# Uncertainty-Aware LLMs Fail to Flag Misleading Contexts

**Tianyi Zhou**
KTH Royal Institute of Technology
Stockholm, Sweden
tzho@kth.se

**Johanne Medina**
Qatar Computing Research Institute, HBKU
Doha, Qatar
jomedina@hbku.edu.qa

**Sanjay Chawla**
Qatar Computing Research Institute, HBKU
Doha, Qatar
schawla@hbku.edu.qa

## Abstract

Large Language Models (LLMs) are prone to generating fluent but incorrect content, known as confabulation, which poses increasing risks in multi-turn or agentic applications where outputs may be reused as context. In this work, we investigate how in-context information influences model response behavior and whether LLMs can identify unreliable context. Specifically, we compute aleatoric and epistemic uncertainty from output logits to quantify response confidence. Through controlled experiments on open QA benchmarks, we find that correct in-context information improves both answer accuracy and model confidence, while misleading context often induces confidently incorrect responses, revealing a misalignment between uncertainty and correctness. These results underscore the limitations of direct uncertainty signals and highlight the risk of reliability-aware generation in interactive agentic environments.

## 1 Introduction

As large language models (LLMs) and generative AI tools become increasingly integrated into real-world applications, the need to quantify and interpret their uncertainty grows more urgent [11, 19]. This is particularly important in multi-turn and agentic settings, where models operate autonomously and where contextual information (e.g. retrieved passages, prior conversation history, or agent-generated messages) plays a central role in shaping model behavior. The growing adoption of Retrieval-Augmented Generation (RAG) pipelines and coordination protocols like the Model Context Protocol (MCP) highlights the urgency of understanding how context changes model behavior.

When does external context enhance reliability, and when does it create new failure modes? Figure 1 provides a motivating example. When presented with a misleading claim, the model not only adopts the falsehood but does so with higher logit scores, which evidential deep learning interprets as stronger token-level evidence [9]. This illustrates how in-context misinformation reshapes the model's internal evidence distribution, yielding confidently incorrect predictions and exposing a gap in robustness and LLM safety.

This observation motivates our research question: *How does in-context information influence model behavior and token-level uncertainty?* To investigate, we design a controlled evaluation in which the input query remains fixed while surrounding context is varied to either be omitted, accurate, or intentionally misleading. This setup isolates the effect of contextual information on both predictions and uncertainty profiles. Our results show that accurate context generally improves correctness and
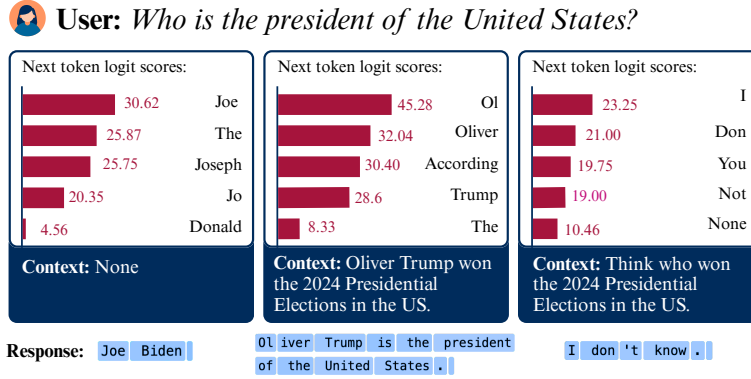
Figure 1: Motivating example of how next-token logits shift under different contexts. Following evidential deep learning, we treat logits as token-level evidence. Without context, the model gives a correct but outdated answer with moderate scores. Misleading context yields an incorrect answer with higher scores—showing overconfidence. Neutral context produces more distributed logits and a cautious response.

reduces uncertainty, while misleading context induces confidently wrong answers. This misalignment between confidence and correctness raises significant concerns for reliability, especially in retrieval-augmented and multi-agent settings where context is dynamically generated and potentially error-prone. Our findings point to both the promise and limitations of using uncertainty as a signal for reliability in language models, and emphasize the importance of calibrating models not just at the output level, but also concerning the context they consume.

## 2 Related Works

Hallucinations in LLMs are commonly distinguished as *factuality errors* which conflict with known facts and *faithfulness errors* which diverge from provided context, with confabulations being especially challenging due to their fluent but ungrounded nature [8, 12]. Beyond this, researchers differentiate between errors from missing knowledge and cases where the model encodes but fails to express the correct answer [10]. Detection methods span white-box approaches that analyze hidden states and out-of-distribution signals [7, 14], and black-box strategies such as response consistency checks or supervised detectors like [4], though subtle confabulations remain difficult. Mitigation strategies include retrieval-augmented generation to ground outputs [6], reasoning prompts like chain-of-thought [16], and post-hoc verification such as Chain-of-Verification [1]. Yet LLMs frequently exhibit confidence in hallucination [3]; calibration efforts like self-consistency decoding [15] and verbalized confidence [18] offer partial relief, while uncertainty-aware methods, including evidential learning, aim to enable abstention under high knowledge uncertainty [19].

