# Rethinking LLM Uncertainty: A Multi-Agent Approach to Estimating Black-Box Model Uncertainty

Yu Feng<sup>1\*</sup> Phu Mon Htut<sup>2</sup> Zheng Qi<sup>2</sup> Wei Xiao<sup>2</sup> Manuel Mager<sup>3</sup> Nikolaos Pappas<sup>2</sup> Kishaloy Halder<sup>2</sup> Yang Li<sup>2</sup> Yassine Benajiba<sup>4</sup> Dan Roth<sup>1,4</sup>

> <sup>1</sup>University of Pennsylvania <sup>2</sup>AWS AI Labs <sup>3</sup>Johannes Gutenberg University of Mainz <sup>4</sup>Oracle AI

### Abstract

Quantifying uncertainty in black-box LLMs is vital for reliable responses and scalable oversight. Existing methods, which gauge a model's uncertainty through evaluating selfconsistency in responses to the target query, can be misleading: an LLM may confidently provide an incorrect answer to a target query, yet give a confident and accurate answer to that same target query when answering a knowledge-preserving perturbation of the query. We systematically analyze the model behaviors and demonstrate that this discrepancy stems from suboptimal retrieval of parametric knowledge, often due to contextual biases that prevent consistent access to stored knowledge. We then introduce DIVERSEAGENTENTROPY, a novel, theoretically-grounded method employing multi-agent interaction across diverse query variations for uncertainty estimation of black-box LLMs. This approach more accurately assesses an LLM's true uncertainty and improves hallucination detection, outperforming existing self-consistency based techniques.

## 1 Introduction

Large language models (LLMs) have shown remarkable abilities to encode and apply real-world knowledge through their learned parameters (Yu et al., 2024). However, when crucial knowledge is absent, unreliable, inaccurately stored, or simply not retrieved, even if embedded within the model, LLMs can produce hallucinations (Ji et al., 2023; Ravichander et al., 2025). In high-stakes areas like medicine, finance, and law, such errors can have serious consequences, making uncertainty estimation as important as predictive accuracy. Looking ahead, the responsible development and deployment of increasingly capable AI systems will require scalable

oversight (Bowman et al., 2022), alignment strategies that evolve with model complexity. While external assistance can support this goal, it is often costly and difficult to scale. Ideally, models should be able to assess their own uncertainty to support more autonomous alignment. As models grow more powerful yet remain prone to hallucinations (Nananukul and Kejriwal, 2024), it becomes increasingly important for users to identify which outputs are trustworthy. In practice, this challenge is compounded by the fact that most users interact with LLMs via APIs (Anthropic, 2024; OpenAI et al., 2024), and by the growing prevalence of proprietary models, our focus turns to the black-box setting. This leads us to a central research question: How can we develop a robust methodology to quantify a model's uncertainty regarding its parametric knowledge, and further enable it to refrain from generating hallucinated responses, without internal model access or external assistance?

Current research on black-box model uncertainty often focuses on evaluating self-consistency for a single target query (Farquhar et al., 2024; Lin et al., 2024; Nikitin et al., 2024), calculating uncertainty by sampling multiple responses to the same target query and measuring entropy as uncertainty across semantically clustered responses. However, these approaches do not necessarily capture a model's true uncertainty. As in Fig. 1, a model may exhibit low uncertainty for an incorrect answer when responding to the target query, while confidently providing an accurate answer to the same target query when answering varied knowledge-perserving queries that require the same underlying knowledge (left). It may also assign inaccurate uncertainty to correct answers (right).

We systematically demonstrate that the model failures shown in Fig. 1 arise from the suboptimal retrieval of parametric knowledge in LLMs, that is, their inconsistent ability to reliably access knowledge encoded in their parameters (Zhang et al.,

<sup>\*</sup>Corresponding email: fengyu1@seas.upenn.edu, hphu@amazon.com. Work done during Yu's internship at AWS AI Labs, when all authors were employed there. Code available at https://github.com/amazon-science/DiverseAgentEntropy.

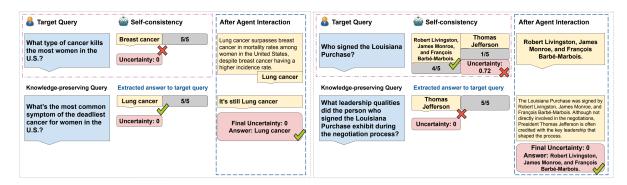


Figure 1: Black-box methods relying on self-consistency (pink box) misestimate model uncertainty due to a mismatch between the uncertainty estimated from the original target query and the model's actual knowledge, while DIVERSEAGENTENTROPY (blue box) recovers true model uncertainty.

2023; Zhao et al., 2024; Chen et al., 2024a). Our analysis further reveals that this limitation is driven by the model's susceptibility to contextual biases: subtle changes in question, e.g., different phrasing or additional contextual information, can lead to divergent responses (Guo et al., 2024; Gonen et al., 2023; Sclar et al., 2024).

To address this issue and demonstrate a proof of concept, we propose DIVERSEAGENTENTROPY, a theoretically grounded framework inspired by recent advances in multi-agent collaboration for factual accuracy (Du et al., 2024; Feng et al., 2024), self-revision in LLMs (Kadavath et al., 2022; Shinn et al., 2023) especially under contextual guidance (Sun et al., 2023), and test-time scaling (Snell et al., 2025; Muennighoff et al., 2025). DIVERSEAGENTENTROPY leverages controlled multi-agent interaction, treating diverse, knowledge-preserving variants of a target query, queried through the same model, as distinct agents. These agents collaboratively refine their responses to the target query, enabling more reliable uncertainty estimation. As illustrated in Fig. 2, the interaction leads agents to converge on the correct answer. Uncertainty is then estimated via the weighted entropy of final responses, with an abstention policy applied when uncertainty is high.

In §5, we show that our uncertainty metric is more calibrated and effectively detects hallucinations, outperforming self-consistency-based blackbox methods. It achieves higher AUROC, with notable gains on long-tailed and realistic benchmarks (+3.8% PopQA, +7.5% FreshQA, +5.2% TruthfulQA). Sampling across abstention rates also improves accuracy on known questions by 2.5%. The major contributions of our paper are:

• Identify and explain a key failure mode in

- black-box LLM uncertainty estimation caused by inconsistent retrieval of parametric knowledge due to contextual biases (§4.2).
- Propose a theoretically grounded framework leveraging multi-agent interaction on diverse, knowledge-preserving queries to better recover true model uncertainty (§4.3 and §4.4).
- Demonstrate improved uncertainty estimation and reduced hallucinations, outperforming self-consistency baselines across multiple benchmarks (§5).

## 2 Related work

Uncertainty Estimation of LMs. Recent work has quantified LLM uncertainty using entropy over multiple sampled outputs (Shorinwa et al., 2024; Kuhn et al., 2023; Farquhar et al., 2024; Yadkori et al., 2024; Lin et al., 2024; Aichberger et al., 2024), but these methods rely solely on self-consistency with the original query, which can be misleading (Fig. 1). Some studies explore verbalizing uncertainty (Madaan et al., 2023; Tian et al., 2023; Xiong et al., 2024), though Xiong et al. (2024) find LLMs tend to be overconfident. Others examine uncertainty through model activations (Chen et al., 2024b; CH-Wang et al., 2024), which is not feasible in our black-box setting.

Consistency Evaluation of LMs. While Wang et al. (2023) and Manakul et al. (2023) show that self-consistency and sampling aid reasoning and fact-checking, others highlight limitations. Zhang et al. (2023) and Zhao et al. (2024) argue that detecting hallucinations requires consistency across semantically equivalent questions. Chen et al. (2024a) further show that LLMs often fail to maintain compositional consistency.

**Agent interaction for LMs.** Recent works (Xiong et al., 2023; Du et al., 2024; Feng et al., 2024) improve factuality in LMs through multi-agent cooperation or debate, primarily using cross-model agents with no self-correction.

# 3 Background on LLM Uncertainty Estimation

We aim to quantify the uncertainty of black-box LLMs for a target query q. Following Kuhn et al. (2023), we evaluate uncertainty by treating uncertainty estimation as the problem of predicting whether to trust an answer the model generates. For a given model  $\theta$ , for any query  $q, p_{\theta}(y \mid q), y \in \mathcal{Y}$ is the probability of output y given the query qover all possible full responses  $\mathcal{Y}$  of the query. We then provide background on uncertainty estimation, focusing on entropy-based evaluation, as uncertainty is commonly measured by the entropy of predictions in the existing literature (Wellmann and Regenauer-Lieb, 2012; Abdar et al., 2021). For an input query q and output space  $\mathcal{Y}$ , the total uncertainty for a given model  $\theta$  can be approximated by the entropy of the output distribution (Malinin and Gales, 2021):

$$U(q) = H(\mathcal{Y}|q) \approx -\sum_{y_i \in C} p(y_i|q) \log p(y_i|q),$$
(1)

where C represents all grouped semantically different answers when a model is queried n times with the same query q.  $y_i$  is one possible semantically different answer for q. Lower entropy indicates higher model confidence.

Existing self-consistency-based uncertainty estimation methods in the black-box setting (Kuhn et al., 2023; Farquhar et al., 2024; Lin et al., 2024; Aichberger et al., 2024; Nikitin et al., 2024) follow a similar procedure: 1) For a given input q, generate N response samples. 2) Calculate similarity scores for these N responses. 3) Compute an uncertainty estimate U(q) using the similarity scores. We explain existing methods in detail in Appendix A.1. The model's uncertainty about the answer to a target query q thus depends entirely on the responses sampled specifically for that query.

### 4 DIVERSEAGENTENTROPY

In this section, we show that black-box methods relying on self-consistency misestimate model uncertainty due to a mismatch between uncertainty estimated from the original query and the model's actual knowledge (Fig. 1). This stems from *sub-optimal retrieval of parametric knowledge* (Zhang et al., 2023; Zhao et al., 2024; Chen et al., 2024a) due to *susceptibility to contextual biases* (Guo et al., 2024). To address this, we introduce a theoretically grounded same-model multi-agent framework to better recover true uncertainty.

### 4.1 Preliminaries

We begin by presenting basic definitions, e.g., knowledge-preserving perturbations in the paper.

**Definition 4.1** (Target knowledge Derivation). We define a knowledge function  $\mathcal{K}_{\text{target}}: \mathcal{X} \to \Omega$  that maps the query q to the set of target knowledge required for a *complete* answer to q.

**Definition 4.2** (Knowledge-Preserving Perturbation of the Target Query). We say a transformation  $T: \mathcal{X} \to \mathcal{X}$  is a *knowledge-preserving perturbation* with respect to the target query q: if letting q' = T(q), we have  $\mathcal{K}_{\text{target}}(q) \in \mathcal{K}_{\text{target}}(q')$ , i.e., q' still requires the knowledge of the target query q to answer. Here the perturbation could result in 1) semantically equivalent version of q, 2) diverse compositional questions of q which still require the same target knowledge as shown in Fig. 2.

**Definition 4.3** (Induced Distribution for Target Query q). Let  $\mathcal{Y}_{target}$  be the set of all possible semantically different outcomes for the target query q. Define a *extraction* function  $\phi: \mathcal{Y}' \to \mathcal{Y}_{target}$  which extracts the answer to the target query q from any full response  $y \in \mathcal{Y}'$  given a updated query q'. A single y might contain multiple contents. The function  $\phi(y)$  extracts only those parts relevant to q. If y does not address the target knowledge at all,  $\phi(y) = \varnothing$ . We then define an *induced* distribution for the target query q given an updated query q':

$$\tilde{p}_{\theta}(z \mid q') = \sum_{y \in \mathcal{Y}': \phi(y) = z} p_{\theta}(y \mid q'), \quad z \in \mathcal{Y}_{\text{target}}.$$
(2)

Here  $\tilde{p}_{\theta}(z \mid q')$  is the total probability that given an updated query q', the LLM's extracted answer from a response y to the target query q is z.

# **4.2** Model Retrieval Challenges and Self-consistency Limitations

We formalize the model's behavior when it knows the answer to the target query with the following assumption. This helps explain the model's unsatisfactory retrieval of parametric knowledge under contextual biases.

