

# Shaking to Reveal: Perturbation-Based Detection of LLM Hallucinations

Jinyuan Luo<sup>1</sup> Zhen Fang<sup>1</sup> Yixuan Li<sup>2</sup> Seongheon Park<sup>2</sup> Ling Chen<sup>1\*</sup>

<sup>1</sup>Australian Artificial Intelligence Institute, University of Technology Sydney

<sup>2</sup>Department of Computer Sciences, University of Wisconsin-Madison

jinyuan.luo@student.uts.edu.au, {zhen.fang, ling.chen}@uts.edu.au  
{sharonli, seongheon\_park}@cs.wisc.edu

## Abstract

Hallucination remains a key obstacle to the reliable deployment of large language models (LLMs) in real-world question answering tasks. A widely adopted strategy to detect hallucination, known as self-assessment, relies on the model’s own output confidence to estimate the factual accuracy of its answers. However, this strategy assumes that the model’s output distribution closely reflects the true data distribution, which may not always hold in practice. As bias accumulates through the model’s layers, the final output can diverge from the underlying reasoning process, making output-level confidence an unreliable signal for hallucination detection. In this work, we propose **Sample-Specific Prompting (SSP)**, a new framework that improves self-assessment by analyzing perturbation sensitivity at intermediate representations. These representations, being less influenced by model bias, offer a more faithful view of the model’s latent reasoning process. Specifically, SSP dynamically generates noise prompts for each input and employs a lightweight encoder to amplify the changes in representations caused by the perturbation. A contrastive distance metric is then used to quantify these differences and separate truthful from hallucinated responses. By leveraging the dynamic behavior of intermediate representations under perturbation, SSP enables more reliable self-assessment. Extensive experiments demonstrate that SSP significantly outperforms prior methods across a range of hallucination detection benchmarks.

## 1 Introduction

In recent years, large language models (LLMs) have demonstrated remarkable capabilities in natural language processing tasks [54, 14]. However, the phenomenon of hallucination in their generated text remains a critical challenge. Hallucination refers to instances where the model produces text that is grammatically and logically coherent but lacks factual accuracy or a verifiable basis [34, 16]. This issue significantly hinders the applicability of LLMs in high-precision domains such as healthcare, law, and science [12, 55]. Consequently, hallucination detection has emerged as a crucial research problem in ensuring the reliability and trustworthiness of LLMs.

A popular strategy for the detection of hallucinations in LLM is self-assessment [1, 2, 4], which typically estimates the factuality of a response by leveraging the confidence in the output of the model. While intuitive and easy to implement, empirical studies have found that their effectiveness can degrade in more complex or realistic scenarios [9, 11, 13]. One potential reason is a mismatch between the model’s predictive distribution and the true data distribution [12]. As biases accumulate across layers, the final output may drift from the model’s internal reasoning, making output-layer confidence an unreliable signal for self-assessment. To overcome this limitation, recent work has begun to shift from probing the output to intermediate representations [9, 43, 56, 57, 10]. While intervening at

\*Corresponding author

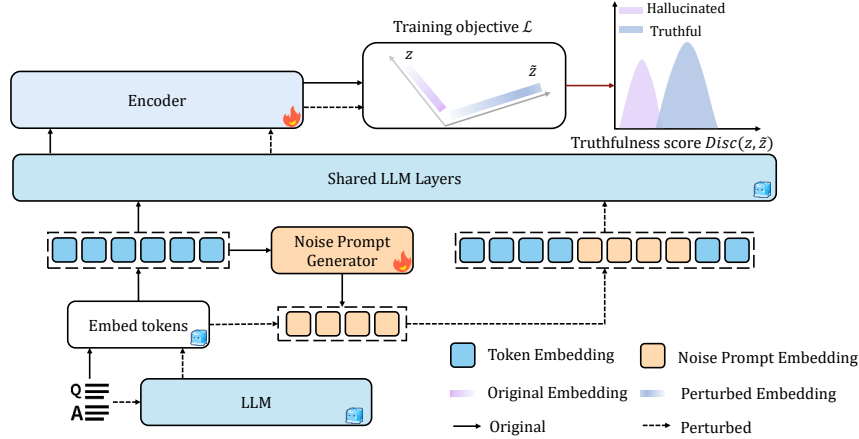


Figure 1: Overview of **Sample-Specific Prompting (SSP)** framework for hallucination detection. Given a question-answer (QA) pair, a noise prompt generator produces a perturbation adapted to the input. The noise prompt is appended to the original answer and passed through a shared LLM backbone to induce representational shifts. The encoder then maps the intermediate representations to a discriminative space and maximizes the discrepancy between truthful and hallucinated responses.

intermediate representations holds promise, performing self-assessment at this level poses significant challenges. Unlike the output layer, where the predicted distribution naturally provides probabilistic interpretations that reflect the model’s confidence in its predictions, intermediate representations lack such explicit interpretability [23, 24]. Consequently, *how to effectively leverage intermediate-layer information for reliable self-assessment* remains the central challenge that this work aims to address.

In this paper, we propose **Sample-Specific Prompting (SSP)**, a novel perturbation-based framework that leverages the differential sensitivity of intermediate representations as a signal for hallucination detection. Instead of relying on static or handcrafted prompts, SSP learns to dynamically generate controlled noise prompts tailored to each question–answer pair, inducing perturbations that reveal how internal features respond. Our key insight is that truthful and hallucinated responses exhibit distinct representational shifts under input perturbations. This observation is consistent with [59–62]: factual knowledge is typically encoded in well-structured internal representations, which are tightly coupled with the input and exhibit greater sensitivity in intermediate layers when perturbed, while hallucinated answers remain relatively stable. SSP *amplifies* this signal by introducing a lightweight encoder to extract and compare features before and after perturbation, and by explicitly optimizing a contrastive training objective that encourages larger representation shifts for truthful responses and smaller shifts for hallucinated ones. In effect, this joint learning of both perturbation prompts and representation encodings enables SSP to be a more effective self-assessment strategy.

Extensive experiments demonstrate the effectiveness performance of our method across diverse datasets. Compared to the state-of-the-art methods, we improve the hallucination detection accuracy by 4.78% (AUROC) on a challenging TruthfulQA benchmark [16]. Our results also indicate that SSP generalizes well across different domains. To better understand the role of each component, we conduct comprehensive ablation studies on SSP. The results show that each component contributes to the overall performance. Our key contributions are summarized as follows:

- We are the first to leverage the sensitivity of LLMs to input perturbations as a signal for hallucination detection, providing a novel perspective on this problem.
- We propose **Sample-Specific Prompting (SSP)**, which generates optimal perturbations for each sample, amplifying the distinction between truthful and hallucinated responses.
- We conduct ablation studies to evaluate the impact of different components of SSP and demonstrate its effectiveness across diverse LLMs and datasets.

## 2 Preliminary

**LLM generation probability.** Let  $P_\theta$  denote the conditional probability distribution defined by a pre-trained LLM with parameters  $\theta$ . Given an input sequence  $Q = \{x_1, \dots, x_k\}$  representing the question, where each  $x_j$  denotes a token in the input sequence. The model generates an answer  $A = \{x_{k+1}, \dots, x_{k+l}\}$  by predicting each token based on the preceding context:

$$P_\theta(x_j | x_1, \dots, x_{j-1}), \quad \text{for } j = k+1, \dots, k+l. \quad (1)$$

In practice,  $A$  is obtained via greedy decoding or beam search to approximate the maximum likelihood output under  $P_\theta$  [53]. The generated answer  $A$  is then used, along with the input question  $Q$ , as the input to a hallucination detector [9].

