

NEXT-TOKEN GRADIENT SENSITIVITY PROBING FOR LLM HALLUCINATION DETECTION

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language models (LLMs) demonstrate remarkable text generation capabilities but often produce hallucinations, *e.g.*, factually inaccurate or unfaithful content, posing significant deployment risks. Detecting such hallucinations remains challenging due to their commonly subtle and localized nature. In this paper, we propose a *gradient*-based paradigm for hallucination detection by probing next-token prediction sensitivity. Specifically, we introduce a statistic called *Next-token Gradient Sensitivity* (NGS), which quantifies the first-order gradient of the maximum log-probability of the next token w.r.t. the current token’s layer embedding, measuring local prediction fragility. We prove that its norm bounds the prediction confidence changes under small perturbations. Building on NGS, we develop *NGS-based Hallucination Detection* (NGS-HD), a method that reframes detection as a *token-level distribution comparison* task. NGS-HD computes the *Maximum Mean Discrepancy* (MMD) between each NGS of test tokens and NGS distributions from referenced truthful and hallucinated tokens, aggregating these MMDs into a global *truthfulness score* for detection. We further derive finite-sample separation bounds for this score, providing theoretical guarantees for its reliability. Extensive experiments demonstrate that NGS-HD outperforms baseline methods, offering a reliable and interpretable solution for detecting LLM hallucinations.

1 INTRODUCTION

The rapid advancement of large language models (LLMs) has led to their widespread deployment in critical applications, from question-answering and content creation to decision-making (Brown et al., 2020; Touvron et al., 2023a;b). Despite their remarkable capabilities, they simultaneously pose critical societal risks through the generation of hallucinated content (*e.g.*, factual inaccuracies in information dissemination (Tonmoy et al., 2024; Huang et al., 2025), limiting the reliable use of LLMs (Zhang et al., 2025; Zhou et al., 2023)). Fundamentally, hallucination represents not merely a technical artifact but a cognitive problem—where models, like humans, make up plausible but unsupported responses (Ji et al., 2023). Consequently, developing effective hallucination detection methods has become a pressing challenge to ensure the safe and trustworthy deployment of LLMs.

A fundamental challenge in hallucination detection lies in capturing the intrinsic uncertainty and instability of the model’s generation process. Intuitively, when an LLM “hallucinates”, its outputs often deviate from the underlying truthful distribution, relying on shallow patterns or weak priors rather than robust, faithful reasoning (Manakul et al., 2023; Huang et al., 2025). This suggests that hallucinated tokens may exhibit higher sensitivity to minor perturbations in their contextual representations, leading to unstable next-token predictions. This observation raises a crucial question:

Can we probe the local sensitivity of next-token predictions to identify hallucinated content reliably?

Two significant difficulties arise: 1) Hallucinations are often localized and context-dependent, appearing amidst otherwise coherent text, which requires detection methods to operate at a fine-grained token level without losing global coherence signals (Manakul et al., 2023). 2) The distortions indicative of hallucinations are often subtle, making them difficult to distinguish from valid variations using surface-level features or prediction probabilities alone (Farquhar et al., 2024).

Existing LLM hallucination detection methods primarily rely on uncertainty estimation (Malinin & Gales, 2021; Manakul et al., 2023; Farquhar et al., 2024) and supervised learning with intermediate

representations (Burns et al., 2022; Du et al., 2024; Park et al., 2025). However, these approaches are often computationally expensive, fail to capture the local instability inherent in hallucinated generations, or aggregate sequence-level information at the cost of token-level granularity, making them potentially less effective for detecting localized hallucinations within longer texts.

In this paper, we propose a *gradient*-based perspective that directly probes the *next-token prediction sensitivity* as a means of detecting hallucinations. As illustrated in Figures 1 and 2, we empirically show that hallucinated tokens exhibit significantly higher sensitivity to small perturbations in their embeddings compared to truthful tokens (see details in Section 3.1).

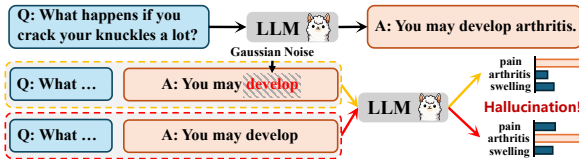


Figure 1: Hallucination detection via next-token sensitivity, where noise shifts predictions, indicating instability.

This motivates the use of *gradients*—which formally capture the first-order sensitivity of the model’s outputs to its inputs—as a natural and discriminative signal for hallucination detection. Inspired by this, we introduce a statistic called **Next-token Gradient Sensitivity (NGS)**, defined as the gradient of the maximum log-probability of the next token w.r.t. the current token’s layer embedding. NGS rigorously quantifies the local instability of the model’s predictions and can be computed efficiently in a single backward pass. We further provide a theoretical bound linking the norm of NGS to the changes in model confidence under perturbation, establishing NGS as a principled measure of prediction fragility (Section 3.2).

Building on NGS, we propose an **NGS-based hallucination detection method (NGS-HD)** in Section 3.3, which detects hallucinations by comparing the distribution of each NGS vector from a test response against pre-collected reference distributions of truthful and hallucinated tokens, as illustrated in Figure 3. To handle variable-length sequences while preserving token-level information, we reframe the detection problem as a *token-level distribution comparison* task using *Maximum Mean Discrepancy* (MMD (Gretton et al., 2012)). Specifically, we compute the MMD between each test token’s NGS vector and each of the two reference distributions, and then aggregate per-token discrepancies into a global *truthfulness score* for detection. This approach effectively detects localized hallucinations without losing fine-grained signals through sequence-level aggregation. We further theoretically derive finite-sample separation bounds for the truthfulness score to guarantee its reliability (Section 3.4). Extensive experiments demonstrate that NGS-HD outperforms existing state-of-the-art methods, validating the effectiveness of gradient signals for hallucination detection.