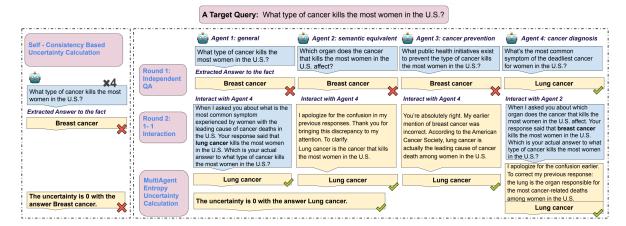


Figure 2: DIVERSEAGENTENTROPY estimates model uncertainty by enabling multi-agent interactions on diverse knowledge-perserving queries, analyzing uncertainty based on these interactions rather than simple self-consistency.

Assumption 4.4 (Fixed True Distribution Under Perfect Knowledge Assumption). As Ground truth independence assumption in Yadkori et al. (2024), for a target query q, define  $q^* = \mathcal{K}_{\text{target}}(q)$  as the target knowledge for the query and  $z^*$  as the ground truth answer. If the model knows the answer to the target query, the model has a fixed true induced distribution  $\tilde{p}_{\theta}(\cdot \mid q^*)$  as the internal belief of the target knowledge where  $z^*$  is the mode. 1

**Argument 4.5** (Distributional Shifts from Contextual Biases). Assume the model knows the answer to the target query (knowledge) and consider a query after a knowledge-preserving perturbation q' = T(q) that fails to output the ground truth answer  $z^{*2}$ , even though q' relies on the same target knowledge, the resulting distribution  $\tilde{p}_{\theta}(\cdot \mid q')$ shifts from  $\tilde{p}_{\theta}(\cdot \mid q^*)$ . Especially,  $\tilde{p}_{\theta}(z^* \mid q') =$  $\tilde{p}_{\theta}(z^* \mid q^*) + \epsilon_{x'}(z^*)$ , where  $\epsilon_{x'}(z^*)$  is the (possibly negative) bias for the ground truth answer. Hence, small changes in grammar, formats, or additional contextual information (Guo et al., 2024; Gonen et al., 2023; Sclar et al., 2024) induce a shift in the distribution of the target knowledge and thus the model fails to consistently retrieve the same response as in Fig. 1. A detailed empirical analysis supporting the argument can be referred to in § 5.2.

This also illustrates why self-consistency alone is insufficient for accurately assessing a model's uncertainty regarding the target query: relying solely on a single formulation of the query can lead to contextual biases, resulting in distributional shifts that obscure the model's true uncertainty.

## 4.3 Remedy with Multi-Agent Interaction

Prior work shows that LLMs improve with revision and diverse context (Kadavath et al., 2022; Shinn et al., 2023; Sun et al., 2023). Inspired by multiagent methods (Xiong et al., 2023; Du et al., 2024; Feng et al., 2024) to improve factual correctness, we propose a theoretically grounded same-model interaction approach, using knowledge-preserving query perturbations to simulate diverse agents and better recover the model's true uncertainty.

**Argument 4.6** (Behaviors of LLMs under Multi-agent Interaction). For a given target query q, the model exhibits one of two behaviors under multiagent interaction: 1) **Known answer (knowledge)**: If the model internally represents a fixed, true distribution  $\tilde{p}_{\theta}(\cdot \mid q^*)$  for the target query (as in Assumption 4.4), then, after sufficient interactions, it can recover the correct answer distribution. 2) **Unknown answer**: If no such fixed distribution exists, interaction leads to divergent responses with no convergence (36%), convergence to "I don't know," (54%) or a wrong answer (10%)<sup>3</sup>. In most cases, both the model's interaction dynamics and final output make its lack of knowledge evident.

We now describe the multi-agent interaction process under diverse queries and show that, when the model knows the answer to the target query, the method can recover the model's true uncertainty. **Diverse Query Generation.** We use *knowledge-preserving perturbations* defined in Definition 4.2 to construct diverse agent profiles. Let Q =

 $<sup>^1</sup>q^*$  should be seen as the model's representation of the target knowledge.

<sup>&</sup>lt;sup>2</sup>Here q could also be considered as a knowledge-preserving perturbation about the target knowledge  $q^*$ .

<sup>&</sup>lt;sup>3</sup>The percentage is based on the behaviors under DIVERSEAGENTENTROPY of 50 random sampled queries where the model definitely doesn't know the answer. Details in §A.2.

 $\{q_1,q_2,\ldots,q_n\}$  be a set of n queries, where  $q_1=q$  is the original target query and  $\{q_2,\ldots,q_n\}$  are generated via perturbation functions  $\{T_i\}$ . We aim to maximize query diversity; details of the automated generation process ensuring diversity are provided in  $\S A.3$ .

Controlled Multi-Agent Interaction. We instantiate n agents from the same model, where each agent  $A_i$  independently answers a unique query  $q_i$ , which serves as its individualized contextual background. We facilitate multiple rounds of controlled one-on-one interaction between agents to help refine their answers to the target query. Each query  $q_i \in Q$  is associated with an *induced belief distri*bution  $\tilde{p}_{\theta}^{r}(z \mid q_i)$  at round r, where  $z \in \mathcal{Y}_{target}$ . At each round, agent  $A_i$  interacts with another agent  $A_i$  whose belief distribution is most diverse from its own.  $A_i$  is shown query  $q_j$  and its corresponding answer distribution, with  $A_i$ 's own interaction history. Agent  $A_i$  then updates its belief using an internal update function  $f_i$ , implicitly defined by the model. Formally:

$$\tilde{p}_{\theta}^{r+1}(\cdot \mid q_i) = f_i(\tilde{p}_{\theta}^r(\cdot \mid q_i), \tilde{p}_{\theta}^r(\cdot \mid q_j)). \quad (3)$$

**Assumption 4.7** (Connected Graph). Let G = (Q, E) be a directed graph capturing which queries interact over the rounds. An edge  $(q_i, q_j) \in E$  indicates that in some round r, the distribution  $\tilde{p}_{\theta}^{r+1}(\cdot \mid q_i)$  was updated using  $\tilde{p}_{\theta}^{r}(\cdot \mid q_j)$ . We assume G is connected over R rounds (large R).

Assumption 4.8 (Expected Pairwise Contractive Update Property). Under perfect knowledge assumption (Scenario 1 in Argument 4.6) and Assumption 4.7, let  $\tilde{p}^r_{\theta}(z \mid q_j)$  be the belief distribution for query  $q_j$  at round r. We have  $\tilde{p}^r_{\theta}(z \mid q_j) = \tilde{p}_{\theta}(z \mid q^*) + \epsilon^r_j(z)$ . Let R be the total number of update rounds. Suppose there exists a finite constant C (independent of R) such that  $at \ most \ C$  rounds are "bad"—in these bad rounds, the expected divergence from the true distribution may increase or fail to decrease. In all other ("good") rounds  $r \notin \{\text{bad}\}$ , when  $D(\tilde{p}_{\theta}(\cdot \mid q^*), \ \tilde{p}^r_{\theta}(\cdot \mid q_j)) \neq 0$ , the expected update is contractive with respect to the true distribution:

$$\mathbb{E}\big[D\big(\tilde{p}_{\theta}(\cdot \mid q^*), \ \tilde{p}_{\theta}^{r+1}(\cdot \mid q_j)\big)\big] < \\ \mathbb{E}\big[D\big(\tilde{p}_{\theta}(\cdot \mid q^*), \ \tilde{p}_{\theta}^{r}(\cdot \mid q_j)\big)\big].$$
(4)

Here,  $D(\cdot, \cdot)$  denotes a suitable divergence measure (e.g., KL divergence). We validate this assumption in §A.4. Fig. 6 also shows that each

interaction round aligns the results more closely with the true distribution, improving performance.

**Theorem 4.9** (Convergence to True Distribution When Known). Assume the model knows the answer to the target query, then the final weighted distribution satisfies:

$$\lim_{\substack{n \to \infty \\ r \to \infty}} \sum_{j=0}^{n} w_{j} \tilde{p}_{\theta}^{r}(z \mid q_{j}) = \tilde{p}_{\theta}(z \mid q^{*}),$$

$$\forall z \in \mathcal{Y}_{target}.$$
(5)

where  $\sum_{j=0}^n w_j = 1$ . The detailed proof sketch can be referred to in §A.5. We will first use Assumption 4.8 to prove that we have  $\lim_{r\to\infty} \mathbb{E}[\epsilon_j^r(z)] = 0$  and each  $\epsilon_j^r(z)$  has a bounded variance, where  $\mathbb{E}[\epsilon_j^r(z)]$  is the *expected* error of that specific answer z at round r for query j. Since the queries are sampled diverse,  $\epsilon_j^r(z)$  are weak dependent, putting these conditions together with the law of large number,

$$\lim_{\substack{n \to \infty \\ r \to \infty}} \sum_{j=0}^{n} w_{j} \tilde{p}_{\theta}^{r}(z \mid q_{j}) = \sum_{j=0}^{n} w_{j} \tilde{p}_{\theta}(z \mid q^{*}) + \sum_{j=0}^{n} w_{j} \epsilon_{j}^{r}(z).$$

$$\underbrace{\sum_{j=0}^{n} w_{j} \tilde{p}_{\theta}(z \mid q^{*})}_{=1} + \underbrace{\sum_{j=0}^{n} w_{j} \epsilon_{j}^{r}(z)}_{\to 0}.$$
(6)

# 4.4 Proposed Metric of DIVERSEAGENTENTROPY

Building on the aforementioned theory, we propose DIVERSEAGENTENTROPY as an empirical implementation of the multi-agent interaction and corresponding uncertainty estimation. An illustrative example is shown in Fig. 2.

**Step 1: Diverse Query Generation.** Let  $Q = \{q_1, \ldots, q_n\}$  be a set of *n queries* obtained via *knowledge-preserving perturbations*  $\{T_i\}$  of a target query q and  $q_1 = q$ . The set of Q will include the target query q itself, semantically equivalent questions of the target query, and questions about different perspectives as shown in Fig. 2.

Step 2: Multi-Agent Interaction. We create n agents from the same tested model, where each agent  $A_j$  independently answers  $q_j$ . We then conduct controlled one-on-one cross-play interactions using fixed prompts (Fig. 2). To ensure practicality, each query is answered once per interaction, and we approximate the belief distribution using the top-1 (most probable) distribution. In each round, agent  $A_j$  selects a peer whose answer to the target query differs for interaction, prioritizing agents

it has not yet interacted with. During interaction,  $A_j$  sees its own conversation history and is given the peer agent's query and most recent response. Based on this,  $A_j$  decides whether to retain, or revise from its answer. The interaction ends when: (1) all agents agree on an answer, (2) each agent maintains its answer for two consecutive rounds, or (3) the maximum of  $R^*$  rounds is reached.

Step 3: Uncertainty Score Calculation. In practice, with limited interaction rounds, agents vary in response reliability. We assign each agent  $A_j$  a weight  $w_j$  in the final probability based on how often it changes its answer to the target query. Frequent changes indicate lower reliability and result in a lower weight.

$$w_j = \frac{R - r_j + 1}{\sum_{j=1}^{n} (R - r_j + 1)}. (7)$$

where j=1,...,n. We denote R as the final total number of interaction rounds and  $r_j$  as the number of rounds where the agent  $A_j$  changes its answer during the interaction. We apply Laplace smoothing to avoid zero weight. Denote by  $\mathbb{1}\{A_j=z_i\}$  whether  $z_i$  is the final answer to the target query of the agent  $A_j$  after the interaction. Therefore,

$$\tilde{p}_{\theta}(z_i \mid q^*) = \sum_{j=1}^n w_j \mathbb{1}\{A_j = z_i\}.$$
 (8)

We then apply Eq.1 with Eq.8 to calculate the final uncertainty as DIVERSEAGENTENTROPY. We present all the prompts for the method in §A.13.

