

Towards Mitigation of Hallucination for LLM-empowered Agents: Progressive Generalization Bound Exploration and Watchdog Monitor

Siyuan Liu^{a,b}, Wenjing Liu^c, Zhiwei Xu^{b,d,*}, Xin Wang^e, Bo Chen^f and Tao Li^{a,**}

^aCollege of Computer Science, Nankai University

^bHaihe Lab of ITAI

^cCollege of intelligent Science and Technology, Inner Mongolia University of Technology

^dInstitute of Computing Technology, Chinese Academy of Sciences

^eDepartment of Electrical and Computer Engineering, Stony Brook University

^fDepartment of Computer Science, Michigan Technological University

Abstract. Empowered by large language models (LLMs), intelligent agents have become a popular paradigm for interacting with open environments to facilitate AI deployment. However, hallucinations generated by LLMs—where outputs are inconsistent with facts—pose a significant challenge, undermining the credibility of intelligent agents. Only if hallucinations can be mitigated, the intelligent agents can be used in real-world without any catastrophic risk. Therefore, effective detection and mitigation of hallucinations are crucial to ensure the dependability of agents. Unfortunately, the related approaches either depend on white-box access to LLMs or fail to accurately identify hallucinations. To address the challenge posed by hallucinations of intelligent agents, we present HalMit, a novel black-box watchdog framework that models the generalization bound of LLM-empowered agents and thus detect hallucinations without requiring internal knowledge of the LLM’s architecture. Specifically, a probabilistic fractal sampling technique is proposed to generate a sufficient number of queries to trigger the incredible responses in parallel, efficiently identifying the generalization bound of the target agent. Experimental evaluations demonstrate that HalMit significantly outperforms existing approaches in hallucination monitoring. Its black-box nature and superior performance make HalMit a promising solution for enhancing the dependability of LLM-powered systems.

1 Introduction

With the rapid proliferation of artificial intelligence in contemporary life, agents empowered by large language models (LLMs) have emerged as pioneers of this technology transformation [31]. However, in conjunction with their widespread deployment, the phenomenon of hallucination has become a major concern in LLM and their agents [13]. LLM hallucination refers to instances in which LLM generated content is inconsistent, unfaithful, contradictory, or unverifiable against established real-world knowledge [13], although it may be presented in a convincing and confident tone. This issue

has been recognized by many academic studies and technical public reports as one of the primary ethical and safety risks associated with LLM agents, along with issues such as bias and toxic content. The hallucination phenomenon becomes thorny when considering the black-box nature of LLMs, and severely undermines the credibility of LLM agents, especially in truth-sensitive fields such as law [9], medicine [7], finance [35], and education [15], where it can have catastrophic cognitive consequences.

Mitigating hallucinations is critical to improving the dependability of LLM agents in real-world applications [13]. Hallucinations typically arise when the generated content significantly exceeds the generalization bounds of the agent [30]. If a generated response lies outside the bound, it is highly likely that this response is hallucinated [34]. Therefore, identifying the generalization bounds is of critical importance in mitigating hallucinations in LLM agents. Although recent efforts have been focused on computing non-vacuous generalization bounds for deep learning models, these bounds tend to become vacuous at the scale of billion-parameter models [19, 34]. In addition, such theoretical bounds are often derived from restrictive statistical assumptions that limit their applicability to models with low generalization capacity. Given the vastness of the semantic space, the tightness of the existing generalization bounds remains difficult to establish [18].

Existing approaches attempt to identify hallucinated responses by analyzing the internal state of the model [14, 10, 38, 11]. This requires full transparency and access to the internal state of the model, known as “white-box access”, and cannot work for large-scale commercial models, which are usually close-sourced. The other approaches [1, 32, 28] are mainly based on cross-checking of the LLM output against external databases. Recent work for modeling generalization bound computes non-vacuous generalization bounds for deep learning models, but these bounds are vacuous for large models at the billion-parameter scale [12, 21] that involve directly asking the LLM to produce its confidence scores in the truthfulness of the statements (referred to as “model confidence” in this paper). Although these black-box/gray-box approaches allow hallucination monitoring through output text or associated confidence scores, they are degraded by limited knowledge of LLMs and often poorly calibrated

* Corresponding Author. Email: xuzhiwei2001@ict.ac.cn.

** Corresponding Author. Email: litao@nankai.edu.cn.

confidence estimates. Ultimately, these score-based approaches may be incorrect, reducing the effectiveness and accuracy of the monitoring. Therefore, there is a strong demand for developing new techniques that can mitigate hallucinations while avoiding the aforementioned limitations.

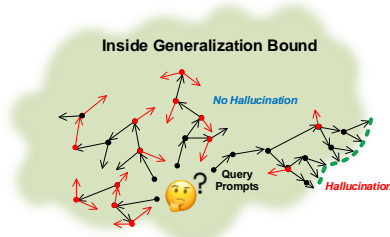


Figure 1: The complexity of generalization bound modeling.

Due to the complexity of the semantic space, deriving a universal generalization bound across all domains is extremely challenging [30]. However, we observe that this bound may be identified much more easily within the domain of each specific agent (Section 2). Building on this insight, we propose HalMit, a fine-grained approach for modeling per-agent generalization bounds to monitor hallucinations that fall outside these boundaries. A major challenge in designing HalMit lies in the need to efficiently identify and model the complex generalization bound of the target agent in a specific domain. As shown in Figure 1, the generalization bound is difficult to pinpoint. The bound exploration process may deviate from the true boundary and become trapped in local loops (red arrows), considering that the bounded space itself is vast.

To accelerate the generalization bound exploration, we propose a probabilistic fractal sampling method to generate a sufficient number of parallel queries so that they can efficiently cover the generalization bound of the targeted agent. Each identified boundary point is stored in a vector database, where the point is represented by its query-response pair along with associated context. During hallucination monitoring, HalMit compares the input query with those retrieved from the vector database. If the query closely resembles the retrieved records in the vector base, it is considered near the boundary, and the response corresponding to the input query is flagged as a potential hallucination. Our major contributions are summarized as follows:

- 1) Through a preliminary study on agent hallucinations, we confirm that hallucinated responses correspond to the agent’s generalization bound and that it is possible to model this bound within specific application domains.
- 2) We propose a novel probabilistic fractal exploration scheme to enable our MAS system to incrementally probe the generalization boundary. In this process, deep reinforcement learning guides the multi-agent exploration by adjusting the probabilities of fractal transformations based on three common semantic patterns. These probabilities are continuously updated to efficiently and effectively model the generalization bound within a specific domain.
- 3) A unique hallucination mitigation technology is provided based on the generalization boundary to enable more dependable monitoring and persistently detecting potential hallucinations.
- 4) We have conducted extensive experimental evaluations and the results are encouraging, demonstrating a significant improvement in hallucination monitoring effectiveness over baseline solutions. A unique hallucination mitigation technology is provided to enable a more dependable monitoring and monitoring of potential hallucinations.

