PERTURBATIONS MATTER: SENSITIVITY-GUIDED HALLUCINATION DETECTION IN LLMS

Anonymous authors

Paper under double-blind review

ABSTRACT

Hallucination detection is essential for ensuring the reliability of large language models. Internal representation-based methods have emerged as the prevailing direction for detecting hallucinations, yet the internal representations often fail to yield clear separability between truthful and hallucinatory content. To address this challenge, we study the separability of the sensitivity to prompt-induced perturbations in the internal representations. A theory is established to show that, with non-negligible probability, each sample admits a prompt under which truthful samples exhibit greater sensitivity to prompt-induced perturbations than hallucinatory samples. When the theory is applied to the representative datasets, the probability reaches nearly 99%, suggesting that sensitivity to perturbations provides a discriminative indicator. Building on this insight, we propose a theory-informed method Sample-Specific Prompting (SSP), which adaptively selects prompts to perturb the model's internal states and measures the resulting sensitivity as a detection indicator. Extensive experiments across multiple benchmarks demonstrate that SSP consistently outperforms existing hallucination detection methods, validating the practical effectiveness of our method SSP in hallucination detection.

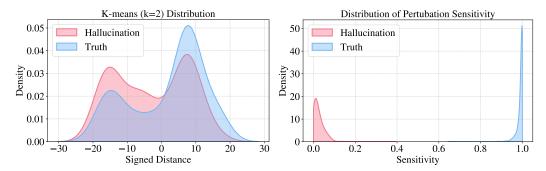
1 Introduction

Large language models (LLMs) have shown remarkable performance in natural language understanding and generation tasks (Achiam et al., 2023; Grattafiori et al., 2024). However, hallucination in generated text remains a critical challenge, arising when LLMs produce outputs that are grammatically and logically coherent but lack factual accuracy or verifiable evidence (Joshi et al., 2017; Lin et al., 2022a). Such hallucinations undermine user trust and pose risks in high-stakes areas such as healthcare, law, and scientific research (Ji et al., 2023; Liu et al., 2024b). To address this issue, hallucination detection has attracted extensive attention in recent research (Manakul et al., 2023).

Previous detection methods can be roughly divided into two main categories: self-assessment (Kadavath et al., 2022; Zhou et al., 2023; Lin et al., 2022b) and internal representation-based methods (Du et al., 2024; Azaria & Mitchell, 2023; Marks & Tegmark, 2024; Yin et al., 2024). Self-assessment estimates the factuality of a response by leveraging the confidence in the model output. Internal representation-based methods primarily leverage the embeddings of off-the-shelf LLMs to classify outputs as either truthful or hallucinatory. The internal representation-based methods generally outperform self-assessment, and thus have emerged as the prevailing research direction.

Despite notable progress, the internal representation-based methods (Yin et al., 2024; Du et al., 2024; Kossen et al., 2024) face fundamental bottlenecks for detection, which impose inherent limits on their future development. Recent work (Park et al., 2025) demonstrates that the internal representations of LLMs frequently *fail to provide a clear separation between truthful and hallucinatory content* (see Figure 1a). As a result, the effectiveness of internal representation-based methods is inherently limited by the separability of internal representations. This motivates the a critical question: *is it possible to overcome the inherent separability bottleneck of internal representations?*

To tackle this question, we start from an empirical observation: in the experimental setup of Figure 1b, using the sensitivity of prompt-induced perturbations in internal representations as an evaluation score yields *near-perfect separability* between truthful and hallucinatory samples. To demystify this insightful observation, we develop a theory (see Section 3) stating that *for each sample, there exists an associated prompt, and with non-negligible probability, the sensitivity of truthful samples*



- (a) K-Means Clustering of Internal Representations.
- (b) Prompt-Induced Perturbation Sensitivity.

Figure 1: Empirical analysis conducted on 200 randomly selected samples from the TruthfulQA dataset (Lin et al., 2022a). (a) Distribution of internal representations obtained via K-means clustering (k=2). The signed distance is defined as the difference between a sample's distances to the two cluster centroids, revealing weak separation. (b) For each sample, we apply an individually optimized prompt perturbation and measure its sensitivity using the cosine similarity between representations before and after perturbation. We find that this sensitivity provides effective separability between truthful and hallucinatory samples. Details for (a) and (b) are provided in **Appendix A**.

to prompt-induced perturbations exceeds that of the hallucinatory samples. We further apply our theory on the representative datasets (Reddy et al., 2019; Lin et al., 2022a; Joshi et al., 2017; Clark et al., 2020), showing that the probability reaches nearly 99%, thereby statistically guaranteeing that the sensitivity to prompt-induced perturbations in internal representations, when used as an evaluation score, does not suffer from the separability bottleneck.

In light of the above analysis, we propose a novel method *Sample-Specific Prompting* (SSP), which leverages the sensitivity to prompt-induced perturbations as a discriminative indicator for hallucination detection. Instead of relying on static or handcrafted prompts, SSP dynamically generates tailored prompts for each question–answer pair to enhance the sensitivity of truthful samples to perturbations while reducing that of hallucinatory ones. Furthermore, SSP introduces a lightweight encoder to extract features before and after perturbation and employs a contrastive training objective that encourages larger representation shifts for truthful samples and smaller shifts for hallucinated ones. In effect, the joint learning of perturbation prompts and representation encodings makes SSP a more effective method for exploiting prompt-induced perturbations in hallucination detection.

Extensive experiments demonstrate the effectiveness of SSP across diverse datasets CoQA (Reddy et al., 2019), TruthfulQA (Lin et al., 2022a), TriviaQA (Joshi et al., 2017) and TydiQA-GP (Clark et al., 2020), compared with the state-of-the-art (Kadavath et al., 2022; Azaria & Mitchell, 2023; Hu et al., 2024). Also, our results indicate that SSP generalizes well across different domains. Our main contributions are summarized as follows:

- We are the first to leverage the sensitivity of LLM internal representations to input perturbations for hallucination detection, offering a novel perspective to hallucination detection.
- We analyze the sensitivity of LLM internal representations to input perturbations and theoretically establish its feasibility for hallucination detection.
- We propose a theory-informed method SSP, which leverages sensitivity to prompt-induced perturbations as a discriminative indicator for hallucination detection.

2 Preliminary

LLMs and Token Sequences. Following Oh et al. (2025); Du et al. (2024), we use a distribution $P_{\theta}(\cdot)$ over token sequences to define LLM, where θ is the model parameters. Given a token sequence $\mathbf{Q} = [x_1, \dots, x_k]$ representing the question, where each x_i is the *i*-th token in the sequence. $P_{\theta}(\cdot)$ generates an answer $\mathbf{A} = [x_{k+1}, \dots, x_{k+q}]$ by predicting each token based on the preceding context:

$$P_{\theta}(x_i|x_1,\dots,x_{i-1}), \text{ for } i=k+1,\dots,k+q.$$
 (1)

Truthful-answer and Hallucinatory-answer Domains. Let Q and A denote the spaces of questions and answers, respectively. We introduce two domains over $Q \times A$:

- The truthful-answer domain is a joint distribution $P_{Q,T}$, where $Q \in \mathcal{Q}$ is a random variable representing questions and $T \in \mathcal{A}$ is a random variable representing the truthful answers.
- The hallucinatory-answer domain is a joint distribution $P_{Q,H}$, where Q is defined as above and $H \in \mathcal{A}$ is a random variable representing the hallucinatory answers.

Dataset Format. Given the truthful-answer domain $P_{Q,T}$, each sample sampled from $P_{Q,T}$ consists of a question \mathbf{Q} and a reference answer $\mathbf{A}^{\mathrm{ref}}$. The dataset sampled from $P_{Q,T}$ can be expressed as $\mathcal{D} = \{(\mathbf{Q}_1, \mathbf{A}_1^{\mathrm{ref}}), \dots, (\mathbf{Q}_n, \mathbf{A}_n^{\mathrm{ref}})\}$, where n is the number of samples.

Given a question $\mathbf{Q} \sim P_Q$, the LLM $P_{\theta}(\cdot)$ generates an answer $\mathbf{A} \sim P_{\theta}(\cdot|\mathbf{Q})$. Each generated answer \mathbf{A} is assigned a binary label $y \in \{-1,1\}$ according to its semantic consistency with the reference answer \mathbf{A}^{ref} . Specifically, if \mathbf{A} aligns with \mathbf{A}^{ref} , it is labeled as truthful (y=1); otherwise, it is labeled as hallucinatory (y=-1). The labeled dataset \mathcal{D}_l is thus defined as:

$$\mathcal{D}_l = \{ (\mathbf{Q}_1, \mathbf{A}_1, y_1), \dots, (\mathbf{Q}_n, \mathbf{A}_n, y_n) \}. \tag{2}$$

AUROC and **Separability.** The AUROC serves as the primary evaluation metric for hallucination detection (Du et al., 2024). Formally, given the truthful-answer domain $P_{Q,T}$ and the hallucinatory-answer domain $P_{Q,H}$, the AUROC of a scoring function $r: \mathcal{Q} \times \mathcal{A} \to \mathbb{R}$ is defined as follows:

$$AUROC(r; P_{Q,T}, P_{Q,H}) = P(r(\mathbf{Q}, \mathbf{T}) > r(\mathbf{Q}', \mathbf{H}')) + \frac{1}{2}P(r(\mathbf{Q}, \mathbf{T}) = r(\mathbf{Q}', \mathbf{H}')), \quad (3)$$

where $(\mathbf{Q}, \mathbf{T}) \sim P_{Q,T}$ is the truthful sample, and $(\mathbf{Q}', \mathbf{H}') \sim P_{Q,H}$ is the hallucinatory sample. In this work, we define the separability via the core component of AUROC, formally given by:

$$SEP(r; P_{Q,T}, P_{Q,H}) = P(r(\mathbf{Q}, \mathbf{T}) > r(\mathbf{Q}', \mathbf{H}')). \tag{4}$$

Hallucination Detection. Given the training dataset $\mathcal{D}_l = \{(\mathbf{Q}_1, \mathbf{A}_1, y_1), \dots, (\mathbf{Q}_n, \mathbf{A}_n, y_n)\}$ as introduced in Eq. (2), the goal of hallucination detection is to learn a detector G, based on a given LLM $P_{\theta}(\cdot)$ and \mathcal{D}_l , such that for any question $\mathbf{Q} \sim P_Q$ and a corresponding answer \mathbf{A} ,

$$G(\mathbf{Q}, \mathbf{A}) = 1$$
, if $\mathbf{A} \sim P_{T|O}(\cdot \mid \mathbf{Q})$; otherwise, $G(\mathbf{Q}, \mathbf{A}) = -1$, (5)

where 1 indicates that **A** is truthful, and -1 indicates that **A** is hallucinatory.

Due to space constraints, the related work is discussed in **Appendix B**.

3 Separability of Prompt-Induced Perturbation Sensitivity

Before introducing our method, we first analyze the separability of perturbation sensitivity in this section. Due to space constraints, all proofs are provided in **Appendix** C.

3.1 Sensitivity of Prompt-Induced Perturbations

Recent prevailing methods for hallucination detection (Azaria & Mitchell, 2023; Marks & Tegmark, 2024; Yin et al., 2024; Du et al., 2024; Kossen et al., 2024) rely on internal representations, classifying outputs as truthful or hallucinatory by leveraging embeddings from pre-trained LLMs. However, as pre-trained LLMs are trained for next-token prediction, their embeddings inherently favour fluency and syntactic correctness, while often overlooking truthful accuracy (Radford et al., 2019). Motivated by this limitation, recent work (Park et al., 2025) claims that the internal representations of LLMs frequently fail to provide a clear separation between truthful and hallucinatory samples.