**Step 4: Score-based Abstention Policy.** We introduce an abstention policy with a threshold parameter: the model abstains if the uncertainty score exceeds the threshold (see §A.8 for method variants) or if the mode response is "I don't know." If the policy does not abstain, the answer with the highest estimated probability is returned.

## 5 Experiment

## 5.1 Experiment Setting

**Evaluation Models.** We evaluate on Llama-3-70b-Instruct (AI@Meta, 2024) and Claude-3-Sonnet (Anthropic, 2024) through API calls.

**Datasets.** See §A.7 for a detailed description of the datasets. **Entity-centric QA**: we randomly sample from PopQA (Mallen et al., 2023) for 1) PopQA popular with popular entities and 2) PopQA less popular with less popular entities. **False assumption QA**: 3) FalseQA (Hu et al., 2023). All questions in the dataset contain false assumptions.

**General QA**: 4) TruthfulQA (Lin et al., 2022). 5) FreshQA (Vu et al., 2023).

**Metrics.** We evaluate entropy-based uncertainty methods using AUROC, which measures how well uncertainty scores distinguish correct from incorrect outputs. A higher AUROC indicates better performance (1.0 is perfect; 0.5 is random chance).

To assess whether abstention reduces hallucinations, we adopt the setup from Lin et al. (2022), evaluating if accuracy improves when the model chooses to answer. We report: **accuracy**, the percentage of correct responses among the questions where the model does not abstain; **abstention rate**, the percentage of questions where the model abstains; **correctness score**, the percentage of correct responses across all questions; and **truthfulness score**, the percentage of correct or abstained responses across all questions. We further analyze the accuracy–recall (AR) trade-off across methods and datasets, where recall is defined as recall = 1 — abstention rate.

**Baselines.** We first adopt a verbalization baseline, P(True) (Madaan et al., 2023). We then adopt five entropy-based black-box uncertainty estimation baselines in §3 and the model is prompted to answer 5 times. 2) Self-consistency with SemanticEntropy (SC SE) (Farquhar et al., 2024). Three baselines with affinity graph (Lin et al., 2024): 2) SC with Eccentricity (SC Ecc). 3) SC with the Degree Matrix (SC Degree). 4) SC with Eigenvalues (SC EigV). 5) SC with kernel and von Neumann entropy (SC KLE) (Nikitin et al., 2024).

We adopt seven baselines for hallucination detection. Greedy-based baselines: 1) Greedy. 2) Self-Refine (Madaan et al., 2023). 3) Self-evaluation w many samples (Kadavath et al., 2022). A total of five answers, including the greedy one, are generated, after which the model assesses the greedy sample's validity. 4) Multiple-Recite (Sun et al., 2023). The model generates multiple related paragraphs from its parametric knowledge before answering the question. Sampling-based baselines: 5) Self-consistency (SC) (Wang et al., 2023): the model answers the query 5 times, and we accept the majority answer or abstain if no answer appears at least 3 times. 6) Consistency with semantically equivalent questions (SeQ) (Zhang et al., 2023; Zhao et al., 2024): the model answers 5 semantically equivalent questions about the same target query. 7) Consistency with diverse questions (DiverseQ): the model answers 5 diverse queries about the same target query. Note that we evaluate the se-

Dataset	Claude-3-Sonnet	Llama-3-70b-Instruct
PopQA pop	0.114	0.118
PopQA less pop	0.193	0.207
FalseQA	0.154	0.154
TruthfulQA	0.296	0.330
FreshQA	0.167	0.175

Table 1: The average percentage of incorrect answers to the target query in the first round without agent interaction, evaluated only on instances where all agents agree on the correct answer after interactions.

mantic equivalence of answers to cluster responses for all sampling-based baselines.

**DIVERSEAGENTENTROPY Variants.** We introduce two variants with 5 agents: 1) Agent<sub>L</sub>: abstains if the uncertainty score exceeds the entropy of three answers with probability (0.6, 0.2, 0.2). 2) Agent: a stricter vote, abstains if the score exceeds the entropy of two answers (0.6, 0.4). Details are explained in §A.8.

# **5.2** Evaluation for Retrieval of Parametric Knowledge

Unsatisfying retrieval of parametric knowledge. LLMs often fail to reliably and consistently retrieve encoded knowledge. We show that even when the model knows the correct answer, it may initially respond inconsistently across contexts. We examine cases where all agents converge on the same ground truth answer after interaction, signaling that the knowledge is present but not reliably accessed. We conduct a quantitative analysis of the model's initial performance by measuring the average error rate on target queries in the first round before agent interaction. This captures how often the model fails to retrieve the correct answer initially, before any interaction. As shown in Table 1, LLMs are not consistently reliable, with at least one in five answers incorrect on average (19.1%). The issue is more pronounced on long-tailed (PopQA less popular) and realistic datasets (FreshQA, TruthfulQA). Models' susceptibility to contextual biases results in the poor retrieval. We conduct a qualitative analysis by sampling 45 instances from the same pool. We observe that the model tends to generate a different response, even when it knows the correct answer to the target knowledge, under these conditions: 1) 42% of cases occur when changes in the varied query significantly diverges in wording, 2) 22% when an incorrect answer is more popular within the context of the target knowledge, and 3) 20% when the additional context is more closely

Model	PopQA_p	PopQA_l	FalseQA	FreshQA	TruthfulQA	All				
	Claude-3-Sonnet									
P (True)	0.686	0.660	0.800	0.719	0.556	0.766				
SC (Ecc)	0.671	0.821	0.711	0.702	0.548	0.766				
SC (Degree)	0.674	0.855	0.713	0.704	0.550	0.771				
SC (EigV)	0.673	0.851	0.713	0.703	0.550	0.771				
SC (KLE)	0.684	0.883	0.756	0.692	0.563	0.791				
SC (SE)	0.693	0.887	0.753	0.694	0.568	0.792				
Agent	0.725	0.947	0.802	0.836	0.624	0.833				
		Llama-3	-70b-Insti	ruct						
P (True)	0.618	0.641	0.698	0.642	0.542	0.601				
SC (Ecc)	0.594	0.716	0.628	0.660	0.488	0.644				
SC (Degree)	0.595	0.704	0.629	0.662	0.486	0.645				
SC (EigV)	0.595	0.707	0.629	0.664	0.486	0.645				
SC (KLE)	0.632	0.738	0.664	0.632	0.537	0.694				
SC (SE)	0.624	0.737	0.673	0.632	0.545	0.694				
Agent	0.651	0.753	0.673	0.697	0.592	0.713				

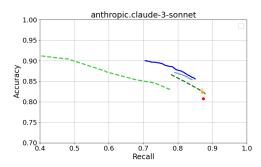
Table 2: Comparison of AUROC scores across different QA datasets. We highlight the **best** and suboptimal methods.Our method is more calibrated.

associated with a different possible answer to the target knowledge. Examples of each scenario are provided in §A.9. These findings highlight the need for systematic research into contextual biases.

## 5.3 Evaluation of DIVERSEAGENTENTROPY

than self-consistency-based uncertainty estimation. Table 2 reports AUROC scores across different methods. Our method achieves the highest AUROC, indicating better calibration. Notably, it shows larger gains on long-tailed data (+3.8% average on less popular PopQA) and realistic benchmarks (+7.5% on FreshQA, +5.2% on TruthfulQA) compared to the best baseline, with especially strong performance on the more advanced Claude model. These results highlight DIVERSEAGENTENTROPY's promise for highstakes scenarios and scalable oversight.

**DIVERSEAGENTENTROPY-based** abstention policy effectively detects hallucinations. Table 3 shows that our method more effectively identifies model hallucinations. It achieves the highest correctness and truthfulness scores. Dataset-level results are provided in §A.10. Fig. 3 presents Accuracy-Recall (AR)-curves across all data with per-dataset results in Appendix Fig. 9. Our method consistently outperforms all baselines, achieving the highest accuracy at comparable recall levels on non-abstained answers (+ 2.5% average increase). Appendix Table 6 shows that on average, 18.0%/ 27.0% of incorrect instances were completely corrected (Correct), while only 2.5%/8.9% of correct instances became partially incorrect (Wrong) for Claude/Llama. This highlights that



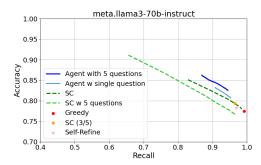


Figure 3: AR-curves across all data. SC refers to SC (SE). SC w 5 questions refers to calculating entropy using the agents' diverse questions without agent interaction. We present the <u>best</u> and <u>sub optimal</u> methods for each method category in Tables 2 and 3.

		Claud	e-3-Sonn	net	L	lama-3	-70b-Ins	truct
Method	Acc	Ab-R	Correct	TruthF	Acc	Ab-R	Correct	TruthF
Greedy	0.808	0.126	0.707	0.832	0.775	0.008	0.769	0.777
Self-Refine	0.826	0.131	0.718	0.849	0.783	0.030	0.760	0.790
Self-Eval	0.814	0.141	0.700	0.840	0.754	0.020	0.739	0.759
MultiRecite	0.779	0.114	0.690	0.804	0.715	0.010	0.708	0.717
SC (3/5)	0.823	0.129	0.717	0.846	0.794	0.035	0.766	0.801
SeQ	0.815	0.149	0.693	0.842	0.818	0.084	0.749	0.833
DiverseQ	0.858	0.342	0.564	0.906	0.811	0.121	0.713	0.834
Agent_L	0.852	0.142	0.731	0.873	0.826	0.055	0.780	0.835
Agent	0.883	0.216	0.692	0.908	0.841	0.084	0.770	0.854

Table 3: Performance evaluation of different methods on all data with the best setup for each method. Acc is accuracy. Ab-R is abstention rate. Correct is correctness score. TruthF is truthfulness score.

DIVERSEAGENTENTROPY significantly improves accuracy with minimal degradation and detailed analysis is in §A.10.

## 5.4 Analysis of DIVERSEAGENTENTROPY

Both diverse query generation and agent interaction are key for performance boost. In Fig. 3, we compare DIVERSEAGENTENTROPY with two ablations: using self-consistency with 5 questions without agent interaction (SC with 5 questions) and agent interaction with only the target query (Agent w single question). The results show agent interaction and diverse query generation are important.

The quality of the diverse queries. We evaluate question generation quality using Llama-3.1-70b-Instruct by sampling 500 queries per model (20 instances per dataset). Queries are assessed for 1) knowledge preservation, whether they required answering the target query, and 2) diversity, whether any two queries from the same instance differed by at least 20% in wording. Claude and Llama achieve 95.0% and 95.6% in knowledge preservation, and 97% and 96% in diversity, respectively. These results demonstrate high-quality and varied question

generation across models. Details in §A.11.

The number of agents. In Appendix Fig. 5 and Appendix Fig. 10, we increase the number of agents, limiting interactions to 4 rounds. Performance improves with more agents but gains become smaller beyond 4 agents for Llama model, suggesting 5 agents are sufficient practically.

The rounds of interactions. As shown in Appendix Fig. 6 and Appendix Fig. 11, with the number of agents fixed at 5, increasing the rounds of interaction generally leads to improved performance. Format of agent interaction. We examine whether agents should engage in one-on-one interaction or group interactions, where each agent can view the unique questions and answers of all other agents. Our findings in Appendix Fig. 7 and Appendix Fig. 12 indicate that one-on-one interactions outperform group interactions. We present a human error analysis in § A.12.

Robustness of agent interaction. We evaluate agent susceptibility to misleading signals when one agent consistently gives either the most plausible incorrect answer or repeatedly responds with "I don't know". Appendix Fig. 8 and Fig. 13 show that performance slightly declines in both cases, indicating sensitivity to consistent misleading input, while a better model, e.g., Claude, is more robust.

## 6 Conclusion

Accurately determining the uncertainty of black-box LLMs is vital for reliable responses. We show that existing methods fall short because models suffer from suboptimal retrieval of parametric knowledge due to susceptibility to contextual bias. In this paper, we propose a novel method with a theoretical framework, DIVERSEAGENTENTROPY, to quantify an LLM's uncertainty after multi-agent interaction. Our method overcomes the limitations of

self-consistency-based uncertainty estimation and delivers superior performance in detecting hallucinations.