These contributions address vital challenges in mitigating hallucination of LLM-empowered agents, paving the way for more trustworthy intelligent systems in real-world applications. To the best of our knowledge, this is the first hallucination monitoring approach that operates without access to internal model knowledge or reliance on cross-verification algorithms. This design enables real-time and persistent hallucination mitigation in deployed intelligent systems. The source code for HalMit will be available on GitHub after the paper is accepted.

2 Motivation

In this section, we present preliminary studies that investigate how hallucinations manifest in LLM-empowered agents across various domains.

- **PS1:** We investigate whether the statistic characteristics of hallucinations vary across different domains.
- **PS2:** We examine the statistic characteristics within each domain to identify potential patterns or regularities in LLM hallucinations.
- **PS3:** We assess whether certain existing characteristics can be directly leveraged to monitor hallucinations.

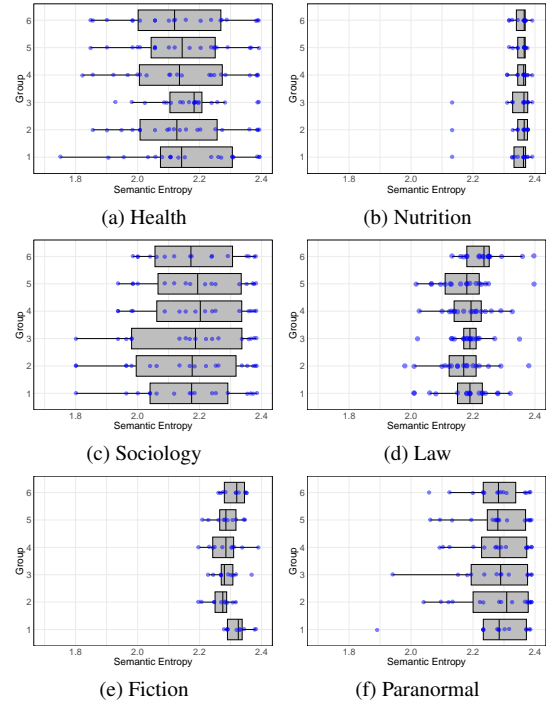


Figure 2: The semantic entropy values for six popular domains.

Experimental Settings: Without loss of generality, we study the response quality of agents powered by Llama3.1-8B across six popular domains, including health, nutrition, sociology, law, fiction, and paranormal. All query-answer pairs and domains used in the study are randomly sampled from a real-world hallucination dataset, TruthfulQA [17], which includes both truthful and hallucinated responses. Semantic entropy [6] is one metric that has been extensively used to assess the uncertainty level of agent responses (see Appendix B), with higher entropy value typically indicating greater uncertainty and a higher likelihood of hallucinations. Unlike general confidence scores of LLM’s response, semantic entropy is specifically designed to characterize the hallucinations of LLMs. Based on semantic entropy, we evaluate response quality across the six domains by measuring semantic entropy, with the results presented in Figure 2.

- **PS1: Significant Variations across Domains.** As shown in Figure 2, the semantic entropy values of agent responses vary substantially across application domains, with noticeable differences in both medians and variances. This indicates that the statistical characteristics of hallucinations differ significantly by domains, suggesting that no universal generalization bound can be established across all domains.
- **PS2: Statistic Stability within The Same Domain.** Further analysis of the boxplots reveals that within each individual domain, the semantic entropy values, while subject to some fluctuation across testing groups, tend to follow a consistent distribution. This internal stability suggests that hallucination patterns are coherent within each domain, making it feasible to identify a domain-specific generalization bound.
- **PS3: Limitations of a Threshold Based on the Existing Metric.** As observed in the boxplots for the Health, Nutrition, Law, Fiction, and Paranormal domains, outliers extending beyond the whiskers in the boxplots are present. This suggests that while semantic entropy has been shown to support hallucination detection [16, 6, 10], relying solely on a fixed threshold is insufficient. The presence of high-entropy yet potentially non-hallucinatory responses (and vice versa) highlights the need for more nuanced detection methods.

Summary of findings: Based on the above observations, our findings reveal that: (1) hallucination patterns vary across application domains but tend to exhibit consistent statistical behavior within the same domain; and (2) hallucinations cannot be effectively detected using a simple threshold on any single existing metric, even one as significant as semantic entropy. These insights suggest that the generalization bound distinguishing hallucinated from non-hallucinated responses is more clearly defined when tailored to a specific agent within a specific domain, rather than derived from a collective set of agents or domains. Consequently, accurate hallucination monitoring requires the identification of these generalization bounds in a fine-grained, domain-aware manner.

Motivated by these observations, we propose *HalMit*, a hallucination mitigation paradigm for LLM-empowered agents.

3 Methodology

To persistently mitigate hallucinations, HalMit functions as a “watchdog” framework for each target agent to monitor hallucinations. Before being used to monitor hallucinations, HalMit first models the agent’s generalization bound based on a proposed multi-agent exploration system. This allows hallucinations to be identified and mitigated based on their deviation from the learned generalization boundary.

3.1 Generalization Exploration Bound with a MAS

Given the black-box nature of LLM-empowered agents, exploring their generalization bounds and identifying hallucinations are critical and challenging [33]. To address this, we introduce a multi-agent bound exploration method that integrates probabilistic fractal sampling into a multi-agent system (MAS) for parallel query generation. To improve the efficiency and relevance of the queries, we propose a reinforcement learning-based scheme to dynamically adjust the fractal probabilities.