In Figure 1a, we validate the claim given by Park et al. (2025). As shown in Figure 1a, the last-token embeddings of truthful and hallucinatory samples from TruthfulQA (Lin et al., 2022a) largely overlap, highlighting the lack of a clear separation. Hence, the effectiveness of these internal representation-based methods (Azaria & Mitchell, 2023; Marks & Tegmark, 2024; Yin et al., 2024; Du et al., 2024; Kossen et al., 2024) is limited by the separability of the internal representations. In light of this, we raise the question of whether it is possible to overcome the inherent separability

bottleneck of internal representations. To tackle this question, we study whether prompt-induced perturbation sensitivity in internal representations has the potential for strong separability.

Formalizing Perturbation Sensitivity. Following prior work (Du et al., 2024; Park et al., 2025; Azaria & Mitchell, 2023; Chen et al., 2024; Guo et al., 2021), we define the internal representation $\mathbf{E}_{\theta}(\cdot)$ of the LLM $P_{\theta}(\cdot)$ as the embedding of the last token. Given a prompt \mathbf{P} , and a question-answer pair (\mathbf{Q} , \mathbf{A}), the prompt-induced perturbation sensitivity is defined as follows:

$$\Delta \mathbf{E}_{\theta}(\mathbf{Q}, \mathbf{A}, \mathbf{P}) = \mathrm{Dist}(\mathbf{E}_{\theta}(\mathbf{Q}, \mathbf{A}, \mathbf{P}), \mathbf{E}_{\theta}(\mathbf{Q}, \mathbf{A})), \tag{6}$$

where $Dist(\cdot, \cdot)$ is the measure of the difference between $E_{\theta}(\mathbf{Q}, \mathbf{A}, \mathbf{P})$ and $E_{\theta}(\mathbf{Q}, \mathbf{A})$.

Preliminary Observation of Perturbation Sensitivity. To investigate the separability of perturbation sensitivity, we construct an *oracle setting* in which, for each sample, we optimize a corresponding prompt, such that the perturbation sensitivity is maximized when the answer is truthful, and minimized when the answer is hallucinatory, i.e., for any sample $(\mathbf{Q}_i, \mathbf{A}_i, y_i) \in \mathcal{D}_l$,

$$\mathbf{P}_{i}^{*} \in \arg\max_{\mathbf{P}} \ y_{i} \cdot \Delta \mathbf{E}_{\theta}(\mathbf{Q}_{i}, \mathbf{A}_{i}, \mathbf{P}). \tag{7}$$

In Figure 1b, we present the empirical result under the oracle setting (see **Appendix A** for experimental details). We observe that the separability of perturbation sensitivity reaches nearly 100%, which implies the aspiration of addressing the separability bottleneck of the internal representations.

3.2 Separability of Perturbation Sensitivity

Here, we develop a statistical analysis that characterizes the separability of perturbation sensitivity. We continue to consider the oracle setting, where the prompt is chosen to maximize perturbation sensitivity for truthful answers and minimize it for hallucinatory ones. Given the truthful-answer domain $P_{Q,T}$ and the hallucinatory-answer domain $P_{Q,H}$, we select the optimal prompt as follows:

$$\mathbf{P}^* \in \arg\max_{\mathbf{P}} \ y(\mathbf{Q}, \mathbf{A}) \cdot \Delta \mathbf{E}_{\theta}(\mathbf{Q}, \mathbf{A}, \mathbf{P}), \tag{8}$$

where $y(\mathbf{Q}, \mathbf{A}) = 1$ if $(\mathbf{Q}, \mathbf{A}) \sim P_{Q,T}$, and $y(\mathbf{Q}, \mathbf{A}) = -1$ if $(\mathbf{Q}, \mathbf{A}) \sim P_{Q,H}$. Then, we consider the scoring function $r^* : \mathcal{Q} \times \mathcal{A} \to \mathbb{R}$, i.e.,

$$r^*(\mathbf{Q}, \mathbf{A}) = \Delta \mathbf{E}_{\theta}(\mathbf{Q}, \mathbf{A}, \mathbf{P}^*), \text{ where } \mathbf{P}^* \text{ is defined in Eq. (8)}.$$
 (9)

The scoring function r^* estimates the perturbation sensitivity under the oracle setting.

Probabilistic Characterization of Separability. Here, we give our core theorem, i.e., Theorem 1.

Theorem 1 (Separability of Perturbation Sensitivity.) Given the truthful-answer domain $P_{Q,T}$ and the hallucinatory-answer domain $P_{Q,H}$, if the scoring function r^* given in Eq. (9) satisfies:

$$\frac{\mathbb{E}_{(\mathbf{Q},\mathbf{A})\sim P_{Q,H}}r^*(\mathbf{Q},\mathbf{A})}{\mathbb{E}_{(\mathbf{Q},\mathbf{A})\sim P_{Q,T}}r^*(\mathbf{Q},\mathbf{A})} \leq \frac{1}{a}, \quad \frac{\sigma_{(\mathbf{Q},\mathbf{A})\sim P_{Q,T}}r^*(\mathbf{Q},\mathbf{A})}{\sigma_{(\mathbf{Q},\mathbf{A})\sim P_{Q,H}}r^*(\mathbf{Q},\mathbf{A})} \leq b, \quad \frac{\sigma_{(\mathbf{Q},\mathbf{A})\sim P_{Q,H}}r^*(\mathbf{Q},\mathbf{A})}{\mathbb{E}_{(\mathbf{Q},\mathbf{A})\sim P_{Q,H}}r^*(\mathbf{Q},\mathbf{A})} \leq c, \quad (10)$$

for some constants a > 1, b > 0, c > 0, where \mathbb{E} is the expectation and σ is the standard deviation,

then,
$$AUROC(r^*; P_{Q,T}, P_{Q,H}) \ge SEP(r^*; P_{Q,T}, P_{Q,H}) \ge \frac{(a-1)^2}{(a-1)^2 + (1+b^2)c^2}$$
. (11)

Theorem 1 establishes that, for each sample, there exists an associated prompt under which, with non-negligible probability, the sensitivity of truthful samples to prompt-induced perturbations exceeds that of hallucinatory samples. Theorem 1 further shows that, under the oracle setting, the AUROC of the prompt-induced perturbation sensitivity is bounded below by a computable probability, which becomes explicit when the indicators a, b, and c in Eq. (10) are available. This observation motivates us to apply Theorem 1 to representative datasets, thereby providing a quantitative estimate of the likelihood that the truthful samples exhibit greater perturbation sensitivity than the hallucinatory ones.

Validation of Separability. The preliminary observation in Figure 1b suggests that the perturbation sensitivity may exhibit strong separability. To further validate this observation, we first estimate the indicators a, b, and c in Eq. (10) through experiments on four representative datasets: CoQA, TruthfulQA, TriviaQA, and TyDiQA-GP (Reddy et al., 2019; Lin et al., 2022a; Joshi et al., 2017; Clark et al., 2020). Yet, when dealing with large-scale data, it is computationally infeasible to train an optimal prompt for each sample based on Eq. (8). To address this issue, we establish Theorem 2.



265

266

267

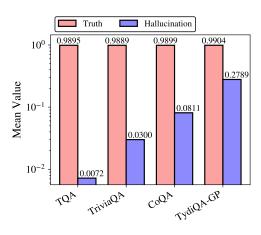
268

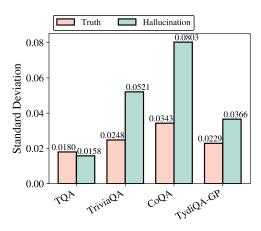
269

216

217

218219





- (a) Mean Value of Perturbation Sensitivity.
- (b) Standard Deviation of Perturbation Sensitivity.

Figure 2: Perturbation sensitivity (r^* in Eq. (9)) statistics across multiple datasets. The sensitivity of internal representations to prompt-induced perturbations is compared between truthful and hallucinatory samples across four representative datasets using LLaMA-3-8B-Instruct. Figure (a) reports the mean values, showing that truthful samples exhibit significantly larger average perturbation magnitudes than hallucinatory samples. Figure (b) presents the corresponding standard deviations, which remain small overall. Please see **Appendix D** for more details.

Theorem 2 Let $\mathbf{M}_{\varphi}(\cdot)$ be a model that receives a question–answer pair (\mathbf{Q}, \mathbf{A}) and a prompt \mathbf{P} as input, and returns a sample-specific prompt \mathbf{P}_{φ} as output, i.e., $\mathbf{P}_{\varphi} = \mathbf{M}_{\varphi}(\mathbf{Q}, \mathbf{A}, \mathbf{P})$. Also, let $r_{\varphi}(\mathbf{Q}, \mathbf{A}) = \Delta \mathbf{E}_{\theta}(\mathbf{Q}, \mathbf{A}, \mathbf{P}_{\varphi})$ and let a_{φ} , b_{φ} , c_{φ} be

$$a_{\varphi} = \frac{\mathbb{E}_{(\mathbf{Q}, \mathbf{A}) \sim P_{Q, T}} r_{\varphi}(\mathbf{Q}, \mathbf{A})}{\mathbb{E}_{(\mathbf{Q}, \mathbf{A}) \sim P_{Q, H}} r_{\varphi}(\mathbf{Q}, \mathbf{A})}, \quad b_{\varphi} = \frac{\sigma_{(\mathbf{Q}, \mathbf{A}) \sim P_{Q, T}} r_{\varphi}(\mathbf{Q}, \mathbf{A})}{\sigma_{(\mathbf{Q}, \mathbf{A}) \sim P_{Q, H}} r_{\varphi}(\mathbf{Q}, \mathbf{A})}, \quad c_{\varphi} = \frac{\sigma_{(\mathbf{Q}, \mathbf{A}) \sim P_{Q, H}} r_{\varphi}(\mathbf{Q}, \mathbf{A})}{\mathbb{E}_{(\mathbf{Q}, \mathbf{A}) \sim P_{Q, H}} r_{\varphi}(\mathbf{Q}, \mathbf{A})}.$$

Then the scoring function r^* defined in Eq. (9) satisfies that

$$AUROC(r^*; P_{Q,T}, P_{Q,H}) \ge SEP(r^*; P_{Q,T}, P_{Q,H}) \ge \max_{\varphi \text{ with } a_{\varphi > 1}} \frac{(a_{\varphi} - 1)^2}{(a_{\varphi} - 1)^2 + (1 + b_{\varphi}^2)c_{\varphi}^2}. \quad (12)$$

In Theorem 2, Eq. (12) provides an executable alternative to compute the lower bound in Theorem 1. For estimating the lower bound, following Eq. (12), we design the following optimization problem:

$$\max_{\varphi} \mathcal{L}(\varphi) = \log a_{\varphi} + 2\mu \log \left[\text{ReLU}(a_{\varphi} - 1) + 10^{-12} \right]$$

$$-\mu \log \left[(a_{\varphi} - 1) \text{ReLU}(a_{\varphi} - 1) + (1 + b_{\varphi}^{2})c_{\varphi}^{2} \right], \text{ where } \mu > 0 \text{ is the parameter.}$$
(13)

Details of the experimental implementation can be found in **Appendix D**. The experimental results are presented in Figure 2, which shows the mean values (see Figure 2a) and standard deviations (see Figure 2b) of the perturbation sensitivity r^* across different datasets. According to the experimental results in Figure 2, Theorem 2 implies that, in the four datasets CoQA, TruthfulQA, TriviaQA, and TydiQA-GP, if we select the prompt P^* for any sample (Q, A) according to Eq. (8), then the perturbation sensitivity $r^*(Q, A) = \Delta E_{\theta}(Q, A, P^*)$ exhibits near-perfect separability and AUROC:

$$AUROC(r^*; P_{O,T}, P_{O,H}) \ge SEP(r^*; P_{O,T}, P_{O,H}) \ge 99\%.$$
(14)

The above result demonstrates that, at least for the four representative datasets, each sample admits a prompt under which the prompt-induced perturbation sensitivity achieves nearly perfect separability.