Our contributions are summarized as follows:

- A gradient-based sensitivity probing for hallucination detection: We propose a statistic *Next-token Gradient Sensitivity* (NGS) that efficiently captures the local instability of LLM generations by measuring the first-order sensitivity of the next-token prediction w.r.t. token layer embeddings. Through theoretical analysis, we show that NGS upper-bounds the changes in prediction confidence under perturbation, establishing it as a principled indicator of hallucination.
- A reliable and explainable hallucination detection method: We introduce NGS-HD that reframes hallucination detection as a token-level distribution comparison task. By computing *Maximum Mean Discrepancy* (MMD) between each NGS of test tokens and reference distributions, and aggregating per-token MMDs into a global *truthfulness score*, NGS-HD detects hallucinations while providing interpretability by offering insights into model behaviors beyond binary classification.
- Theoretical guarantees and empirical validation: We derive finite-sample separation bounds of NGS-HD for the proposed truthfulness score, ensuring its reliability under reasonable assumptions. Extensive results on various benchmarks show that NGS-HD outperforms existing state-of-the-art methods, validating the effectiveness of gradient signals for hallucination detection.

2 RELATED WORK

Hallucination Detection. Although “truthful” is a natural requirement of language generation, large language models (LLMs) may still produce factually incorrect or contextually inconsistent outputs, termed hallucinations. Detecting hallucinations (Ren et al., 2023; Kuhn et al., 2023; Lin et al., 2024; Manakul et al., 2023; Chen et al., 2024a; Lin et al., 2022a; Du et al., 2024; Park et al., 2025; Zhang et al., 2023b) has therefore become a critical research focus for safe and reliable deployments.

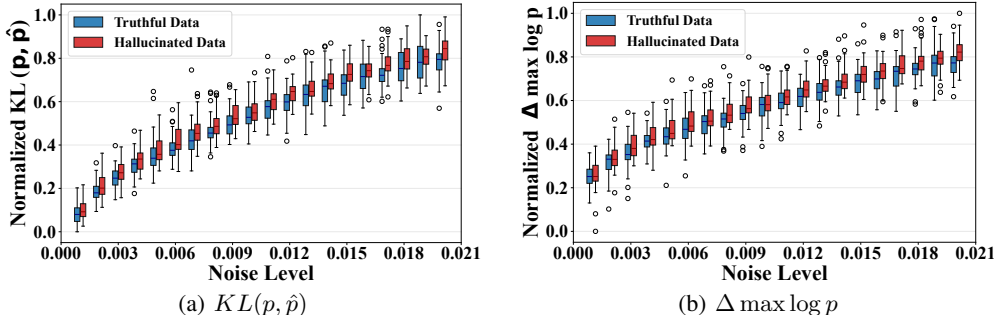


Figure 2: Comparisons between hallucinated data (red) and truthful data (blue) across noise levels in (a) normalized $KL(p, \hat{p})$ and (b) normalized $\Delta \max \log p$ on TruthfulQA. Both capture output discrepancies of LLaMA-3.1-8b before and after perturbing the previous answer token’s embedding, showing that hallucinated tokens exhibit stronger local output-sensitivity than truthful ones.

One perspective attributes hallucinations to the uncertainty in LLM predictions. Logits-based methods employ perplexity (Ren et al., 2023), LN-entropy (Malinin & Gales, 2021), and Semantic Entropy (Kuhn et al., 2023) to measure language-invariant uncertainty. Consistency-based methods assess similarity across multiple samples using metrics such as ROUGE (Lin et al., 2024), BERTScore (Zhang et al., 2019), natural language inference, and prompt-based comparisons (Manakul et al., 2023), or via eigenvalues of the response covariance matrix (Chen et al., 2024a). Verbalized-based methods (Lin et al., 2022a) and self-evaluation (Kadavath et al., 2022) further allow models to express uncertainty in natural language through verbalized confidence predictions. Nonetheless, hallucinations still arise in high-confidence generations, limiting the reliability of uncertainty-based approaches. Multi-sampling may alleviate this issue, but incurs high computational cost.

A complementary perspective leverages LLM internal states to infer truthfulness. CCS (Burns et al., 2022) mines latent knowledge from activations, SAPLMA (Azaria & Mitchell, 2023) trains classifiers on hidden states, HaloScope (Du et al., 2024) identifies hallucination subspaces via singular value decomposition, and TSV (Park et al., 2025) introduces learnable steering vectors to reshape latent features for improved separability. While effective, these methods often overlook the local instability present in hallucinated outputs or summarize sequence-level information at the expense of token-level details, limiting their ability to detect localized hallucinations within longer texts.

3 PROBING GRADIENTS SIGNALS FOR LLM HALLUCINATION DETECTION

LLM Hallucinations. LLM hallucinations refer to the generated content is either factually incorrect or contextually unfaithful to the provided source or internal knowledge (Huang et al., 2025). They often arise when the model generates content it is inherently uncertain about, due to knowledge gaps or insufficient reasoning, making the probing of such uncertainty a central focus of detection.

LLM Hallucination Detection. Hallucination detection aims to identify whether a generated text contains non-factual or unfaithful information (Ji et al., 2023; Farquhar et al., 2024). Formally, we define the truthful distribution \mathbb{P}_{true} as the joint distribution over prompt-truthful response pairs. Given a prompt $\mathbf{x}_{\text{prompt}} = (\mathbf{x}_1, \dots, \mathbf{x}_m)$ and a model-generated continuation $\tilde{\mathbf{x}} = (\mathbf{x}_{m+1}, \dots, \mathbf{x}_{m+n})$, we aim to determine whether the full sequence $\mathbf{x}_{\text{prompt}} \oplus \tilde{\mathbf{x}}$ is from \mathbb{P}_{true} .

Challenges for Hallucination Detection. Detecting hallucinations is challenging as they often appear as locally plausible but globally inconsistent outputs, deeply embedded in the model’s representations (Huang et al., 2025). This requires capturing both fine-grained token-level uncertainty and higher-level semantic coherence. Moreover, the local nature of many hallucinations demands detection that preserves token-level signals without being misled by generally fluent context. Furthermore, the sensitivity of LLMs to input variations can amplify output instability, challenging detection methods to robustly quantify such uncertainty without excessive computational overhead.