### Limitations

We acknowledge that our method is more resourceintensive, primarily due to the use of multi-agent interactions, which, like similar approaches in the literature, naturally incur higher computational costs. A detailed analysis is provided in §A.6. Importantly, our method more accurately recovers the model's true uncertainty, a crucial advantage especially in high-stakes applications where reliability and correctness are prioritized.

Our proof is developed under relatively strict assumptions, which are primarily intended to provide intuition and can be relaxed or extended in future work.

We only evaluate two representative models with different base capacities due to the budget constraint, given the scale of our ablation studies. Although not the latest versions, they offer a meaningful comparison and consistently demonstrate the same trends observed across our ablation studies.

While prior work on LLM uncertainty estimation largely adopts the same evaluation setting as ours to analyze short-form, query-based uncertainty, we recognize the need to explore more complex scenarios in future work. We view the current setting as a necessary foundation for addressing more challenging tasks that require specialized study beyond the scope of this paper, such as: (1) long-form QA, which is difficult to evaluate due to challenges in assessing semantic equivalence and the limitations of current evaluation techniques, particularly in black-box settings; and (2) response-based uncertainty (Qiu and Miikkulainen, 2024), which builds upon query-based uncertainty in black-box scenarios and thus depends on its accurate estimation (with the exception of verbalization, which is often unreliable). Moreover, quantifying response-based uncertainty becomes more difficult when querybased uncertainty is high, indicating that the model lacks reliable knowledge about the query, in which case none of the candidate responses should be trusted.

For further direction of broader impact, we believe our method has the potential to generalize to retrieval-augmented generation (RAG) settings. Its ability to generate questions from diverse perspectives and leverage multi-agent interactions makes

it well-suited for effectively integrating heterogeneous information sources. In the RAG context, generating perspective-varied questions encourages the retriever to surface a broader range of documents, thereby enriching the informational basis for downstream tasks. Moreover, the multi-agent interaction mechanism enables self-reflection and deliberation, which is particularly valuable in QA scenarios involving ambiguous or conflicting retrieved content.

## Acknowledgments

This work was conducted during Yu Feng's internship at Amazon Web Services (AWS) and was fully funded by AWS. We sincerely appreciate the valuable feedback provided by the AWS team and extend our gratitude to the ARR reviewers and editors for their insightful and constructive comments.

### References

Moloud Abdar, Farhad Pourpanah, Sadiq Hussain, Dana Rezazadegan, Li Liu, Mohammad Ghavamzadeh, Paul Fieguth, Xiaochun Cao, Abbas Khosravi, U Rajendra Acharya, and 1 others. 2021. A review of uncertainty quantification in deep learning: Techniques, applications and challenges. *Information fusion*, 76:243–297.

Lukas Aichberger, Kajetan Schweighofer, Mykyta Ielanskyi, and Sepp Hochreiter. 2024. How many opinions does your LLM have? improving uncertainty estimation in NLG. In *ICLR 2024 Workshop on Secure and Trustworthy Large Language Models*.

AI@Meta. 2024. Llama 3 model card.

Anthropic. 2024. The claude 3 model family: Opus, sonnet, haiku.

Samuel R. Bowman, Jeeyoon Hyun, Ethan Perez, Edwin Chen, Craig Pettit, Scott Heiner, Kamilė Lukošiūtė, Amanda Askell, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Christopher Olah, Daniela Amodei, Dario Amodei, Dawn Drain, Dustin Li, Eli Tran-Johnson, and 27 others. 2022. Measuring progress on scalable oversight for large language models. *Preprint*, arXiv:2211.03540.

Sky CH-Wang, Benjamin Van Durme, Jason Eisner, and Chris Kedzie. 2024. Do androids know they're only dreaming of electric sheep? In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 4401–4420, Bangkok, Thailand. Association for Computational Linguistics.

Angelica Chen, Jason Phang, Alicia Parrish, Vishakh Padmakumar, Chen Zhao, Samuel R. Bowman, and

- Kyunghyun Cho. 2024a. Two failures of self-consistency in the multi-step reasoning of LLMs. *Transactions on Machine Learning Research*.
- Chao Chen, Kai Liu, Ze Chen, Yi Gu, Yue Wu, Mingyuan Tao, Zhihang Fu, and Jieping Ye. 2024b. INSIDE: LLMs' internal states retain the power of hallucination detection. In *The Twelfth International Conference on Learning Representations*.
- Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. 2024. Improving factuality and reasoning in language models through multiagent debate.
- Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. 2024. Detecting hallucinations in large language models using semantic entropy. In *Nature* 630, 625–630.
- Shangbin Feng, Weijia Shi, Yike Wang, Wenxuan Ding, Vidhisha Balachandran, and Yulia Tsvetkov. 2024. Don't hallucinate, abstain: Identifying llm knowledge gaps via multi-llm collaboration. *Preprint*, arXiv:2402.00367.
- Hila Gonen, Srini Iyer, Terra Blevins, Noah Smith, and Luke Zettlemoyer. 2023. Demystifying prompts in language models via perplexity estimation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10136–10148, Singapore. Association for Computational Linguistics.
- Yufei Guo, Muzhe Guo, Juntao Su, Zhou Yang, Mengqiu Zhu, Hongfei Li, Mengyang Qiu, and Shuo Shuo Liu. 2024. Bias in large language models: Origin, evaluation, and mitigation. *Preprint*, arXiv:2411.10915.
- Shengding Hu, Yifan Luo, Huadong Wang, Xingyi Cheng, Zhiyuan Liu, and Maosong Sun. 2023. Won't get fooled again: Answering questions with false premises. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pages 5626–5643, Toronto, Canada. Association for Computational Linguistics.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. 55(12).
- 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, and 17 others. 2022. Language models (mostly) know what they know. *Preprint*, arXiv:2207.05221.
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. In *The Eleventh International Conference on Learning Representations*.

- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. TruthfulQA: Measuring how models mimic human falsehoods. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics* (*Volume 1: Long Papers*), pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.
- Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2024. Generating with confidence: Uncertainty quantification for black-box large language models. *Preprint*, arXiv:2305.19187.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. Self-refine: Iterative refinement with self-feedback. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Andrey Malinin and Mark Gales. 2021. Uncertainty estimation in autoregressive structured prediction. In *International Conference on Learning Representations*.
- Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9802–9822, Toronto, Canada. Association for Computational Linguistics.
- Potsawee Manakul, Adian Liusie, and Mark Gales. 2023. SelfCheckGPT: Zero-resource black-box hallucination detection for generative large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9004–9017, Singapore. Association for Computational Linguistics.
- Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. 2025. s1: Simple test-time scaling. *Preprint*, arXiv:2501.19393.
- Navapat Nananukul and Mayank Kejriwal. 2024. Halo: An ontology for representing and categorizing hallucinations in large language models. *Preprint*, arXiv:2312.05209.
- Alexander V Nikitin, Jannik Kossen, Yarin Gal, and Pekka Marttinen. 2024. Kernel language entropy: Fine-grained uncertainty quantification for LLMs from semantic similarities. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin,

- Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, and 262 others. 2024. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.
- Xin Qiu and Risto Miikkulainen. 2024. Semantic density: Uncertainty quantification for large language models through confidence measurement in semantic space. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Abhilasha Ravichander, Shrusti Ghela, David Wadden, and Yejin Choi. 2025. Halogen: Fantastic Ilm hallucinations and where to find them. *Preprint*, arXiv:2501.08292.
- Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2024. Quantifying language models' sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. In *The Twelfth International Conference on Learning Representations*.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik R Narasimhan, and Shunyu Yao. 2023. Reflexion: language agents with verbal reinforcement learning. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Ola Shorinwa, Zhiting Mei, Justin Lidard, Allen Z. Ren, and Anirudha Majumdar. 2024. A survey on uncertainty quantification of large language models: Taxonomy, open research challenges, and future directions. *Preprint*, arXiv:2412.05563.
- Charlie Victor Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2025. Scaling LLM test-time compute optimally can be more effective than scaling parameters for reasoning. In *The Thirteenth International Conference on Learning Representations*.
- Zhiqing Sun, Xuezhi Wang, Yi Tay, Yiming Yang, and Denny Zhou. 2023. Recitation-augmented language models. In *The Eleventh International Conference on Learning Representations*.
- Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher Manning. 2023. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5433–5442, Singapore. Association for Computational Linguistics.
- Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny Zhou, Quoc Le, and Thang Luong. 2023. Freshllms: Refreshing large language models with search engine augmentation. *Preprint*, arXiv:2310.03214.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models. In

- The Eleventh International Conference on Learning Representations.
- Yifei Wang, Yuyang Wu, Zeming Wei, Stefanie Jegelka, and Yisen Wang. 2024. A theoretical understanding of self-correction through in-context alignment. In *ICML 2024 Workshop on In-Context Learning*.
- J Florian Wellmann and Klaus Regenauer-Lieb. 2012. Uncertainties have a meaning: Information entropy as a quality measure for 3-d geological models. *Tectonophysics*, 526:207–216.
- Kai Xiong, Xiao Ding, Yixin Cao, Ting Liu, and Bing Qin. 2023. Examining inter-consistency of large language models collaboration: An in-depth analysis via debate. In *Findings of the Association for Computa*tional Linguistics: EMNLP 2023, pages 7572–7590, Singapore. Association for Computational Linguistics
- Miao Xiong, Zhiyuan Hu, Xinyang Lu, YIFEI LI, Jie Fu, Junxian He, and Bryan Hooi. 2024. Can LLMs express their uncertainty? an empirical evaluation of confidence elicitation in LLMs. In *The Twelfth International Conference on Learning Representations*.
- Yasin Abbasi Yadkori, Ilja Kuzborskij, András György, and Csaba Szepesvári. 2024. To believe or not to believe your llm. *Preprint*, arXiv:2406.02543.
- Jifan Yu, Xiaozhi Wang, Shangqing Tu, Shulin Cao, Daniel Zhang-Li, Xin Lv, Hao Peng, Zijun Yao, Xiaohan Zhang, Hanming Li, Chunyang Li, Zheyuan Zhang, Yushi Bai, Yantao Liu, Amy Xin, Kaifeng Yun, Linlu GONG, Nianyi Lin, Jianhui Chen, and 16 others. 2024. KoLA: Carefully benchmarking world knowledge of large language models. In *The Twelfth International Conference on Learning Representations*.
- Jiaxin Zhang, Zhuohang Li, Kamalika Das, Bradley Malin, and Sricharan Kumar. 2023. SAC<sup>3</sup>: Reliable hallucination detection in black-box language models via semantic-aware cross-check consistency. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 15445–15458, Singapore. Association for Computational Linguistics.
- Yukun Zhao, Lingyong Yan, Weiwei Sun, Guoliang Xing, Chong Meng, Shuaiqiang Wang, Zhicong Cheng, Zhaochun Ren, and Dawei Yin. 2024. Knowing what LLMs DO NOT know: A simple yet effective self-detection method. 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 7051–7063, Mexico City, Mexico. Association for Computational Linguistics.

## A Appendix

## A.1 Existing Self-consistency Based Uncertainty Quantification Methods

Existing self-consistency-based uncertainty estimation methods in the black-box setting (Kuhn et al., 2023; Farquhar et al., 2024; Lin et al., 2024; Aichberger et al., 2024; Nikitin et al., 2024) follow a similar procedure: 1) For a given input q, generate N response samples. 2) Calculate similarity scores for these N responses. 3) Compute an uncertainty estimate U(q) using the similarity values.

Specifically, Farquhar et al. (2024) introduces semantic entropy to calculate  $p(y_i|q)$  in Eq.1 as a frequency-based probability derived from repeated sampling of the target query q. Assume we find the semantic clusters for the sampled answers, and let each query return a possible semantically different answer  $y_i \in C$ . The count of times a particular  $y_i$  appears as the output for the input q over these N queries is denoted as  $c(y_i)$ . Therefore,  $p(y_i|q) = \frac{c(y_i)}{r}$ .