As illustrated in Figure 3, the proposed MAS consists of three specialized agent types: core agent (CA), query generation agent (QGA),

and evaluation agent (EA). Among them, the CA coordinates interactions between QGAs and the target LLM-powered agent. Meanwhile, the EA, guided by HalluBench criteria [36], evaluates the quality of the response from the target and provides essential feedback to refine the query generation process.

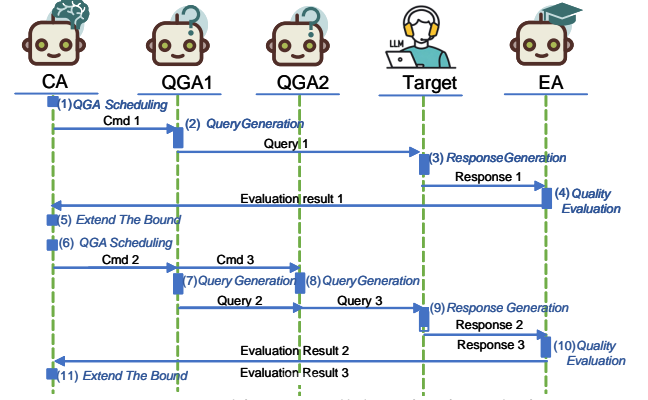


Figure 3: Multi-agent collaboration in HalMit.

3.1.1 Probabilistic Fractal-based Query Generation

Leveraging the self-similarity feature of the fractal in natural language, we propose a novel probabilistic fractal-based query generation method for use in QGAs. This method iteratively constructs increasingly complex query structures that progressively approach the generalization bound of the target agent. Unlike conventional fractal systems that apply all affine transformations in each iteration step, our method proposes an Iterated Function System with Probabilities (IFSP) (See appendix A) that enough queries can be generated quickly to cover the generalization bound of the target agent τ . More specifically, according to semantic theory [23], three semantic extension patterns, *induction*, *deduction*, and *analogy*, are used as fractal affine transformations to extend specific queries and navigate the semantic space. The execution probability of each transformation is dynamically adjusted based on the IFSP system:

$$\mathcal{F} = \{FT_i : \mathcal{P}_{t-1}^{\tau} \xrightarrow{p_i} \mathcal{P}_t^{\tau}, \sum_{i=1}^3 p_i = 1\}, \quad (1)$$

where P_t^{τ} are queries used in round t , functions $FT1, FT2, FT3$ correspond to three fractal affine transformations, p_i is the execution probability for FT_i . These three types of fractal affine transformations are introduced as follows:

- **FT1: Semantic Deduction.** This transformation generates more specific queries by deriving them from general rules or concepts presented in the previous iteration. For example, given a query, “Did humans really land on the moon in 1969?”, a deductive transformation would produce a more focused follow-up such as “What were the technological advancements that enabled humans to land on the moon in 1969?”.
- **FT2: Semantic Analog.** This transformation broadens the scope of the original query by leveraging semantic associations such as synonyms, antonyms, or functional analogies. In this way, a new query, “What historical events in space exploration paralleled the significance of the moon landing in 1969?”, can be generated. This helps the system probe parallel narratives and conceptual similarities in the semantic space.

- **FT3: Semantic Induction.** This transformation generates broader, more abstract queries by generalizing from specific instances and inferring underlying linguistic or conceptual patterns. As another example, “How can we assess the reliability of commonly accepted events in the history of space exploration?”, is obtained by following FT3. This supports exploration of overarching themes and epistemological questions within the domain.

These three affine transformations are applied in each iteration of the exploration process. In this way, new queries are generated through the iterative process that starts from basic concepts or principles. The core agent can order multiple query generation agents to generate queries, so the iteration can be realized through parallel processes, to significantly increase the speed in identifying the generalization bound. This iteration process includes four steps:

- 1) To cover a broader semantic space within the bound, the CA randomly initializes multiple queries in a domain and sends each query to the target agent τ as initial questions.
- 2) The target agent responds to queries with responses that may contain hallucinations. Therefore, each QA pair is sent to an EA. After receiving the QA pair, the EA assesses whether the response in the received QA pair contains a hallucination, and sends a report back to the CA, including the QA pair and the corresponding evaluation results.
- 3) Depending on the evaluation result, QGA will perform query generation in two ways. In case a hallucination is reported, the CA embeds the QA pair and the context information into a vector database as a point of the generalization bound of agent τ (detailed in Section 3.2). To expand the exploration range and cover more of the generalization bound, the CA schedules the QGAs to generate new queries through the fractal affine transformations, FT1 and FT2, where the probability of each is determined through reinforcement learning (details in Section 3.1.2). Otherwise, in case no hallucination is reported, new queries will be randomly generated, and sent back to the CA.
- 4) The new queries generated by the QGAs are sent back to the target agent. In this way, a new round of fractal exploration with a pipeline is scheduled by CA, where multiple new iterations are generated in parallel in response to each exploration path. More EAs are scheduled to assess all QA pairs and the evaluation reports will be sent to the CA. The bound search speed can be exponentially increased in this way.

During this iterative process, the ratio of the hallucinations among all QA pairs, γ , is incrementally updated by the core agent. Once γ becomes larger than an empirical threshold ϵ , it indicates that the generalization bound of agent τ can be identified by the vector database, and this iterative process of \mathcal{F} ends. These parameters are evaluated in ablation studies in Section 4.5.

3.1.2 Reinforced Determination of Fractal Probabilities

To further increase the efficiency in identifying the generalization bound, we use deep reinforcement learning [26] to determine the probability of each transformation function in IFSP \mathcal{F} to go for in the next step. This probability is adjusted in each iteration so that the exploration process can more efficiently converge towards the bound. Deep reinforcement learning is a framework in which a policy network is trained to maximize long-term objectives. In our design, the policy network is trained to efficiently select the appropriate probabilities of fractal affine transformations in each iteration, driving the fast convergence to the bound as its long-term objective, as shown in Figure 4. This design significantly accelerates the exploration process.