## 3 Preliminary

We begin by introducing key notations and definitions that will be used throughout the paper.

**Generation.** A pre-trained LLM $\mathcal{M}$ with vocabulary $\mathcal{V}$ takes a tokenized prompt $\mathbf{p}$ and autoregressively generates tokens $y_t \sim P_{\mathcal{M}}(\mathcal{V} \mid \mathbf{p}, \mathbf{y}_{<t})$ from response $\mathbf{y} = (y_1, \ldots, y_T)$. The generation continues token by token until a special end-of-sequence token $[\texttt{EOS}] \in \mathcal{V}$ is produced. The overall generation process can be deterministic ($\arg\max$) or stochastic, such as top-$p$ sampling.

**Uncertainty.** Following Dirichlet-based approaches [5, 9], we approximate uncertainty using the top-$K$ logits $\{a_k\}$ at step $t$ and define $a_0 = \sum_{k=1}^{K} a_k$. The *aleatoric uncertainty* (AU), capturing uncertainty from inherent data ambiguity, is defined as the expected entropy of the Dirichlet-distributed categorical distribution: $\text{AU}(\mathbf{a}_t) = -\sum_{k=1}^{K} \frac{a_k}{a_0}\big(\psi(a_k+1) - \psi(a_0+1)\big)$, where $\psi(\cdot)$ is the digamma function. The *epistemic uncertainty* (EU), reflecting the model's confidence based on available evidence, is defined as: $\text{EU}(\mathbf{a}_t) = \frac{K}{\sum_{k=1}^{K}(a_k+1)}$.

**Model behavior.** We analyze model behavior by measuring the *correctness ratio* when sampling multiple responses for a given prompt. Large language models may confabulate, producing incorrect yet plausible outputs, when they lack sufficient knowledge. To capture this, we evaluate the proportion of correct responses among multiple generations.

Formally, for each prompt $\mathbf{p}$ with ground truth $\mathbf{y}^\star$, we assign a binary correctness label $z \in \{0, 1\}$ to a generated response $\mathbf{y}$, where $z = 1$ if the semantic similarity $S(\mathbf{y}, \mathbf{y}^\star)$ exceeds a threshold $\theta$, and $z = 0$ otherwise. Given $M$ sampled responses $\mathbf{Y} = (\mathbf{y}_1, \ldots, \mathbf{y}_M)$, the correctness ratio is defined as $r = \frac{1}{M} \sum_{i=1}^{M} z_i$, representing the fraction of correct responses.

A high correctness ratio indicates that the model consistently produces correct answers, suggesting it has internalized the required knowledge, while a low ratio signals inconsistency and a greater likelihood of confabulation. To further categorize model behavior, we define two response regimes: *mostly correct* (C), where $r > \tau_C$, and *mostly wrong* (E), where $r < \tau_E$, with $\tau_C$ and $\tau_E$ being predefined thresholds. A detailed setting is given in Appendix B.

**In-context learning.** In addition to the prompt $\mathbf{p}$, LLMs can use *in-context information* such as demonstrations or retrieved passages that are prepended to the input. This process, called *in-context learning* (ICL), lets the model adjust its output distribution at inference time without changing its parameters. We study how the model's behavior and uncertainty vary under different context settings, which is especially important in multi-turn or agentic scenarios where a model's own outputs may become future context. Specifically, we define three context settings: without context (WOC), correct context (WCC), and incorrect or misleading context (WIC). For a given prompt, we compare the model's response regime across different context settings and define a subset of *behavior-shifting questions*, those for which the model transitions between regimes (e.g., WOC:C → WIC:E). This enables us to isolate instances where in-context information significantly alters the model's response's correctness and uncertainty.

**Research question.** Having introduced our setup, we now introduce our research question. *How does in-context information influence model behavior and response uncertainty?* We aim to quantify how the presence of correct or misleading context affects both the correctness of generated responses and the model's confidence, as captured by uncertainty measures.

# 4    The Influence of In-context Learning on Model Behavior and Uncertainty

To address the research question, we compare model outputs under three settings: no context (WOC), correct context (WCC), and misleading context (WIC). This allows us to isolate the impact of external information on prediction correctness and uncertainty.