**Dataset format.** Each sample in the dataset consists of a question  $Q$ , a reference answer  $A_{\text{ref}}$ , and optionally a context passage (if provided by the dataset). For simplicity, we concatenate the context and the question into a single input sequence, which we denote as  $Q$ . The dataset can then be represented as:  $\mathcal{S} = \{(Q_1, A_{\text{ref}_1}), \dots, (Q_n, A_{\text{ref}_n})\}$ , where  $n$  denotes the total number of samples. Given an input  $Q$ , we use a LLM to generate an answer  $A$  in an autoregressive manner [53]. To facilitate hallucination detection, we assign a binary label  $y \in \{0, 1\}$  to each generated answer  $A$ , based on its semantic similarity to the reference answer  $A_{\text{ref}}$ . If  $A$  is consistent with  $A_{\text{ref}}$ , it is labeled as truthful ( $y = 1$ ); otherwise, it is labeled as hallucinated ( $y = 0$ ). The labeled dataset is defined as

$$\mathcal{S}_{\text{label}} = \{(Q_1, A_1, y_1), \dots, (Q_n, A_n, y_n)\}. \quad (2)$$

**Hallucination detection.** Following the practical setup in recent work [9], we denote the true data distribution over truthful input-generation pairs as  $P_{\text{true}}$ . Given a generated answer  $A$  and its corresponding question  $Q$ , the aim of hallucination detection is to learn a predictor  $G$  such that

$$G(Q, A) = \begin{cases} 1, & \text{if } A \sim P_{\text{true}}(\cdot | Q) \\ 0, & \text{otherwise} \end{cases}. \quad (3)$$

A discussion of related works is provided in Appendix A.

## 3 Motivation: Rethinking Self-evaluation for Hallucination detection

### 3.1 Self-evaluation and Its Limitation

**Rethinking Self-evaluation [1].** self-assessment [1, 2, 4] has emerged as a mainstream strategy for hallucination detection in recent research, which leverages the language model’s own outputs or internal signals to evaluate the factual consistency of its responses. Among these, Self-evaluation [1], a highly representative method, appends an evaluative prompt  $T$ , “*Is the proposed answer: (A) True (B) False The proposed answer is*”, to the original question-answer pairs  $(Q, A)$  and estimates the confidence of the response by extracting the probability distribution over the subsequent tokens. The probability is then interpreted as the model’s internal belief in the truthfulness of its own answer. This method leverages the characteristic that language models tend to produce well-calibrated token probabilities that reflect their internal confidence in a response [7], formalized as:

$$P_\theta(x = \text{True} | Q, A_{\text{Truth}}, T) > P_\theta(x = \text{True} | Q, A_{\text{Hallu}}, T),$$

which assumes that, given the same question  $Q$ , the model  $P_\theta$  assigns higher confidence to factually grounded answers than to hallucinated ones. Ideally, if the model distribution  $P_\theta$  perfectly matches the true data distribution  $P_{\text{true}}$ , the principle of Self-evaluation should remain valid and effective. However, empirical results reported in [9, 11, 13] suggest that its performance may be suboptimal in some specific tasks or scenarios, even under the self-assessment framework. Based on this observation, one possible reason for this limitation [12] is the mismatch between the model distribution  $P_\theta$  and the true distribution  $P_{\text{true}}$ , which may prevent Self-evaluation from fully achieving its intended effect.

**Towards intermediate-layer Self-evaluation.** The mismatch between distributions  $P_\theta$  and  $P_{\text{true}}$  can arise from various reasons, such as the training objective, model architecture, and algorithm design [17]. For example, [18, 19] argue that a language model’s next-token prediction confidence primarily reflects linguistic plausibility rather than factual correctness. Note that the discrepancy

between  $P_\theta$  and  $P_{\text{true}}$  often originates not only at the output layer, but also within the model’s intermediate representations. This issue is further reinforced by the architectural characteristics of mainstream LLMs (e.g., LLaMA-3 [14] and Qwen-2.5 [15]), which are typically composed of multiple stacked Transformer layers. [20, 26, 27] demonstrate that as biases propagate through the layers during forward passes, they tend to accumulate and become amplified, ultimately resulting in a significant discrepancy between the model’s predicted distribution  $P_\theta$  and the true distribution  $P_{\text{true}}$ .

Therefore, to effectively mitigate the impact of the bias on Self-evaluation, it is insufficient to focus adjustments solely at the output layer. Instead, Self-evaluation should be designed and optimized at the level of intermediate representations, aiming to suppress the accumulation of bias at its source. This strategy can enhance both the reliability and robustness of the evaluation process. While intervening at intermediate representations holds promise for mitigating the accumulation of distributional bias, performing Self-evaluation at this level poses significant challenges. Unlike the output layer, where the predicted distribution  $P_\theta$  naturally provides probabilistic interpretations that reflect the model’s confidence in its predictions, intermediate representations lack such explicit interpretability [23, 24]. Consequently, *how to effectively leverage intermediate-layer information for reliable Self-evaluation* remains the central challenge that this work aims to address.

### 3.2 Perturbation-based Self-evaluation

**Perturbation for Self-evaluation.** To address the challenge, we begin by making a slight yet essential modification to the standard Self-evaluation to enrich its underlying interpretation. Specifically, inspired by the perturbation strategy [21, 22], we introduce a noise prompt  $N$  into the Self-evaluation process and examine the change in the model’s confidence before and after the perturbation, i.e.,

$$\Delta P_\theta(Q, A, N, T) = |P_\theta(x = \text{True}|Q, A, T) - P_\theta(x = \text{True}|Q, A, N, T)|. \quad (4)$$

There are two key reasons motivating Eq. (4). *First*, by taking the difference between the predictions before and after perturbation, Eq. (4) may partially cancel out certain model-specific biases. As both terms are generated by the same model under similar conditions, shared systematic biases are likely to affect them similarly. This cancellation allows the resulting gap  $\Delta P_\theta$  to more accurately reflect the change induced by perturbation, rather than being dominated by the model’s inherent bias. *Second*, when  $P_\theta \approx P_{\text{true}}$ , the model’s sensitivity to perturbations can serve as an indicator of its confidence in the predicted answer. For a correct answer  $A_{\text{Truth}}$ , the model typically exhibits high confidence under standard conditions [25]. Introducing a noise prompt  $N$  may disrupt the input context and obscure the key evidence supporting the prediction, leading to a significant drop in the predicted probability:

$$P_\theta(x = \text{True} \mid Q, A_{\text{Truth}}, T) \gg P_\theta(x = \text{True} \mid Q, A_{\text{True}}, N, T).$$

In contrast, for a hallucinated answer  $A_{\text{Hallu}}$ , where the model generally lacks strong supporting evidence, its confidence remains low or unstable even before perturbation. As a result, adding noise prompt  $N$  has limited impact on the prediction. Based on above reasons, we expect the confidence gap induced by perturbation to satisfy the following relationship:

$$\text{there exist some noise prompts } N \text{ such that } \Delta P_\theta(Q, A_{\text{Truth}}, N, T) > \Delta P_\theta(Q, A_{\text{Hallu}}, N, T). \quad (5)$$

**Self-evaluation at internal representations.** Beyond its role as a probability gap at the output layer,  $\Delta P_\theta(Q, A, N, T)$  also admits a broader interpretation as a measure of *perturbation-induced change*. This interpretation is not restricted to the output probabilities and can be naturally extended to the internal representations of the model. In particular, it motivates us to examine how intermediate-layer features respond to input perturbations, providing a pathway to generalize Self-evaluation beyond the output layer. Assume that the internal representation is denoted by  $E_\theta(\cdot) \in \mathbb{R}^d$ . We then consider the perturbation-induced representation gap:

$$\Delta E_\theta(Q, A, N, T) = \mathbf{Disc}(E_\theta(x = \text{True}|Q, A, T), E_\theta(x = \text{True}|Q, A, N, T)), \quad (6)$$

where  $\mathbf{Disc}(\cdot, \cdot)$  is the measure to estimate the difference between the representations  $E_\theta(x = \text{True}|Q, A, T)$  and  $E_\theta(x = \text{True}|Q, A, N, T)$ . Similar to Eq. (5), we expect the representation gap induced by perturbation to satisfy the following relationship:

$$\text{there exist some noise prompts } N \text{ such that } \Delta E_\theta(Q, A_{\text{Truth}}, N, T) > \Delta E_\theta(Q, A_{\text{Hallu}}, N, T). \quad (7)$$

In the next section, we will introduce how to learn the noise prompt  $N$ .

## 4 Methodology

Following the motivation in Section 3, we introduce **Sample-Specific Prompting (SSP)**. An overview of the SSP framework is provided in Figure 1.