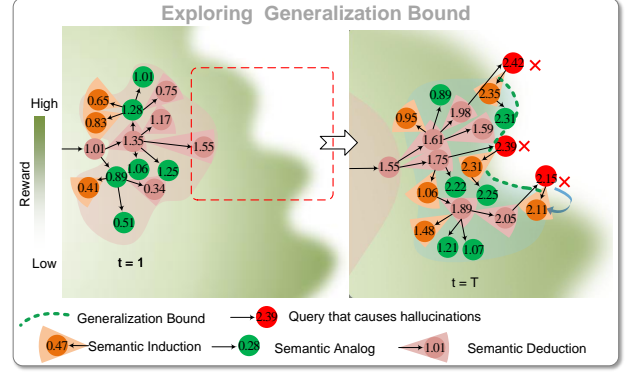


Figure 4: Process of exploring generalization bound. The numbers in the circles are the corresponding semantic entropies.

Training Data Prepared for Policy Network:

To evaluate the effectiveness of a fractal affine transformation in state s_i , we repetitively send the corresponding query, \mathcal{P}_i^τ , to the target agent K times, collect K responses, $\{a_i^\tau\}_K$, and according to [6], calculate the semantic entropy for the target agent in state s_i , H_i^τ . More specifically, every item in $\{a_i^\tau\}_K$, $a_i^\tau(k)$, is included in the token sequence group for calculating H_i^τ . Define a function $\text{sig}(a_i^\tau(k))$, which is equal to 0 if the answer of the target agent contains hallucinations; otherwise, $\text{sig}(a_i^\tau(k))$ is equal to 1. Hence, the reward for the triple $\{\mathcal{P}_i^\tau, \{a_i^\tau\}_K, H_{i-1}^\tau\}$ can be given as:

$$R_i^\tau(\mathcal{P}_i^\tau, \{a_i^\tau\}_K, H_{i-1}^\tau, H_i^\tau) = \begin{cases} \Delta H_i^\tau & \text{if } \prod_{k=1}^K \text{sig}(a_i^\tau(k))^\tau \neq 0 \\ \frac{1}{R_{i-1}^\tau} & \text{if } \prod_{k=1}^K \text{sig}(a_i^\tau(k))^\tau = 0 \end{cases} \quad (2)$$

We generate queries based on three fractal affine transformations, collect their queries, answers, and semantic entropy, and calculate their rewards. Considering that each of these rewards represents the degree to which the enhancement in the exploration of generalization bounds is achieved through the corresponding affine transformation, we configure the probability for an affine transformation j ($j \in \{1, 2, 3\}$) in our IFSP as:

$$p_j = R_j^\tau / \sum_{k=1}^3 R_k^\tau. \quad (3)$$

Each p_j , along with the corresponding input query, response, and semantic entropy, is included in the training database of the policy network, which consists of a quadruple, $\{\mathcal{P}_j^\tau, \{a_j^\tau\}_K, H_{j-1}^\tau, p_j\}$. Through the iteration process of fractal sampling on the target agent, more quadruples can be collected until the scale of this dataset is sufficient for the policy network training.

Policy Network Training: A popular design of the policy network multilayer perceptron (MLP) [26] is introduced to capture the non-linear representation of state features and predict the probability distribution of three affine transformations $\{p_j\}$ in IFSP \mathcal{F} . States in the generalization space of the target agent are defined as:

$$s_i^\tau = \left\lfloor \frac{\Delta \text{Sim}(\mathcal{P}_0^\tau, \mathcal{P}_i^\tau) \times e^{H_i^\tau}}{\omega} \right\rfloor, \quad (4)$$

where \mathcal{P}_0^τ is the initial query used to explore the generalization bound. Its semantic similarity to \mathcal{P}_i^τ is involved as a scaling parameter and ω is a normalization parameter.

For each state s_i^τ , the action-value function $Q()$ for affine transformations f_i^τ at the time step i is $Q(s_i^\tau, f_i^\tau) = \mathbb{E}[R_i^\tau | s_i^\tau, f_i]$, where R_i^τ is the reward defined in Formula (2). The objective $Q(s_i^\tau, f_i; \theta_i)$

is to train the policy network, parameterized by θ denoted. The network is optimized to align with the optimal action-value $Q^*(s_i^\tau, f_i)$, which corresponds to the highest possible reward. To achieve this, a loss function is defined using the L2 norm:

$$\mathcal{L}(\theta) = \sum_{i=1}^N \|Q^*(s_i^\tau, f_i) - Q(s_i^\tau, f_i; \theta_i)\|, \quad (5)$$

where N is the mini-batch of training samples. Finally, the optimal parameters θ^* of the policy network can be obtained by:

$$\theta^* = \arg \min_{\theta_i} \sum_{i=1}^N \mathcal{L}(Q^*(s_i^\tau, f_i), Q(s_i^\tau, f_i; \theta_i)). \quad (6)$$

The convergence study of the exploration is provided in the experimental section (Section 4.3).

3.2 Hallucination Monitoring

In this section, we describe how the generalization bound can be leveraged to monitor if there exist hallucinations in the response from a target agent. This is achieved by comparing the response with the information retrieved from the vector database that represents the generalization bound of the target agent. As the bound often has an irregular shape (as illustrated in Figure 1), hallucination monitoring can be very difficult.

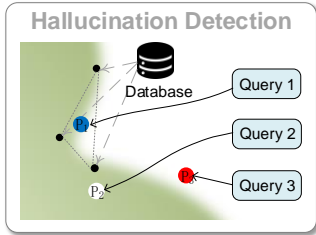


Figure 5: Hallucination monitoring

Hallucinated responses tend to diverge significantly, making it difficult to achieve a meaningful comparison with the bound. Instead, we choose the use of input query to compare against the bound. To start hallucination monitoring, the input query P_q^τ is compared with all related items in the vector database built in Section 3. This is achieved by evaluating the cosine similarity of the query vector Q_v and each related vector v_i denoted by S_i . Specifically, we consider three cases:

- When there are more than three similar items in the database that exceed a threshold ϵ , we calculate the centroid of three most similar items $P_{v_1}^\tau, P_{v_2}^\tau, P_{v_3}^\tau$:

$$C = \frac{\sum_{i=1}^3 S_i \cdot v_i}{\sum_{i=1}^3 S_i}, \quad (7)$$

where S_i is the similarity score of the i -th vector. Next, we calculate the cosine similarity S_C between the query vector Q_v and the normalized centroid C . If $S_C > \epsilon$, the input query is considered to be beyond the generalization bound of the target agent (corresponding to the blue point in Figure 5).