3.1 MOTIVATIONS AND METHOD OVERVIEW

Motivations. Hallucinations in LLMs often occur when generated outputs deviate from the truthful distribution \mathbb{P}_{true} , producing plausible but unfaithful responses. We empirically observe that

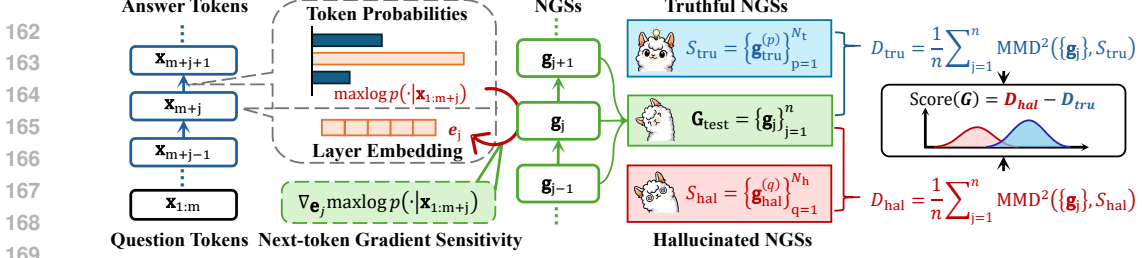


Figure 3: Overview of the proposed NGS-HD. Given reference sets of Next-token Gradient Sensitivity (NGS) vectors S_{tru} and S_{hal} from truthful tokens and hallucinated tokens, respectively, for a test prompt-answer sequence $(\mathbf{x}_{1:m}, \mathbf{x}_{m+1}, \dots, \mathbf{x}_{m+n})$, we first compute NGS vectors for all answer tokens $(\mathbf{x}_{m+1}, \dots, \mathbf{x}_{m+n})$ through a single backward pass of the LLM, then calculate the MMD between the test NGS distribution and each reference distribution, and finally aggregate the per-token MMD values into a global truthfulness score for detection.

hallucinated tokens tend to exhibit stronger local output-sensitivity than truthful ones, i.e., small perturbations to a token’s embedding induce significant changes in next-token predictions. Intuitively, hallucinated tokens are supported mainly by weak priors or shallow contextual cues, so even minor representation changes can greatly alter model behaviors.

To quantify this, we perturb the embedding \mathbf{e} of token \mathbf{x}_{m+t} with noise $\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$ and measure changes in the next-token distribution $p(\mathbf{x}_{m+t+1} | \mathbf{x}_{1:m+t})$. We evaluate both the KL divergence $KL(p, \hat{p})$, capturing overall distribution shift, and the change in max log probability $\Delta \max \log p = \log \hat{p}_l - \log p_l$ (where $l = \arg \max_l p_l$), reflecting top-prediction stability. In Figures 2 (a)–(b), both metrics are markedly larger for hallucinated tokens—particularly at multiple noises (AUROCs > 0.70 , see details in Appendix D.4)—confirming their higher sensitivity. This observed instability, though *variable under random perturbation*, motivates the use of gradients—which precisely capture first-order sensitivity—as a principled and efficient signal for hallucination detection.

Method Overview. Motivated by the observation that hallucinated tokens exhibit unstable predictive behaviors, we propose to quantify this instability via a statistic termed **Next-token Gradient Sensitivity (NGS)**, which captures the first-order sensitivity of the model’s next-token prediction to its layer embeddings (Section 3.2). Based on this statistic, we propose a NGS-based hallucination detection method, **NGS-HD**, as shown in Figure 3. Particularly, for a prompt-answer test sample, NGS-HD computes NGSs for all answer tokens in one backward pass, and calculates *Maximum Mean Discrepancy* (MMD) between their empirical distribution and pre-constructed reference distributions of truthful and hallucinated tokens, then aggregates per-token MMDs into a global *truthfulness score* for detection (Section 3.3). We further provide theoretical analyses, establishing finite-sample separation bounds for the truthfulness score to guarantee its reliability (Section 3.4).

3.2 PROBING NEXT-TOKEN GRADIENT SENSITIVITY

Perturbation-based experiments in Section 3.1 provide intuitive evidence of instability, which, however, are often computationally expensive since they require multiple forward passes with different noise samples to obtain stable estimates. To overcome this, we directly probe the model’s local behaviors through gradients, which offer an exact first-order measure of output sensitivity to embedding perturbations without requiring explicit perturbation samplings. To this end, we propose a statistic, **next-token gradient sensitivity (NGS)** that captures such sensitivity in a rigorous manner.

Next-token Gradient Sensitivity. Formally, consider a prompt-answer sequence $\mathbf{x}_{1:m+n} = (\mathbf{x}_1, \dots, \mathbf{x}_m, \mathbf{x}_{m+1}, \dots, \mathbf{x}_{m+n})$, with token embeddings $\mathbf{e}_i \in \mathbb{R}^d$ from a specific layer of the LLM, where $i = m+1, \dots, m+n$. Let $p(\cdot | \mathbf{x}_{1:i})$ denote the model’s conditional probability distribution for the next-token prediction at position i . We then introduce the following definition:

Definition 1. (Next-token Gradient Sensitivity (NGS)) For a token \mathbf{x}_{m+j} with its layer embedding $\mathbf{e}_j \in \mathbb{R}^d$ ($j = 1, \dots, n$), the Next-token Gradient Sensitivity (NGS) of \mathbf{x}_{m+j} is defined as the gradient of the maximum log-probability for the next token w.r.t. the layer embedding \mathbf{e}_j :

$$\mathbf{g}_j = \nabla_{\mathbf{e}_j} \max \log p(\cdot | \mathbf{x}_{1:m+j}) \in \mathbb{R}^d, \quad j = 1, \dots, n. \quad (1)$$

Remark. NGS specifically computes the gradient w.r.t. the current token’s layer embedding \mathbf{e}_j , not all previous tokens’. This design offers significant computational advantages, reducing complex-

ity from quadratic to linear in sequence length, while maintaining theoretical focus by specifically measuring the local sensitivity of a token’s own representation. Importantly, since $p(\cdot|\mathbf{x}_{1:m+j})$ is a function of all prior context, its gradient implicitly captures aggregated sensitivity signals from the entire history, making NGS an efficient and informative statistic of generation instability.