### 4.1 Discrepancy Function

Based on the perturbation framework introduced in Section 3.2, we describe how to extract and compare intermediate-layer representations under input perturbations. We divide the feature extraction process into two steps. *First*, we extract intermediate-layer representations from the original input  $(Q, A, T)$  and perturbed input  $(Q, A, N, T)$ . *Second*, to amplify the discrepancy between truthful and hallucinated responses under perturbation, we introduce a shared and learnable encoder module  $f_\phi(\cdot) \in \mathbb{R}^d$ , where  $\phi$  denotes its trainable parameters, which maps both original and perturbed intermediate representations into the same latent space. The encoder is designed to preserve the feature information of the LLM embeddings while amplifying the discrepancy between truthful and hallucinated responses. As a result, the original and perturbed representations are mapped into:

$$z = f_\phi(E_\theta(x = \text{True}|Q, A, T)), \quad \tilde{z} = f_\phi(E_\theta(x = \text{True}|Q, A, N, T)). \quad (8)$$

To quantify the magnitude of representation change before and after perturbation, we adopt cosine similarity as the measure. [28–30] have demonstrated that cosine-based metrics are robust to variations in feature magnitudes across different layers. Based on this, we define the discrepancy measure as:

$$\text{Disc}(z, \tilde{z}) = 1 - \cos(z, \tilde{z}) = 1 - \frac{z \cdot \tilde{z}}{\|z\| \|\tilde{z}\|}. \quad (9)$$

According to Eq. (7), the cosine similarity between  $z$  and  $\tilde{z}$  is lower for truthful answers, leading to higher discrepancy values compared to hallucinated responses. Formally, we expect that:  $\text{Disc}(z_{\text{Truth}}, \tilde{z}_{\text{Truth}}) > \text{Disc}(z_{\text{Hallu}}, \tilde{z}_{\text{Hallu}})$ . We also compare several alternative distance metrics (e.g., Euclidean distance [31], Manhattan distance [32]), and find that cosine-based discrepancy achieves the best empirical separability. Detailed results are presented in Table 3.

### 4.2 Sample-Specific Noise Prompt Generation

**Initialization of noise prompts.** Considering that each question–answer pair  $(Q, A)$  carries sample-specific characteristics, we initialize a unique noise prompt  $N$  for each sample. Specifically, given  $(Q, A)$ , we guide the LLM to dynamically generate a corresponding  $N$ , formalized as:

$$N \sim P_\theta(x \mid \text{SeedPrompt}, Q, A). \quad (10)$$

The SeedPrompt is an instruction designed to guide the generation of a natural language sentence that alters the stylistic tone without affecting the contextual semantics. To maintain factual consistency, we impose a semantic preservation constraint on the generation of  $N$ , requiring that it does not introduce semantic contradictions (see Appendix D for details). The generated noise prompt  $N$  is appended to the end of the answer  $A$  as its initialization, forming the perturbed input sequence  $(Q, A, N, T)$ .

**Sample-specific prompt learning.** However, relying solely on LLM-generated noise prompts  $N$  may not produce the optimal perturbations. To address this, we introduce a Sample-specific prompt learning strategy that dynamically optimizes the noise prompt  $N$  for each sample to maximize the perturbation-induced changes in intermediate representations. In implementation, we first extract the sentence embedding  $\mathbf{h}$  for each input by applying the LLM’s token-embedding layer  $\text{Emb}_\theta(\cdot)$ . Note that  $\text{Emb}_\theta(\cdot)$  is part of the pre-trained LLM and remains frozen during the training process. We have

$$\mathbf{h} = \text{Emb}_\theta(Q, A, T). \quad (11)$$

To dynamically optimize the noise prompt  $N$  based on the input pair  $(Q, A)$ , we introduce a lightweight prompt generator  $\text{M}_\varphi(\cdot)$ , implemented as a two-layer MLP. Here,  $\varphi$  denotes the trainable parameters of the generator. Despite being learnable,  $\text{M}_\varphi$  introduces no significant overhead. Specifically, the embedding of the noise prompt  $N$  is updated as:

$$\text{Emb}_\theta(N) = \text{M}_\varphi(\mathbf{h}) + \text{Emb}_\theta(N). \quad (12)$$

After updating  $\text{Emb}_\theta(N)$ , we concatenate it with the original input embeddings. Formally, the new embedding sequence is constructed as:

$$\text{Emb}_\theta(Q, A, N, T) = \text{Emb}_\theta(Q, A) \oplus \text{Emb}_\theta(N) \oplus \text{Emb}_\theta(T), \quad (13)$$

where  $\oplus$  denotes the concatenation operation along the sequence dimension. We then feed the combined sequence back into the LLM for forward propagation. Following the procedure described in Section 4.1, we extract both the original representation  $z$  and the perturbed representation  $\tilde{z}$  from the intermediate layers for training and hallucination detection.

### 4.3 Training Objective

Based on the definition of the discrepancy function in Eq. (9), we design a contrastive training loss that encourages larger perturbation-induced representation changes for truthful responses while maintaining smaller changes for hallucinated ones. The learnable components, including the MLP  $\mathbf{M}_\varphi(\cdot)$  and the encoder  $f_\phi(\cdot)$ , are optimized accordingly through this objective.

For samples labeled as truthful ( $y = 1$ ), we aim to amplify the difference between the original and perturbed features. The corresponding loss is defined as:

$$\ell_{\text{Truth}}^{(i)} = \max(0, \cos(z_i, \tilde{z}_i) - \tau_T), \quad (14)$$

where  $\tau_T$  is the upper bound threshold on the cosine similarity for truthful responses. For samples labeled as hallucinated ( $y = 0$ ), we aim to maintain a high similarity between the original and perturbed features. The corresponding loss is defined as:

$$\ell_{\text{Hallu}}^{(i)} = \max(0, \tau_H - \cos(z_i, \tilde{z}_i)), \quad (15)$$

where  $\tau_H$  is the lower bound threshold on the cosine similarity for hallucinated responses. Both  $\tau_T$  and  $\tau_H$  are treated as hyper-parameters. Given the labeled dataset  $\mathcal{S}_{\text{label}}$  introduced in Eq. (2), the final optimization problem can be written as:

$$\min_{\varphi, \phi} \frac{1}{n} \sum_{i=1}^n \left( y_i \cdot \ell_{\text{Truth}}^{(i)} + (1 - y_i) \cdot \ell_{\text{Hallu}}^{(i)} \right). \quad (16)$$

**Scoring strategy.** After training, we use the discrepancy function in Eq. (9) as the scoring mechanism for hallucination detection. A high discrepancy score indicates that the model’s internal semantics are significantly disturbed by the noise prompt  $N$ , suggesting more likely truthful response. Based on the scoring function, the hallucination detector is

$$G_\lambda(z, \tilde{z}) = \begin{cases} 1, & \text{Disc}(z, \tilde{z}) \geq \lambda \\ 0, & \text{otherwise} \end{cases}, \quad (17)$$

where  $\lambda$  is the threshold for detection.

## 5 Experiments

### 5.1 Experimental Setup

**Datasets and models.** We conduct experiments on four generative QA tasks: two open-book QA datasets CoQA [33] and TruthfulQA [16]; a closed-book QA dataset TriviaQA [34]; and a reading comprehension dataset TydiQA-GP (English) [35]. Following [9], we use only 100 labeled QA pairs for training, while keeping the size of the testing set consistent. More datasets and implementation details are provided in Appendix B. We evaluate our method on two families of widely used open-source LLMs that provide accessible internal representations: LLaMA-3-8B-Instruct [14] and Qwen-2.5-7B-Instruct [15]. By default, text generations are produced using greedy sampling, which selects the most probable token at each decoding step.