- Otherwise, we compare the semantic entropy of the query $H_{Q_v}^\tau$ with the semantic entropy of the most similar vector $H_{v_i}^\tau$ in the vector database. If $H_{Q_v}^\tau$ is larger, the input query is likely to be outside the generalization bound and may cause a hallucination (corresponding to the red point in Figure 5e).

- If $H_{Q_v}^\tau < H_{v_i}^\tau$, the input query is within the generalization bound, and will obtain a rational response (corresponding to the white point in Figure 5).

The details of the monitoring process are given in Algorithm 1.

Algorithm 1 Hallucination monitoring algorithm

```

1: Input: Input query  $P_q^\tau$ , vector database  $\mathcal{V}^\tau$ , similarity threshold  $\epsilon$  of
   hallucination monitoring;
2: Output: Hallucination status of  $P_q^\tau$ ;
3: for each vector  $v_i \in \mathcal{V}^\tau$  do
4:   Normalize  $v_i$ :  $v_i \leftarrow \frac{v_i}{\|v_i\|}$ ;
5: end for
6: Represent  $P_q^\tau$  with a vector:  $Q_v^\tau \leftarrow \text{Embedding}(P_q^\tau)$ ;
7: Normalize  $Q_v^\tau$ :  $Q_v^\tau \leftarrow \frac{Q_v^\tau}{\|Q_v^\tau\|}$ ;
8: Initialize an empty list of results;
9: for each vector  $v_i \in \mathcal{V}^\tau$  do
10:   Compute cosine similarity  $S_i$ ;
11:   Append  $(v_i, S_i)$  to results;
12:   Sort results by similarity score  $S_i$  in descending order;
13: end for
14: if results[3] >  $\epsilon$  then
15:   Compute the centroid by Formula (7);
16:   Normalize the centroid:  $C \leftarrow \frac{C}{\|C\|}$ ;
17: end if
18: Compute the similarity  $S_C$  between  $Q_v^\tau$  and the centroid  $C$ ;
19: if  $S_C \geq \epsilon$  then
20:   Report  $P_q^\tau$  may cause a hallucination;
21: else
22:   Compute the semantic entropy  $H(Q_v^\tau)$  of the query;
23:   if  $H(Q_v^\tau) > \max(H(v_i^\tau))$  then
24:     Report  $P_q^\tau$  may cause a hallucination;
25:   else
26:     Return the response of the agent for  $P_q^\tau$ ;
27:   end if
28: end if

```

4 Experimental Evaluation

4.1 Datasets

To evaluate the performance of HalMit, two popular public Query-Answer (QA) datasets, MedQuAD [2] and SQuAD [22], are used.

- MedQuAD is a collection of question-answer pairs meticulously curated from 12 trusted National Institute of Health (NIH) websites and covers various medical topics including diseases, medications, and diagnostic tests.
- SQuAD consists of questions posed by crowd-workers on a set of Wikipedia articles, where the answer to every question is a segment of text or span, and each QA pair is paired with a title.

Without loss of generality, we randomly select four domains from each of these two datasets, including "Treatment", "Inheritance", "New York City", and "Modern History", to construct different types of agents. Since the response from each target agent may not be completely matched with the correct response in the QA pair from the database, it is necessary to determine whether the response is a hallucination or not. We follow existing work [25] to use the GQA metric to identify the responses that include hallucinations. More specifically, a binary label is assigned according to the average of unigram F1 and ROUGE-L. A hallucination is labeled if this average is less than 0.5 [4].

4.2 Evaluation Setup

To evaluate the performance of HalMit, six LLMs, including Llama (Llama2-7B-Instruct and Llama3.1-8B-Instruct), Mistral-7b, qwen2-1.5b, Falcon-7b and Vicuna-7b, are used to support agent inference.

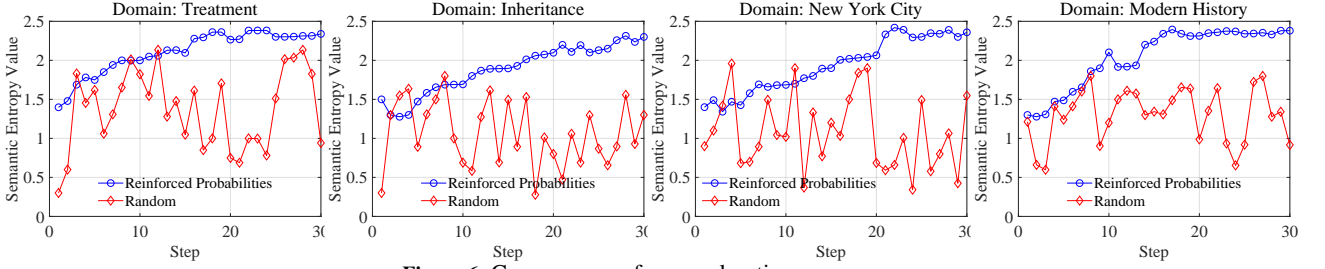


Figure 6: Convergence of our exploration process.

All agents in specific domains are implemented using RAG technology. Specifically, this RAG pipeline is constructed with Elasticsearch [5] served as the vector database, which is used to record generalized bounds in the vector format to serve as references for identifying conditions where hallucinations occur. In addition, the m3e-base [27] is used as the embedding model to vectorize the information of a generalized bound, including the query and the corresponding responses stored in the database that represent the bounds, as well as the input query in the hallucination monitoring process. We also utilize a repository within Agentscope V0.1.0 [8] to enable exploration of the generalization bound and monitoring of hallucinations of LLM-empowered agents.

During modeling the generalization bound, Qwen-max is used to generate queries, while GPT 4 is used to judge whether each response of the target LLM has hallucinations. In addition, we also incorporate a supervised method. When using GPT 4 for assistance in judgment, we evaluate the confidence in the inference results. If confidence is less than 60%, we perform a manual review to ensure the accuracy of the judgment before proceeding [3]. During training the policy network of reinforcement determination on fractal probability, the learning rate is set to 10^{-4} and the batch size is 64. It converges in 300 epochs. We set the ratio γ to 0.6 to guide the search for bounds, randomly initialize ten queries of the bound searching for each domain, and the similarity threshold ϵ to 0.8 for the monitoring of hallucinations. The impacts of setting the two parameters to different values are studied in Section 4.5.