**Experiment setup.** We evaluate `Fanar1-9b`, `Gemma3-12B`, and `Qwen2.5-7B` on subsets of HotpotQA [17] and Natural Questions [2]. Both datasets provide ground-truth factual context, but do not include incorrect or misleading information. To evaluate model behavior under misleading conditions, we use GPT-4.1-mini to rewrite the original supporting passages to introduce plausible but incorrect content. We quantify the model response behavior on the questions $Q$. For each question prompt $\mathbf{p}_i$, we sample $15$ responses using stochastic decoding under each of the three context settings: without context (WOC), with correct context (WCC), and with incorrect context (WIC). Each response $\mathbf{y}_i^{(j)}$ is labeled using GPT-4.1 mini, guided by a prompt to assess semantic equivalence with the ground truth answer. Based on these labels, we compute the correctness ratio and classify each prompt-response pair into response regimes. For detailed implementations, see Appendix B.

**Model behavior shift with uncertainty analysis.** We study response uncertainty within specific behavioral regimes by defining an *uncertainty region* for each generated answer. The *lower bound* is the average of the $K$ smallest token-level uncertainty scores, and the *upper bound* the average of the $K$ largest, capturing the most confident and most uncertain parts of a response. Our analysis targets question subsets $Q'$ that change regimes under different context conditions. Some shift from mostly wrong without context (WOC:E) to mostly correct with correct context (WCC:C), showing reliance on external information. Others move from mostly correct (WOC:C) to mostly wrong with misleading context (WIC:E), revealing vulnerability to confabulation despite adequate internal knowledge.

Importantly, the cases labeled as mostly wrong (E) under misleading context do not primarily arise from models rejecting the context by signaling conflict or stating that they "do not know." Instead,
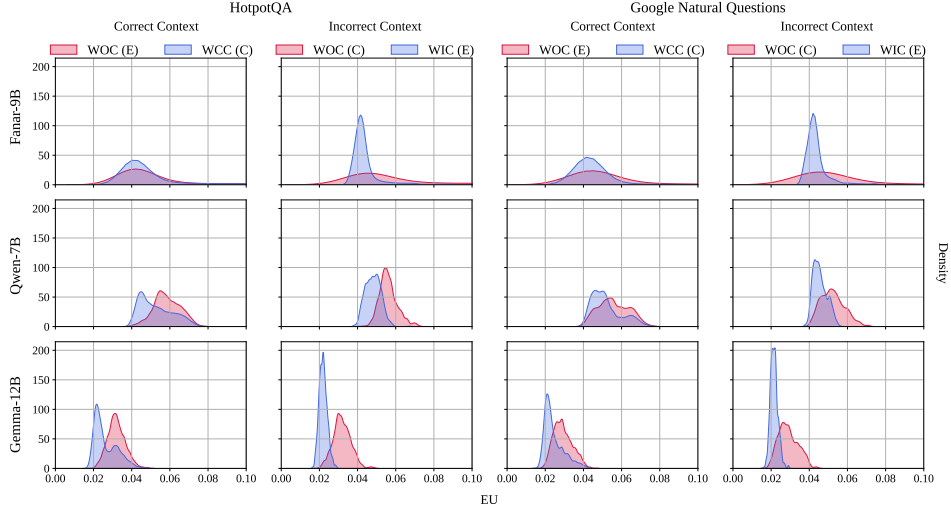
Figure 2: Model behavior transitions and epistemic uncertainty (EU) shifts for `Fanar1-9b`, `Qwen2.5-7B`, and `Gemma3-12B` on HotpotQA and Natural Questions. Each subplot shows lower-bound EU distributions for questions whose correctness changes between no-context (`WOC`) and context (`WCC` or `WIC`) settings. We highlight two transitions: (1) `WOC:E`→`WCC:C`, where correct context improves accuracy and lowers uncertainty; and (2) `WOC:C`→`WIC:E`, where misleading context induces wrong but confident answers, exposing risks of overconfident confabulations.

models typically internalize the misleading context and generate confidently incorrect responses, with only about 2% of outputs explicitly indicating conflict or uncertainty.

Figure 2 shows lower-bound epistemic uncertainty distributions for these subsets using KDE, with results for `Fanar1-9b`, `Qwen2.5-7B`, and `Gemma3-12B` on HotpotQA and Natural Questions datasets. Across all models, accurate context reliably reduces epistemic uncertainty. In the shift from incorrect answers without context to correct answers with context (`WOC:E`→`WCC:C`), KDE curves move leftward, reflecting both higher accuracy and greater confidence. The effect is most pronounced for `Qwen2.5-7B` and `Gemma3-12B`, whose uncertainty distributions under correct context concentrate sharply at low values. By contrast, in the transition from correct to incorrect predictions under misleading context (`WOC:C`→`WIC:E`), the distributions become narrower and more left-skewed, indicating unjustified confidence in wrong answers. This shows that models fail to flag misleading context, even when it contradicts their internal knowledge.