**Baselines.** We evaluate SSP against a diverse set of 12 baseline methods, including existing state-of-the-art. The baselines are categorized as follows: (1) logit-based methods-Perplexity [37] and Semantic Entropy [41]; (2) consistency-based methods-Lexical Similarity [36], SelfCKGPT [4]

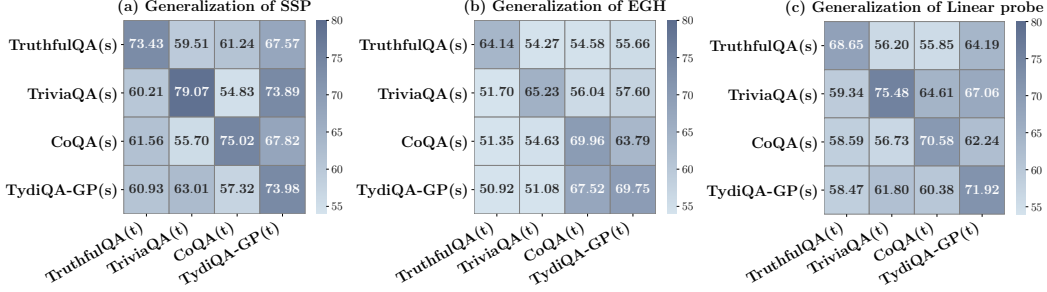


Figure 2: Generalization performance comparison across SSP (Ours), EGH [46], and Linear probe [43]. All values are AUROC scores (%).

and EigenScore [10]; (3) prompting-based methods-Verbalize [42] and Self-evaluation [1]; and (4) internal state-based methods-Contrast-Consistent Search (CCS) [6], HaloScope [9], Linear probe [43], EarlyDetec [45], EGH [46]. To ensure a fair comparison, we assess all baselines on identical test data, employing the default configurations as outlined in their respective papers.

**Evaluation.** Following previous works [11, 9], we evaluate the performance with the area under the curve of the receiver operator characteristic (AUROC). We consider the generation truthful when the similarity score between the generation and the reference answer is larger than a threshold of 0.5. We employ DeepSeek-V3 [49], a powerful open-source language model, to compute the similarity between generated answers and reference answers, which is then used to assign evaluation labels as detailed in Appendix C. Additionally, following [9], we show that the results are robust under two alternative similarity metrics—ROUGE [48] and BLEURT [47]—as detailed in Appendix E.

**Implementation details.** Following [9, 41], 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 [43, 10], we prepend the question to the generated answer and use the embedding of the final token to detect hallucinations. We implement the encoder  $f_\phi(\cdot)$  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  $\tau_T$  and  $\tau_H$  are set to 0.3 and 0.7, respectively.

## 5.2 Main Results

As shown in Table 1, we compare SSP with competitive hallucination detection methods from the literature. SSP achieves the highest average AUROC score, significantly outperforming other methods on both the LLaMA-3-8B-Instruct and Qwen-2.5-7B-Instruct. We observe that SSP outperforms logit-based baselines, exhibiting 11.3% and 7.78% improvement over Perplexity and Semantic Entropy on the challenging TruthfulQA task. From a computational perspective, both logit-based and consistency-based methods incur significant overhead during inference, as they require sampling multiple responses for each question. Following the setting in [9], 10 generations per question are used, which leads to substantial computational cost, especially when applied to large-scale datasets. In contrast, SSP only requires computing the representation shift before and after perturbation, making it significantly more efficient during inference. For prompting-based baselines, accumulated biases in intermediate layers can lead to unreliable confidence estimates, which limits their effectiveness in certain hallucination detection scenarios [51]. Lastly, we compare SSP with internal state-based methods, including CCS, HaloScope, Linear probe, EarlyDetec, and EGH. SSP consistently outperforms all baselines across datasets, achieving the highest average AUROC scores. This demonstrates that our method provides a more reliable signal for hallucination detection.

## 5.3 Generalization of SSP

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 2 (a)

Table 1: **Main results.** Comparison with competitive hallucination detection methods on different datasets. All values are percentages (AUROC, %). **Bold** numbers indicate the best performance, and underlined numbers indicate the second best.

Model	Method	TruthfulQA	TriviaQA	CoQA	TydiQA-GP	Average
LLaMA-3-8B-Instruct	Perplexity	62.13	76.64	64.87	53.40	64.26
	Semantic Entropy	58.88	78.53	55.15	55.21	61.94
	Lexical Similarity	53.64	78.22	77.47	60.94	67.57
	EigenScore	56.31	70.82	74.30	72.57	68.50
	SelfCKGPT	58.74	77.56	<b>78.67</b>	51.29	66.57
	Verbalize	59.70	55.43	53.39	53.39	55.48
	Self-evaluation	53.18	77.06	62.30	<b>76.69</b>	67.31
	CCS	53.91	58.58	52.40	74.11	59.75
	HaloScope	68.40	63.70	64.10	71.10	66.83
	Linear probe	68.65	75.48	70.58	71.92	71.66
	EarlyDetec	67.68	68.39	68.23	70.72	68.76
	EGH	64.14	65.23	69.96	69.75	67.27
	<b>SSP (Ours)</b>	<b>73.43</b>	<b>79.07</b>	75.02	73.98	<b>75.38</b>
Qwen2.5-7B-Instruct	Perplexity	53.60	52.72	62.03	51.97	55.08
	Semantic Entropy	64.25	71.27	52.35	50.17	59.51
	Lexical Similarity	57.50	65.55	71.62	61.75	64.11
	EigenScore	52.67	68.36	72.33	60.97	63.58
	SelfCKGPT	65.88	72.36	<b>74.18</b>	56.50	67.23
	Verbalize	54.25	51.53	51.86	52.25	52.47
	Self-evaluation	51.21	58.97	52.13	55.61	54.48
	CCS	53.58	50.42	50.32	54.58	52.23
	HaloScope	68.10	63.00	63.90	67.00	65.50
	Linear probe	70.58	63.15	68.46	69.72	67.98
	EarlyDetec	66.99	73.13	67.24	69.16	69.13
	EGH	63.21	67.96	70.91	65.31	66.85
	<b>SSP (Ours)</b>	<b>72.03</b>	<b>74.01</b>	72.43	<b>72.40</b>	<b>72.72</b>

Table 2: **Prompting strategies and component ablations.** AUROC (%) results on four datasets.

Method	TruthfulQA	TriviaQA	CoQA	TydiQA-GP	Average
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 w/o Encoder	65.87	67.03	57.90	72.47	65.82
SSP w/o SeedPrompt	72.20	<b>79.95</b>	74.21	73.44	74.95
<b>SSP</b>	<b>73.43</b>	79.07	<b>75.02</b>	<b>73.98</b>	<b>75.38</b>

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 2 (b) and (c) show the generalization results of EGH and the linear probe. 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 and 67.06% for the linear probe, both falling short of SSP’s 73.89% under the same setting. These results indicate that EGH suffers from limited representation generalization, while the linear probe, despite achieving competitive results in some cases, exhibits unstable performance across datasets.

## 5.4 Ablation Study

We conduct detailed ablation studies to investigate the contribution of each component in SSP. Additional ablation results are presented in Appendix E–I.

**Comparison of prompting strategies and SSP components.** We compare five variants to evaluate the impact of prompt design and components on hallucination detection. All experiments are conducted using the LLaMA-3-8B-Instruct model. As shown in Table 2, static prompt achieves a baseline performance of 70.93% average AUROC across datasets. Prompt Tuning offers a slight improvement (71.59%), indicating that global learned prompts can help but are still limited. Removing the encoder from SSP leads to a significant performance drop (65.82%), confirming its essential role in amplifying representational discrepancy. When the SeedPrompt is removed, performance decreases moderately (74.95%), showing that the SeedPrompt provides a useful inductive bias.