The term $\max \log p(\cdot|\mathbf{x}_{1:m+j})$ represents the model’s confidence in its top prediction at step j . The gradient \mathbf{g}_j thus quantifies the first-order sensitivity of this confidence to infinitesimal perturbations of the token \mathbf{x}_{m+j} . To formally show that \mathbf{g}_j serves as a reliable measure of local prediction instability, we present the following theorem which bounds the changes in confidence under perturbations.

Theorem 1. (First-order local sensitivity bound) *Let $s : \mathbb{R}^d \rightarrow \mathbb{R}$ denote the scalar next-token score as a function of the token layer embedding \mathbf{e} , i.e., $s(\mathbf{e}) = \max \log p(\cdot|\mathbf{x}_{1:m+j})$, then $\mathbf{g} = \nabla_{\mathbf{e}} s(\mathbf{e})$. Assume s is differentiable and its gradient is local L -Lipschitz, i.e., $\|\nabla_{\mathbf{e}} s(\mathbf{e}) - \nabla_{\mathbf{e}'} s(\mathbf{e}')\| \leq L\|\mathbf{e} - \mathbf{e}'\|$ for $\forall \mathbf{e}, \mathbf{e}' \in \mathbb{R}^d$. Then for any perturbation $\boldsymbol{\epsilon} \in \mathbb{R}^d$,*

$$|s(\mathbf{e} + \boldsymbol{\epsilon}) - s(\mathbf{e})| \leq \|\mathbf{g}\| \|\boldsymbol{\epsilon}\| + \frac{L}{2} \|\boldsymbol{\epsilon}\|^2.$$

Theorem 1 provides the theoretical foundation for using NGS as a probe for model instability. It shows that the norm of \mathbf{g}_j upper bounds the confidence changes of the model under small perturbations of the embedding \mathbf{e}_j . Consequently, a large $\|\mathbf{g}\|$ indicates that the model’s top prediction is highly fragile—even small perturbations to the token’s representation can lead to significant shifts in confidence. This formulation captures the precise notion of unstable reasoning that characterizes hallucinations, justifying NGS as a principled and interpretable metric for detection.

Constructing Sentence Gradient Sensitivity. Given a prompt $\mathbf{x}_{prompt} = (\mathbf{x}_1, \dots, \mathbf{x}_m)$ and a generated answer sequence $\tilde{\mathbf{x}} = (\mathbf{x}_{m+1}, \dots, \mathbf{x}_{m+n})$, we compute the NGS for every answer token. This results in a sequence of gradient vectors that encapsulates the evolution of the model’s sensitivity throughout the generation process. Stacking these vectors leads to:

$$\mathbf{G} = \{\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_n\}. \quad (2)$$

In practice, the NGSs \mathbf{G} can be computed in *a single backward pass*. After the autoregressive forward pass, we compute the next-token confidence scores s_{m+1}, \dots, s_{m+n} (i.e., the max log probabilities) for the generated positions and stack them into a vector $\mathbf{s} \in \mathbb{R}^n$. We then perform a single backward/autograd call on \mathbf{s} (i.e., jacobian function in `torch.autograd.functional`) to obtain the sensitivity of these scalar scores w.r.t. the token embeddings. This operation efficiently yields all required gradients $\mathbf{g}_1, \dots, \mathbf{g}_n$ simultaneously. Since each entry of \mathbf{s} is a scalar value, this computation is significantly more efficient than computing full Jacobians of vector-valued outputs and requires only standard gradient operations. The resulting gradients $\mathbf{G} = \{\mathbf{g}_1, \dots, \mathbf{g}_n\}$ thus serve as a computationally efficient, token-level feature representation for our truthfulness detector.

3.3 PROBING NGS FOR DETECTING LLM HALLUCINATIONS

While the norm of the NGS vector $\|\mathbf{g}_j\|$ provides a simple measure of local instability, it discards the rich directional information in the gradients, which may encode the specific nature of the model’s uncertainty. Therefore, a core challenge arises: how to develop a detector that can effectively utilize the full vector information of the NGS sequence $\mathbf{G} = \{\mathbf{g}_1, \dots, \mathbf{g}_n\}$ for answers of variable lengths n ? Common approaches (Du et al., 2024; Park et al., 2025) to handling variable-length sequences often involve aggregation techniques like averaging or pooling (Kim, 2014), or the use of recurrent networks (Sutskever et al., 2014) or transformers (Vaswani et al., 2017) to compress the sequence into a fixed-dimensional representation. However, these methods can inevitably lose fine-grained, token-level information that is crucial for identifying localized hallucinations within longer sequences.

To address the above issue, we reframe the binary classification problem as a *per-token distribution comparison* task by assessing whether each token NGSs of a test sample’s answer more closely resemble the distribution of truthful or hallucinated tokens. This transforms the challenge from classifying a single sequence representation to comparing sets of token-level features against referenced distributions. The idea is motivated by that hallucinations often appear as localized inconsistencies within otherwise coherent contexts (Bender et al., 2021; Manakul et al., 2023). By operating at the token distribution level, this approach preserves sensitivity to such fine-grained deviations without the information loss from aggregation. To achieve this, we propose a NGS-based hallucination detection method, called **NGS-HD**. The core of NGS-HD involves comparing the empirical distribution of NGSs from a test sample’s answer against two reference distributions from truthful and

hallucinated tokens. Intuitively, a truthful answer should statistically align closer to truthful tokens, while a hallucinated one aligns more with hallucinated tokens.