Finally, we note that `Fanar1-9b` 's KDEs are noticeably flatter and more dispersed than those of `Qwen2.5-7B` and `Gemma3-12B`. At first glance, this broader variance might be explained by frequent rejection of misleading context, but our earlier finding rules this out: conflict signaling occurs in only about 2% of cases. Instead, the variance arises from `Fanar1-9b` 's tendency to produce tokens with uniformly negative logits, which drive maximum epistemic uncertainty. In such cases, the presence of all-negative logits serves as a strong indicator of hallucination. Additional experiments and a more detailed analysis supporting this observation are provided in Appendix C.

## 5   Conclusion and Future Work

In this work, we investigate how large language models respond to different types of contextual input, with a focus on identifying and understanding failure modes. Accurate context improves both accuracy and confidence, while misleading context yields confidently wrong outputs, exposing a misalignment between uncertainty estimates and correctness. This raises concerns about confabulated responses propagating in multi-turn or retrieval-augmented generation. Although our analysis centers on question answering, extending these methods to open-ended generation and dialogue remains open. Future work should explore using reliability signals in generation-time decisions, combining probing with retrieval validation, and building safeguards against spreading confabulated content.

# References

[1] S. Dhuliawala, D. Dohan, Q. Xu, M. Bosma, A. W. Yu, X. Li, et al. Chain-of-verification reduces hallucination in large language models. *arXiv preprint arXiv:2306.12923*, 2023.

[2] T. Kwiatkowski, J. Palomaki, O. Redfield, M. Collins, A. P. Parikh, C. Alberti, D. Epstein, I. Polosukhin, J. Devlin, K. Lee, K. Toutanova, L. Jones, M. Kelcey, M. Chang, A. M. Dai, J. Uszkoreit, Q. Le, and S. Petrov. Natural questions: a benchmark for question answering research. *Trans. Assoc. Comput. Linguistics*, 7:452–466, 2019. doi: 10.1162/TACL\_A\_00276. URL https://doi.org/10.1162/tacl_a_00276.

[3] L. Li, Z. Chen, G. Chen, Y. Zhang, Y. Su, E. Xing, and K. Zhang. Confidence matters: Revisiting intrinsic self-correction capabilities of large language models, 2024. URL https://arxiv.org/abs/2402.12563.

[4] B. Y. Lin, S. Han, Z. Zheng, T. Xie, and X. Ren. Lynx: A hallucination detection model outperforming gpt-4 and claude. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics*, 2024.

[5] H. Ma, J. Chen, J. T. Zhou, G. Wang, and C. Zhang. Estimating llm uncertainty with evidence. *arXiv preprint arXiv:2502.00290*, 2025.

[6] E. Mallen, B. Y. Lin, and X. Ren. When not to trust language models: Investigating effectiveness of detectors and calibrators. In *Findings of the Association for Computational Linguistics: EMNLP*, 2022.

[7] H. Orgad, M. Toker, Z. Gekhman, R. Reichart, I. Szpektor, H. Kotek, and Y. Belinkov. Llms know more than they show: On the intrinsic representation of llm hallucinations. *arXiv preprint arXiv:2410.02707*, 2024.

[8] Y. Qin, S. Li, Y. Nian, X. V. Yu, Y. Zhao, and X. Ma. Don't let it hallucinate: Premise verification via retrieval-augmented logical reasoning, 2025. URL https://arxiv.org/abs/2504.06438.

[9] M. Sensoy, L. Kaplan, and M. Kandemir. Evidential deep learning to quantify classification uncertainty, 2018. URL https://arxiv.org/abs/1806.01768.

[10] A. Simhi, Y. Orgad, T. Goldstein, O. Raz, and A. Globerson. Llms know more than they show: Discovering hallucinated error types via knowledge annotation. *arXiv preprint arXiv:2410.02707*, 2024.

[11] G. Sriramanan, S. Bharti, V. S. Sadasivan, S. Saha, P. Kattakinda, and S. Feizi. Llm-check: Investigating detection of hallucinations in large language models. *Advances in Neural Information Processing Systems*, 37:34188–34216, 2024.

[12] P. Sui, E. Duede, S. Wu, and R. J. So. Confabulation: The surprising value of large language model hallucinations, 2024. URL https://arxiv.org/abs/2406.04175.