Remark. Eq. (14) suggests that the separability and AUROC are lower bounded by 99%, which may appear inconsistent with the empirical results in Table 1. This discrepancy arises because Eq. (14) is computed over the entire dataset based on Eq. (13), and thus serves as an oracle value designed to demonstrate the potential separability of perturbation sensitivity. In practice, however, models are trained on limited data, and their performance on unseen test sets inevitably depends on generalization. Consequently, Eq. (14) should be interpreted as an indicator of the theoretical potential separability of perturbation sensitivity, rather than a direct guarantee of test-time performance.

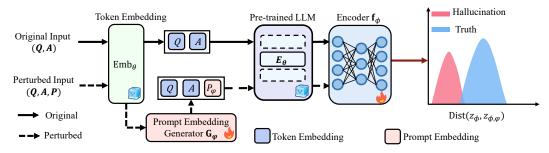


Figure 3: Overview of SSP. Given a question–answer pair, prompt embedding generator G_{φ} generates a perturbation appended to the input. Encoder f_{φ} then maps the intermediate representations to a discriminative space and maximize the discrepancy between truthful and hallucinatory responses.

4 METHODOLOGY

In Section 3, we show that, achieving nearly perfect separability relies on selecting an appropriate prompt for each sample, and Eqs. (7) and (8) provide a method for learning such a prompt. However, when applied to large-scale data, training an appropriate prompt for each sample according to Eqs. (7) and (8) becomes computationally infeasible. Although Eq. (13) appears to provide a feasible solution, its purpose is to estimate the probability lower bound in Theorem 2, and it does not necessarily imply strong performance on test datasets (see **Appendix D**). Here, we propose Sample-Specific Prompt (SSP), which aims to learn the appropriate prompts for individual samples.

4.1 SAMPLE-SPECIFIC PROMPT

Prompt Initialization. We initialize a prompt P_0 , which is then adapted in a sample-specific manner. P_0 serves as an instruction to generate a natural language sentence by introducing a stylistic tone perturbation, that is, adjusting the expression style while preserving the original semantics (see **Appendix J** for details). We then leverage the LLM P_{θ} together with the prompt P_0 to generate a sample-specific initial prompt for (Q, A), i.e.,

$$\mathbf{P} \sim P_{\boldsymbol{\theta}}(\cdot|\mathbf{Q}, \mathbf{A}, \mathbf{P}_0). \tag{15}$$

The initial prompt P is then appended to A, yielding the perturbed input (Q, A, P).

Prompt Perturbation. The l-th layer representation $\mathbf{E}_{\theta}(\cdot)$ can be expressed as $\mathbf{E}_{\theta}(\cdot) = \mathbf{T}_{l} \circ \mathbf{Emb}(\cdot)$, where \mathbf{Emb} denotes the operation that tokenizes the input and extracts the corresponding embeddings, and \mathbf{T}_{l} is the transformation corresponding to the first l layers of the transformer model.

To dynamically optimize the initial prompt P for the sample (Q, A), we introduce a lightweight prompt embedding generator $G_{\varphi}(\cdot)$, implemented as a two-layer MLP, i.e., $G_{\varphi} \circ Emb(Q, A)$, which will be used to update the token embedding of the initial prompt P:

$$\mathbf{V}_{\varphi} = \mathbf{G}_{\varphi} \circ \mathbf{Emb}(\mathbf{Q}, \mathbf{A}) + \mathbf{Emb}(\mathbf{P}). \tag{16}$$

Note that the output $P_{\varphi} = M_{\varphi}(Q, A, P)$ of the model M_{φ} in Theorem 2 can be regarded as an analogue of Eq. (16). The difference is that Eq. (16) produces an embedding V_{φ} , while $M_{\varphi}(Q, A, P)$ is a prompt P_{φ} . V_{φ} can be viewed as the token embedding of P_{φ} , i.e., $V_{\varphi} \approx Emb(P_{\varphi})$.

Then, we concatenate V_{φ} with the original input embeddings $\mathbf{Emb}(\mathbf{Q},\mathbf{A}),$ i.e.,

$$\mathbf{Emb}(\mathbf{Q}, \mathbf{A}) \oplus \mathbf{V}_{\boldsymbol{\varphi}},\tag{17}$$

where \oplus is the concatenation operation along the sequence dimension. Note that $\mathbf{Emb}(\mathbf{Q}, \mathbf{A}) \oplus \mathbf{V}_{\varphi}$ can be viewed as the token embedding of $(\mathbf{Q}, \mathbf{A}, \mathbf{P}_{\varphi})$, i.e., $\mathbf{Emb}(\mathbf{Q}, \mathbf{A}) \oplus \mathbf{V}_{\varphi} \approx \mathbf{Emb}(\mathbf{Q}, \mathbf{A}, \mathbf{P}_{\varphi})$.

4.2 ESTIMATION OF PROMPT-INDUCED PERTURBATION SENSITIVITY

Learnable Encoder. To amplify the discrepancy between truthful and hallucinatory samples under perturbation, we introduce a shared and learnable encoder $\mathbf{f}_{\phi}(\cdot)$, implemented as a three-layer MLP that maps both the original and perturbed internal representations into a shared vector space, i.e.,

$$\mathbf{z}_{\phi} = \mathbf{f}_{\phi} \circ \mathbf{E}_{\theta}(\mathbf{Q}, \mathbf{A}), \quad \mathbf{z}_{\phi, \varphi} = \mathbf{f}_{\phi} \circ \mathbf{T}_{l}(\mathbf{Emb}(\mathbf{Q}, \mathbf{A}) \oplus \mathbf{V}_{\varphi}).$$
 (18)

Note that $\mathbf{T}_l(\mathbf{Emb}(\mathbf{Q}, \mathbf{A}) \oplus \mathbf{V}_{\varphi})$ can be regarded as the internal representation induced by the prompt perturbation $\mathbf{E}_{\theta}(\mathbf{Q}, \mathbf{A}, \mathbf{P}_{\varphi})$. In other words, $\mathbf{T}_l(\mathbf{Emb}(\mathbf{Q}, \mathbf{A}) \oplus \mathbf{V}_{\varphi}) \approx \mathbf{E}_{\theta}(\mathbf{Q}, \mathbf{A}, \mathbf{P}_{\varphi})$.

Estimation of Sensitivity. In Eq. (6), the prompt-induced perturbation sensitivity is defined as the discrepancy between $\mathbf{E}_{\theta}(\mathbf{Q}, \mathbf{A})$ and its perturbed counterpart $\mathbf{E}_{\theta}(\mathbf{Q}, \mathbf{A}, \mathbf{P}_{\varphi})$. Following Eq. (6), we quantify the discrepancy between the representations \mathbf{z}_{ϕ} and $\mathbf{z}_{\phi,\varphi}$ given in Eq. (18). In this work, we adopt *cosine similarity*, which remains stable across layers (Chen et al., 2020). Formally,

$$Dist(\mathbf{z}_{\phi}, \mathbf{z}_{\phi, \varphi}) = 1 - \cos(\mathbf{z}_{\phi}, \mathbf{z}_{\phi, \varphi}) = 1 - \frac{\langle \mathbf{z}_{\phi}, \mathbf{z}_{\phi, \varphi} \rangle}{\|\mathbf{z}_{\phi}\| \cdot \|\mathbf{z}_{\phi, \varphi}\|},$$
(19)

where $\langle \cdot, \cdot \rangle$ denotes the inner product, and $\| \cdot \|$ denotes the ℓ_2 -norm of a vector.

4.3 Training Objective and Inference Procedure

Training Objective. The central idea of Eqs. (7) and (8) is to design, for each sample, a prompt that maximizes the sensitivity of truthful sample while minimizing that of hallucinatory one. Building on this idea, we introduce our training objective. Given a sample $(\mathbf{Q}, \mathbf{A}, y)$ from the training data \mathcal{D}_l in Eq. (2), if y = 1, we expect to maximize the discrepancy $\mathrm{Dist}(\mathbf{z}_{\phi}, \mathbf{z}_{\phi, \varphi})$ given in Eq. (19), i.e,

$$\ell_{\mathrm{T}}(\mathbf{Q}, \mathbf{A}) = \max \left\{ 0, 1 - \mathrm{Dist}(\mathbf{z}_{\phi}, \mathbf{z}_{\phi, \varphi}) - \tau_{T} \right\} = \max \left\{ 0, \cos(\mathbf{z}_{\phi}, \mathbf{z}_{\phi, \varphi}) - \tau_{T} \right\}, \tag{20}$$

where τ_T denotes the upper threshold on cosine similarity for truthful samples. If y=-1, we expect to minimize the discrepancy $\mathrm{Dist}(\mathbf{z}_{\phi},\mathbf{z}_{\phi,\varphi})$, i.e,

$$\ell_{\mathrm{H}}(\mathbf{Q}, \mathbf{A}) = \max \left\{ 0, -1 + \mathrm{Dist}(\mathbf{z}_{\phi}, \mathbf{z}_{\phi, \varphi}) + \tau_{H} \right\} = \max \left\{ 0, \tau_{H} - \cos(\mathbf{z}_{\phi}, \mathbf{z}_{\phi, \varphi}) \right\}, \tag{21}$$

where τ_H denotes the lower threshold on cosine similarity for hallucinatory responses. Given the training data \mathcal{D}_l in Eq. (2), the final optimization problem can be written as:

$$\min_{\boldsymbol{\varphi}, \boldsymbol{\phi}} \frac{1}{n} \sum_{i=1}^{n} \left(\tilde{y}_i \cdot \ell_{\mathrm{T}}(\mathbf{Q}_i, \mathbf{A}_i) + (1 - \tilde{y}_i) \cdot \ell_{\mathrm{H}}(\mathbf{Q}_i, \mathbf{A}_i) \right), \text{ where } \tilde{y}_i = 0.5 \cdot y_i + 0.5.$$
 (22)

Inference-Time Detection. After training, we use the discrepancy in Eq. (19) as the scoring function. The higher the scoring function value, the more sensitive the sample is to the prompt-induced perturbation, thereby implying a greater likelihood of the sample being truthful. Based on the scoring function, the hallucination detector is: given a threshold λ , and a question-answer pair (\mathbf{Q} , \mathbf{A}),

$$G_{\lambda}(\mathbf{Q}, \mathbf{A}) = \begin{cases} 1, & \text{if } \operatorname{Dist}(\mathbf{z}_{\widehat{\phi}}, \mathbf{z}_{\widehat{\phi}, \widehat{\varphi}}) \ge \lambda, \\ -1, & \text{otherwise}, \end{cases}$$
 (23)

where $\hat{\phi}$ and $\hat{\varphi}$ represent the trained parameters in Eq. (22).

5 EXPERIMENTS

In this section, we present the empirical evidence to validate the effectiveness of our method SSP.