Three popular hallucination detection methods are included as baselines: 1) Predictive Probability (PP) [20], 2) In-Context-Learning Prompt (ICL) [29], and 3) SelfCheckGPT (SCG) [20]. In addition, a comprehensive set of metrics is used in our evaluation, including: 1) the area under the receiver operator characteristic curve (AUROC), 2) the area under the precision recall curve (AUC-PR), 3) the F1 score, and 4) the accuracy. In addition, we also recorded the semantic entropy defined in Formula (2) in the Appendix of the output of the target agent to illustrate the uncertainty metric.

4.3 Convergence Study of The Exploration

Our probabilistic fractal exploration method is reinforced to explore towards the generalization bound. Through a reward mechanism based on increases in semantic entropy, our fractal exploration process is directed toward generating statements with higher uncertainty—indicating proximity to the generalization bound. To evaluate the completeness of fractal exploration, we compare the performance of using the reinforced determination of fractal probabilities against a baseline that uses randomly assigned probabilities.

As shown in Figure 6, we present the semantic entropy values over the final 30 steps of the exploration process. This can also be observed to investigate how semantic entropy evolves with the exploration of the generalization bound through fractal transformations. The results demonstrate that with reinforced fractal probability se-

lection, each exploration step consistently increases or maintains semantic entropy, signaling effective converge on the generalization bound. In contrast, the random probability strategy yields volatile entropy values, with no clear trend, indicating an unreliable and less directed search process. These findings suggest that reinforcement-guided fractal exploration offers a more robust and targeted approach to identifying generalization boundaries.

4.4 Effectiveness of Hallucination Monitoring

We evaluated the effectiveness of HalMit, and the results are presented in Table 1. HalMit achieves the best performance in Inheritance and Modern History domains, demonstrating both its superiority and adaptability to various types of agent. Specifically, our method improves the AUROC and AUC-PR metrics up to 8% over the best baseline, highlighting its effectiveness in distinguishing qualified output from hallucinations. The only exception is the New York City topic, where SelfCheckGPT outperforms our method in monitoring hallucinations. This may be due to the miscellaneous slangy dialogues on this topic, which SelfCheckGPT appears to be better equipped to handle.

Compared to baselines that use a fixed threshold for hallucination detecting, HalMit shows greater effectiveness across agent hallucinations in different domains. The Treatment and Inheritance domains primarily involve scientific knowledge, while New York City and Modern History consist of miscellaneous questions. HalMit performs particularly well in domains that allow for divergent responses, probably because it better adapts to the semantic diversity and complexity inherent in such queries. Finally, as a basic scheme only relying on LLM to detect hallucinations, ICL performs poorly for hallucination identification, while the used LLMs struggle to reliably detect their own hallucinations without external guidance.

To evaluate the effectiveness of HalMit in monitoring hallucinations among agents empowered by other LLM models, Table 2 presents the hallucination monitoring performance for agents with the other four models, Mistral-7b, Qwen2-1.5b, Falcon-7b and Vicuna-7b. Without loss of generalization, the agents in Treatment domain are selected to construct experiments. The experimental results demonstrate that HalMit consistently outperforms the baselines, achieving the highest accuracy and F1 scores. In particular, HalMit shows the most significant improvement for agents empowered by Qwen2-1.5b, reaching the highest accuracy of 0.85. PP and ICL show lower monitoring accuracy in agents with Vicuna, SelfCheckGPT, and HalMit experiences a marked increase in F1 score. These results demonstrate the superiority of HalMit as the most robust approach, enhancing its performance across a variety of LLM architectures.

4.5 Ablation Study

We investigate the impact of parameter γ and ϵ on the monitoring accuracy in the Inheritance domain. As shown in Figures 7a and 7b,

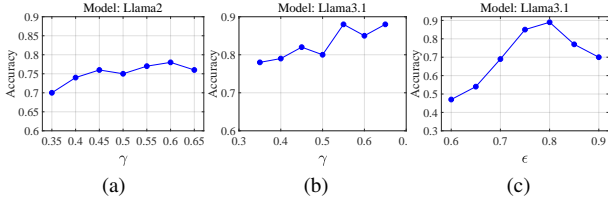
Table 1: Hallucination monitoring performance on different agents

	Backbone	AUROC↑	AUC-PR↑	F1↑	Acc↑	Backbone	AUROC↑	AUC-PR↑	F1↑	Acc↑
Treatment										
PP	Llama2	0.56	0.54	0.56	0.66	Llama3.1	0.56	0.69	0.6	0.65
ICL		0.59	0.56	0.60	0.55		0.48	0.71	0.7	0.61
SCG		0.71	0.72	0.69	0.7		0.74	0.84	0.7	0.75
HalMit		0.76	0.80	0.79	0.73		0.80	0.86	0.82	0.88
Inheritance										
PP	Llama2	0.68	0.60	0.77	0.75	Llama3.1	0.71	0.79	0.72	0.71
ICL		0.69	0.67	0.59	0.56		0.61	0.73	0.65	0.65
SCG		0.70	0.73	0.68	0.75		0.85	0.84	0.85	0.85
HalMit		0.70	0.79	0.8	0.78		0.90	0.86	0.82	0.88
New York City										
PP	Llama2	0.54	0.12	0.21	0.74	Llama3.1	0.60	0.58	0.52	0.53
ICL		0.59	0.22	0.32	0.75		0.58	0.64	0.63	0.45
SCG		0.82	0.72	0.85	0.82		0.86	0.84	0.87	0.86
HalMit		0.88	0.77	0.75	0.89		0.89	0.82	0.88	0.84
Modern History										
PP	Llama2	0.74	0.76	0.73	0.79	Llama3.1	0.48	0.45	0.46	0.56
ICL		0.69	0.67	0.64	0.67		0.59	0.55	0.54	0.61
SCG		0.78	0.79	0.81	0.82		0.78	0.77	0.6	0.76
HalMit		0.84	0.80	0.77	0.89		0.84	0.84	0.67	0.89