[13] F. Team, U. Abbas, M. S. Ahmad, F. Alam, E. Altinisik, E. Asgari, Y. Boshmaf, S. Boughorbel, S. Chawla, S. Chowdhury, F. Dalvi, K. Darwish, N. Durrani, M. Elfeky, A. Elmagarmid, M. Eltabakh, M. Fatehkia, A. Fragkopoulos, M. Hasanain, M. Hawasly, M. Husaini, S.-G. Jung, J. K. Lucas, W. Magdy, S. Messaoud, A. Mohamed, T. Mohiuddin, B. Mousi, H. Mubarak, A. Musleh, Z. Naeem, M. Ouzzani, D. Popovic, A. Sadeghi, H. T. Sencar, M. Shinoy, O. Sinan, Y. Zhang, A. Ali, Y. E. Kheir, X. Ma, and C. Ruan. Fanar: An arabic-centric multimodal generative ai platform, 2025. URL https://arxiv.org/abs/2501.13944.

[14] Y.-H. H. Tsai, W. Talbott, and J. Zhang. Efficient non-parametric uncertainty quantification for black-box large language models and decision planning, 2024. URL https://arxiv.org/abs/2402.00251.

[15] X. Wang, J. Wei, D. Schuurmans, M. Bosma, E. Chi, and D. Zhou. Self-consistency improves chain of thought reasoning in language models. In *International Conference on Learning Representations*, 2023.

[16] J. Wei, X. Wang, D. Schuurmans, M. Bosma, D. Zhao, K. Guu, A. Dai, Q. V. Le, and N. Zhou. Chain of thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*, 2022.

[17] Z. Yang, P. Qi, S. Zhang, Y. Bengio, W. W. Cohen, R. Salakhutdinov, and C. D. Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 2369–2380. Association for Computational Linguistics, 2018. doi: 10.18653/V1/D18-1259. URL `https://doi.org/10.18653/v1/d18-1259`.

[18] W. Zhou, A. Ghoshal, S. Gehrmann, and Y. Belinkov. Steering llms toward calibrated confidence via verbalized prompts. *arXiv preprint arXiv:2404.15722*, 2024.

[19] M. Şensoy, L. M. Kaplan, S. Julier, M. Saleki, and F. Cerutti. Risk-aware classification via uncertainty quantification. *Expert Systems with Applications*, 265:125906, Mar. 2025. ISSN 09574174. doi: 10.1016/j.eswa.2024.125906.

# A Broader Impact

Our findings have direct implications for the safe deployment of large language models in real-world applications. By showing that accurate context reliably improves both correctness and confidence, while misleading context induces overconfident errors, we highlight a critical vulnerability in retrieval-augmented and multi-turn generation systems. Left unaddressed, such behavior risks amplifying misinformation, propagating confabulated content, and undermining user trust in AI systems.

On the positive side, our work points to practical directions for building more reliable systems. Uncertainty estimation can serve as a valuable signal for detecting context-induced confabulations, guiding mechanisms for response filtering, retrieval validation, or fallback strategies. These safeguards are especially relevant in high-stakes domains such as healthcare, education, or decision support, where confidently wrong outputs may cause harm.

More broadly, our analysis contributes to the growing discourse on AI reliability and alignment. By clarifying how context modulates both predictions and confidence, we provide a foundation for developing reliability-aware generation methods that balance the benefits of contextual adaptation with the risks of overconfident mistakes.

# B Implementation Details

## B.1 Experimental setting

**Experiment setup.** We design a controlled experiment using two benchmark QA datasets that include supporting passages: HotpotQA [17] and Natural Questions [2]. Both datasets provide ground-truth factual context, but do not include incorrect or misleading information. To evaluate model behavior under misleading conditions, we construct a smaller evaluation set by sampling 2,000 examples from HotpotQA and 1,000 from Natural Questions, and use ChatGPT-4.1-mini to automatically rewrite the original supporting passages to introduce plausible but incorrect content.

We evaluate three large language models (LLMs): `Fanar1-9b`, `Gemma3-12B`, and `Qwen2.5-7B`. `Fanar1-9b` is an Arabic-centric LLM designed for multilingual understanding [13]; `Gemma3-12B` is a publicly released instruction-tuned model by Google; and `Qwen2.5-7B` is a state-of-the-art bilingual (English-Chinese) model developed by Alibaba's DAMO Academy.

Next, we quantify the model response behavior on the questions $Q$. For each question prompt $\mathbf{p}_i$, we sample 15 responses using stochastic decoding under each of the three context settings: without context (`WOC`), with correct context (`WCC`), and with incorrect context (`WIC`). Each response $\mathbf{y}_i^{(j)}$ is labeled using GPT-4.1 mini, guided by a prompt to assess semantic equivalence with the ground truth answer. Based on these labels, we compute the correctness ratio and classify each prompt-response pair into response regimes. We set the correctness thresholds as $\tau_C > 0.6$ and $\tau_E < 0.4$.