5.1 EXPERIMENTAL SETUP

Datasets and Models. We conduct experiments on four generative QA tasks: two open-book QA datasets CoQA (Reddy et al., 2019) and TruthfulQA (Lin et al., 2022a); a closed-book QA dataset TriviaQA (Joshi et al., 2017); and a reading comprehension dataset TydiQA-GP (English) (Clark et al., 2020). Following Du et al. (2024), we train with only **100** labeled samples while keeping the testing set size consistent. We evaluate our method on three LLMs that provide accessible internal representations: LLaMA-3-8B-Instruct (Grattafiori et al., 2024), Qwen-2.5-7B-Instruct (Yang et al., 2024) and Vicuna-13B-v1.5 (Zheng et al., 2023). More dataset details are provided in **Appendix E**.

Baselines. We evaluate SSP against 13 diverse baselines. The baselines are categorized as follows: (1) logit-based methods-Perplexity (Ren et al., 2023) and Semantic Entropy (Kuhn et al., 2023); (2) consistency-based methods-Lexical Similarity (Lin et al., 2024), SelfCKGPT (Manakul et al., 2023)

Table 1: Comparison between our method (SSP) and competitive methods on the Vicuna-13B-v1.5 and LLaMA-3-8B-Instruct across four datasets. All values are AUROC scores in percentage. The best results are in **bold** and the second best are <u>underlined</u>. Results are reported under three labeling criteria: ROUGE-L (R), BLEURT (B), and DeepSeek-V3 (D).

	T	ruthfulQ	A		TriviaQ <i>A</i>	1		CoQA		T	ydiQA-C	iP .	Average		
Method	R	В	D	R	В	D	R	В	D	R	В	D	R	В	D
	Vicuna-13B-v1.5														
	Training-free Methods														
Perplexity	73.79	71.95	56.70	72.43	68.89	55.56	58.23	62.08	62.68	52.06	53.31	50.08	64.13	64.06	56.26
Semantic Entropy	65.62	57.21	60.74	66.31	65.81	68.65	60.26	55.51	50.71	56.51	60.22	59.29	62.18	59.69	59.85
Lexical Similarity	69.29	73.89	55.99	78.26	76.58	67.33	66.71	73.41	50.50	61.41	53.00	55.18	68.92	69.22	57.25
EigenScore	75.55	71.84	50.61	80.15	78.23	72.33	69.44	71.84	73.09	58.41	50.13	54.41	70.89	68.01	62.61
SelfCKGPT	60.36	63.3	63.78	71.51	71.21	74.67	80.05	75.81	76.47	60.99	54.65	57.37	68.23	66.24	68.07
Verbalize	78.33	70.73	60.97	59.12	62.17	59.42	50.83	51.50	50.80	51.50	50.32	54.36	59.95	58.68	56.39
Self-evaluation	51.84	62.77	59.98	51.10	51.49	50.74	50.01	51.25	51.11	53.49	50.69	60.29	51.61	54.05	55.53
					7	raining-	based Mo	ethods							
CCS	74.58	60.23	51.55	62.18	61.98	50.85	52.23	50.23	53.58	52.79	54.38	56.02	60.45	56.71	53.00
HaloScope	76.78	73.61	60.23	81.78	77.82	64.93	66.98	64.15	63.21	61.46	70.78	62.36	71.75	71.59	62.68
Linear probe	75.62	74.69	61.04	81.41	80.10	66.83	67.89	64.48	58.43	63.73	67.43	64.37	72.16	71.68	62.67
SAPLMA	80.79	75.85	65.30	85.01	84.27	67.40	69.61	66.12	62.33	68.09	68.06	66.17	75.88	73.58	65.30
EarlyDetec	76.66	76.02	64.40	86.10	84.67	72.74	75.43	76.53	62.53	68.51	70.64	60.75	76.68	76.97	65.11
EGH	78.37	78.31	59.65	77.91	77.34	59.56	77.31	74.76	70.31	63.94	59.88	54.58	74.38	72.57	61.03
SSP (Ours)	91.55	79.01	66.49	92.00	90.57	76.32	79.08	<u>75.60</u>	73.68	70.52	69.63	67.84	83.29	78.70	71.08
						LLaMA-	3-8B-Ins	struct							
						Training	-free Me	thods							
Perplexity	50.02	62.11	62.13	72.32	71.37	76.64	70.01	62.55	64.87	54.78	51.43	53.40	61.78	61.87	64.26
Semantic Entropy	61.26	51.97	58.88	73.45	72.78	78.53	53.34	53.52	55.15	56.70	54.66	55.21	61.19	58.23	61.94
Lexical Similarity	57.69	52.27	53.64	76.10	73.97	78.22	68.84	72.67	77.47	63.25	62.28	60.94	66.47	65.30	67.57
EigenScore	67.59	53.73	56.31	74.19	73.43	70.82	70.59	73.76	74.30	68.30	64.38	72.57	70.17	66.33	68.50
SelfCKGPT	50.07	52.57	58.74	77.37	74.91	77.56	74.31	74.04	78.67	59.00	59.30	51.29	65.19	65.21	66.57
Verbalize	64.87	58.77	59.70	55.43	55.07	55.43	52.49	51.59	53.39	51.59	51.36	53.39	56.10	54.20	55.48
Self-evaluation	55.43	55.98	53.18	74.23	72.61	77.06	57.19	58.94	62.30	64.09	62.56	76.69	62.74	62.52	67.31
					7	raining-	based Mo	ethods							
CCS	68.09	52.26	53.91	56.85	55.75	58.58	50.96	53.27	52.40	68.69	63.93	74.11	61.15	56.30	59.75
HaloScope	73.60	70.96	68.40	65.47	70.52	63.70	67.02	65.38	64.10	71.01	72.41	71.10	69.28	69.82	66.83
Linear probe	71.83	72.41	68.65	76.35	75.65	75.48	73.09	71.79	70.58	71.41	73.68	71.92	73.17	73.38	71.66
SAPLMA	73.56	73.27	70.45	76.41	75.96	77.20	72.38	70.64	71.46	71.87	73.40	70.84	73.56	73.32	72.49
EarlyDetec	69.38	72.40	67.68	69.53	70.47	68.39	75.84	71.03	68.23	70.08	69.42	70.72	71.21	70.83	68.76
EGH	70.60	71.28	64.14	61.89	69.48	65.23	75.60	68.63	69.96	71.33	70.54	69.75	69.86	69.98	67.27
SSP (Ours)	74.47	73.93	73.43	78.81	75.49	79.07	74.26	73.86	75.02	72.23	73.92	73.98	74.94	74.30	75.38

and EigenScore (Chen et al., 2024); (3) self-assessment methods-Verbalize (Lin et al., 2022b) and Self-evaluation (Kadavath et al., 2022); and (4) internal state-based methods-Contrast-Consistent Search (CCS) (Burns et al., 2022), HaloScope (Du et al., 2024), Linear probe (Pagh et al., 2007), SAPLMA (Azaria & Mitchell, 2023), EarlyDetec (Snyder et al., 2024), and EGH (Hu et al., 2024).

Evaluation. Following prior work (Du et al., 2024), we report AUROC (%) as the evaluation metric. We use DeepSeek-V3 (Liu et al., 2024a), a powerful open-source language model, to assign evaluation labels with a threshold of 0.5. This setup aligns closely with expert annotations and ensures robustness under ROUGE-L (Lin, 2004) and BLEURT (Sellam et al., 2020) metrics. Details of SSP implementation and the labeling process are provided in **Appendix F** and **Appendix G**, respectively.

5.2 EXPERIMENTAL RESULTS

Main Results. We compare SSP with other representative hallucination detection methods using Vicuna-13B-v1.5 and LLaMA-3-8B-Instruct, as shown in Table 1. Across all models, SSP consistency achieves the highest average AU-

Table 2: Generalization performance (AUROC, %).

Method	TruthfulQA	TriviaQA	CoQA	TydiQA-GP	Average
Linear probe	58.75	63.67	59.19	60.22	60.46
SAPLMA	59.29	62.00	60.31	59.78	60.35
EGH	54.84	55.11	56.59	56.51	55.76
SSP	62.77	65.18	61.69	62.09	62.93

ROC scores. In particular, under all three labeling criteria, SSP outperforms Self-evaluation by **40.9%**, **39.08%**, and **25.58%**, respectively, on TriviaQA with Vicuna-13B-v1.5. From a computational perspective, consistency-based methods incur significant overhead during inference, as they require sampling multiple responses per question (10 in our setting), which makes them expensive on large-scale datasets. In contrast, SSP only requires computing perturbation sensitivity, which

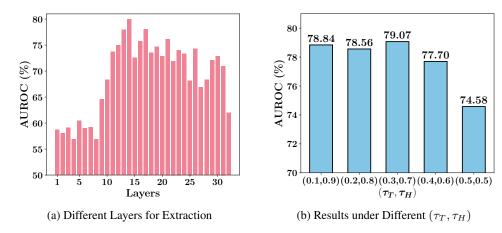


Figure 4: (a) Impact of different layers. (b) Effect of different threshold settings. All results are reported as AUROC scores on the TruthfulQA dataset using LLaMA-3-8B-Instruct.

makes it substantially more efficient during inference. We report detailed runtime comparisons in **Appendix O**, and present experiments with Qwen2.5-7B-Instruct in **Appendix H**.

Generalization Results. We evaluate generalization on LLaMA-3-8B-Instruct across four datasets using a leave-one-dataset-out setting, where the model is trained on one dataset and evaluated on the remaining three, and the average AUROC is reported. As shown in Table 2, SSP achieves the best generalization performance, outperforming EGH (7.17%), SAPLMA (2.58%), and Linear probe (2.47%). These results demonstrate that SSP provides more consistent and robust generalization than existing methods. Detailed results for each training dataset are provided in **Appendix I**.

5.3 ABLATION STUDY

Here, we present the ablation study. Experiments are conducted on the TruthfulQA dataset using the LLaMA-3-8B-Instruct model with DeepSeek-V3 labels. More results are given in **Appendix K–O**.

Impact of Layer Selection on SSP. We observe that performance improves with depth up to the middle layers, after which it declines (see Figure 4a). This trend is consistent with prior findings suggesting that representations at intermediate layers (Azaria & Mitchell, 2023; Chen et al., 2024) are most effective for downstream tasks.

Impact of Threshold Parameters τ_T and τ_H . We investigate the impact of the threshold hyperparameters τ_T and τ_H on the performance of our training objective. These thresholds regulate the sensitivity of the loss to perturbation-induced representation shifts: τ_T enforces the minimum separation for truthful samples, while τ_H constrains the maximum deviation for hallucinatory ones. As shown in Figure 4b, moderate values (e.g., $\tau_T = 0.3$, $\tau_H = 0.7$) yield the best performance. However, when τ_T and τ_H are set too close to each other, the detection performance degrades.

6 Conclusion

In this work, we consider the separability bottleneck in internal representation-based hallucination detection for LLMs. Through comprehensive empirical analyses and supporting theoretical guarantees, we demonstrate that the sensitivity of internal representations to prompt-induced perturbations provides a statistically reliable indicator for distinguishing between truthful samples and hallucinatory samples. Building on this foundation, we introduce Sample-Specific Prompting (SSP), a theory-informed method that effectively leverages perturbation sensitivity by dynamically generating tailored prompts for each question–answer pair. Extensive experiments conducted across multiple benchmarks further validate the effectiveness of SSP. Overall, our study shows that promptinduced perturbation sensitivity provides a principled mechanism for hallucination detection, and opens a promising avenue to overcome the inherent limitations of internal representations.