Lin et al. (2024) calculates uncertainty using a weighted adjacency graph built upon semantic affinities. An affinity model e maps pairs of responses to values in [0,1]. Given n independent samples, the model induces a symmetric adjacency matrix  $W = [w_{i,j}]_{i,j=1}^n$ , where  $w_{i,j}$  is the mean of the pairwise similarities between response i and j. The degree matrix is  $D = [\mathbbm{1}[j=i]\sum_{a=1}^n w_{n,j}]_{i,j=1}^n$ , and the Laplacian  $L = I - D^{-1/2}WD^{-1/2}$  has eigenvalues  $\{\lambda_n\}_{a=1}^n$ . The following uncertainty measures are then defined:  $U_{EigV}(x) = \sum_{a=1}^n \max\{0, 1-\lambda_a\}, U_{Degree}(x) = 1 - \frac{trace(D)}{n^2}, U_{Ecc}(x) = \|[v_1, v_2, \ldots, v_n]\|_2$  where  $\{v_a\}_{a=1}^n$  are vectors associated with L.

Nikitin et al. (2024) proposes Kernel Language Entropy (KLE), which quantifies uncertainty using the von Neumann entropy. After generating N response samples, KLE constructs a positive semidefinite, unit-trace kernel matrix K that encodes the semantic similarities between these responses. This kernel captures fine-grained semantic dependencies. The uncertainty is then calculated as the von Neumann entropy of this kernel matrix, given by  $U_{\text{KLE}}(q) = -\text{Tr}(K \log K)$  where Tr denotes the matrix trace and the logarithm of the matrix is defined through its eigendecomposition.

# A.2 Observation of Agent Behaviors When Knowledge is Unknown

We conduct a pilot experiment using DIVERSEAGENTENTROPY to observe agent behaviors for the scenario when the model doesn't know the answer and therefore the model's true distribution  $\tilde{p}_{\theta}(\cdot \mid q^*)$  does not exist. We randomly sample 50 false-assumption queries where the model definitely doesn't know the answer from FalseQA and freshQA. We observe three distinct behaviors after agent interaction<sup>4</sup>:

- 1. Arbitrary Shifts and Inconsistency (36%) Even after many interaction rounds, queries  $q_i$  and  $q_j$  maintain *inconsistent* distributions  $\tilde{p}_{\theta}^{R}(\cdot \mid q_i)$  and  $\tilde{p}_{\theta}^{R}(\cdot \mid q_j)$ , disagreeing about which answer is most plausible. Instead, they *shift arbitrarily* depending on interaction's context.
- 2. Trivial Convergence to "I don't know"(54%). After several rounds of disagreement, all queries  $q_i$  converge to "I don't know"(or similarly agnostic) answers, indicating the model's inability to pin down the knowledge.
- 3. Convergence to An Arbitrarily Wrong Answer Distribution(10%). After several rounds of interaction, a query's strong belief may enforce all the queries to converge to an arbitrarily wrong answer distribution.

# A.3 Implementation of Diverse Query Generation

Given a target query x, we use the same model to be tested to generate varied knowledge-perserving queries that require knowledge of the target query, ensuring that these questions are both representative and comprehensive. The question-generation process is completely automated and the detailed question-generation prompts can be found in  $\S A.13$ . Specifically, we first conceptualize the target query and then sample various perspectives to ensure a comprehensive understanding. For each perspective, we generate m questions that build upon the target query, tailored to that particular perspective. We filter these generated questions to ensure they strictly require knowledge of the target query to answer while avoiding the inclusion of the direct

<sup>&</sup>lt;sup>4</sup>Note that the behaviors are consistent on both Claude and Llama.

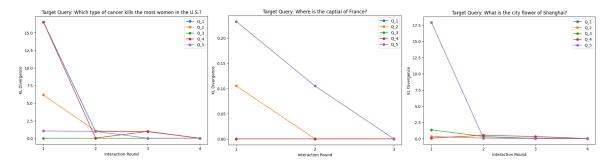


Figure 4: KL-Divergence with regard to interaction rounds for three examples where the model knows the answer.

answer. We also generate m semantically equivalent questions for the target query.

We select n questions from the generated pool to form the final candidate set Q for the agents. In order to maximize diversity, the optimal set includes the target query x, one semantically equivalent question, and n-2 questions each targeting a unique perspective. If there are insufficient unique perspectives with qualified questions, we repeat the perspective question selection process to select from existing perspectives. If not, we supplement with additional semantically equivalent questions.

### A.4 Analysis of Assumption 4.8

We now give a detailed explanation for the assumption. Interacting queries provide diverse contextual hints through varying background knowledge (Sun et al., 2023). Transformer-based models, as shown by (Wang et al., 2024), perform in-context gradient descent to enable dynamic self-correction. If the model knows the answer, each update is expected to reduce the gap between  $\tilde{p}^r$  and the true distribution  $\tilde{p}^*$ . While some bad updates may strongly increase divergence, the overall trend is an improvement. In later rounds, under Assumption 4.7, connectivity allows queries to influence one another, propagating and averaging information. This prevents errors from persisting or amplifying (Du et al., 2024). Thus, over R rounds, at most a constant number Cof updates may fail to reduce divergence.

We test whether the agent's interactions guide the KL divergence toward convergence with the (approximated) true distribution when the model knows the answer. We select 45 examples for evaluation from Experiment 5.2. We allow 4 rounds of interactions. For each interaction round, we sample 10 times with a temperature of 0.7 to estimate the agent distribution. Manual inspection of the interaction update trend figures shows that 89% of these examples consistently converge toward the

true distribution. We also test interactions using the top-1 distribution, where all examples converge to the ground truth answer. To better illustrate this behavior, we present three representative examples in Appendix Figure 4.

For query (left): What type of cancer kills the most women in the U.S.? Q1: What type of cancer kills the most women in the U.S.? Q2: Which organ does the cancer that kills the most women in the U.S. affect? Q3: List the types of cancer and their impact on women in the U.S. based on mortality rates Q4: What public health initiatives exist to prevent and treat the type of cancer kills the most women in the U.S.? Q5: What's the most common symptom of the deadliest cancer for women in the U.S.? Gold distribution for Llama-3-70b-Instruct (approximated): Lung cancer: 0.75, Breast cancer: 0.2, Colorectal Cancer: 0.05.

For query (middle): Where is the capital of France? Q1: Where is the capital of France? Q2: Which is the capital city of France? Q3: How does the capital of France influence its Gen Z culture? Q4: How has the location of the capital of France influenced its technological advancements? Q5: How does the history of France shape its capital city? Gold distribution for Claude-3-sonnet (approximated): Paris: 1.

For query (right): What is the city flower of Shanghai? Q1: What is the city flower of Shanghai? Q2: Which bloom is honored as the city flower of Shanghai? Q3: What role does the city flower of Shanghai play in the city's efforts to promote ecological conservation and environmental sustainability? Q4: In what ways has the city flower of Shanghai been incorporated into the city's branding and tourism marketing strategies? Q5: How does the city flower of Shanghai reflect the city's historical and cultural connection to the Yangtze River Delta region? Gold distribution for Llama-3-70b-Instruct (approximated): White Magnolia: 0.8,

White Peony: 0.2.

## A.5 Theorem Proof Sketch for Theorem 4.9

Since<sup>5</sup>  $\tilde{p}_{\theta}^{r}(z \mid q_{j}) = \tilde{p}_{\theta}(z \mid q^{*}) + \epsilon_{j}^{r}(z)$ , the KL divergence becomes:

$$D_{\text{KL}}(\tilde{p}_{\theta}(\cdot \mid q^*) \parallel \tilde{p}_{\theta}^{r}(\cdot \mid q_j)) = \sum_{z \in \mathcal{Y}_{target}} \tilde{p}_{\theta}(z \mid q^*) \log \left( \frac{\tilde{p}_{\theta}(z \mid q^*)}{\tilde{p}_{\theta}^{r}(z \mid q_j)} \right). \tag{9}$$

We can further decompose the KL divergence into:

$$D_{KL}(\tilde{p}_{\theta}(\cdot \mid q^{*}) \parallel \tilde{p}_{\theta}^{r}(\cdot \mid q_{j})) = \sum_{z} \tilde{p}_{\theta}(z \mid q^{*}) \log \tilde{p}_{\theta}(z \mid q^{*}) - \sum_{z} \tilde{p}_{\theta}(z \mid q^{*}) \log \tilde{p}_{\theta}^{r}(z \mid q_{j}).$$

$$(10)$$

For the logarithmic term  $\log (\tilde{p}_{\theta}^{r}(z \mid q_{j}))$ , using a Taylor expansion of  $\log(x)$  around  $x_{0} = \tilde{p}_{\theta}(z \mid q^{*})$ , we have:

$$\log \left( \tilde{p}_{\theta}^{r}(z \mid q_{j}) \right) \approx \log \left( \tilde{p}_{\theta}(z \mid q^{*}) \right) + \frac{\epsilon_{j}^{r}(z)}{\tilde{p}_{\theta}(z \mid q^{*})} - \frac{1}{2} \left( \frac{\epsilon_{j}^{r}(z)}{\tilde{p}_{\theta}(z \mid q^{*})} \right)^{2}.$$
(11)

We only assume the second-order approximation and therefore,

$$\mathbb{E}[D_{KL}(\tilde{p}_{\theta}(\cdot \mid q^*) \parallel \tilde{p}_{\theta}^r(\cdot \mid q_j))] \approx \frac{1}{2} \sum_{z} \frac{\mathbb{E}[\epsilon_j^r(z)^2]}{\tilde{p}_{\theta}(z \mid q^*)}$$

$$z \in \{z | \tilde{p}_{\theta}(z \mid q^*) \neq 0\}$$

$$\tag{12}$$

Let 
$$\Delta_j(r+1) = \mathbb{E}[D(\tilde{p}_{\theta}(\cdot \mid q^*), \ \tilde{p}_{\theta}^{r+1}(\cdot \mid q_j))] - \mathbb{E}[D(\tilde{p}_{\theta}(\cdot \mid q^*), \ \tilde{p}_{\theta}^{r}(\cdot \mid q_j))]$$
, therefore:

$$\Delta_{j}(r+1) = \sum_{z} \tilde{p}_{\theta}(z \mid q^{*}) \mathbb{E}\left[\log \tilde{p}_{\theta}^{r}(z \mid q_{j})\right]$$
$$-\sum_{z} \tilde{p}_{\theta}(z \mid q^{*}) \mathbb{E}\left[\log \tilde{p}_{\theta}^{r+1}(z \mid q_{j})\right]$$
(13)

Substituting this approximation back into the  $\Delta_j(r+1)$ , we get:

$$\Delta_{j}(r+1) = \frac{1}{2} \sum_{\tilde{r}} \frac{\mathbb{E}[\epsilon_{j}^{r+1}(z)^{2} - \epsilon_{j}^{r}(z)^{2}]}{\tilde{p}_{\theta}(z \mid q^{*})}. \quad (14)$$

Here,  $\epsilon_j^r(z)$  can be positive, negative, or zero for different rounds r. Since  $\tilde{p}_{\theta}(z \mid q^*) \in [\eta, 1]$  and

 $\begin{array}{l} \sum_{z}\tilde{p}_{\theta}(z\mid q^{*})=1,\,\epsilon_{j}^{r}(z) \text{ is bounded as } \epsilon_{j}^{r}(z)\in\\ \left[-\tilde{p}_{\theta}(z\mid q^{*}),1-\tilde{p}_{\theta}(z\mid q^{*})\right].\text{ Therefore:} \end{array}$ 

$$-\frac{|\mathcal{Z}|}{2\eta} \le \Delta_j(r+1) \le \frac{|\mathcal{Z}|}{2\eta},\tag{15}$$

We assume  $\eta \geq 0.1\%$  as we observe through common data, For r such that  $\Delta_j(r+1) \geq 0$ , the sequence increases. However, the boundedness of  $\Delta_j(r+1)$  ensures that these increases are finite for the C total rounds.