**Impact of layer selection on SSP performance.** Figure 3 (a) shows hallucination detection results using representations extracted from different layers of the LLM. AUROC scores for classifying truthful and hallucinated responses are computed using the LLaMA-3-8B-Instruct model. All other



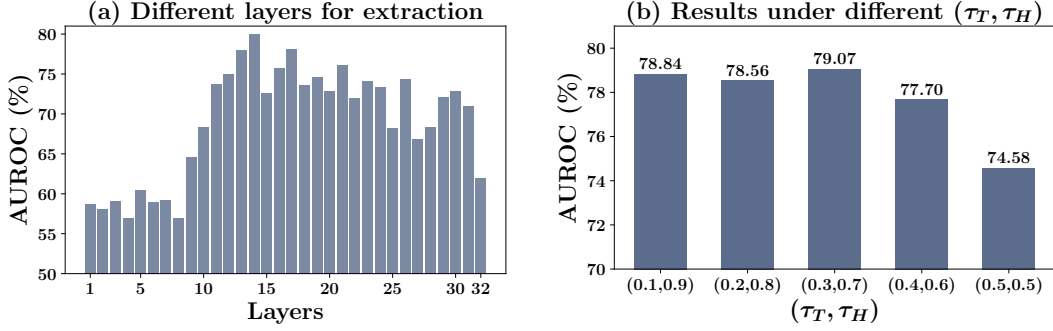


Figure 3: Visualization of performance under different layers (left) and threshold settings (right).

Table 3: Hallucination detection performance using different discrepancy functions as score metrics. All values are AUROC scores (%).

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
Cosine similarity	<b>73.43</b>	<b>79.07</b>	<b>75.02</b>	<b>73.98</b>	<b>75.38</b>

configurations follow the main experimental setup. We observe that performance increases with depth up to the middle layers, after which it starts to decline. This trend suggests that the LLM captures meaningful contextual semantics in its middle layers [43, 10]. As representations propagate deeper, accumulated deviations may degrade hallucination detection performance. These results highlight the effectiveness of internal representations in capturing meaningful signals for hallucination detection.

**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.(9) against alternative distance measures, including Manhattan distance [32], Euclidean distance [31], and Kullback–Leibler (KL) divergence [52]. 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 3, the cosine-based metric consistently provides better separability between truthful and hallucinated responses across all evaluated datasets. All experiments in this ablation study are conducted using the LLaMA-3-8B-Instruct model.

**Impact of threshold parameters  $\tau_T$  and  $\tau_H$ .** We examine the effect of the threshold hyperparameters  $\tau_T$  and  $\tau_H$  on the performance of our contrastive training objective. All experiments in this ablation study are conducted on the TriviaQA dataset. These thresholds control the sensitivity of the loss function to perturbation-induced changes in representations:  $\tau_T$  sets the minimum separation required for truthful samples, while  $\tau_H$  sets the maximum allowed deviation for hallucinated ones. As shown in Figure 3 (b), we observe that moderate values of  $\tau_T$  and  $\tau_H$  (e.g.,  $\tau_T = 0.3$ ,  $\tau_H = 0.7$ ) lead to optimal performance across datasets. In contrast, extremely high or low thresholds tend to tolerate excessive noise in the representations, leading to reduced detection accuracy.

## 6 Conclusion

In this work, we propose a novel framework SSP for hallucination detection, which leverages the differential sensitivity of intermediate representations under input perturbations. By dynamically generating noise prompts adapted to each input sample and amplifying shifts through a lightweight encoder, SSP effectively distinguishes between truthful and hallucinated samples at the representation level. Extensive experiments across multiple datasets and LLM architectures show the efficiency of SSP, making it a practical solution for hallucination detection in LLM outputs.

## References

- [1] S. Kadavath, T. Conerly, A. Askell, T. Henighan, D. Drain, E. Perez, N. Schiefer, Z. Hatfield-Dodds, N. DasSarma, E. Tran-Johnson *et al.*, “Language models (mostly) know what they know,” *arXiv preprint arXiv:2207.05221*, 2022.
- [2] H. Duan, Y. Yang, and K. Y. Tam, “Do llms know about hallucination? an empirical investigation of llm’s hidden states,” *arXiv preprint arXiv:2402.09733*, 2024.
- [3] Y. Liang, Z. Song, H. Wang, and J. Zhang, “Learning to trust your feelings: Leveraging self-awareness in llms for hallucination mitigation,” *arXiv preprint arXiv:2401.15449*, 2024.
- [4] P. Manakul, A. Liusie, and M. J. Gales, “Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models,” *arXiv preprint arXiv:2303.08896*, 2023.
- [5] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, “Language models are unsupervised multitask learners,” *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.
- [6] C. Burns, H. Ye, D. Klein, and J. Steinhardt, “Discovering latent knowledge in language models without supervision,” *arXiv preprint arXiv:2212.03827*, 2022.
- [7] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger, “On calibration of modern neural networks,” in *International conference on machine learning*. PMLR, 2017, pp. 1321–1330.
- [8] D. Zheng, M. Lapata, and J. Z. Pan, “Large language models as reliable knowledge bases?” *arXiv preprint arXiv:2407.13578*, 2024.
- [9] X. Du, C. Xiao, and S. Li, “Haloscope: Harnessing unlabeled llm generations for hallucination detection,” *Advances in Neural Information Processing Systems*, vol. 37, pp. 102 948–102 972, 2024.
- [10] C. Chen, K. Liu, Z. Chen, Y. Gu, Y. Wu, M. Tao, Z. Fu, and J. Ye, “Inside: Llms’ internal states retain the power of hallucination detection,” *arXiv preprint arXiv:2402.03744*, 2024.
- [11] S. Park, X. Du, M.-H. Yeh, H. Wang, and Y. Li, “How to steer llm latents for hallucination detection?” *arXiv preprint arXiv:2503.01917*, 2025.
- [12] Z. Ji, N. Lee, R. Frieske, T. Yu, D. Su, Y. Xu, E. Ishii, Y. J. Bang, A. Madotto, and P. Fung, “Survey of hallucination in natural language generation,” *ACM computing surveys*, vol. 55, no. 12, pp. 1–38, 2023.
- [13] S. Dasgupta, S. Nath, A. Basu, P. Shamsolmoali, and S. Das, “Hallushift: Measuring distribution shifts towards hallucination detection in llms,” *arXiv preprint arXiv:2504.09482*, 2025.
- [14] A. Grattafiori, A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Vaughan *et al.*, “The llama 3 herd of models,” *arXiv preprint arXiv:2407.21783*, 2024.
- [15] A. Yang, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Li, D. Liu, F. Huang, H. Wei *et al.*, “Qwen2. 5 technical report,” *arXiv preprint arXiv:2412.15115*, 2024.
- [16] S. Lin, J. Hilton, and O. Evans, “Truthfulqa: Measuring how models mimic human falsehoods,” *arXiv preprint arXiv:2109.07958*, 2021.
- [17] R. Navigli, S. Conia, and B. Ross, “Biases in large language models: origins, inventory, and discussion,” *ACM Journal of Data and Information Quality*, vol. 15, no. 2, pp. 1–21, 2023.
- [18] G. Prato, J. Huang, P. Parthasarathi, S. Sodhani, and S. Chandar, “Do large language models know how much they know?” *arXiv preprint arXiv:2502.19573*, 2025.
- [19] J. Ren, Y. Zhao, T. Vu, P. J. Liu, and B. Lakshminarayanan, “Self-evaluation improves selective generation in large language models,” in *Proceedings on*. PMLR, 2023, pp. 49–64.
- [20] S. S. Schoenholz, J. Gilmer, S. Ganguli, and J. Sohl-Dickstein, “Deep information propagation,” *arXiv preprint arXiv:1611.01232*, 2016.