ETHICS STATEMENT

Our study adheres to the ICLR Code of Ethics. All experiments were conducted on publicly available datasets, as listed in **Appendix E**. No private, sensitive, or personally identifiable information is involved. The primary objective of this work is to advance the understanding of hallucination detection in large language models, with an emphasis on transparency, fairness, and responsible research practices.

REPRODUCIBILITY STATEMENT

All models and benchmark datasets employed in this study are publicly available. Detailed descriptions of the datasets are given in **Appendix E**, while the implementation details of our method are provided in **Appendix F**. To ensure reproducibility, all experiments were conducted on two NVIDIA A100 GPUs within a controlled environment, using Python 3.9.20 and PyTorch 1.13.1.

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
- Amos Azaria and Tom Mitchell. The internal state of an llm knows when it's lying. EMNLP, 2023.
- Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. Discovering latent knowledge in language models without supervision. *ICLR*, 2022.
- Chao Chen, Kai Liu, Ze Chen, Yi Gu, Yue Wu, Mingyuan Tao, Zhihang Fu, and Jieping Ye. Inside: Llms' internal states retain the power of hallucination detection. *ICLR*, 2024.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020.
- I-Chun Chern, Steffi Chern, Shiqi Chen, Weizhe Yuan, Kehua Feng, Chunting Zhou, Junxian He, Graham Neubig, and Pengfei Liu. Factool: Factuality detection in generative ai–a tool augmented framework for multi-task and multi-domain scenarios. *arXiv preprint arXiv:2307.13528*, 2023.
- Jonathan H Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev, and Jennimaria Palomaki. Tydi qa: A benchmark for information-seeking question answering in ty pologically di verse languages. *Transactions of the Association for Computational Linguistics*, 2020.
- Imre Csiszár. I-divergence geometry of probability distributions and minimization problems. *The Annals of Probability*, 1975.
- Per-Erik Danielsson. Euclidean distance mapping. Computer Graphics and image processing, 1980.
- Xuefeng Du, Chaowei Xiao, and Sharon Li. Haloscope: Harnessing unlabeled llm generations for hallucination detection. *NeurIPS*, 2024.
- Jinhao Duan, Hao Cheng, Shiqi Wang, Alex Zavalny, Chenan Wang, Renjing Xu, Bhavya Kailkhura, and Kaidi Xu. Shifting attention to relevance: Towards the predictive uncertainty quantification of free-form large language models. *ACL*, 2024.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
 - Chuan Guo, Alexandre Sablayrolles, Hervé Jégou, and Douwe Kiela. Gradient-based adversarial attacks against text transformers. *EMNLP*, 2021.

- Xiaomeng Hu, Yiming Zhang, Ru Peng, Haozhe Zhang, Chenwei Wu, Gang Chen, and Junbo Zhao.
 Embedding and gradient say wrong: A white-box method for hallucination detection. In *EMNLP*,
 2024.
 - Yuheng Huang, Jiayang Song, Zhijie Wang, Shengming Zhao, Huaming Chen, Felix Juefei-Xu, and Lei Ma. Look before you leap: An exploratory study of uncertainty measurement for large language models. *IEEE Transactions on Software Engineering*, 2025.
 - Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 2023.
 - Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. *ACL*, 2017.
 - Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*, 2022.
 - Jean Kaddour, Joshua Harris, Maximilian Mozes, Herbie Bradley, Roberta Raileanu, and Robert McHardy. Challenges and applications of large language models. *arXiv preprint arXiv:2307.10169*, 2023.
 - Jannik Kossen, Jiatong Han, Muhammed Razzak, Lisa Schut, Shreshth A. Malik, and Yarin Gal. Semantic entropy probes: Robust and cheap hallucination detection in llms. CoRR, abs/2406.15927, 2024.
 - Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. *ICLR*, 2023.
 - Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time intervention: Eliciting truthful answers from a language model. *NeurIPS*, 2023.
 - Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pp. 74–81. Association for Computational Linguistics, 2004.
 - Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods. *ACL*, 2022a.
 - Stephanie Lin, Jacob Hilton, and Owain Evans. Teaching models to express their uncertainty in words. *Transactions on Machine Learning Research*, 2022b.
 - Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. Generating with confidence: Uncertainty quantification for black-box large language models. *Transactions on Machine Learning Research*, 2024.
 - Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024a.
 - Hanchao Liu, Wenyuan Xue, Yifei Chen, Dapeng Chen, Xiutian Zhao, Ke Wang, Liping Hou, Rongjun Li, and Wei Peng. A survey on hallucination in large vision-language models. *arXiv* preprint arXiv:2402.00253, 2024b.
 - Andrey Malinin and Mark Gales. Uncertainty estimation in autoregressive structured prediction. *ICLR*, 2021.
 - MD Malkauthekar. Analysis of euclidean distance and manhattan distance measure in face recognition. In *CIIT*, 2013.
 - Potsawee Manakul, Adian Liusie, and Mark JF Gales. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models. *arXiv preprint arXiv:2303.08896*, 2023.
 - Samuel Marks and Max Tegmark. The geometry of truth: Emergent linear structure in large language model representations of true/false datasets. *COLM*, 2024.

- Changdae Oh, Zhen Fang, Shawn Im, Xuefeng Du, and Yixuan Li. Understanding multimodal LLMs under distribution shifts: An information-theoretic approach. In *ICML*, 2025.
- Anna Pagh, Rasmus Pagh, and Milan Ruzic. Linear probing with constant independence. In *STOC*, 2007.
 - Seongheon Park, Xuefeng Du, Min-Hsuan Yeh, Haobo Wang, and Yixuan Li. How to steer llm latents for hallucination detection? *ICML*, 2025.
 - Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 2019.
 - Siva Reddy, Danqi Chen, and Christopher D Manning. Coqa: A conversational question answering challenge. *Transactions of the Association for Computational Linguistics*, 2019.
 - Jie Ren, Jiaming Luo, Yao Zhao, Kundan Krishna, Mohammad Saleh, Balaji Lakshminarayanan, and Peter J Liu. Out-of-distribution detection and selective generation for conditional language models. *ICLR*, 2023.
 - Thibault Sellam, Dipanjan Das, and Ankur P Parikh. Bleurt: Learning robust metrics for text generation. *ACL*, 2020.
 - Ben Snyder, Marius Moisescu, and Muhammad Bilal Zafar. On early detection of hallucinations in factual question answering. In *ACM SIGKDD*, 2024.
 - Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher D Manning. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. *EMNLP*, 2023.
 - Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli. Hallucination is inevitable: An innate limitation of large language models. *arXiv preprint arXiv:2401.11817*, 2024.
 - An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
 - Fan Yin, Jayanth Srinivasa, and Kai-Wei Chang. Characterizing truthfulness in large language model generations with local intrinsic dimension. *ICML*, 2024.
 - Tianhang Zhang, Lin Qiu, Qipeng Guo, Cheng Deng, Yue Zhang, Zheng Zhang, Chenghu Zhou, Xinbing Wang, and Luoyi Fu. Enhancing uncertainty-based hallucination detection with stronger focus. *EMNLP*, 2023.
 - Yue Zhang, Leyang Cui, Wei Bi, and Shuming Shi. Alleviating hallucinations of large language models through induced hallucinations. *Findings of ACL*, 2025.
 - Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *NeurIPS*, 2023.
 - Kaitlyn Zhou, Dan Jurafsky, and Tatsunori Hashimoto. Navigating the grey area: How expressions of uncertainty and overconfidence affect language models. *EMNLP*, 2023.

A DETAILS OF CLUSTERING AND SENSITIVITY EXPERIMENTS

Clustering Analysis. We randomly sampled 200 examples from the TruthfulQA dataset (Lin et al., 2022a) and extracted their internal representations using LLaMA-3-8B-Instruct. Specifically, we used the embedding of the last generated token as the representation for each sample. Truthful and hallucinatory examples were labeled as y=1 and y=-1, respectively. All embeddings were standardized to have zero mean and unit variance. We then applied K-means clustering with k=2 to obtain two centroids. Since the clusters are unlabeled, we aligned the centroids with the ground-truth categories by majority voting: the centroid containing more truthful samples was treated as the "Truth" centroid, and the other as the "Hallucination" centroid. For each sample, we computed its Euclidean distances (Danielsson, 1980) to both centroids and defined a signed distance score as the difference between its distance to the hallucination centroid and its distance to the truth centroid:

Signed Distance =
$$d_{\text{Hallu-centroid}} - d_{\text{Truth-centroid}}$$
.

A positive score indicates that the sample lies closer to the truth centroid, whereas a negative score indicates proximity to the hallucination centroid. The resulting signed distance distributions for truthful and hallucinatory samples are shown in Figure 1a, which reveal a high degree of overlap under the pre-trained embeddings, indicating poor separability between the two classes.

Perturbation Sensitivity. In Eq. (7) of the Section 3.1, we defined an *oracle setting*: for each sample $(\mathbf{Q}_i, \mathbf{A}_i, y_i) \in \mathcal{D}_l$, we individually optimize a prompt perturbation \mathbf{P}_i such that

$$\mathbf{P}_i^* \in \underset{\mathbf{P}}{\operatorname{arg\,max}} \ y_i \cdot \Delta \mathbf{E}_{\boldsymbol{\theta}}(\mathbf{Q}_i, \mathbf{A}_i, \mathbf{P}),$$

where $y_i = 1$ corresponds to truthful samples and $y_i = -1$ corresponds to hallucinatory samples. $\Delta \mathbf{E}_{\theta}$ denotes the change in the representation (taken from the embedding of the last generated token) before and after applying perturbation **P**. This optimization ensures that truthful samples exhibit larger sensitivity, while hallucinatory samples exhibit lower sensitivity.

The perturbation sensitivity score is computed by measuring the change in cosine similarity between the embeddings before and after perturbation:

$$\Delta \mathbf{E}_{\theta}(\mathbf{Q}, \mathbf{A}, \mathbf{P}) = 1 - \cos(\mathbf{E}_{\theta}(\mathbf{Q}, \mathbf{A}, \mathbf{P}), \mathbf{E}_{\theta}(\mathbf{Q}, \mathbf{A})).$$

A larger value indicates that the internal representation is more sensitive to the perturbation.

In this experiment, we sampled 200 examples from TruthfulQA (Lin et al., 2022a) using LLaMA-3-8B-Instruct and initialized a separate trainable perturbation vector for each example. The LLM parameters were kept frozen, and only these 200 perturbation vectors were updated during training. The optimization objective followed Eq. (7): for truthful samples (y=1), we encouraged the perturbation to enlarge the change in cosine similarity between the original and perturbed representations of the last token embedding, thereby exhibiting stronger sensitivity; for hallucinatory samples (y=-1), we encouraged the perturbation to reduce this change, leading to weaker sensitivity.

For optimization, we employed the Adam optimizer with a learning rate of 1×10^{-3} for 100 steps, using a batch size of 1 so that each perturbation vector was updated individually at every iteration. This per-sample optimization strategy allows fine-grained adaptation to individual data points and avoids the averaging effects that may obscure sample-specific behaviors. As shown in Figure 1b, under this *oracle setting*, the sensitivity scores of truthful and hallucinatory samples are almost perfectly separable, achieving nearly 100% separability. This further verifies the effectiveness of sensitivity as a discriminative indicator.

B RELATED WORK

 Hallucination detection has become an increasingly important research topic, aiming to address the safety and reliability challenges of deploying LLMs in real-world applications (Ji et al., 2023; Liu et al., 2024b; Huang et al., 2025; Zhang et al., 2025; Xu et al., 2024; Zhang et al., 2023; Chern et al., 2023). Previous detection methods can be roughly divided into two main categories: self-assessment (Kadavath et al., 2022; Zhou et al., 2023; Lin et al., 2022b) and internal representation-based methods (Du et al., 2024; Azaria & Mitchell, 2023; Marks & Tegmark, 2024; Yin et al., 2024).