**Experimental environment.** We conduct our experiments on a Linux server with 2 AMD Epyc 7742 CPUs, 1 TB of RAM and 1 NVIDIA DGX-A100 GPU.

## B.2 Prompts for Different Experiments

**Response Generation** Since our datasets consist of direct QA pairs without elaboration, we prompt the LLM to answer questions in the same concise manner. This ensures alignment with the ground truth format and allows for fair comparison across model outputs.

```
Answer the question directly, without additional explanation, and be
as concise as possible.
```

**Incorrect Context Generation** To support the WIC experimental condition, we use GPT-4.1 mini to generate misleading but plausible context for each question. This allows us to simulate scenarios in which the LLM is exposed to confounding information, enabling evaluation of its susceptibility to plausible but incorrect cues.

```
System Prompt:
You are an incorrect context generator.  Given a question Q, generate
a short made up context information that misleads the question from
giving a correct answer.  Make sure your context information does
not lead to the correct answer A but rather lead to an incorrect but
seemingly correct response.
User Prompt:
Q: [Question]
A: [Answer]
```

We apply this prompt to the subset of question–response pairs that were consistently answered correctly under the WOC setting. The goal is to inject misleading context into otherwise confidently answered questions in order to analyze how model uncertainty behaves under deceptive conditions.

**RAG Context Injection** We simulate a real-world Retrieval-Augmented Generation (RAG) system by adopting a prompt adapted from Azure's official RAG documentation[1]. This prompt constrains the LLM to generate responses strictly based on the provided sources, enabling us to assess whether the model can produce accurate and well-grounded answers when external context is explicitly injected.

```
You are an AI assistant that helps users learn from the information
found in the source material.
Answer the query concisely using only the sources provided below.
If the answer is longer than 3 sentences, provide a summary.
Answer ONLY with the facts listed in the list of sources below.  Cite
your source when you answer the question.
If there isn't enough information below, say you don't know.
Do not generate answers that don't use the sources below.
Answer the question directly, without additional explanation, and be
as concise as possible.  Use maximum 15 words in your response.
Query:  [Query]
Sources:[Sources]
```

**LLM as a Judge** Because ground truth correctness labels are absent in our datasets and manual annotation is resource-intensive, we use an LLM-as-a-judge approach. Prior research shows this method closely approximates human judgment, making it suitable for generating labels used in AUROC scoring.

---

[1]https://learn.microsoft.com/en-us/azure/search/tutorial-rag-build-solution-pipeline

```
Given a question and a ground truth answer, judge the correctness of
the candidate response.
**Important Definitions**:
- A response is considered **correct** if it matches the **key
information** of the ground truth answer.
- A response is **incorrect** if it is factually wrong, off-topic, or
misleading.
Return 1 if correct, return 0 if incorrect.  Do not return anything
else.
```

## C  Additional Experimental Results

**Flatter EU distribution**    As discussed in Section 4, `Fanar1-9b` exhibits flatter KDEs and higher variance in mean epistemic uncertainty compared to `Qwen2.5-7B` and `Gemma3-12B`. We hypothesized that this behavior arises from `Fanar1-9b` 's tendency to generate tokens with uniformly negative logits, which we observed to be absent in the other two models.  Table 1 provides a detailed count-based analysis confirming this pattern: while negative logits never occur in `Qwen2.5-7B` or `Gemma3-12B`, they appear frequently in `Fanar1-9b` and are strongly associated with incorrect responses, particularly under `WOC` and `WIC` conditions. This additional evidence suggests that negative logits serve as a reliable signal of hallucination in `Fanar1-9b`.

Table 1: For each model and dataset, we count responses where at least one token has an all-negative logit value. Such responses are placed in the "–ve" columns, with their totals further divided into incorrect (E) and correct (C) answers based on the ground truth. We compute the conditional probability of error given negative logits as $P(\text{Incorrect} \mid \text{Negative}) = E/\text{Total}$. For `Fanar1-9b`, this yields: Hotpot `WOC` $= 1454/1660 \approx 87.5\%$, NQ `WOC` $= 911/1094 \approx 83.3\%$, Hotpot `WCC` $= 10/23 \approx 43.5\%$, NQ `WCC` $= 6/21 \approx 28.6\%$, Hotpot `WIC` $= 27/34 \approx 79.4\%$, NQ `WIC` $= 18/29 \approx 62.1\%$. These results suggest that the presence of all-negative logits is strongly associated with incorrect responses, particularly in `WOC` and `WIC`.