Table 2: Additional results for evaluating hallucination monitoring performance

Model	Method	AUROC↑	AUC-PR↑	Acc↑	F1↑
Mistral-7b	PP	0.43	0.49	0.56	0.37
	ICL	0.58	0.52	0.59	0.49
	SCG	0.56	0.67	0.69	0.64
	HalMit	0.69	0.70	0.79	0.77
qwen2-1.5b	PP	0.49	0.50	0.50	0.42
	ICL	0.52	0.39	0.54	0.44
	SCG	0.72	0.80	0.78	0.68
	HalMit	0.79	0.80	0.85	0.81
Falcon-7b	PP	0.64	0.59	0.49	0.44
	ICL	0.57	0.60	0.51	0.55
	SCG	0.68	0.71	0.75	0.79
	HalMit	0.73	0.75	0.79	0.80
Vicuna-7b	PP	0.39	0.41	0.45	0.67
	ICL	0.60	0.58	0.56	0.79
	SCG	0.68	0.72	0.71	0.81
	HalMit	0.75	0.71	0.84	0.89

when γ varies from 0.35 to 0.65, the monitoring accuracy remains consistently high, maintaining a level between 0.78 and 0.88. The curve exhibits remarkable stability across different γ values, with only a slight fluctuation. The stable performance across a wide range of γ values demonstrates that our proposed monitoring method is relatively insensitive to the choice of γ . This robustness is particularly valuable for practical applications, as it suggests that the method can maintain reliable performance without requiring precise fine-tuning of γ . Similarly, the impact of the parameter ϵ on the accuracy of the monitoring is evaluated, and the tested values of ϵ vary between 0.6 and 0.9. As illustrated in Figure 7c, ϵ increases from 0.6 to 0.8, the accuracy improves, and the highest performance is achieved at ϵ of 0.8. These observations suggest that $\epsilon = 0.8$ strikes an optimal trade-off, yielding the best performance.


Figure 7: Ablation Study of γ and ϵ .

5 Related Work

Most of the related approaches take LLM as a white box to detect hallucinations. Ji et al. [14] used a mutual information-based feature selection method to select sensitive neurons from the last activation

layer and trained a classifier using Llama MLP. Han et al. [10] proposed a method to approximate semantic entropy from the hidden state of the model, converting the semantic entropy into binary labels, and trained a logistic regression classifier to predict hallucinations. Zhu et al. [38] used PCA to reduce the dimensionality of the hidden layer embedding and adopted interval partitioning or GMM clustering to establish abstract states. In addition, they also used Markov models and hidden Markov models to capture state transitions, and used a small amount of annotated reference data to link internal state transitions with hallucination/factual output behaviors. He et al. [11] combined the static features within the model with the dynamic features and used Siamese networks to identify situations where the answers of the large language model deviate from the facts. Gaurang et al. [24] proposes a hallucination detection method that analyze internal model signals. Xiaoling et al. [37] propose HADEMIF for detecting hallucinations in LLMs by leveraging a Deep Dynamic Decision Tree and an MLP to calibrate model predictions. Requiring the access to internal states of LLMs to detect hallucinations, these methods not only suffer from high complexity and computational demand but also may not be feasible for commercial LLM software. This highlights the need for research on the detection of hallucinations without accessing the internal states of the LLMs.

The other solutions detect hallucinations whereas considering the black-box nature of LLMs. Hou et al. [12] proposed using the belief tree, a probabilistic framework, to detect hallucinations using the logical consistency between model beliefs. Quevedo et al. [21] extracted features using two LLMs and used these features to train logistic regression and simple neural networks to detect hallucinations. Although these methods detect hallucination through output features or associated confidence scores, the limited understanding of the generalization bound of LLMs leads to poorly calibrated confidence estimates.

6 Conclusion

In this work, we present an in-depth study of the hallucination phenomenon in LLM-empowered agents through the lens of fine-grained domains. Based on our findings, we propose an effective and efficient hallucination monitor HalMit, according to a few key observations from our study. Our research reveals that LLMs exhibit similar generalization bounds within the same domain, providing a foundation for accurately monitoring hallucinations in specific domains. To take advantage of this key insight, we have designed a reinforced probabilistic fractal exploration method that efficiently identifies the general-

ization bound of an LLM-empowered agent within a domain. Despite the black-box nature of LLM-empowered agents, this approach significantly accelerates the boundary identification process while improving both the accuracy and efficiency of hallucination monitoring based on the generalization bound. Extensive experimental evaluations demonstrate that our method outperforms existing mainstream hallucination monitoring techniques across multi-topic datasets and different foundation LLMs. This work not only offers a novel technical pathway for monitoring hallucinations for agents in a specific domain, but also provides robust theoretical and practical support to enhance their security and dependability in critical applications.