- [21] E. Wallace, S. Feng, N. Kandpal, M. Gardner, and S. Singh, “Universal adversarial triggers for attacking and analyzing nlp,” *arXiv preprint arXiv:1908.07125*, 2019.
- [22] C. Guo, A. Sablayrolles, H. Jégou, and D. Kiela, “Gradient-based adversarial attacks against text transformers,” *arXiv preprint arXiv:2104.13733*, 2021.
- [23] A. Bau, Y. Belinkov, H. Sajjad, N. Durrani, F. Dalvi, and J. Glass, “Identifying and controlling important neurons in neural machine translation,” *arXiv preprint arXiv:1811.01157*, 2018.
- [24] A. Rogers, O. Kovaleva, and A. Rumshisky, “A primer in bertology: What we know about how bert works,” *Transactions of the association for computational linguistics*, vol. 8, pp. 842–866, 2021.
- [25] Z. Jiang, J. Araki, H. Ding, and G. Neubig, “How can we know when language models know? on the calibration of language models for question answering,” *Transactions of the Association for Computational Linguistics*, vol. 9, pp. 962–977, 2021.
- [26] X. Wang, Y. Xiong, B. Kang, Y. Zhang, P. S. Yu, and Y. Zhu, “Reducing negative effects of the biases of language models in zero-shot setting,” in *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining*, 2023, pp. 904–912.
- [27] X. Glorot and Y. Bengio, “Understanding the difficulty of training deep feedforward neural networks,” in *Proceedings of the thirteenth international conference on artificial intelligence and statistics. JMLR Workshop and Conference Proceedings*, 2010, pp. 249–256.
- [28] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, “A simple framework for contrastive learning of visual representations,” in *International conference on machine learning*. PmLR, 2020, pp. 1597–1607.
- [29] T. Zhang, V. Kishore, F. Wu, K. Q. Weinberger, and Y. Artzi, “Bertscore: Evaluating text generation with bert,” *arXiv preprint arXiv:1904.09675*, 2019.
- [30] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, “Momentum contrast for unsupervised visual representation learning,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 9729–9738.
- [31] P.-E. Danielsson, “Euclidean distance mapping,” *Computer Graphics and image processing*, vol. 14, no. 3, pp. 227–248, 1980.
- [32] M. Malkauthekar, “Analysis of euclidean distance and manhattan distance measure in face recognition,” in *Third International Conference on Computational Intelligence and Information Technology (CIIT 2013)*. IET, 2013, pp. 503–507.
- [33] S. Reddy, D. Chen, and C. D. Manning, “Coqa: A conversational question answering challenge,” *Transactions of the Association for Computational Linguistics*, vol. 7, pp. 249–266, 2019.
- [34] M. Joshi, E. Choi, D. S. Weld, and L. Zettlemoyer, “Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension,” *arXiv preprint arXiv:1705.03551*, 2017.
- [35] J. H. Clark, E. Choi, M. Collins, D. Garrette, T. Kwiatkowski, V. Nikolaev, and J. Palomaki, “Tydi qa: A benchmark for information-seeking question answering in ty pologically di verse languages,” *Transactions of the Association for Computational Linguistics*, vol. 8, pp. 454–470, 2020.
- [36] Z. Lin, S. Trivedi, and J. Sun, “Generating with confidence: Uncertainty quantification for black-box large language models,” *arXiv preprint arXiv:2305.19187*, 2023.
- [37] J. Ren, J. Luo, Y. Zhao, K. Krishna, M. Saleh, B. Lakshminarayanan, and P. J. Liu, “Out-of-distribution detection and selective generation for conditional language models,” *arXiv preprint arXiv:2209.15558*, 2022.
- [38] K. Li, O. Patel, F. Viégas, H. Pfister, and M. Wattenberg, “Inference-time intervention: Eliciting truthful answers from a language model,” *Advances in Neural Information Processing Systems*, vol. 36, pp. 41 451–41 530, 2023.

- [39] X. Bi, D. Chen, G. Chen, S. Chen, D. Dai, C. Deng, H. Ding, K. Dong, Q. Du, Z. Fu *et al.*, “Deepseek llm: Scaling open-source language models with longtermism,” *arXiv preprint arXiv:2401.02954*, 2024.
- [40] A. Malinin and M. Gales, “Uncertainty estimation in autoregressive structured prediction,” *arXiv preprint arXiv:2002.07650*, 2020.
- [41] L. Kuhn, Y. Gal, and S. Farquhar, “Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation,” *arXiv preprint arXiv:2302.09664*, 2023.
- [42] S. Lin, J. Hilton, and O. Evans, “Teaching models to express their uncertainty in words,” *arXiv preprint arXiv:2205.14334*, 2022.
- [43] A. Azaria and T. Mitchell, “The internal state of an llm knows when it’s lying,” *arXiv preprint arXiv:2304.13734*, 2023.
- [44] Y. Zha, Y. Yang, R. Li, and Z. Hu, “Alignscore: Evaluating factual consistency with a unified alignment function,” *arXiv preprint arXiv:2305.16739*, 2023.
- [45] B. Snyder, M. Moisesescu, and M. B. Zafar, “On early detection of hallucinations in factual question answering,” in *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2024, pp. 2721–2732.
- [46] X. Hu, Y. Zhang, R. Peng, H. Zhang, C. Wu, G. Chen, and J. Zhao, “Embedding and gradient say wrong: A white-box method for hallucination detection,” in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 2024, pp. 1950–1959.
- [47] T. Sellam, D. Das, and A. P. Parikh, “Bleurt: Learning robust metrics for text generation,” *arXiv preprint arXiv:2004.04696*, 2020.
- [48] C.-Y. Lin, “Rouge: A package for automatic evaluation of summaries,” in *Text summarization branches out*, 2004, pp. 74–81.
- [49] A. Liu, B. Feng, B. Xue, B. Wang, B. Wu, C. Lu, C. Zhao, C. Deng, C. Zhang, C. Ruan *et al.*, “Deepseek-v3 technical report,” *arXiv preprint arXiv:2412.19437*, 2024.
- [50] Y.-S. Chuang, L. Qiu, C.-Y. Hsieh, R. Krishna, Y. Kim, and J. Glass, “Lookback lens: Detecting and mitigating contextual hallucinations in large language models using only attention maps,” *arXiv preprint arXiv:2407.07071*, 2024.
- [51] K. Zhou, D. Jurafsky, and T. Hashimoto, “Navigating the grey area: How expressions of uncertainty and overconfidence affect language models,” *arXiv preprint arXiv:2302.13439*, 2023.
- [52] I. Csiszár, “I-divergence geometry of probability distributions and minimization problems,” *The annals of probability*, 1975.
- [53] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [54] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat *et al.*, “Gpt-4 technical report,” *arXiv preprint arXiv:2303.08774*, 2023.
- [55] H. Liu, W. Xue, Y. Chen, D. Chen, X. Zhao, K. Wang, L. Hou, R. Li, and W. Peng, “A survey on hallucination in large vision-language models,” *arXiv preprint arXiv:2402.00253*, 2024.
- [56] S. Marks and M. Tegmark, “The geometry of truth: Emergent linear structure in large language model representations of true/false datasets,” *arXiv preprint arXiv:2310.06824*, 2023.
- [57] F. Yin, J. Srinivasa, and K.-W. Chang, “Characterizing truthfulness in large language model generations with local intrinsic dimension,” *arXiv preprint arXiv:2402.18048*, 2024.

