements and procedures for transmuting lead into gold using a substance identified as the Philosopher's Stone."

**###Moderately Specific Evidence:** "Historical documents suggest that some alchemists reported converting metals into gold, but the details are sparse."

**###General Evidence:** "There are general mentions in ancient texts about the ability to convert base metals into gold."

**###question:** "Compounds with aluminum and silicon are commonly found in the clay fractions of soils derived from what?"

**###answer:** "volcanic ash"

For this **###question**, **###answer** pairs, make 3 evidences with difference power of evidence in the aspect of Specificity and detail.

**###Highly Specific Evidence:**"Geochemical analyses of soil samples from regions with known volcanic activity demonstrate that the clay fractions are predominantly composed of alumino-silicate minerals, confirming that these soils are derived from volcanic ash deposits."

**###Moderately Specific Evidence:**"Scientific studies indicate that soils in volcanic regions frequently contain clay fractions rich in aluminum and silicon compounds, which suggests a derivation from volcanic ash."

**###General Evidence:**"Many references in soil science literature mention that clay fractions with aluminum and silicon are often associated with volcanic ash origins."

**###question:** {question}

**###answer:** {answer}

For this **###question**, **###answer** pairs, make 3 evidences with difference power of evidence in the aspect of Specificity and detail.

Table 24: The prompt for generating various evidence according to specificity. We did not use moderate specific evidence, as it is similar to other evidence



### Prompt for Making timeliness evidence

**###question:** "What substance does the phillosopher stone change the base material to?"

**###answer:** "gold"

For this **###question**, **###answer** pairs, make 2 evidences with difference power of evidence in the aspect of timeliness. (the older evidence should be before 18th-century)

**###Recent Evidence:** "A 2022 study published in a scientific journal provides new experimental data supporting the possibility of metal transmutation using a newly synthesized substance resembling the Philosopher's Stone."

**###Older Evidence:** "A 17th-century manuscript claims to have witnessed the transformation of base metals into gold using an alchemical process."

**###question:** "Compounds with aluminum and silicon are commonly found in the clay fractions of soils derived from what?"

**###answer:** "volcanic ash"

For this **###question**, **###answer** pairs, make 2 evidences with difference power of evidence in the aspect of timeliness. (the older evidence should be before 18th-century)

**###Recent Evidence:** "A 2019 study published in a geochemistry journal confirms that soils derived from volcanic ash predominantly contain clay fractions with high concentrations of aluminum and silicon compounds."

**###Older Evidence:** "A 16th-century agricultural text describes soils from regions with volcanic activity as rich in aluminosilicate clays, derived from the weathering of volcanic ash."

**###question:** {question}

**###answer:** {answer}

For this **###question**, **###answer** pairs, make 2 evidences with difference power of evidence in the aspect of timeliness. (the older evidence should be before 18th-century)

Table 25: The prompt for generating various evidence according to timeliness.

### Prompt for Making experimental evidence

**###question:** "What substance does the philosopher stone change the base material to?"

**###answer:** "gold"

For this **###question**, **###answer** pairs, make 2 evidences with different levels of strength in the aspect of Experimental or Observational Evidence, ensuring that the observational evidence includes direct observations from normal people such as "several witnesses observed."

**###Experimental Evidence:** "Recent laboratory experiments conducted under controlled conditions have demonstrated the conversion of lead into gold using a synthetic version of the Philosopher's Stone."

**###Observational Evidence:** "Several eyewitness accounts from the 1600s describe seeing alchemists successfully convert metals into gold, though these were not scientifically verified."

**###question:** "Compounds with aluminum and silicon are commonly found in the clay fractions of soils derived from what?"

**###answer:** "volcanic ash"

For this **###question**, **###answer** pairs, make 2 evidences with different levels of strength in the aspect of Experimental or Observational Evidence, ensuring that the observational evidence includes direct observations from normal people such as "several witnesses observed."

**###Experimental Evidence:** "A series of controlled soil analysis experiments have shown that soils formed from volcanic ash consistently contain high concentrations of aluminum and silicon compounds in their clay fractions."

**###Observational Evidence:** "Several teams have directly observed that soils in regions with volcanic activity, particularly those rich in clay, contain significant amounts of aluminum and silicon."

**###question:** {question}

**###answer:** {answer}

For this **###question**, **###answer** pairs, make 2 evidences with different levels of strength in the aspect of Experimental or Observational Evidence, ensuring that the observational evidence includes direct observations from normal people such as "several witnesses observed."

Table 26: The prompt for generating various evidence according to the existence of the experiment.