Measuring NGS Proximity with MMD. Formally, let $S_{\text{tru}} = \{\mathbf{g}_{\text{tru}}^{(p)}\}_{p=1}^{N_t} \sim P_{\text{tru}}$ and $S_{\text{hal}} = \{\mathbf{g}_{\text{hal}}^{(q)}\}_{q=1}^{N_h} \sim P_{\text{hal}}$ denote two reference sets of NGSs from truthful and hallucinated tokens, respectively. For a test sample with answers’ NGSs $\mathbf{G}_{\text{test}} = \{\mathbf{g}_j\}_{j=1}^n$, we calculate the MMD via Eqn. (4) for each NGS in \mathbf{G}_{test} and aggregate these per-token distances across the answer by averaging for the hallucination set and truthful set, respectively, yielding:

$$D_{\text{hal}} = \frac{1}{n} \sum_{j=1}^n \text{MMD}^2(\{\mathbf{g}_j\}, S_{\text{hal}}), \quad D_{\text{tru}} = \frac{1}{n} \sum_{j=1}^n \text{MMD}^2(\{\mathbf{g}_j\}, S_{\text{tru}}). \quad (3)$$

The *Maximum Mean Discrepancy* (MMD) (Gretton et al., 2012) measures the distance between distributions, assessing whether two sets of examples are drawn from the same distribution, providing a robust distance between sets of vectors (Zhang et al., 2024). For a single NGS \mathbf{g} and a reference set $S = \{\mathbf{g}_\ell\}_{\ell=1}^N$ (i.e., S_{tru} or S_{hal}), the empirical MMD² between $\{\mathbf{g}\}$ and S is defined as:

$$\widehat{\text{MMD}}_b^2[\{\mathbf{g}\}, S; \mathcal{H}_k] = \frac{1}{N^2} \sum_{l, l'=1}^N k(\mathbf{g}_l, \mathbf{g}_{l'}) - \frac{2}{N} \sum_{l=1}^N k(\mathbf{g}_l, \mathbf{g}) + k(\mathbf{g}, \mathbf{g}), \quad (4)$$

where the kernel $k : \mathcal{G} \times \mathcal{G} \rightarrow \mathbb{R}$ maps NGS features to a reproducing kernel Hilbert space (RKHS) \mathcal{H}_k , such as the Gaussian kernel $k(\mathbf{a}, \mathbf{b}) = \exp(-\|\mathbf{a} - \mathbf{b}\|^2 / (2\sigma^2))$. In practice, we collect two sets of NGSs of truthful and hallucinated answer tokens from all training data, respectively, and train a deep kernel MMD following Liu et al. (2020). We refer readers to Appendix C.5 for more details.

Truthfulness Score. The final truthfulness score for the test answer is defined as the difference between its distance to the hallucinatory distribution and to the truthful distribution as:

$$\widehat{\text{Score}}(\mathbf{G}_{\text{test}}) = D_{\text{hal}} - D_{\text{tru}}. \quad (5)$$

A high score (e.g., $\widehat{\text{Score}}(\mathbf{G}_{\text{test}}) > 0$) indicates that the test answer’s NGS distribution is, on average, closer to the truthful distribution than to the hallucinatory one, suggesting a truthful generation. Conversely, a low score (e.g., $\widehat{\text{Score}}(\mathbf{G}_{\text{test}}) < 0$) indicates the answer is more similar to the hallucinatory distribution and is thus classified as a potential hallucination.

Remarks. 1) While MMD is classically a tool for comparing two distributions, recent studies (Zhang et al., 2023a; 2024) validate its efficacy in single-sample detection by quantifying discrepancies from references. 2) Crucially, the benefit of our pipeline does not come from MMD alone but from operating on token-level NGSs: NGS encodes inherent local instability of model predictions, and comparing NGSs rather than pooled answers preserves the fine-grained signals to expose local hallucinations. 3) We do not perform a direct two-sample MMD between the test-token set (often very small) and each reference distribution—which can be statistically weak when the test set contains few tokens (Liu et al., 2021)—but instead compute per-token MMD scores (treating each test token as a singleton) and aggregate them. This strategy combines the statistical strength of prebuilt reference sets with the per-token granularity required for reliable detection.

Interpretation and Advantages of NGS-HD. The NGS-HD leverages the next-token gradient sensitivity, i.e., NGS, to capture the *inherent local instability* of hallucinated tokens. Unlike methods that compress sequence information through some aggregation techniques, our approach preserves fine-grained, token-level information by assessing the distributional similarity of NGS vectors, enabling effective detection of localized hallucinations within longer responses. Moreover, our NGS-HD is model-agnostic, requiring only gradient signals from any differentiable language model without architectural modifications. Furthermore, the per-token MMD scores provide intrinsic interpretability by allowing practitioners to localize potential hallucinations within the generated text, offering valuable insights into model behaviors beyond simple binary classification.

3.4 THEORETICAL GUARANTEES FOR NGS-HD

To provide statistical guarantees for the reliability of NGS-HD, we present a concentration bound for the empirical truthfulness score. Let $k(\cdot, \cdot)$ be a positive-definite kernel with feature map $\varphi :$

$\mathcal{G} \rightarrow \mathcal{H}_k$ such that $k(\mathbf{a}, \mathbf{b}) = \langle \varphi(\mathbf{a}), \varphi(\mathbf{b}) \rangle_{\mathcal{H}}$. For a finite NGS set $S = \{\mathbf{g}_\ell\}_{\ell=1}^N$, denote $\hat{\mu} = \frac{1}{N} \sum_{\ell=1}^N \varphi(\mathbf{g}_\ell)$. Then the empirical score can be expressed as $\widehat{\text{Score}}(\mathbf{G}_{\text{test}}) = \frac{1}{n} \sum_{j=1}^n \left(\|\varphi(\mathbf{g}_j) - \hat{\mu}_{\text{hal}}\|_{\mathcal{H}}^2 - \|\varphi(\mathbf{g}_j) - \hat{\mu}_{\text{tru}}\|_{\mathcal{H}}^2 \right)$. We now establish the following finite-sample separation bounds:

Theorem 2. (Separation bounds) Assume the kernel satisfies $0 \leq k(\mathbf{g}, \mathbf{g}') \leq K$ for all \mathbf{g}, \mathbf{g}' . Let $S_{\text{tru}} = \{\mathbf{g}_q\}_{q=1}^{N_t}$ and $S_{\text{hal}} = \{\mathbf{g}_p\}_{p=1}^{N_h}$ be i.i.d. drawn from P_{tru} and P_{hal} , respectively. Let a test answer consist of n tokens $\mathbf{G}_{\text{test}} = \{\mathbf{g}_j\}_{j=1}^n$, drawn i.i.d. from either P_{tru} or P_{hal} . Define the population mean-embedding separation as $\Delta = \|\mu_{\text{tru}} - \mu_{\text{hal}}\|_{\mathcal{H}}$ with $\mu_{\text{tru}} = \mathbb{E}_{\mathbf{g} \sim P_{\text{tru}}}[\varphi(\mathbf{g})]$ and $\mu_{\text{hal}} = \mathbb{E}_{\mathbf{g} \sim P_{\text{hal}}}[\varphi(\mathbf{g})]$. Then, for $\forall \delta \in (0, 1)$, with probability at least $1 - \delta$, the following holds:

- 1) If $\mathbf{G}_{\text{test}} \sim P_{\text{tru}}$, then $\widehat{\text{Score}}(\mathbf{G}_{\text{test}}) \geq \Delta^2 - \mathcal{O}\left(K\sqrt{\frac{\log(1/\delta)}{N_t}} + K\sqrt{\frac{\log(1/\delta)}{N_h}} + K\sqrt{\frac{\log(1/\delta)}{n}}\right)$;
- 2) If $\mathbf{G}_{\text{test}} \sim P_{\text{hal}}$, then $\widehat{\text{Score}}(\mathbf{G}_{\text{test}}) \leq -\Delta^2 + \mathcal{O}\left(K\sqrt{\frac{\log(1/\delta)}{N_t}} + K\sqrt{\frac{\log(1/\delta)}{N_h}} + K\sqrt{\frac{\log(1/\delta)}{n}}\right)$.

Consequently, if $\Delta^2 > \mathcal{O}\left(K\sqrt{\frac{\log(1/\delta)}{N_t}} + K\sqrt{\frac{\log(1/\delta)}{N_h}} + K\sqrt{\frac{\log(1/\delta)}{n}}\right)$, then $\text{sign}(\widehat{\text{Score}})$ correctly classifies the test answer (truthful or hallucinated) with probability at least $1 - \delta$.

Theorem 2 reveals that our detector will correctly classify an answer with high probability provided that the inherent distributional separation Δ is relatively large to overcome the estimation error introduced by finite-sized reference sets and test samples. Crucially, our empirical results in Section 4.2 show that NGS-HD achieves satisfactory performance even with moderate reference set sizes (e.g., $N_t = N_h = 200$). This indicates that the natural distributional separation Δ in practice is substantial enough that the estimation error terms become negligible compared to Δ^2 . This aligns with our empirical findings: the core effectiveness of NGS-HD stems from the inherent separability of NGS distributions, which is reliably detectable even with limited but representative data.

4 EXPERIMENTS

Datasets. We evaluate on four QA benchmarks: TruthfulQA (Lin et al., 2022b), SciQ (Welbl et al., 2017), TriviaQA (Joshi et al., 2017), and NQ Open (Kwiatkowski et al., 2019). For TruthfulQA, we use all 817 samples. From SciQ, we sample 3,000 training pairs for stable evaluation. TriviaQA uses its full test set. Each dataset is split 3 : 1 into training and test sets, with results averaged over 10 random splits. All responses are generated via greedy decoding. Additional benchmarks, e.g., Wikipedia (Foundation, 2022) and NarrativeQA (Kočíský et al., 2018), are in Appendix D.11.

Models. We compare our method on four LLMs, i.e., LLaMA-3.1-8b (Dubey et al., 2024), Qwen-2.5-7b (Yang et al., 2024a), Qwen-3-8b (Yang et al., 2025), Qwen-3-14b (Yang et al., 2025).

Baselines. 1) Logit-based methods: Perplexity (Ren et al., 2023), LengthNormalized Entropy (LN-entropy) (Malinin & Gales, 2021) and Semantic Entropy (Kuhn et al., 2023); 2) Consistency-based methods: Lexical Similarity (Lin et al., 2024), SelfCKGPT (Manakul et al., 2023) and EigenScore (Chen et al., 2024a); 3) Verbalized methods: Self-evaluation (Kadavath et al., 2022); 4) Internal state-based methods: Contrast-Consistent Search (CCS) (Burns et al., 2022), SAPLMA (Azaria & Mitchell, 2023), HaloScope (Du et al., 2024) and TSV (Park et al., 2025).

Evaluation. Following Farquhar et al. (2024), we evaluate with Area Under the Receiver Operating Characteristic curve (AUROC). We adopt DeepSeek-V3 (Liu et al., 2024) as an automatic judge to determine if an answer is correct by checking consistency with the gold reference. Additionally, our method is also robust using the BLEURT score (Sellam et al., 2020) (see Appendix D.1).

4.1 COMPARISONS WITH HALLUCINATION DETECTION BASELINES

We compare our NGS-HD with hallucination detection baselines on LLaMA-3.1-8b and Qwen-3-8b in Table 1, respectively. We defer more results of larger LLMs like Qwen-3-14b in Appendix D.3.

Results on LLaMA-3.1-8b. As shown in Table 1, existing methods exhibit clear limitations: uncertainty-based approaches such as LN-Entropy (54.01% Avg.) and Lexical Similarity (55.07% Avg.) struggle on most datasets, indicating their inability to capture nuanced hallucination patterns.

Table 1: Comparisons with hallucination detection baselines on different datasets for LLaMA-3.1-8b and Qwen-3-8b in terms of AUROC (%), where † denotes methods trained on fully labeled datasets.