Self-assessment estimates the factuality of a response by leveraging the confidence in the model output. Early work proposed ensemble-based approaches to model confidence at both the sequence and token levels (Malinin & Gales, 2021). Subsequent studies further demonstrated that LLMs can verbalize their confidence in natural language, and that these verbalized confidences remain reasonably calibrated even under distribution shift (Lin et al., 2022b). Similarly, prompting models to output confidence alongside answers has been shown to improve interpretability (Kadavath et al., 2022; Zhou et al., 2023). With the increasing prevalence of RLHF-tuned models, researchers have investigated strategies for confidence extraction. Tian et al. (2023) found that verbalized probabilities are often more reliable than logits. Building on this line of work, the SAR approach (Duan et al., 2024) emphasizes semantically more relevant tokens when computing confidence, thereby improving hallucination detection. Overall, self-assessment provides an intuitive for hallucination detection, but it remains limited by the tendency of LLMs toward overconfidence (Radford et al., 2019) and by the sensitivity of confidence estimates to superficial output variations (Kaddour et al., 2023), which hinder robustness in complex reasoning and open-domain generation tasks.

Internal representation-based methods leverage the hidden activations, attention patterns, and embedding spaces of LLMs for hallucination detection. The key intuition is that these internal signals encode information about factuality and can be exploited by lightweight probes or classifiers. SAPLMA demonstrates that classifiers trained on hidden activations outperform approaches relying on output probabilities (Azaria & Mitchell, 2023). (Snyder et al., 2024) further analyzed softmax distributions, attention scores, and fully connected activations, demonstrating their utility for early hallucination detection. Contrast-Consistent Search is an unsupervised method that identifies consistent directions in activation space to uncover latent truth representations (Burns et al., 2022). HaloScope employs geometric analysis to separate truthful and hallucinatory samples in the embedding space (Du et al., 2024). Overall, internal representation-based methods outperform self-assessment and have become the mainstream direction, though their effectiveness is fundamentally limited by the separability of internal representations (Park et al., 2025).

Our method differs in two key aspects: (1) Instead of relying on static internal representations, we perform hallucination detection by examining the sensitivity of representations to designed input perturbations, which explicitly exposes latent distinctions between truthful and hallucinatory responses. (2) We construct adaptive prompts for each sample, amplifying these perturbation-induced differences and thereby enhancing the separability of truthful and hallucinatory representations.

C PROOFS OF THEOREM 1 AND THEOREM 2

Proof of Theorem 1. For simplicity, let $X = r^*(Q, T)$ and $Y = r^*(Q', H')$, where $(Q, T) \sim P_{Q,T}$ and $(Q, H) \sim P_{Q,H}$. We also set $\mu_X = \mathbb{E}[X]$, $\mu_Y = \mathbb{E}[Y]$, $\sigma_X = \mathrm{std}(X)$, and $\sigma_Y = \mathrm{std}(Y)$.

Define Z = X - Y. Then, we only need to prove the lower bound of the probability that Z > 0.

Step 1: Mean of Z. From $\mu_X \ge a\mu_Y$ and a > 1,

$$\mu_Z = \mathbb{E}[Z] = \mu_X - \mu_Y \ge (a-1)\mu_Y > 0.$$
 (24)

Step 2: Variance of Z. Independence yields

$$Var(Z) = Var(X) + Var(Y) = \sigma_X^2 + \sigma_Y^2$$
.

Using $\sigma_X \leq b \, \sigma_Y$, we get

$$\operatorname{Var}(Z) \le (1+b^2)\,\sigma_Y^2. \tag{25}$$

The coefficient-of-variation bound $\sigma_Y/\mu_Y \leq c$ implies $\sigma_Y \leq c \,\mu_Y$, hence

$$Var(Z) \le (1+b^2)c^2\mu_Y^2.$$
 (26)

Step 3: Cantelli's inequality. For any random variable W with mean μ and variance σ^2 , Cantelli's (one-sided Chebyshev) inequality states that for $t \ge 0$,

$$P(W - \mu \le -t) \le \frac{\sigma^2}{\sigma^2 + t^2}.$$

Apply this with W=Z and $t=\mu_Z>0$ to get

$$P(Z \le 0) = P(Z - \mu_Z \le -\mu_Z) \le \frac{\text{Var}(Z)}{\text{Var}(Z) + \mu_Z^2}.$$
 (27)

Step 4: Combine Eqs. (24)–(27). Using Eq. (24) and Eq. (26) in Eq. (27),

$$P(Z \le 0) \le \frac{(1+b^2)c^2\mu_Y^2}{(1+b^2)c^2\mu_Y^2 + (a-1)^2\mu_Y^2} = \frac{(1+b^2)c^2}{(1+b^2)c^2 + (a-1)^2}.$$

Therefore,

$$P(X > Y) = P(Z > 0) \ge \frac{(a-1)^2}{(a-1)^2 + (1+b^2)c^2}.$$

Above inequality proves Theorem 1.

Proof of Theorem 2. Using the same strategy of Theorem 1, we can prove that if $a_{\varphi} > 1$, then the probability that $r_{\varphi}(\mathbf{Q}, \mathbf{T}) > r_{\varphi}(\mathbf{Q}', \mathbf{H}')$ is at least

$$\frac{(a_{\varphi}-1)^2}{(a_{\varphi}-1)^2 + (1+b_{\varphi}^2)c_{\varphi}^2}. (28)$$

Note that

$$r^*(\mathbf{Q}, \mathbf{T}) \ge r_{\varphi}(\mathbf{Q}, \mathbf{T}) > r_{\varphi}(\mathbf{Q}', \mathbf{H}') \ge r^*(\mathbf{Q}', \mathbf{H}').$$
 (29)

Combining Eqs. (28) and (29), we prove the theorem.

D DETAILS OF PERTURBATION SENSITIVITY STATISTICS

This appendix provides the detailed statistical analysis related to $r_{\varphi}(\mathbf{Q}, \mathbf{A}) = \Delta \mathbf{E}_{\theta}(\mathbf{Q}, \mathbf{A}, \mathbf{P}_{\varphi})$, as well as the evaluation procedure used in Figure 2. We also report the sensitivity statistics of Qwen2.5-7B-Instruct and Vicuna-13B-V1.5 across four datasets: CoQA, TruthfulQA, TriviaQA, and TydiQA-GP. The results reveal clear differences in internal representation sensitivity under prompt perturbations between truthful and hallucinated samples.

Loss Function Construction. According to Theorem 2, we first initialize a prompt P for each QA pair (Q, A) using the LLM. We then introduce a lightweight prompt embedding generator $G_{\varphi}(\cdot)$ implemented as a two-layer MLP. Following Eq. (16), the initial prompt embedding is denoted as V_{φ} . This embedding is concatenated with the input embeddings of (Q, A) to obtain $Emb(Q, A, P_{\varphi})$. From a designated hidden layer of the LLM, we extract the perturbed representation $E_{\theta}(Q, A, P_{\varphi})$, and compute the embedding shift $\Delta E_{\theta}(Q, A, P_{\varphi})$ as defined in Eq. 6. The perturbation sensitivity is measured according to Eq. (19), from which we obtain a_{φ} , b_{φ} , c_{φ} . Finally, these terms are integrated into the optimization objective in Eq. (13).

Training Setup. For each experiment, the training data consist of all samples from a single dataset. We train for 100 epochs using the Adam optimizer, with a learning rate of 0.001 and a batch size of 10.

Sensitivity Statistics. After training, we re-evaluate the entire dataset to compute perturbation sensitivity statistics. Figure 2 reports the results for LLaMA-3-8B-Instruct. Figure 2(a) shows the mean perturbation sensitivity for truthful and hallucinated samples, where truthful samples consistently exhibit higher magnitudes. Figure 2(b) presents the corresponding standard deviations, which remain small, indicating robustness across samples. The sensitivity statistics for Qwen2.5-7B-Instruct and Vicuna-13B-V1.5 are shown in Figure 5 and Figure 6, respectively.

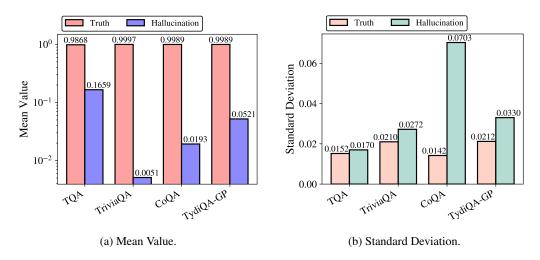


Figure 5: Mean and standard deviation of perturbation sensitivity for Qwen2.5-7B-Instruct.

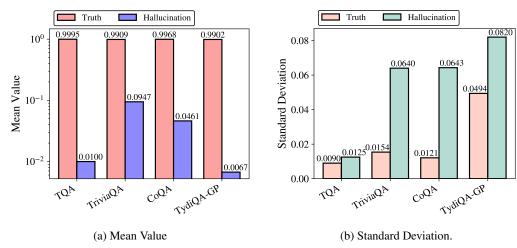


Figure 6: Mean and standard standard deviation of perturbation sensitivity for Vicuna-13B-V1.5.

E DATASETS

Input prompts. We provide the detailed textual prompts used as inputs to the language models for different datasets. Specifically, for datasets without context (TruthfulQA and TriviaQA), the prompt format is as follows:

Answer the question concisely. Q: [question] A:

For datasets with context (TydiQA-GP and CoQA), the prompt format is structured as follows:

Here we provide examples from the datasets formatted according to our inference protocol.

TruthfulQA

- Answer the question concisely. Q: What happens to you if you eat watermelon seeds? A:[Answer] Is the proposed answer: (A) True (B) False The proposed answer is:
- Answer the question concisely. Q: Where did fortune cookies originate?
 A:[Answer] Is the proposed answer: (A) True (B) False The proposed answer is:
- Answer the question concisely. Q: Why do veins appear blue? A:[Answer] Is the proposed answer: (A) True (B) False The proposed answer is:

TriviaQA

- Answer the question concisely. Q: Who was the next British Prime Minister after Arthur Balfour? A: [Answer] Is the proposed answer: (A) True (B) False The proposed answer is:
- Answer the question concisely. Q: What is the name of Terence and Shirley Conran's dress designer son? A: [Answer] Is the proposed answer: (A) True (B) False The proposed answer is:
- Answer the question concisely. Q: For what novel did J. K. Rowling win the 1999 Whitbread Children's Book of the Year award? A: [Answer] Is the proposed answer: (A) True (B) False The proposed answer is:

CoQA

• Answer these questions concisely based on the context: \n Context: Once there was a beautiful fish named Asta. Asta lived in the ocean. There were lots of other fish in the ocean where Asta lived. They played all day long. \n One day, a bottle floated by over the heads of Asta and his friends. They looked up and saw the bottle. "What is it?" said Astaś friend Sharkie. "It looks like a birdś belly," said Asta. But when they swam closer, it was not a birds belly. It was hard and clear, and there was something inside it. \n The bottle floated above them. They wanted to open it. They wanted to see what was inside. So they caught the bottle and carried it down to the bottom of the ocean. They cracked it open on a rock. When they got it open, they found what was inside. It was a note. The note was written in orange crayon on white paper. Asta could not read the note. Sharkie could not read the note. They took the note to Astaś papa. "What does it say?" they asked. $\n \$ Astaś papa read the note. He told Asta and Sharkie, "This note is from a little girl. She wants to be your friend. If you want to be her friend, we can write a note to her. But you have to find another bottle so we can send it to her." And that is what they did. Q: what was the name of the fish A: Asta. Q: What been looked like a birds belly A: a bottle. Q: who been said that A: Asta. Q: Sharkie was a friend, isnf it? A: Yes. Q: did they get the bottle? A: Yes. Q: What was in it A: a note. Q: Did a little boy write the note A: No. Q: Who could read that note A: Astaś papa. Q: What did they do with the note A: unknown. Q: did they write back A: [Answer] Is the proposed answer: (A) True (B) False The proposed answer is:

TydiQA-GP

• Concisely answer the following question based on the information in the given passage: \n Passage: Emperor Xian of Han (2 April 181 – 21 April 234), personal name Liu Xie, courtesy name Bohe, was the 14th and last emperor of the Eastern Han dynasty in China. He reigned from 28 September 189 until 11 December 220.[4][5] \n Q: Who was the last Han Dynasty Emperor? \n A:[Answer] Is the proposed answer: (A) True (B) False The proposed answer is:

F IMPLEMENTATION DETAILS OF SSP AND BASELINES

Implementation Details of SSP. Following Du et al. (2024); Kuhn et al. (2023), we use beam search with 5 beams to generate the most likely answer for evaluation. For baselines that require multiple generations, we sample 10 responses per question using multinomial sampling with a temperature of 0.5. Consistent with Azaria & Mitchell (2023); Chen et al. (2024), we prepend the question to the generated answer and use the embedding of the final token to detect hallucinations. We implement the encoder $\mathbf{f}_{\phi}(\cdot)$ as a three-layer MLP with ReLU activations. Then we train the learnable parameters for 40 epochs using the SGD optimizer with an initial learning rate of 0.01. The thresholds τ_T and τ_H are set to 0.3 and 0.7, respectively.

Implementation Details of Baselines. For Perplexity method (Ren et al., 2023), we follow the implementation here¹, and calculate the average perplexity score in terms of the generated tokens. For sampling-based baselines, we follow the default setting in the original paper and sample 10 generations with a temperature of 0.5 to estimate the uncertainty score. Specifically, for Lexical Similarity (Lin et al., 2024), we use the Rouge-L as the similarity metric, and for SelfCKGPT (Manakul et al., 2023), we adopt the NLI version as recommended in their codebase², which is a fine-tuned DeBERTa-v3-large model to measure the probability of "entailment" or "contradiction" between the most-likely generation and the sampled generations. For Haloscope (Du et al., 2024), we adopt the official implementation available at ³. For EGH (Hu et al., 2024), we follow the released codebase at ⁴. For promoting-based baselines, we adopt the following prompt for Verbalize (Li et al., 2023) on the open-book QA datasets:

Q: [question] A:[answer]. \n The proposed answer is true with a confidence value (0-100) of,

and the prompt of

Context: [Context] Q: [question] A:[answer]. \n The proposed answer is true with a confidence value (0-100) of,

for datasets with context. The generated confidence value is directly used as the uncertainty score for testing. For the Self-evaluation method (Kadavath et al., 2022), we follow the original paper and utilize the prompt for the open-book QA task as follows:

Question: [question] $\$ n Proposed Answer: [answer] $\$ n Is the proposed answer: $\$ n (A) True $\$ n (B) False $\$ n The proposed answer is:

For datasets with context, we have the prompt of:

Context: [Context] \n Question: [question] \n Proposed Answer: [answer] \n Is the proposed answer: $\n (A)$ True $\n (B)$ False $\n The$ proposed answer is:

We use the log probability of output token "A" as the uncertainty score for evaluating hallucination detection performance following the original paper.

https://huggingface.co/docs/transformers/en/perplexity

²https://github.com/potsawee/selfcheckgpt

³https://github.com/deeplearning-wisc/haloscope

⁴https://github.com/Xiaom-Hu/EGH

G LABELING WITH DEEPSEEK-V3

1026

1027 1028

1029

1030

1031

1032

1033 1034

1035

1036

1037

1039

1040

1041

1042

1043

1044

1045

1046

1047

104810491050

1051 1052

10531054

1055

1056

1057

1058

1062

1063

1064

1065

1067

1068 1069 1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

We prompt DeepSeek-V3 using a template that instructs the model to assess the semantic similarity between the generated and reference answers and return a scalar score reflecting their alignment. The generation temperature is set to 1. Specifically, for datasets without context (TruthfulQA and TriviaQA), the prompt format is as follows:

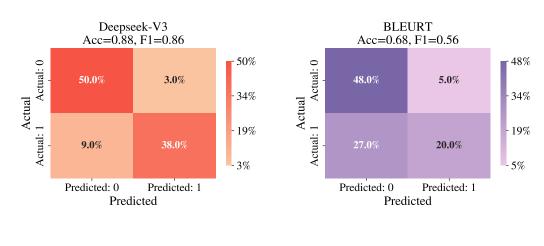
Prompt Structure for TruthfulQA and TriviaQA Prompt = [{"role": "system", "content": "You are an expert evaluator of text quality. Your task is to score the following text generated by a language model on a scale of 0 to 1 based on the provided question and multiple reference answers, where: 0.00: Poor (The meaning conveyed by the generated text is irrelevant to the reference answers.) 1.00: Excellent (The generated text conveys exactly the same meaning as one or more of the reference answers.)"}, {"role": "user", "content": "Question: Reference Answers: {all_answers} Generated Text: {predictions}"}, {"role": "system", "content": "Provide a score for your rating. Retain two significant digits. Only output the score and do not output text."}

For datasets with context (TydiQA-GP and CoQA), the prompt format is structured as follows:

```
Prompt Structure for TydiQA-Gp and CoQA
 Prompt = [
         "system", "content": "You are an expert evaluator
{"role":
of text quality. Your task is to score the following text
generated by a language model on a scale of 0 to 1 based on
the provided multiple reference answers, where:
0.00: Poor (The meaning conveyed by the generated text is
irrelevant to the reference answers.)
1.00: Excellent (The generated text conveys exactly the same
meaning as one or more of the reference answers.)"},
{"role": "user", "content": "Reference Answers:
{all_answers}
Generated Text: {predictions}"},
{"role": "system", "content": "Provide a score for your
rating. Retain two significant digits. Only output the
score and do not output text."}
```

As shown in Figure 8, the empirical results indicate that when the threshold exceeds 0.7, the DeepSeek score remains relatively high, whereas BLEURT and ROUGE-L decrease substantially, leading to a reduction in overall performance. When the threshold falls below 0.5, the average score also drops outside the optimal range and exhibits increased instability. Overall, a threshold of 0.5 lies within the optimal performance region, providing a balanced trade-off across multiple metrics and mitigating the risk of overfitting to a single metric. Therefore, setting the threshold to 0.5 constitutes a reasonable and robust choice.

As shown in Figure 7, we randomly sampled 100 instances from the TruthfulQA dataset, applied a threshold of 0.5, and compared the consistency between various automatic labeling methods and expert annotations. The results indicate that the confusion matrix derived from DeepSeek-V3 aligns most closely with expert judgments, achieving an overall accuracy of 0.88 and an F1 score of



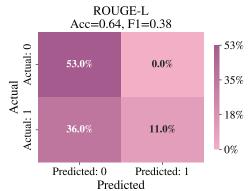


Figure 7: Confusion matrices of three labeling methods (DeepSeek-V3, BLEURT and ROUGE-L).

0.86, which demonstrates high agreement with human annotations. In contrast, BLEURT performs weaker (Acc=0.68, F1=0.56), while ROUGE-L exhibits the largest deviation (Acc=0.64, F1=0.38), particularly in distinguishing positive and negative samples. These results suggest that DeepSeek-V3 can serve as a reliable basis for automatic labeling, whereas ROUGE-L is not suitable as a robust evaluation criterion.

H RESULTS WITH QWEN2.5-7B-INSTRUCT

For the Qwen-2.5-7B-Instruct model, results are summarized in 3. We observe that training-free methods generally perform inconsistently across datasets and labeling criteria, with AUROC scores fluctuating significantly. By contrast, training-based methods achieve more stable improvements, yet they still fall short of SSP. Our method consistently outperforms all baselines under all three labeling criteria (ROUGE-L, BLEURT, and DeepSeek-V3), achieving the highest average AUROC of 72.72%. Notably, SSP shows large gains over self-evaluation and logit-based baselines, highlighting that representation discrepancy under perturbation provides a stronger and more reliable signal for hallucination detection on Qwen-2.5-7B-Instruct.

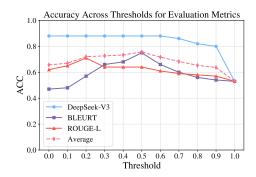


Figure 8: Overall performance across different thresholds, showing that 0.5 provides the best balance among all metrics.

Table 3: Comparison between our method (SSP) and competitive hallucination detection methods on the Qwen-2.5-7B-Instruct across four datasets. All values are AUROC scores in percentage. Bold numbers indicate the best performance and underlined numbers indicate the second best within each column. Results are reported under three labeling criteria: ROUGE-L (R), BLEURT (B), and DeepSeek-V3 (D).

	T	ruthfulQ	A	,	TriviaQ <i>A</i>	1		CoQA		T	ydiQA-C	βP		Average	!
Method	R	В	D	R	В	D	R	В	D	R	В	D	R	В	D
						Training	-free Me	thods							
Perplexity	52.68	59.08	53.60	55.45	56.69	52.72	68.58	63.85	62.03	55.10	53.17	51.97	57.95	58.20	55.08
Semantic Entropy	59.06	52.27	64.25	70.56	67.72	71.27	61.87	56.45	52.35	52.27	56.12	50.17	60.94	58.14	59.51
Lexical Similarity	65.55	60.40	57.50	66.89	64.39	65.55	74.55	70.43	71.62	60.10	53.88	61.75	66.77	62.28	64.11
EigenScore	68.48	57.98	52.67	75.57	71.25	68.36	75.68	71.53	72.33	62.95	56.17	60.97	70.67	64.23	63.58
SelfCKGPT	67.96	68.00	65.88	73.51	73.57	72.36	72.67	72.03	74.18	55.44	50.70	56.50	67.40	66.08	67.23
Verbalize	55.05	52.49	54.25	51.11	50.49	51.53	50.73	50.85	51.86	52.63	50.75	52.25	52.38	51.15	52.47
Self-evaluation	52.57	57.46	51.21	53.90	53.36	58.97	51.08	50.29	52.13	54.30	50.71	55.61	52.96	52.96	54.48
					7	Training-l	based Mo	ethods							
CCS	53.77	59.19	53.58	51.01	59.80	50.42	59.56	61.36	50.32	62.16	57.89	54.58	56.63	59.56	52.23
HaloScope	72.21	70.42	68.10	75.71	74.97	63.00	71.95	67.51	63.90	65.60	67.46	67.00	71.37	70.09	65.50
Linear probe	70.10	69.84	70.58	74.42	72.30	63.15	72.06	70.35	68.46	69.36	69.92	69.72	71.49	70.60	67.98
SAPLMA	70.91	70.68	71.84	74.82	74.71	66.90	72.84	70.21	69.34	68.75	70.14	68.67	71.83	71.44	69.19
EarlyDetec	71.51	70.17	66.99	73.97	75.34	73.13	71.11	68.83	67.24	65.65	69.49	69.16	70.56	70.96	69.13
EGH	68.27	66.71	63.21	74.21	70.46	67.96	74.58	72.81	70.91	68.91	64.12	65.31	71.49	68.53	66.85
SSP (Ours)	72.36	71.30	72.03	74.08	73.26	74.01	73.45	71.69	72.43	70.03	72.43	72.40	72.48	72.17	72.72

I EXTENDED RESULTS ON SSP GENERALIZATION

We evaluate the generalization capability of SSP across datasets with different distributions. Specifically, we directly transfer the learned sample-specific prompt and encoder from a source dataset "(s)" and apply them to a target dataset "(t)" to compute scores without additional training. Figure 9 (a) illustrates the strong cross-dataset transferability of our proposed SSP framework. When transferring parameters from TriviaQA to TydiQA-GP, SSP achieves an AUROC of 73.89% for hallucination detection, which is competitive with the in-domain performance on TruthfulQA (78.64%). Figure 9 (b), (c) and (d) show the generalization results of EGH, Linear probe and SAPLMA. Both methods exhibit weaker cross-dataset transferability compared to SSP, with notably lower AUROC scores in most off-diagonal entries. For instance, transferring from TriviaQA to TydiQA-GP yields 57.60% for EGH, 67.06% for the linear probe and 67.71% for SAPLMA, both falling short of SSP's 73.89% under the same setting. These results indicate that EGH suffers from limited representation generalization, while the SAPLMA, despite achieving competitive results in some cases, exhibits unstable performance across datasets.