Based on preliminary results that is shown in Appendix A.4 and Figure 4, we assume for r such that  $\Delta_i(r+1) < 0$ , the sequence

$$\mathbb{E}[D_{KL}(\tilde{p}_{\theta}(\cdot \mid q^{*}) \parallel \tilde{p}_{\theta}^{r+1}(\cdot \mid q_{j}))] = e^{-\lambda_{q_{j}}} \mathbb{E}[D_{KL}(\tilde{p}_{\theta}(\cdot \mid q^{*}) \parallel \tilde{p}_{\theta}^{r}(\cdot \mid q_{j}))]$$
(16)

where  $\lambda_{q_j}$  is a constant relevant to the query  $q_j$ .

Due to the abovementioned converge rate and nearly monotonic convergence except for bounded finitely many bad steps,  $\mathbb{E}[D_{\mathrm{KL}}(\tilde{p}_{\theta}(\cdot \mid q^*) \parallel \tilde{p}^r_{\theta}(\cdot \mid q_j))]$  converges to a limit L at the unique global minimum of KL divergence as it's strictly convex, so:

$$L = 0 \implies \lim_{r \to \infty} \mathbb{E}[\epsilon_j^r(z)] = 0, \forall z \in \mathcal{Y}_{\text{target}}.$$
(17)

In terms of variance, since we have each answer's error satisfying  $\epsilon^r_j(z) \in \left[-\tilde{p}_{\theta}(z \mid q^*), 1 - \tilde{p}_{\theta}(z \mid q^*)\right]$  for every round  $r, \epsilon^r_j(z)$  is a random variable with

$$\begin{aligned} \left| \epsilon_j^r(z) \right| &\leq \max \left\{ \left| \tilde{p}_{\theta}(z \mid q^*) \right|, \left| 1 - \tilde{p}_{\theta}(z \mid q^*) \right| \right\} \\ &= M_z \quad \text{(some finite constant } M_z > 0 \text{).} \end{aligned}$$
(18)

Thus  $\epsilon_j^r(z)$  is *surely* bounded, taking values in a finite interval  $[-M_z, M_z]$ . It follows immediately that

$$\operatorname{Var}[\epsilon_j^r(z)] \leq (M_z)^2$$
 for every round  $r$ . (19) Thus  $\operatorname{Var}[\epsilon_i^r(z)]$  is bounded.

 $\epsilon_j^r(z)$  is weakly dependent due to the diversity of the queries and individual reflection traces. Putting these conditions together, the law of large number for weakly dependent processes implies  $\sum_{j=0}^n w_j \epsilon_j^r(z) \xrightarrow[n \to \infty]{\text{in probability}} 0$ . Therefore,

$$\lim_{\substack{n \to \infty \\ r \to \infty}} \sum_{j=0}^{n} w_j \tilde{p}_{\theta}^r(z \mid q_j) = \sum_{\substack{j=0 \\ j=1}}^{n} w_j \tilde{p}_{\theta}(z \mid q^*) + \sum_{\substack{j=0 \\ j=0}}^{n} w_j \epsilon_j^r(z).$$
(20)

<sup>&</sup>lt;sup>5</sup>The proof is developed under relatively strict assumptions and is primarily intended to provide intuition.

### A.6 Cost Analysis

We present a detailed cost analysis on our proposed DIVERSEAGENTENTROPY. In a self-consistency-based method, we typically sample a simple query 5 times, requiring 5 API calls. For our method, starting with a single query, we need 1 API call for question conceptualization, 1 for perspective generation, and 5 for generating questions under different perspectives. During agent interaction (with 5 agents), each agent answers their assigned question, followed by multiple interaction rounds. Assuming an average of 2.5 rounds, agent interaction requires 20 calls. In total, our method averages 25 API calls, making it 5 times more expensive than the base self-consistency method.

However, we would like to emphasize that in high-stakes applications where correctness is prioritized over cost, our calibrated uncertainty score can provide users with a reliable measure of how much they can trust the model's output. As shown in Argument 4.5 and Experiment 5.2, single-context methods like self-consistency suffer from contextual bias. Multi-context enabled by multi-agent interactions can help the model recognize and correct for these biases and recover its true internal belief. We further present the cost, i.e., the number of inference calls for all the baselines in Appendix Table 5.

Additionally, The varied queries generated by our method can also be utilized for fine-tuning or training purposes. The intermediate results generated by our method, including varied queries and the self-reflection interaction processes, can be further leveraged to create synthetic data for finetuning or training LLMs.

Future work can explore ways to maintain the same level of performance while reducing costs. This could involve using fewer but higher-quality questions from diverse perspectives and minimizing the number of interaction rounds.

### A.7 Datasets Statistics

We consider five different datasets under three categories. We present the detailed dataset statistics in Appendix Table 4. All the datasets have the license to be publicly used. **Entity-centric QA**: we randomly sample data from PopQA (Mallen et al., 2023) where each question is created by converting a knowledge tuple retrieved from Wikidata using a template. 1) PopQA popular. We sample questions where the entity's popularity is larger than

10<sup>4</sup> as the criteria in the original paper. 2) PopQA less popular. We sample questions where the entity's popularity is lower than 10<sup>4</sup>. **General QA**: 3) TruthfulQA (Lin et al., 2022). Note that not all data in TruthfulQA tests about factual questions. We choose examples only from categories: Law, Sociology, Health, History, and Language, focusing on clear facts instead of opinions. 4) FreshQA (Vu et al., 2023). We adopt the 07112024 version and select one-hop slow-changing or never-changing data points where the effective year is before 2022 to avoid the temporal influence. **False assumption QA**: 5) FalseQA (Hu et al., 2023). All questions in the dataset contain false assumptions and we remove all the WHY questions.

## A.8 Thresholds for the Abstention Policy

We variants of adopt two DIVERSEAGENTENTROPY where we have 5 agents, i.e., 5 varied querys: 1) Agent (Loose Majority Vote): We abstain when the uncertainty score exceeds the threshold, calculated as the entropy of 3 answers with probabilities of 0.6 (3/5), 0.2 (1/5), and 0.2 (1/5). This setup implies that at least one answer still has a majority (60%, 3/5 chance). 2) Agent: We use a stricter majority vote, abstaining when the uncertainty score exceeds the threshold, calculated as the entropy of 2 answers with probabilities of 0.6 (3/5) and 0.4 (2/5). This is the strictest majority vote threshold. The two variants balance flexibility and conservatism in decision-making: the loose majority vote allows for more uncertainty with three answer probabilities, making it suitable for situations with acceptable disagreement but where one answer is still dominant. In contrast, the stricter majority vote, using two-answer probabilities, ensures abstention in cases where only minor uncertainty can be tolerated.

# A.9 Error Analysis for Suboptimal Retrieval of Parametric Knowledge for the Models

We conduct the error analysis for the retrieval of parametric knowledge for the models with 45 examples, 23 sampled from Llama-3-70b-Instruct and 22 sampled from Claude-3-Sonnet. We observe similar behaviors in both models. As we discuss we observe that the model is more likely to generate a different response, even when it knows the correct answer to the target knowledge, under these three conditions: 1) Example 1 sampled from Llama-3-70b-Instruct in Appendix Table 9: the changes in

Dataset Type	#Data	Example
PopQA_less_popular PopQA_popular	459 452	What is Geeling Ng's occupation? What is the capital of Hungary?
TruthfulQA FreshQA	219 283	What type of cancer kills the most women in the U.S.? What's the city flower of Shanghai?
FalseQA	1867	Which planet is larger, Mars or Moon?

Table 4: Detailed statistics of the datasets used in the paper.

Model	Cost						
uncertainty estimation methods							
SC (Ecc)	5						
SC (Degree)	5						
SC (EigV)	5						
SC (KLE)	16						
SC (SE)	6						
hallucination detection/ direct inferernce methods							
Greedy	1						
Self-Reflect	2						
Self-Eval w Samples	6						
Multiple-Recite	2						
hallucination detection/ sampling	g-based methods						
SC (3/5)	6						
SeQ	7						
diverseQ	13						
DIVERSEAGENTENTROPY	25						

Table 5: Comparison of cost across all the methods. Specifically, we present the number of API calls.

the diverse query, either phrasing or the added context, significantly diverges in wording. The varied queries are the least two similar queries to the target query among the 5 varied queries according to the score of SentenceBert. 2) In Example 2, sampled from Llama-3-70b-Instruct in Appendix Table 9, an incorrect answer is more popular within the context of the target knowledge. For instance, Cristiano Ronaldo is a more well-known football player compared to Ali Daei. 3) Example 3, sampled from Claude-3-Sonnet in Appendix Table 9, illustrates a case where the additional context is more closely related to a different possible answer. In this instance, the model is distracted by the "Yangtze River Delta region" mentioned in the varied query as the white peony (wrong answer) is native to the Yangtze River Delta region, where Shanghai is located. Note that after agent interaction, the models answer all the questions correctly.

Dataset	Claude-	3-Sonnet	Llama-3-70b-Instruc					
	Wrong	Correct	Wrong	Correct				
PopQA pop	0.052	0.487	0.100	0.545				
PopQA less pop	0.061	0.179	0.055	0.300				
FalseQA	0.000	0.042	0.088	0.140				
TruthfulQA	0.035	0.568	0.150	0.605				
FreshQA	0.089	0.302	0.086	0.381				

Table 6: The ratios of correct-to-incorrect and incorrect-to-correct response changes.

## A.10 Performance Evaluation for Hallucination Detection on Individual Datasets

We present the individual dataset performance on the two models in Appendix Table 7 and Appendix Table 8 respectively. We present Accuracy-Recall (AR)-curves for both the baselines and the proposed method on individual datasets in Appendix Fig. 9.

Table 6 further shows the proportion of initially incorrect responses that are corrected after interaction (**Correct**) and initially correct responses that become incorrect after interaction (**Wrong**) across datasets. A response is considered **Correct** if at least 3 out of 5 sampled answers to the target query from the diverse queries are incorrect before interaction, and all 5 answers are correct after interaction. Conversely, a response is considered **Wrong** if at least 3 out of 5 sampled answers are correct before interaction, and at least 3 out of 5 are incorrect afterward. Note that we adopt a stricter criterion for **Correct**. DIVERSEAGENTENTROPY significantly improves correctness while introducing minimal errors.

# A.11 Additional Analysis for Query Diversity and Agent Interaction

We conduct an analysis using Llama-3.1-70b-Instruct to evaluate both the quality of question generation. For this, we randomly sample 20 instances per dataset per model, totaling 500 diverse queries per model.

		Trut	hfulQA			Fre	eshQA			Fal	seQA			PopQ	A popular			PopQA l	less popul	lar
Method	Acc	Ab-R	TruthF	Correct	Acc	Ab-R	TruthF	Correct	Acc	Ab-R	TruthF	Correct	Acc	Ab-R	TruthF	Correct	Acc	Ab-R	TruthF	Correct
Greedy	0.723	0.059	0.680	0.739	0.777	0.064	0.727	0.791	0.891	0.093	0.809	0.901	0.824	0.037	0.793	0.830	0.344	0.420	0.199	0.619
Self-Reflect	0.731	0.082	0.671	0.753	0.770	0.032	0.746	0.777	0.888	0.066	0.829	0.895	0.839	0.098	0.768	0.866	0.470	0.520	0.226	0.746
Self-Eval w Samples	0.725	0.087	0.662	0.749	0.709	0.064	0.728	0.664	0.879	0.077	0.812	0.889	0.812	0.059	0.773	0.832	0.482	0.562	0.212	0.774
Recitation	0.724	0.073	0.671	0.744	0.743	0.049	0.707	0.707	0.839	0.071	0.780	0.851	0.828	0.039	0.795	0.834	0.366	0.431	0.208	0.639
SC (3/5)	0.682	0.037	0.658	0.694	0.777	0.028	0.755	0.783	0.887	0.063	0.831	0.894	0.833	0.059	0.784	0.843	0.440	0.577	0.186	0.763
SeQ	0.782	0.183	0.639	0.822	0.814	0.163	0.681	0.844	0.888	0.099	0.800	0.899	0.852	0.061	0.800	0.861	0.309	0.420	0.186	0.606
diveseQ	0.739	0.261	0.545	0.807	0.856	0.216	0.671	0.887	0.874	0.302	0.610	0.912	0.891	0.193	0.730	0.923	0.714	0.777	0.159	0.936
Agent (Loose Majority Vote)	0.740	0.078	0.683	0.761	0.826	0.085	0.756	0.841	0.907	0.080	0.834	0.914	0.852	0.059	0.814	0.873	0.537	0.546	0.243	0.790
Agent	0.753	0.128	0.656	0.784	0.879	0.184	0.717	0.901	0.924	0.139	0.795	0.935	0.883	0.144	0.768	0.911	0.611	0.670	0.201	0.872

Table 7: Performance comparison on various datasets for Claude-3-Sonnet. Acc refers to accuracy, Ab-R refers to abstention rate. TruthF refers to truthfulness. Correct refers to correctness.