Model	Method	TruthfulQA	TriviaQA	SciQ	NQ Open	Avg.
Llama-3.1-8b	LN-Entropy (Malinin & Gales, 2021)	59.06 \pm 1.4	51.70 \pm 1.2	51.96 \pm 0.9	53.31 \pm 2.3	54.01 \pm 0.77
	Semantic Entropy (Kuhn et al., 2023)	54.25 \pm 1.5	52.06 \pm 0.1	53.75 \pm 0.9	56.06 \pm 2.1	54.03 \pm 0.68
	Lexical Similarity (Lin et al., 2024)	57.82 \pm 2.7	50.97 \pm 1.1	58.96 \pm 1.0	52.51 \pm 1.5	55.07 \pm 0.86
	EigenScore (Chen et al., 2024a)	53.83 \pm 0.6	56.78 \pm 1.0	49.46 \pm 1.2	61.75 \pm 2.3	55.46 \pm 0.71
	SelfCKGPT (Manakul et al., 2023)	54.95 \pm 0.4	73.20 \pm 1.0	51.82 \pm 0.6	64.92 \pm 0.7	61.22 \pm 0.35
	Perplexity (Ren et al., 2023)	58.99 \pm 1.9	57.27 \pm 1.4	58.10 \pm 2.0	51.87 \pm 1.3	56.56 \pm 0.84
	Self-evaluation (Kadavath et al., 2022)	54.96 \pm 1.5	50.26 \pm 3.3	50.07 \pm 2.3	52.03 \pm 2.9	51.83 \pm 1.30
	CCS (Burns et al., 2022)	58.69 \pm 0.4	51.23 \pm 0.4	68.60 \pm 2.7	57.96 \pm 2.5	59.12 \pm 0.93
	SAPLMA† (Azaria & Mitchell, 2023)	68.09 \pm 2.8	69.96 \pm 0.7	62.89 \pm 2.7	63.77 \pm 1.1	66.18 \pm 1.03
	Haloscope† (Du et al., 2024)	65.95 \pm 4.3	58.82 \pm 4.1	64.29 \pm 0.5	53.06 \pm 0.4	60.53 \pm 1.50
	TSV (Park et al., 2025)	63.95 \pm 3.8	62.70 \pm 2.3	63.33 \pm 3.1	58.09 \pm 2.4	62.02 \pm 1.48
	TSV† (Park et al., 2025)	80.50 \pm 4.7	84.31 \pm 0.5	72.85 \pm 1.1	71.56 \pm 1.2	77.30 \pm 1.25
	NGS-HD (Ours)		82.18\pm3.0	85.65\pm0.7	80.67\pm2.1	74.84\pm0.9
Qwen-3-8b	LN-Entropy (Malinin & Gales, 2021)	55.14 \pm 1.0	59.05 \pm 2.4	66.26 \pm 1.6	59.76 \pm 1.9	60.05 \pm 0.90
	Semantic Entropy (Kuhn et al., 2023)	54.64 \pm 1.3	53.85 \pm 0.8	52.63 \pm 0.2	54.74 \pm 1.9	53.97 \pm 0.61
	Lexical Similarity (Lin et al., 2024)	57.44 \pm 0.8	53.05 \pm 2.1	52.87 \pm 0.3	54.48 \pm 0.7	54.46 \pm 0.59
	EigenScore (Chen et al., 2024a)	64.79 \pm 1.4	69.01 \pm 3.1	49.70 \pm 1.1	69.67 \pm 1.4	61.17 \pm 1.19
	SelfCKGPT (Manakul et al., 2023)	62.69 \pm 0.3	71.28 \pm 3.4	48.7 \pm 3.2	74.73 \pm 3.8	64.35 \pm 1.43
	Perplexity (Ren et al., 2023)	63.75 \pm 1.0	66.32 \pm 2.0	66.17 \pm 1.9	54.76 \pm 2.0	62.75 \pm 0.89
	Self-evaluation (Kadavath et al., 2022)	51.25 \pm 2.1	53.90 \pm 2.3	50.93 \pm 1.4	53.85 \pm 1.0	52.48 \pm 0.89
	CCS (Burns et al., 2022)	54.53 \pm 0.1	56.63 \pm 0.1	56.86 \pm 0.6	53.83 \pm 0.3	55.46 \pm 0.17
	SAPLMA† (Azaria & Mitchell, 2023)	57.08 \pm 5.1	77.58 \pm 1.3	69.15 \pm 2.2	77.04 \pm 1.2	70.21 \pm 1.46
	Haloscope† (Du et al., 2024)	59.92 \pm 0.8	56.85 \pm 2.1	64.02 \pm 4.0	57.26 \pm 1.0	59.51 \pm 1.16
	TSV (Park et al., 2025)	54.84 \pm 2.8	57.35 \pm 1.5	61.86 \pm 2.7	54.84 \pm 3.0	57.22 \pm 1.28
	TSV† (Park et al., 2025)	65.77 \pm 2.8	80.76 \pm 0.7	76.18 \pm 2.4	73.39 \pm 2.6	74.03 \pm 1.14
	NGS-HD (Ours)		77.76\pm1.9	82.38\pm0.6	78.77\pm2.7	82.30\pm1.1

Table 2: Impact of per-token distribution comparison vs. token averaging for LLaMA-3.1-8b.

Method	TruthfulQA	SciQ	NQ Open
Lexical Similarity	57.82 \pm 2.7	58.96 \pm 1.0	47.49 \pm 1.5
Perplexity	58.99 \pm 1.9	58.10 \pm 2.0	51.87 \pm 1.3
Haloscope†	65.95 \pm 4.3	64.29 \pm 0.5	53.06 \pm 0.4
TSV†	80.50 \pm 4.7	72.85 \pm 1.1	71.56 \pm 1.2
NGS-Norm	73.50 \pm 0.3	56.63 \pm 0.1	60.41 \pm 0.1
NGS-AVG+CE	80.05 \pm 2.9	76.27 \pm 1.5	69.67 \pm 0.8
NGS-CLS+CE	79.92 \pm 3.3	77.26 \pm 1.7	69.59 \pm 1.2
NGS-AVG+MMD	80.09 \pm 2.2	74.48 \pm 2.3	68.23 \pm 1.5
NGS-CLS+MMD	80.80 \pm 2.5	75.74 \pm 1.6	69.24 \pm 0.6
NGS-HD	82.18\pm3.0	80.67\pm2.1	74.84\pm0.9

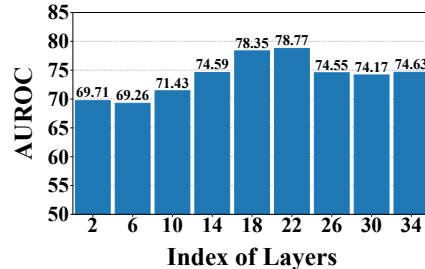


Figure 4: AUROC values for different layers on SciQ for Qwen-3-8b.

Training-dependent detectors like TSV† perform better (77.30% Avg.) yet still fall short on benchmarks such as SciQ (72.85% Avg.), revealing their limited generalizability. In contrast, our NGS-HD achieves best performance with an average AUROC of 80.84%, outperforming TSV† by 3.54% \uparrow and showing strong gains on SciQ by 7.82% \uparrow (80.67% vs. 72.85%) and NQ Open by 3.28% \uparrow (74.84% vs. 71.56%). These consistent improvements stem from our token-level gradient design, which directly probes localized instability in prediction, enabling finer hallucination detection.