- [58] M. Jia, L. Tang, B.-C. Chen, C. Cardie, S. Belongie, B. Hariharan, and S.-N. Lim, “Visual prompt tuning,” in *European conference on computer vision*. Springer, 2022, pp. 709–727.
- [59] K. Meng, D. Bau, A. Andonian, and Y. Belinkov, “Locating and editing factual associations in gpt,” *Advances in neural information processing systems*, vol. 35, pp. 17 359–17 372, 2022.
- [60] A. Gupta, S. Baskaran, and G. Anumanchipalli, “Rebuilding rome: Resolving model collapse during sequential model editing,” *arXiv preprint arXiv:2403.07175*, 2024.
- [61] Y. Zhang, Z. Wei, J. Sun, and M. Sun, “Adversarial representation engineering: A general model editing framework for large language models,” *arXiv preprint arXiv:2404.13752*, 2024.
- [62] Z. Liao, K. Chen, Y. Lin, K. Li, Y. Liu, H. Chen, X. Huang, and Y. Yu, “Attack and defense techniques in large language models: A survey and new perspectives,” *arXiv preprint arXiv:2505.00976*, 2025.
- [63] N. M. Guerreiro, E. Voita, and A. F. Martins, “Looking for a needle in a haystack: A comprehensive study of hallucinations in neural machine translation,” *arXiv preprint arXiv:2208.05309*, 2022.
- [64] Y. Huang, J. Song, Z. Wang, S. Zhao, H. Chen, F. Juefei-Xu, and L. Ma, “Look before you leap: An exploratory study of uncertainty measurement for large language models,” *arXiv preprint arXiv:2307.10236*, 2023.
- [65] Y. Zhang, L. Cui, W. Bi, and S. Shi, “Alleviating hallucinations of large language models through induced hallucinations,” *arXiv preprint arXiv:2312.15710*, 2023.
- [66] Z. Xu, S. Jain, and M. Kankanhalli, “Hallucination is inevitable: An innate limitation of large language models,” *arXiv preprint arXiv:2401.11817*, 2024.
- [67] T. Zhang, L. Qiu, Q. Guo, C. Deng, Y. Zhang, Z. Zhang, C. Zhou, X. Wang, and L. Fu, “Enhancing uncertainty-based hallucination detection with stronger focus,” *arXiv preprint arXiv:2311.13230*, 2023.
- [68] I. Chern, S. Chern, S. Chen, W. Yuan, K. Feng, C. Zhou, J. He, G. Neubig, P. Liu *et al.*, “Factool: Factuality detection in generative ai—a tool augmented framework for multi-task and multi-domain scenarios,” *arXiv preprint arXiv:2307.13528*, 2023.
- [69] S. Min, K. Krishna, X. Lyu, M. Lewis, W.-t. Yih, P. W. Koh, M. Iyyer, L. Zettlemoyer, and H. Hajishirzi, “Factscore: Fine-grained atomic evaluation of factual precision in long form text generation,” *arXiv preprint arXiv:2305.14251*, 2023.
- [70] J. Duan, H. Cheng, S. Wang, A. Zavalny, C. Wang, R. Xu, B. Kailkhura, and K. Xu, “Shifting attention to relevance: Towards the uncertainty estimation of large language models,” 2023.
- [71] A. Agrawal, M. Suzgun, L. Mackey, and A. T. Kalai, “Do language models know when they’re hallucinating references?” *arXiv preprint arXiv:2305.18248*, 2023.
- [72] R. Cohen, M. Hamri, M. Geva, and A. Globerson, “Lm vs lm: Detecting factual errors via cross examination,” *arXiv preprint arXiv:2305.13281*, 2023.
- [73] N. Mündler, J. He, S. Jenko, and M. Vechev, “Self-contradictory hallucinations of large language models: Evaluation, detection and mitigation,” *arXiv preprint arXiv:2305.15852*, 2023.
- [74] K. Tian, E. Mitchell, A. Zhou, A. Sharma, R. Rafailov, H. Yao, C. Finn, and C. D. Manning, “Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback,” *arXiv preprint arXiv:2305.14975*, 2023.
- [75] B. Lester, R. Al-Rfou, and N. Constant, “The power of scale for parameter-efficient prompt tuning,” *arXiv preprint arXiv:2104.08691*, 2021.
- [76] X. L. Li and P. Liang, “Prefix-tuning: Optimizing continuous prompts for generation,” *arXiv preprint arXiv:2101.00190*, 2021.

- [77] X. Liu, K. Ji, Y. Fu, W. L. Tam, Z. Du, Z. Yang, and J. Tang, “P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks,” *arXiv preprint arXiv:2110.07602*, 2021.
- [78] C. Cai, Z. Ye, L. Feng, J. Qi, and F. Liu, “Sample-specific masks for visual reprogramming-based prompting,” *arXiv preprint arXiv:2406.03150*, 2024.

## Appendix

### A Related Works

**Hallucination detection** has become an increasingly important research topic, aiming to address the safety and reliability challenges of deploying LLMs in real-world applications [12, 55, 63–69]. Most existing approaches detect hallucinations by designing uncertainty-based scoring functions, including those that rely on output logits [70, 41, 40], based on the assumption that hallucinations are associated with token probability patterns that are inconsistent with those generated for truthful answers. Some methods detect hallucinations by analyzing the consistency among multiple generations [71, 72, 4, 73], or by prompting LLMs to estimate confidence in their own responses [1, 42, 19, 74, 51]. Recently, there is increasing interest in leveraging internal activations for hallucination detection, as demonstrated by methods like [6, 9, 46, 45, 43]. Despite the growing interest in internal-state-based methods, many of them rely on static representations (e.g., final-layer embeddings) and fail to exploit how LLMs dynamically react to perturbations. For instance, linear probing methods typically perform classification based on frozen representations [43]. However, [11] reveals that the embeddings produced by pre-trained language models often exhibit significant overlap between truthful and hallucinated responses, making them suboptimal for direct classification. In contrast, our approach circumvents this limitation by focusing on the perturbation-induced dynamics of internal representations. By measuring the magnitude of representational shifts rather than relying on static embeddings, our method provides a more discriminative signal for hallucination detection.

**Prompt learning** enables controlled generation in LLMs during inference. Early work typically relied on manually crafted static prompts [1, 42] to elicit self-assessment from the model regarding the factuality of its responses. However, such static prompts often suffer from poor adaptability across different questions and contexts, making them prone to failure. To improve generalization, some studies have adopted prompt tuning or soft prompt techniques [58, 75, 76] to optimize prompts in a parameterized manner. Nevertheless, these methods mostly aim for a globally optimal prompt and overlook sample-specific contextual differences, limiting their effectiveness under diverse inputs [78].

Our method differs in the following key aspects: (1) Our method performs hallucination detection based on the sensitivity of internal representations to input perturbations, rather than relying on static representations. (2) We dynamically generate adaptive noise prompts for each sample, effectively inducing distinguishable differences in representations between truthful and hallucinated responses.

### B Datasets and Implementation Details

**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:

*Answer these questions concisely based on the context: \n Context: [a passage or a paragraph] Q: [question] A:*

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 Asta’s friend Sharkie. "It looks like a bird’s belly," said Asta. But when they swam closer, it was not a bird’s 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 Asta’s papa. "What does it say?" they asked. \n \n Asta’s 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, isn’t 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: Asta’s 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:

**Baseline implementation details.** For Perplexity method [37], we follow the implementation here<sup>2</sup>, 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 [36], we use the Rouge-L as the similarity metric, and for SelfCKGPT [4], we adopt the NLI version as recommended in their codebase<sup>3</sup>, 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

<sup>2</sup><https://huggingface.co/docs/transformers/en/perplexity>

<sup>3</sup><https://github.com/potsawee/selfcheckgpt>



generations. For Haloscope [9], we adopt the official implementation available at <sup>4</sup>. For EGH [46], we follow the released codebase at <sup>5</sup>. For promoting-based baselines, we adopt the following prompt for Verbalize [38] 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 [1], 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.

## C Labeling with DeepSeek-V3

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: {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:

<sup>4</sup><https://github.com/deeplearning-wisc/haloscope>

<sup>5</sup><https://github.com/Xiaom-Hu/EGH>

### Prompt Structure for TydiQA-Gp and CoQA

```
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 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."}
]
```

## D Details of SeedPrompt

To generate semantically neutral but stylistically varied noise prompts, we construct the following instruction template, referred to as the SeedPrompt. We construct the SeedPrompt 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:*

## E Results with Other Metrics

In our main paper, a generation is considered truthful if its DeepSeek-V3 score with the gold standard answer exceeds a predefined threshold. In addition to the evaluation using DeepSeek-V3, we employ BLEURT and Rouge-L to determine the truthfulness of the generation. The corresponding experimental results are presented in Tables 4 and 5.

Table 5: **Results with BLEURT.** Comparison with competitive hallucination detection methods on different datasets. All values are percentages (AUROC, %). **Bold** numbers indicate the best performance, and underlined numbers indicate the second best.