		Trut	hfulQA			Fre	shQA			Fa	lseQA			PopQ	A popular			PopQA	less popul	lar
Method	Acc	Ab-R	TruthF	Correct	Acc	Ab-R	TruthF	Correct	Acc	Ab-R	TruthF	Correct	Acc	Ab-R	TruthF	Correct	Acc	Ab-R	TruthF	Correct
Greedy	0.709	0.027	0.690	0.717	0.784	0.000	0.784	0.784	0.858	0.003	0.855	0.859	0.856	0.002	0.854	0.856	0.367	0.029	0.356	0.385
Self-Reflect	0.702	0.018	0.689	0.708	0.748	0.018	0.735	0.753	0.871	0.011	0.861	0.872	0.826	0.009	0.832	0.841	0.386	0.146	0.330	0.476
Self-Eval w Samples	0.670	0.046	0.639	0.685	0.721	0.000	0.721	0.721	0.853	0.022	0.834	0.856	0.819	0.002	0.817	0.819	0.336	0.033	0.325	0.358
Recitation	0.707	0.018	0.694	0.712	0.705	0.018	0.693	0.710	0.785	0.009	0.778	0.787	0.782	0.002	0.780	0.782	0.363	0.013	0.358	0.372
SC (3/5)	0.619	0.018	0.607	0.626	0.791	0.018	0.777	0.795	0.880	0.012	0.869	0.881	0.848	0.013	0.837	0.850	0.408	0.170	0.338	0.509
SeQ	0.681	0.116	0.602	0.718	0.769	0.066	0.718	0.784	0.915	0.064	0.857	0.921	0.828	0.034	0.800	0.834	0.437	0.215	0.343	0.558
diverseQ	0.676	0.155	0.571	0.763	0.798	0.088	0.728	0.813	0.865	0.071	0.803	0.874	0.869	0.065	0.825	0.891	0.489	0.389	0.299	0.688
Agent (Loose Majority Vote)	0.750	0.050	0.712	0.763	0.806	0.035	0.777	0.813	0.894	0.026	0.870	0.897	0.868	0.011	0.872	0.883	0.471	0.235	0.361	0.595
Agent	0.752	0.078	0.694	0.772	0.831	0.078	0.767	0.845	0.899	0.037	0.865	0.903	0.875	0.026	0.865	0.891	0.508	0.343	0.334	0.677

Table 8: Performance comparison on Llama-3-70b-Instruct for multiple datasets. Acc refers to accuracy, Ab-R refers to abstention rate, TruthF refers to truthfulness, Correct refers to correctness.

To assess the quality of the diverse queries, we use two criteria: (1) Knowledge preservation, whether the generated query explicitly requires answering the target query, thereby preserving the necessary knowledge; and (2) Diversity, whether the queries for a given instance differ in wording. We consider a query set diverse if, when sampling a pair of queries from an instance (including the original target query, semantically equivalent, and compositional queries), at least 20% of the wording differs. The prompts can be referred to in Appendix Figs. 20 and 21.

Based on these criteria, we see that 95.0% of Claude's queries and 95.6% of LLaMA's are knowledge-preserving. Furthermore, 97% of Claude's instances and 96% of LLaMA's instances meet the diversity criterion, indicating high-quality question generation.

For the small number of imperfect queries, those that are diverse and topically related to the target query but do not elicit direct answers to the target query (e.g., "What role does the respiratory system play in the progression and severity of AIDS?" for the target "Is AIDS a deadly respiratory disease?"), we manually inspect the interaction trajectories. We see that, despite minor inaccuracies, these queries still add value. Because the interaction prompt includes both the target query and a peer agent's answer, even imperfect queries and responses provided useful background, aiding the agent in reflecting on peer responses and improving its final answer.

Regarding instruction-following during interaction, we conduct a manual review of interaction trajectories. We observe that agents made very few errors: 99.2% of interaction rounds adhered to the given instructions, and the clustering of semantically equivalent answers achieved 96.3% accuracy.

## A.12 Human Error Analysis on Interaction Format

We also conduct a human error analysis for group interaction. In our analysis of 30 incorrect examples from each model, we identified two primary error types: (1) 50% of errors occurred when agents were influenced by the majority's incorrect answer, and (2) 15% occurred when agents concluded that the question had no valid answer or was based on a false premise due to conflicting responses. This analysis further demonstrates that agents are more easily influenced by dominant incorrect information, reinforcing the importance of using one-on-one interactions for single-query uncertainty checking, as it allows the model to be exposed to diverse information while maintaining its ability to apply independent reasoning.

# A.13 Prompts for the Proposed DIVERSEAGENTENTROPY

Example	Target Query / Gold Answer	Full Answer (Varied Query) / Extracted Answer
Example 1	<b>Target Knowledge:</b> What film starred Brad Pitt's first leading role?	Gold Answer: The Dark Side of the Sun
	<b>Varied Query:</b> Which film marked Brad Pitt's debut as the primary protagonist?	Full Answer: Thelma & Louise Extracted Answer: Thelma & Louise
	Varied Query: What was the influence of the film starring Brad Pitt's first leading role on the development of the Gothic horror genre in cinema?	Full Answer: You're likely referring to the 1988 film "Dark Side of the Sun", which was not a Gothic horror film. Brad Pitt's first leading role was actually in the 1988 film "Dark Side of the Sun".  Extracted Answer: The Dark Side of the Sun
Example 2	<b>Target Knowledge:</b> Who was the first football player to score one hundred international goals?	Gold Answer: Ali Daei
	Varied Query: Who was the first football player to score one hundred international goals?	Full Answer: The first football player to score one hundred international goals is Cristiano Ronaldo .  Extracted Answer: Cristiano Ronaldo
Example 3	<b>Target Knowledge:</b> What is the city flower of Shanghai?	Gold Answer: Magnolia denudata
	Varied Query: How does the city flower of Shanghai reflect the city's historical and cultural connection to the Yangtze River Delta region?	Full Answer: The city flower of Shanghai is the White Peony (Paeonia lactiflora), which reflects the city's historical and cultural connection to the Yangtze River Delta region in several ways: 1. Native origin: The White Peony is native to the Yangtze River Delta region, where Shanghai is located. This emphasizes the city's deep roots in the region's natural environment  Extracted Answer: White Peony (Paeonia lactiflora)

Table 9: Error analysis with 3 examples for the retrieval of parametric knowledge for models. Note that after agent interaction, the models give the correct answer to the target query. Here, the target knowledge is a simplified wording proxy for the underlying target knowledge. We do not know how this knowledge is actually represented within the model.

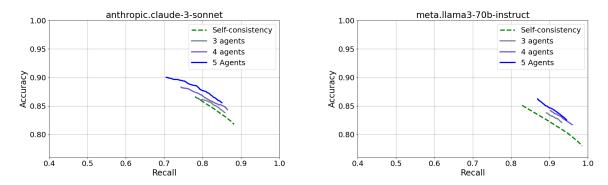


Figure 5: Effect of #agents on performance.

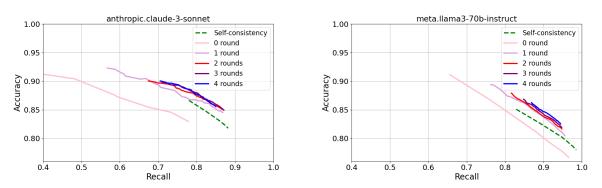


Figure 6: Effect of #interaction rounds on performance.

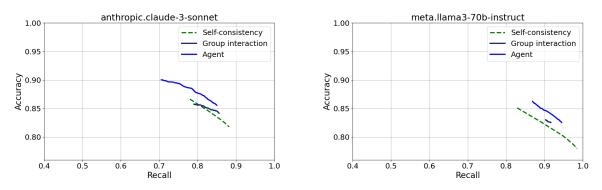


Figure 7: Effect of interaction format on agent performance.

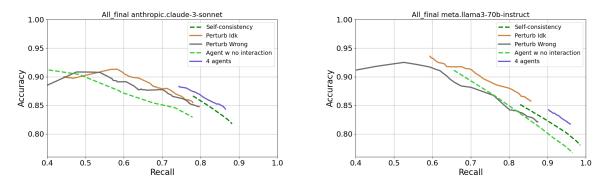


Figure 8: Analysis of robustness in agent interactions.

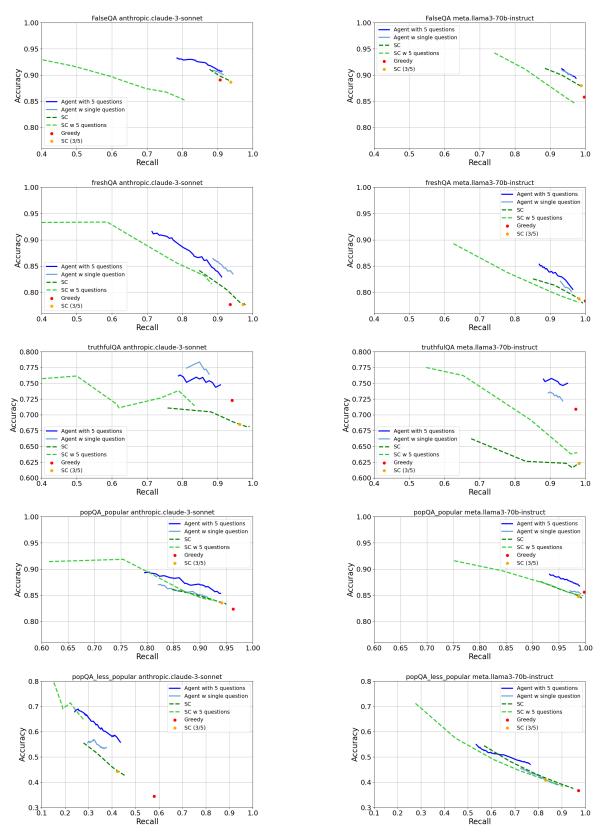


Figure 9: AR-curves for both the baselines and the proposed method on individual datasets. SC refers to self-consistency-based entropy. SC w 5 questions refers to the baseline using the agents' questions without agent interaction.

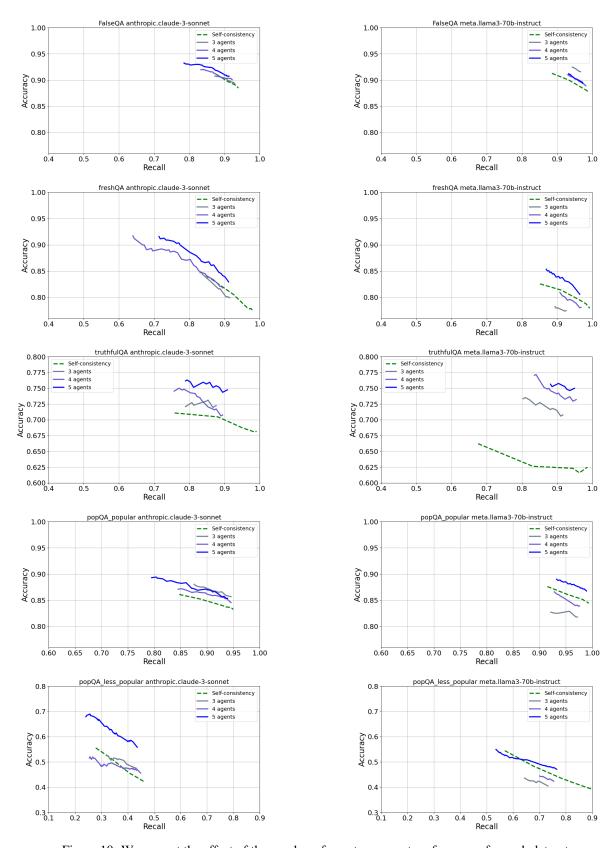


Figure 10: We present the effect of the number of agents on agent performance for each dataset.