References

- [1] K. Andriopoulos and J. Pouwelse. Augmenting llms with knowledge: A survey on hallucination prevention. *arXiv preprint arXiv:2309.16459*, 2023.
- [2] A. Ben Abacha and D. Demner-Fushman. A question-entailment approach to question answering. *BMC Bioinf*, 20, 2019.
- [3] S. Casper, X. Davies, C. Shi, T. K. Gilbert, J. Scheurer, J. Rando, R. Freedman, T. Korbak, D. Lindner, P. Freire, et al. Open problems and fundamental limitations of reinforcement learning from human feedback. *arXiv preprint arXiv:2307.15217*, 2023.
- [4] Y. Chen, Q. Fu, Y. Yuan, Z. Wen, G. Fan, D. Liu, D. Zhang, Z. Li, and Y. Xiao. Hallucination detection: Robustly discerning reliable answers in large language models. In *CIKM*, 2023.
- [5] Elastic. Elasticsearch. <https://www.elastic.co/guide/en/elasticsearch/reference/current/elasticsearch-intro-what-is-es.html>, 2023. Accessed: 2023-10-10.
- [6] S. Farquhar, J. Kossen, L. Kuhn, and Y. Gal. Detecting hallucinations in large language models using semantic entropy. *Nature*, 630(8017), 2024.
- [7] O. Freyer, I. C. Wiest, J. N. Kather, and S. Gilbert. A future role for health applications of large language models depends on regulators enforcing safety standards. *Lancet Digit Health*, 6(9), 2024.
- [8] D. Gao, Z. Li, X. Pan, W. Kuang, Z. Ma, B. Qian, F. Wei, W. Zhang, Y. Xie, D. Chen, L. Yao, H. Peng, Z. Y. Zhang, L. Zhu, C. Cheng, H. Shi, Y. Li, B. Ding, and J. Zhou. Agentscope: A flexible yet robust multi-agent platform. *CoRR*, abs/2402.14034, 2024.
- [9] C. M. Greco and A. Tagarelli. Bringing order into the realm of transformer-based language models for artificial intelligence and law. *Artif Intell Law*, 2023.
- [10] J. Han, J. Kossen, M. Razzak, L. Schut, S. A. Malik, and Y. Gal. Semantic entropy probes: Robust and cheap hallucination detection in llms. In *ICML*, 2024.
- [11] J. He, Y. Gong, Z. Lin, C. Wei, Y. Zhao, and K. Chen. Llm factoscope: Uncovering llms’ factual discernment through measuring inner states. In *ACL Findings*, 2024.
- [12] B. Hou, Y. Zhang, J. Andreas, and S. Chang. A probabilistic framework for llm hallucination detection via belief tree propagation. *arXiv preprint arXiv:2406.06950*, 2024.
- [13] L. Huang, W. Yu, W. Ma, W. Zhong, Z. Feng, H. Wang, Q. Chen, W. Peng, X. Feng, B. Qin, et al. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Trans Inf Syst*, 2023.
- [14] Z. Ji, D. Chen, E. Ishii, S. Cahyawijaya, Y. Bang, B. Wilie, and P. Fung. Llm internal states reveal hallucination risk faced with a query. *arXiv preprint arXiv:2407.03282*, 2024.
- [15] E. Kasneci, K. Seßler, S. Küchemann, M. Bannert, D. Dementieva, F. Fischer, U. Gasser, G. Groh, S. Günemann, E. Hüllermeier, et al. Chatgpt for good? on opportunities and challenges of large language models for education. *Learn Individ Differ*, 103, 2023.
- [16] L. Kuhn, Y. Gal, and S. Farquhar. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. *arXiv preprint arXiv:2302.09664*, 2023.
- [17] S. Lin, J. Hilton, and O. Evans. Truthfulqa: Measuring how models mimic human falsehoods. In *ACL*, 2022.
- [18] S. Lotfi, M. Finzi, Y. Kuang, T. G. Rudner, M. Goldblum, and A. G. Wilson. Non-vacuous generalization bounds for large language models. *arXiv preprint arXiv:2312.17173*, 2023.
- [19] S. Lotfi, Y. Kuang, B. Amos, M. Goldblum, M. Finzi, and A. G. Wilson. Unlocking tokens as data points for generalization bounds on larger language models. *arXiv preprint arXiv:2407.18158*, 2024.
- [20] P. Manakul, A. Liusie, and M. Gales. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models. In *EMNLP*, 2023.
- [21] E. Quevedo, J. Yero, R. Koerner, P. Rivas, and T. Cerny. Detecting hallucinations in large language model generation: A token probability approach. *arXiv preprint arXiv:2405.19648*, 2024.
- [22] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang. SQuAD: 100,000+ questions for machine comprehension of text. In *EMNLP*, 2016.
- [23] N. Riemer. Remetonymizing metaphor: Hypercategories in semantic extension. 2002.
- [24] G. Sriramanan, S. Bharti, V. S. Sadasivan, S. Saha, P. Kattakinda, and S. Feizi. Llm-check: Investigating detection of hallucinations in large language models. *Adv Neural Inf Process Syst*, 37:34188–34216, 2024.
- [25] D. Su, X. Li, J. Zhang, L. Shang, X. Jiang, Q. Liu, and P. Fung. Read before generate! faithful long form question answering with machine reading. In *ACL Foundings*, 2022.
- [26] D. Udekwe, O. ofe Ajayi, O. Ubadike, K. Ter, and E. Okafor. Comparing actor-critic deep reinforcement learning controllers for enhanced performance on a ball-and-plate system. *Expert Syst Appl*, 245, 2024. ISSN 0957-4174.
- [27] Y. Wang, Q. Sun, and S. He. M3E: Moka Massive Mixed Embedding Model, 2023.
- [28] J. Wei, Y. Yao, J.-F. Ton, H. Guo, A. Estornell, and Y. Liu. Measuring and reducing llm hallucination without gold-standard answers via expertise-weighting. *arXiv preprint arXiv:2402.10412*, 2024.
- [29] S. Xu and C. Zhang. Misconfidence-based demonstration selection for llm in-context learning. *arXiv preprint arXiv:2401.06301*, 2024.
- [30] H. Yang, H. Lu, W. Lam, and D. Cai. Exploring compositional generalization of large language models. In *NAACL-HLT*, pages 16–24, 2024.
- [31] J. Yang, H. Jin, R. Tang, X. Han, Q. Feng, H. Jiang, S. Zhong, B. Yin, and X. Hu. Harnessing the power of llms in practice: A survey on chatgpt and beyond. *ACM Trans Knowl Discov Data*, 18(6), 2024.
- [32] C. Zhang. User-controlled knowledge fusion in large language models: Balancing creativity and hallucination. *arXiv preprint arXiv:2307.16139*, 2023.
- [33] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.
- [34] Y. Zhang, S. Li, J. Liu, P. Yu, Y. R. Fung, J. Li, M. Li, and H. Ji. Knowledge overshadowing causes amalgamated hallucination in large language models. *arXiv preprint arXiv:2407.08039*, 2024.
- [35] H. Zhao, Z. Liu, Z. Wu, Y. Li, T. Yang, P. Shu, S. Xu, H. Dai, L. Zhao, G. Mai, et al. Revolutionizing finance with llms: An overview of applications and insights. *arXiv preprint arXiv:2401.11641*, 2024.
- [36] Z. Zhao, B. Wang, L. Ouyang, X. Dong, J. Wang, and C. He. Beyond hallucinations: Enhancing lvlms through hallucination-aware direct preference optimization. *arXiv preprint arXiv:2311.16839*, 2023.
- [37] X. Zhou, M. Zhang, Z. Lee, W. Ye, and S. Zhang. Hademif: Hallucination detection and mitigation in large language models. In *ICLR*.
- [38] D. Zhu, D. Chen, Q. Li, Z. Chen, L. Ma, J. Grossklags, and M. Fritz. Pollmgraph: Unraveling hallucinations in large language models via state transition dynamics. *arXiv preprint arXiv:2404.04722*, 2024.