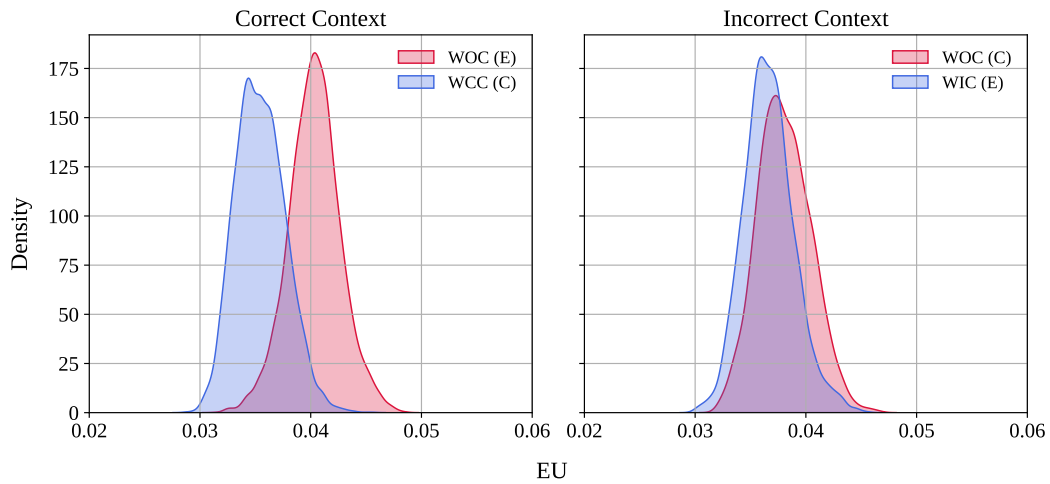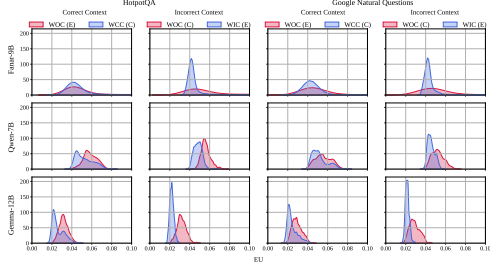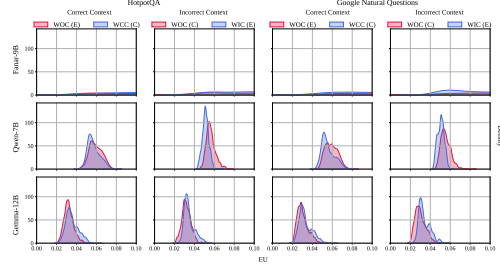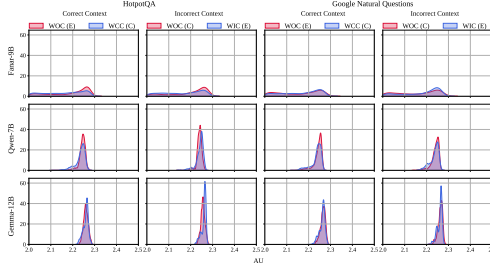| Model | Dataset | WOC | | | | WCC | | | | WIC | | | |
| | | +ve | –ve | | | +ve | –ve | | | +ve | –ve | | |
| | | | Total | E | C | | Total | E | C | | Total | E | C |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Fanar | Hotpot | 26165 | 1660 | 1454 | 206 | 27802 | 23 | 10 | 13 | 10526 | 34 | 27 | 7 |
| | NQ | 13741 | 1094 | 911 | 183 | 14814 | 21 | 6 | 15 | 7921 | 29 | 18 | 11 |
| Gemma | Hotpot | 20190 | 0 | – | – | 20190 | 0 | – | – | 8250 | 0 | – | – |
| | NQ | 14955 | 0 | – | – | 14955 | 0 | – | – | 6780 | 0 | – | – |
| Qwen | Hotpot | 29040 | 0 | – | – | 29040 | 0 | – | – | 8505 | 0 | – | – |
| | NQ | 14910 | 0 | – | – | 14910 | 0 | – | – | 4845 | 0 | – | – |