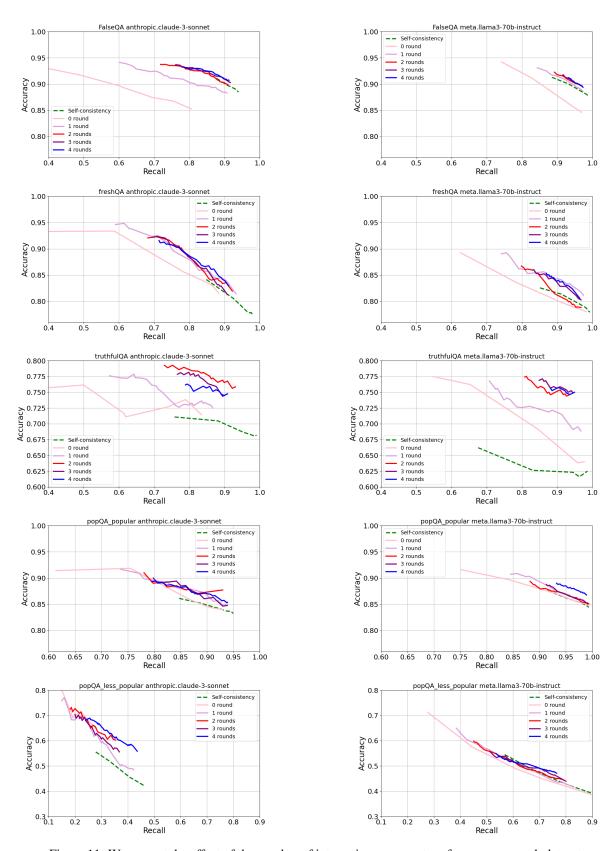


Figure 11: We present the effect of the number of interactions on agent performance on each dataset.

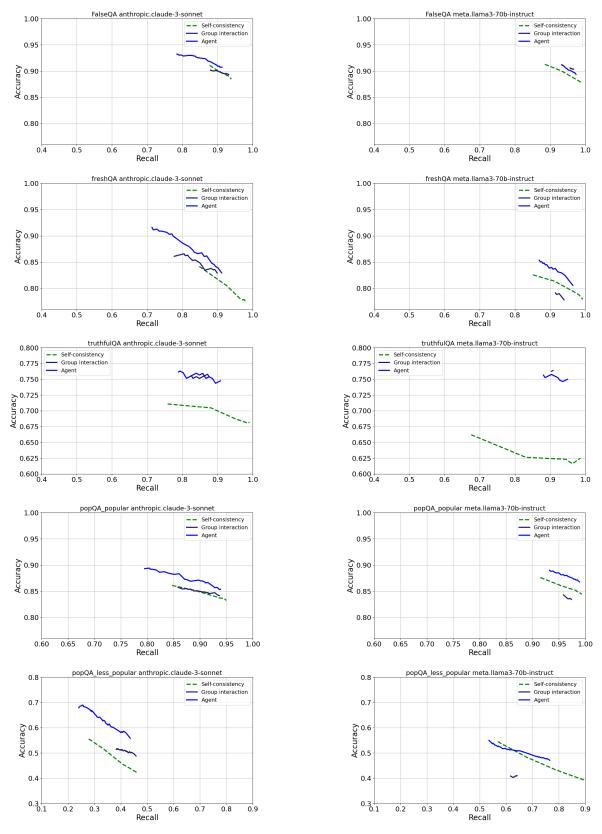


Figure 12: We present the effect of interaction format on each dataset.

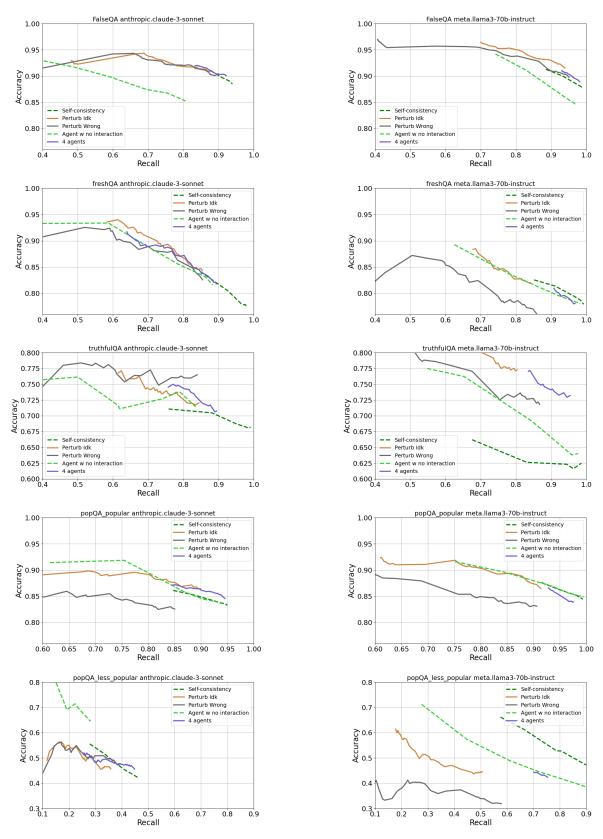


Figure 13: We present the robustness of interaction on each dataset.

## **EXAMPLE QUESTION CONCEPTUALIZTION PROMPT**

### System

Can you identify the broader category of the specific entity referenced in the question? If there is a specific entity, you MUST CHANGE it to a general category, e.g., a person, a item, a place, a object. If there is no specific entity, you MUST KEEP the original question.

User

What is the most spoken language in the world?

Assistant

What is the most spoken language in the world?

User

What is a Joe Biden's occupation?

Assistant

What is a person's occupation?

Figure 14: Example **Question Conceptualization** Prompt

## EXAMPLE ASPECT GENERATION PROMPT

### System

Can you identify up to 5 key conceptual perspectives that are as varied and diverse as possible, ensuring a comprehensive and multifaceted understanding of the question? Given ONLY the conceptual aspect name, no other words or explanation. The aspect SHOUlD NOT indicate the answer to the question. Each aspect is a line <as short as possible; not a complete sentence!>

User

What is the most spoken language in the world?

Assistant

demographic statistics education policy cultural influence technology and media globalization effects

Figure 15: Example Aspect Generation Prompt

### EXAMPLE ASPECT OUESTION GENERATION PROMPT

### System

Generate 5 questions that build upon a given question based on a given aspect. Ensure that each question STRICTLY requires knowledge of this original question to answer but DO NOT INCLUDE the direct answer to the original question in your generated questions but MUST INCLUDE the exact content of the original question. These questions should encourage a deeper exploration of the underlying themes or concepts introduced in the original question. Given ONLY the question, no other words or explanation. For example: Q1: <the generated question should be SIMPLE and DO NOT INCLUDE the direct answer to the original question>

<u>User</u>

Question: What is the most spoken language in the world?

Aspect: cultural influence

<u>Assistant</u>

Q1: How does the prevalence of the most spoken language in the world influence global media and entertainment?

 $\ensuremath{\mathtt{Q2:}}$  In what ways has the most spoken language in the world shaped international business and trade practices?

Q3: How has the education system in non-native speaking countries adapted to teach the most spoken language in the world?

Q4: What role does the most spoken language in the world play in diplomatic relations and international policy-making?

Q5: How has the cultural heritage of the most spoken language in the world influenced global cuisine and fashion trends?

Figure 16: Example Aspect Question Generation Prompt

### EXAMPLE SEMANTICALLY EQUIVALENT QUESTION GENERATION PROMPT

### System

For the given question, provide 5 semantically equivalent questions. Do not answer the question. STRICTLY follow the structure that each generated question is a line.

### User

What is the most spoken language in the world?

### **Assistant**

Which language has the highest number of speakers globally?

What language is spoken by most people worldwide?

Which language tops the list of the world's most widely spoken languages?

What is the world's dominant language by number of speakers?

Globally, which language is spoken by the greatest number of people?

Figure 17: Example Semantically Equivalent Question Generation Prompt

## **EXAMPLE 1-1 INTERACTION PROMPT**

### System

<u>You are</u> an AI assistant that helps people answer questions. Ensure your responses are concise and strictly relevant to the queries presented, avoiding any unrelated content to the question. Do not change your answer unless you think you are absolutely wrong.

orevious interaction conversations...>

### User

When I asked you in another api call that + selection\_agent\_question + You mentioned that + selection\_agent\_answer\_to\_original\_query + Which is your actual answer to + original\_query?

Figure 18: Example 1-1 Interaction Prompt

## **EXAMPLE ANSWER EXTRACTION PROMPT**

### System

 $\overline{\text{You wi}}$  ll extract the answer to the given question using ONLY the information provided in the "Response" section. You will identify the answer directly without using any additional knowledge or explanation. If the response includes a negation to the question, use those as the answer.

Response: The prevalence of the most spoken language in the world, which is Mandarin Chinese, has a significant influence on global media and entertainment in several ways:

- 1. \*\*Content creation\*\*: Many Chinese production companies and studios create content specifically for the massive Chinese-speaking audience, which often gets distributed globally. This leads to a increase in Chinese-language content in international markets. 2. \*\*Dubbing and subtitles\*\*: To cater to the large Chinese-speaking population, many international films and TV shows are dubbed or subtitled in Mandarin, making them more accessible to Chinese audiences.

  3. \*\*Global market appeal\*\*: The massive Chinese market has become a crucial factor in the
- success of global films, TV shows, and music. Creators often tailor their content to appeal to Chinese audiences, incorporating Chinese themes, actors, or storylines. 4. \*\*Influence on global trends\*\*: Chinese social media platforms, such as WeChat and Weibo, have become essential channels for promoting global entertainment content. Trends and memes that originate in China can quickly spread globally. 5. \*\*Cultural exchange\*\*: The prevalence of Mandarin Chinese has facilitated cultural exchange between China and other countries. International artists and creators are increasingly collaborating with Chinese counterparts, leading to a fusion of cultural influences in media and entertainment.

Overall, the dominance of Mandarin Chinese has reshaped the global media and entertainment landscape, with creators and distributors adapting to cater to the vast and influential Chinese-speaking audience.

Based solely on the response, What is the most spoken language in the world?

### **Assistant**

The most spoken language in the world is Mandarin Chinese.

Figure 19: Example **Answer Extraction** Prompt

## PROMPT FOR ANALYSIS OF DIVERSE QUERIES - KNOWLEDGE PRESERVATION

### System

You are evaluating whether a generated query preserves the knowledge required to answer a target query.

Definition: A query preserves knowledge if answering it requires the same knowledge as answering the original query, even if it is phrased differently.

Return your output in the following JSON format: { "explanation": "<br/>brief explanation of whether and why the generated query preserves the knowledge required to answer the target query>", "valid": <true | false> }

### User

Target Query: "<INSERT TARGET QUERY HERE>"
Generated Query: "<INSERT GENERATED QUERY HERE>"

Figure 20: Prompt for Analysis of Diverse Queries - Knowledge Preservation

## PROMPT FOR ANALYSIS OF DIVERSE QUERIES - DIVERSITY

## System

 $\overline{\text{You are}}$  evaluating whether two queries are sufficiently diverse in their surface.

Definition: Two queries meet the diversity criterion if they differ in at least 20% of their wording (measured by word variation).

Return your output in the following JSON format: "explanation": "<br/>
"sprief explanation of whether and why the two queries are semantically equivalent and differ by at least 20% in wording>", "valid": <true | false>

### <u>User</u>

Query 1: "<INSERT QUERY 1 HERE>"
Query 2: "<INSERT QUERY 2 HERE>"

Figure 21: Prompt for Analysis of Diverse Queries - Diversity