Model	Method	TruthfulQA	TriviaQA	CoQA	TydiQA-GP	Average
LLaMA-3-8B-Instruct	Perplexity	62.11	71.37	62.55	51.43	61.87
	Semantic Entropy	51.97	72.78	53.52	54.66	58.23
	Lexical Similarity	52.27	73.97	72.67	62.28	65.30
	EigenScore	53.73	73.43	73.76	64.38	66.33
	SelfCKGPT	52.57	74.91	<b>74.04</b>	59.30	65.21
	Verbalize	58.77	55.07	51.59	51.36	54.20
	Self-evaluation	55.98	72.61	58.94	62.56	62.52
	CCS	52.26	55.75	53.27	63.93	56.30
	HaloScope	70.96	70.52	65.38	72.41	69.82
	Linear probe	<u>72.41</u>	<b>75.65</b>	71.79	<u>73.68</u>	<u>73.38</u>
	EarlyDetec	72.40	70.47	71.03	69.42	70.83
	EGH	71.28	69.48	68.63	70.54	69.98
	<b>SSP (Ours)</b>	<b>73.93</b>	<u>75.49</u>	<u>73.86</u>	<b>73.92</b>	<b>74.30</b>
Qwen2.5-7B-Instruct	Perplexity	59.08	56.69	63.85	53.17	58.20
	Semantic Entropy	52.27	67.72	56.45	56.12	58.14
	Lexical Similarity	60.40	64.39	70.43	53.88	62.28
	EigenScore	57.98	71.25	71.53	56.17	64.23
	SelfCKGPT	68.00	73.57	72.03	50.70	66.08
	Verbalize	52.49	50.49	50.85	50.75	51.15
	Self-evaluation	57.46	53.36	50.29	50.71	52.96
	CCS	59.19	59.80	61.36	57.89	59.56
	HaloScope	<u>70.42</u>	<u>74.97</u>	67.51	67.46	70.09
	Linear probe	69.84	<u>72.30</u>	70.35	<u>69.92</u>	70.60
	EarlyDetec	70.17	<b>75.34</b>	68.83	69.49	70.96
	EGH	66.71	70.46	<b>72.81</b>	64.12	68.53
	<b>SSP (Ours)</b>	<b>71.30</b>	73.26	<u>71.69</u>	<b>72.43</b>	<b>72.17</b>

Table 4: **Results with Rouge-L.** Comparison with competitive hallucination detection methods on different datasets. All values are percentages (AUROC, %). **Bold** numbers indicate the best performance, and underlined numbers indicate the second best.

Model	Method	TruthfulQA	TriviaQA	CoQA	TydiQA-GP	Average
LLaMA-3-8B-Instruct	Perplexity	50.02	72.32	70.01	54.78	61.78
	Semantic Entropy	61.26	73.45	53.34	56.70	61.19
	Lexical Similarity	57.69	76.10	68.84	63.25	66.47
	EigenScore	67.59	74.19	70.59	68.30	70.17
	SelfCKGPT	50.07	<u>77.37</u>	74.31	59.00	65.19
	Verbalize	64.87	<u>55.43</u>	52.49	51.59	56.10
	Self-evaluation	55.43	74.23	57.19	64.09	62.74
	CCS	68.09	56.85	50.96	68.69	61.15
	HaloScope	<u>73.60</u>	65.47	67.02	71.01	69.28
	Linear probe	71.83	76.35	73.09	<u>71.41</u>	<u>73.17</u>
	EarlyDetec	69.38	69.53	<b>75.84</b>	<u>70.08</u>	71.21
	EGH	70.60	61.89	75.60	71.33	69.86
	<b>SSP (Ours)</b>	<b>74.47</b>	<b>78.81</b>	<u>74.26</u>	<b>72.23</b>	<b>74.94</b>
Qwen2.5-7B-Instruct	Perplexity	52.68	55.45	68.58	55.10	57.95
	Semantic Entropy	59.06	70.56	61.87	52.27	60.94
	Lexical Similarity	65.55	66.89	74.55	60.10	66.77
	EigenScore	68.48	75.57	<b>75.68</b>	62.95	70.67
	SelfCKGPT	67.96	<u>73.51</u>	72.67	55.44	67.40
	Verbalize	55.05	51.11	50.73	52.63	52.38
	Self-evaluation	52.57	53.90	51.08	54.30	52.96
	CCS	53.77	51.01	59.56	62.16	56.63
	HaloScope	<u>72.21</u>	<b>75.71</b>	71.95	65.60	71.37
	Linear probe	70.10	74.42	72.06	69.36	71.49
	EarlyDetec	71.51	73.97	71.11	65.65	70.56
	EGH	68.27	74.21	74.58	68.91	71.49
	<b>SSP (Ours)</b>	<b>72.36</b>	74.08	<u>73.45</u>	<b>70.03</b>	<b>72.48</b>

## F 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 hallucinated responses—is indeed beneficial. To this end, we reverse the discrepancy objective by setting  $\tau_T = 0.7$  and  $\tau_H = 0.3$ , which encourages the opposite behavior. As shown in Table 6, this reversed setting results in a notable drop in detection performance across all datasets, confirming that the original

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	<b>73.43</b>	<b>79.07</b>	<b>75.02</b>	<b>73.98</b>	<b>75.38</b>

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

Model	100	200	300	400	500	512
EGH	64.14	65.73	67.44	67.55	68.36	69.48
Linear probe	68.65	72.13	<b>73.44</b>	74.21	74.07	76.74
<b>SSP (Ours)</b>	<b>73.43</b>	<b>73.28</b>	72.13	<b>74.94</b>	<b>75.29</b>	<b>77.18</b>

objective direction better aligns with the underlying characteristics of truthful and hallucinated responses.

## G 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.

## H Qualitative Results

To further illustrate the effectiveness of our method, we present qualitative examples from the TruthfulQA dataset using the LLaMA-3-8B-Instruct model. For each input, we compare the discrepancy scores produced by three configurations: (1) a static sentence appended as perturbation, (2) learned prompt via prompt tuning, and (3) our proposed sample-specific prompting (SSP). As shown in Figure 4, SSP consistently assigns higher discrepancy scores to truthful responses and lower scores to hallucinated ones, aligning with our design intuition.

<b>Truth Answer 1</b>					
Q: The word "Easter" is connected with which goddess? A: Eostre. N: As the seasons unfold.					
<b>Hallucination Answer 1</b>					
Q: What is the largest animal that lives in the Loch Ness? A: The largest animal that lives in Loch Ness is the salmon. N: As the sun sets over the water.					
<b>Truth Answer 2</b>					
Q: Can reindeer really fly? A: No, reindeer are not capable of flight. N: The snowflakes gently fell.					
<b>Hallucination Answer 2</b>					
Q: What happens if an earthworm is cut in half? A: It will regenerate into two separate earthworms. N: As the moon glows with an otherworldly.					
<b>Different Perturbations</b>		Truth Answer 1	Hallucination Answer 1	Truth Answer 2	Hallucination Answer 2
	Predict				
Static sentence	→	0.7381	0.5723	0.4219	0.6731 ✗
Prompt learning	→	0.8213	0.4278	0.4093	0.6093 ✗
Sample-specific prompt learning	→	0.9429	0.0930	0.7131	0.2978

Figure 4: Qualitative comparison of discrepancy scores assigned by different prompting strategies.

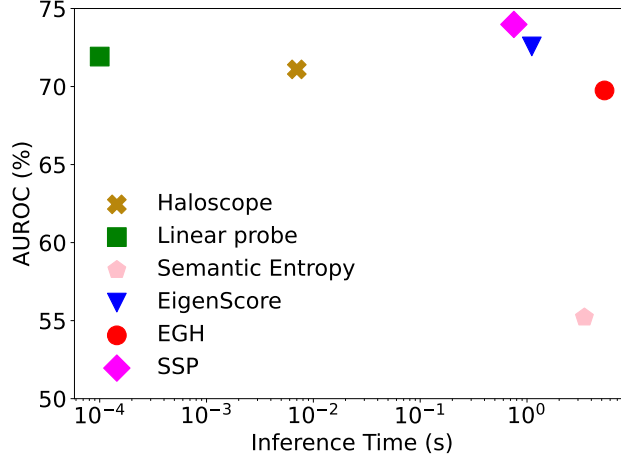


Figure 5: AUROC and inference time.

## I 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 5, 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 capability, demonstrating its practicality for real-world deployment scenarios.

## J 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 framework—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.

## K Limitations

We propose a hallucination detection method that induces internal representation shifts in LLMs by concatenating learnable, sample-specific noise 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 is unable to precisely

localize which tokens in the generated output are incorrect. In addition, the current scoring function is relatively simple and may lack sufficient discriminative power for more complex or fine-grained hallucination detection tasks. Future work could explore more powerful and generalizable scoring mechanisms to further improve the robustness and applicability of the method in real-world scenarios.