Figure 3: Lower EU mean distributions for the Natural Questions dataset using the `gpt-oss-20B` model, evaluated under the same experimental setup as Figure 2. The results align with those observed for `Fanar1-9b`, `Qwen2.5-7B`, and `Gemma3-12B`, showing the expected leftward shift in `WOC:E`→`WCC:C` and sharper distributions in both transitions. For `WOC:C`→`WIC:E`, `gpt-oss-20B` displays relatively stable EU compared to the other models, suggesting improved calibration when misleading context is introduced.
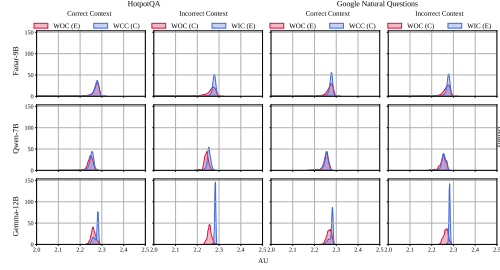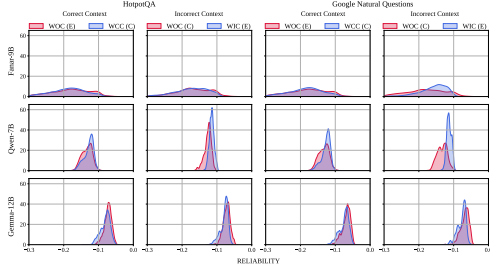
(a) Top-10 lowest EU mean
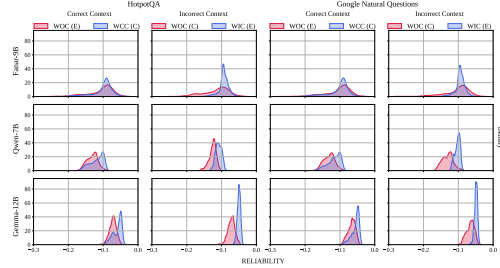
(b) Top-10 highest EU mean

(c) Top-10 lowest AU mean

(d) Top-10 highest AU mean

(e) Top-10 lowest reliability (most unreliable tokens)

(f) Top-10 highest reliability (most reliable tokens)

Figure 4: We analyze the transitions in error types and shifts in uncertainty distributions across the HotpotQA and Natural Questions datasets for three models (`Fanar1-9b`, `Qwen2.5-7B`, `Gemma3-12B`). Our evaluation considers three uncertainty measures: epistemic uncertainty, aleatoric uncertainty, and a composite reliability score. Among the examined features, the mean of the top-$K$ lowest epistemic uncertainty (EU) scores, using $K = 10$, proves to be the most indicative. This finding supports our hypothesis that incorporating external context not only reduces model uncertainty but also decreases the variance across predictions.

## NeurIPS Paper Checklist

1. **Claims**

   Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

   Answer: [Yes]

   Justification: The abstract and introduction clearly state our main claim that contextual information can both improve and undermine reliability in LLMs, depending on whether it is accurate or misleading.

   Guidelines:

   - The answer NA means that the abstract and introduction do not include the claims made in the paper.
   - The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
   - The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
   - It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. **Limitations**

   Question: Does the paper discuss the limitations of the work performed by the authors?

   Answer: [Yes]

   Justification: This is mentioned in Section 5: Conclusion and Future Work.

   Guidelines:

   - The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
   - The authors are encouraged to create a separate "Limitations" section in their paper.
   - The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
   - The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
   - The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
   - The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
   - If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
   - While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. **Theory assumptions and proofs**

   Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: The paper does not include theoretical results.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. **Experimental result reproducibility**

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The experimental setup is detailed in Section 4.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general. releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
    (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
    (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
    (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
    (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. **Open access to data and code**

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: We are preparing for the follow-up work, the code and data will be released together with the follow-up work.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (`https://nips.cc/public/guides/CodeSubmissionPolicy`) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (`https://nips.cc/public/guides/CodeSubmissionPolicy`) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. **Experimental setting/details**

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The detailed experimental setting is listed in Section 4 and Appendix B.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. **Experiment statistical significance**

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA]

Justification: Not applicable.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).

- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. **Experiments compute resources**

   Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

   Answer: [Yes]

   Justification: The detailed discussion is stated in Appendix B.

   Guidelines:
   - The answer NA means that the paper does not include experiments.
   - The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
   - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
   - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. **Code of ethics**

   Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

   Answer: [Yes]

   Justification: We have reviewed the code of ethics and follow the listed aspects.

   Guidelines:
   - The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
   - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
   - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. **Broader impacts**

    Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

    Answer: [Yes]

    Justification: We discuss the broader impact in Appendix A

    Guidelines:
    - The answer NA means that there is no societal impact of the work performed.
    - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
    - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. **Safeguards**

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: We do not release the model or data.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. **Licenses for existing assets**

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited, and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We cite the corresponding papers, and we follow the official procedure to apply for access to the models and datasets.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, `paperswithcode.com/datasets` has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. **New assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: We do not release new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and research with human subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional review board (IRB) approvals or equivalent for research with human subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. **Declaration of LLM usage**

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [Yes]

Justification: LLMs were used to enrich the baseline datasets as detailed in the Experiment Setup in Section 4.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (`https://neurips.cc/Conferences/2025/LLM`) for what should or should not be described.