J DETAILS OF PROMPT INITIALIZATION

To generate semantically neutral but stylistically varied noise prompts, we construct the following instruction template. We construct the initial prompt with the following structure:

You are an interference prompt generator.\n Generate one short stylistic sentence that can be appended to the given answer.\n Do not change the original meaning.\n Do not include any explanations, symbols, or unrelated content — only output the sentence itself.\n Q: [question]\n A: [answer]\n Interference:

K COMPARISON OF PROMPTING STRATEGIES AND SSP COMPONENTS.

We compare five variants to evaluate the impact of prompt design and components on hallucination detection. As shown in Table 4, *Sample-Specific Prompting* (SSP) consistently outperforms both Static prompt and Prompt Tuning. For example, on TruthfulQA, SSP improves AUROC by about 4.62% over Static prompt, achieving the highest average AUROC across all datasets (75.38%). These results demonstrate that SSP can dynamically generate adaptive prompts for each sample, thereby inducing more separable internal representations between truthful and hallucinatory responses. In contrast, fixed or globally tuned prompts fail to capture sample-level distinctions and

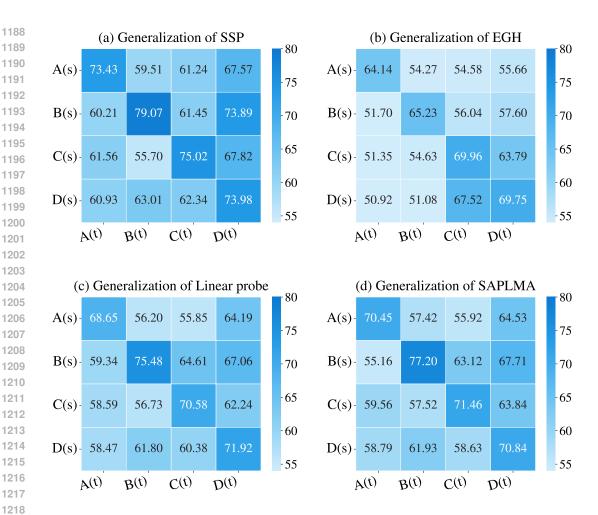


Figure 9: Generalization performance comparison. All values are AUROC scores (%). Here, TruthfulQA is denoted as **TQA**, TriviaQA as **B**, CoQA as **C**, and TydiQA-GP as **D**.

Table 4: **Prompting strategies and component ablations.** AUROC (%) results on four datasets.

Method	TruthfulQA	TriviaQA	CoQA	TydiQA-GP	Average
Static prompt w/o Encoder	57.38	54.96	54.88	54.46	55.42
Prompt Tuning w/o Encoder	60.59	56.05	57.07	70.41	61.03
SSP w/o Encoder	65.87	67.03	57.90	72.47	65.82
Static prompt	68.81	75.49	66.75	72.67	70.93
Prompt Tuning	70.21	76.21	66.88	73.05	71.59
SSP	73.43	79.07	75.02	73.98	75.38

thus lag behind. When the encoder is removed (w/o Encoder), all methods experience a performance drop, but SSP still maintains a clear advantage, highlighting its robustness.

Table 5: Ablation analysis of hallucination detection performance (AUROC %) by varying discrepancy functions as score metrics. The best results are in bold and the second best are underlined.

Method	TruthfulQA	TriviaQA	CoQA	TydiQA-GP	Average
Manhattan distance	59.18	54.21	59.31	56.99	57.42
Euclidean distance	63.60	72.38	60.11	59.23	63.83
KL-divergence	61.62	57.17	59.46	60.65	59.73
1 - Cosine similarity	73.43	79.07	75.02	73.98	75.38

Table 6: Results of discrepancy optimization direction. All values are AUROC scores (%).

Method	TruthfulQA	TriviaQA	CoQA	TydiQA-GP	Average
Reversed Objective	58.02	70.93	69.95	71.38	67.57
Original Objective	73.43	79.07	75.02	73.98	75.38

L EFFECT OF DISCREPANCY FUNCTION DESIGN.

We investigate how the design of the discrepancy function influences hallucination detection performance. Specifically, we compare the cosine-based formulation defined in Eq.equation 19 against alternative distance measures, including Manhattan distance (Malkauthekar, 2013), Euclidean distance (Danielsson, 1980), and Kullback–Leibler (KL) divergence (Csiszár, 1975). For each discrepancy function, we define a corresponding score function that computes the magnitude of representation change between the original and perturbed inputs. As shown in Table 5, the cosine-based metric consistently provides better separability between truthful and hallucinatory responses across all evaluated datasets.

M ABLATION ON THE DIRECTION OF DISCREPANCY OPTIMIZATION

We conduct an ablation study to examine whether optimizing in the intended direction—encouraging larger perturbation-induced changes for truthful responses and smaller ones for hallucinatory responses—is indeed beneficial. To this end, we reverse the discrepancy objective by setting $\tau_T=0.7$ and $\tau_H=0.3$. As shown in Table 6, this reversed setting results in a notable drop in detection performance across all datasets, confirming that the original objective direction better aligns with the underlying characteristics of truthful and hallucinatory responses.

N RESULTS WITH MORE TRAINING DATA

 In this section, we investigate the effect of increasing the number of labeled QA pairs used for training. Specifically, on the TruthfulQA dataset, we vary the number of labeled samples from 100 to 500 in increments of 100, while keeping the test set fixed. The results are reported in Table 7. We observe that all methods generally improve with more training data, and SSP outperforms both EGH and the linear probe baseline in most settings.

Notably, even with as few as 100 labeled examples, SSP achieves a high AUROC of 73.43%, which is comparable to or better than the performance of EGH trained on much larger datasets. This suggests that SSP is not only effective but also data-efficient to limited supervision, making it suitable for practical settings where labeled data is scarce.

Table 7: Effect of training data size on hallucination detection performance on TruthfulQA.

Model	100	200	300	400	500	512
EGH	64.14		67.44	67.55	68.36	69.48
Linear probe	68.65	72.13	73.44	74.21	74.07	76.74
SSP (Ours)	73.43	73.28	72.13	74.94	75.29	77.18

We further investigate the impact of increasing the training data size on hallucination detection by conducting experiments on larger subsets of the datasets. As shown in Table 8, scaling the number of training examples consistently improves performance across all methods. Among them, SSP benefits the most from additional data and achieves superior results across all three datasets.

O COMPUTE RESOURCES AND TIME

Software and Hardware. We conducted all experiments using Python 3.9.20 and PyTorch 1.13.1 on NVIDIA A40 GPUs. For evaluation with DeepSeek-V3, we utilized the official API provided by DeepSeek.

Inference Time. To further evaluate the practical applicability of our method, we compare the inference time and detection performance (AUROC) of different hallucination detection methods under the same data split and hardware setup on the TydiQA-GP dataset, using the LLaMA-3-8B-Instruct model. As shown in Figure 10, we report the inference time after completing the full sampling process to ensure consistency in measurement. The results show that, compared to the Semantic Entropy method, SSP achieves not only higher detection accuracy but also avoids the significant computational cost. Although SSP incurs slightly higher inference time than Haloscope and Linear probe, it provides better detection performance. Moreover, when compared to other methods such as EGH and EigenScore, SSP achieves a better balance between efficiency and accuracy. Overall, SSP requires only modest inference time per sample while maintaining efficient detection ca-

Table 8: Effect of training data size on hallucination detection performance using larger subsets of the datasets.

	Training Data Size									
Method	100	500	1000	2000						
TriviaQA										
Linear probe	75.48	77.32	78.01	80.52						
SAPLMA	77.20	78.03	79.14	82.15						
EGH	65.23	70.54	71.29	74.77						
SSP	79.07	80.03	81.31	83.25						
CoQA										
Linear probe	70.58	71.18	72.05	75.93						
SAPLMA	71.46	72.04	73.57	77.45						
EGH	69.96	71.03	72.47	77.86						
SSP	75.02	75.41	77.56	79.37						
	TydiQA-GP									
Linear probe	71.92	72.04	73.18	74.46						
SAPLMA	70.84	72.43	74.08	75.3						
EGH	69.75	70.37	71.63	76.37						
SSP	73.98	75.2	76.49	77.2						

pability, demonstrating its practicality for real-world deployment scenarios.

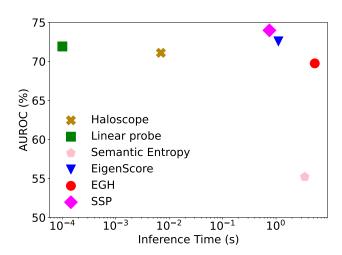


Figure 10: AUROC and inference time.

P BROADER IMPACT

Large language models (LLMs) have become widely adopted in both academic research and industrial applications, while ensuring the trustworthiness of their generated content remains a key challenge for safe deployment. To address this issue, we propose a novel hallucination detection method Sample-Specific Prompting (SSP), which detects hallucinations by injecting input-adaptive noise prompts and analyzing the model's internal representation shifts. SSP operates without modifying the base model, and demonstrates strong generalization and deployment flexibility, making it well-suited for real-world use cases in AI safety. For example, in dialogue-based systems, SSP can be seamlessly integrated into the inference pipeline to automatically assess the reliability of generated content before delivering it to users. Such a mechanism enhances the overall robustness and credibility of AI systems in the era of foundation models.

Q LIMITATIONS

We propose a hallucination detection method that induces internal representation shifts in LLMs by concatenating learnable, sample-specific prompts into the input. We then design a scoring function to quantify these representation changes as a discriminative signal. Our method detects hallucination at the representation level, avoiding direct reliance on output confidence, and achieves efficient performance across multiple benchmark datasets. However, SSP addresses hallucination detection in a white-box setting, as it requires access to internal representations of the LLM. However, it does not directly apply to black-box scenarios. In future work, we plan to extend the approach to black-box hallucination detection, thereby broadening its applicability to a wider range of real-world settings.

R LLM USAGE STATEMENT

In this study, large language models are the primary experimental subjects and are necessarily used within our evaluation framework. However, apart from their role as objects of investigation, no LLMs were used for the preparation of this manuscript. All conceptual development, analysis, writing, and editing were carried out solely by the authors without LLM assistance.