Results on Qwen-3-8b. Similar trends are observed on Qwen-3-8b: common uncertainty metrics like Perplexity (62.75% Avg.) and Lexical Similarity (54.46% Avg.) perform poorly, while trained baselines like TSV† (74.03% Avg.) still vary substantially across domains (*e.g.*, only 65.77% on TruthfulQA). In contrast, NGS-HD again delivers superior and stable results, averaging 80.30%, *i.e.*, a 6.27% \uparrow improvement over TSV†, with notable performance on TruthfulQA by 11.99% \uparrow (77.76% vs. 65.77%) and NQ Open by 8.91% \uparrow (82.30% vs. 73.39%). The robust generalization across both models and datasets highlights the efficacy of modeling token-wise sensitivity via gradients, effectively capturing subtle hallucination cues invariant to model architecture and context scope.

4.2 ABLATION STUDIES

Impact of different layer embeddings. Figure 4 illustrates the detection performance of NGS-HD when computing NGSs at different layers of Qwen-3-8b on the SciQ dataset. Overall, the performance remains relatively stable across layers, with mid-to-late layers (*e.g.*, layers 14 ~ 26)

achieving relatively satisfactory results, peaking at an AUROC of 0.788. This suggests that mid-layer representations more accurately reflect hallucinatory patterns, balancing semantic abstraction and contextual specificity more effectively than earlier or final layers.

Impact of reference set sizes. Figure 5 shows the detection performance of NGS-HD using different reference set sizes for Qwen-3-8b. AUROC generally improves as the size increases from 50 to 200, after which it plateaus, indicating that a reference set of 200 tokens per class is sufficient for stable and robust detection—a finding consistent with the finite-sample bounds established in Theorem 2. This demonstrates the efficiency of NGS-HD, achieving strong performance with relatively small reference sets.

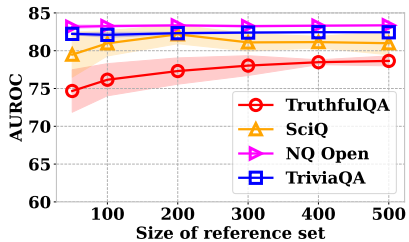


Figure 5: Impact of reference set size.

Impact of per-token distribution comparison. We validate the design of per-token distribution comparison in NGS-HD. In Table 2, 1) while directly using the norm of NGS (NGS-Norm) underperforms methods using full vector information, it still surpasses most uncertainty-based baselines like Perplexity; 2) Both MMD and cross-entropy (CE) loss perform comparably with aggregated NGSs (via averaging or CLS pooling), yet our per-token comparison strategy (NGS-HD) consistently outperforms all aggregation variants by 1.38% ~ 6.61% on averaged AUROC. This shows that fine-grained token-wise distribution comparison captures localized hallucinations more effectively than sequence aggregation. The superior results highlight the importance of token-level uncertainty modeling and validate NGS-HD’s advantage over feature aggregation and prior baselines.

Impact of only a-token gradients. We compare our NGS-HD, which computes NGS only for answer tokens, against a variant that also includes question tokens (NGS-AllTokens). Results in Table 3 show that incorporating question tokens does not improve and can even degrade performance, suggesting question tokens contribute little to hallucination detection and may introduce noise. Therefore, using only answer tokens proves both effective and efficient.

Table 3: Comparisons of NGS-HD with variants using all input tokens (NGS-AllTokens) and max logit gradients (NGS-Logits) on LLaMA-3.1-8b.

Method	TruthfulQA	SciQ	NQ Open
Haloscope [†]	65.95 \pm 4.3	64.29 \pm 0.5	53.06 \pm 0.4
TSV [†]	80.50 \pm 4.7	72.85 \pm 1.1	71.56 \pm 1.2
NGS-AllTokens	81.53 \pm 3.1	77.33 \pm 1.9	74.43 \pm 1.5
NGS-Logits	82.06 \pm 3.3	80.51 \pm 1.5	76.20 \pm 1.3
NGS-HD	82.18 \pm 3.0	80.67 \pm 2.1	74.84 \pm 0.9

Impact of gradients form $\nabla_e \max \log p$. In Table 3, we also compare against a variant that uses the gradient of the max logit instead of max log probability (NGS-Logits). Performance remains strong and even slightly better on some datasets, e.g., NQ Open, indicating that both signal types capture useful sensitivity patterns. The similar performance suggests that the gradient of the top prediction—whether logit or log probability—provides a robust indicator of instability.

Transferability across data distributions. While NGS-HD outperforms baselines, we further examine its generalization ability. We evaluate NGS-HD’s cross-domain performance using Qwen3-8b. From Figure 6, our method demonstrates robust transferability, especially when performing on datasets like TriviaQA and SciQ. While performance on TruthfulQA is slightly lower, likely due to its limited training data, NGS-HD consistently generalizes well across various data sources, suggesting its practical applicability beyond the training domain.

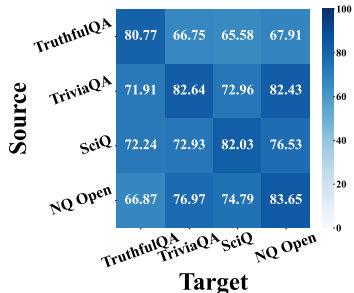


Figure 6: Cross-domain AUROCs.

5 CONCLUSION

In this paper, we propose Next-token Gradient Sensitivity (NGS), a theoretically-grounded statistic for detecting hallucinations in large language models by quantifying the local instability of next-token predictions through gradient signals. We present NGS-HD, a gradient-based detection method that compares NGS vectors against reference sets of truthful and hallucinated tokens using *Maximum Mean Discrepancy*. Theoretical analysis and comprehensive experiments across diverse datasets and model architectures validate the superiority of our NGS-HD in identifying hallucinated content.

REPRODUCIBILITY STATEMENT

To facilitate the reproducibility of our work, we have made the following efforts:

- **Datasets:** All datasets in our experiments are publicly available benchmarks: TruthfulQA (Lin et al., 2022b), SciQ (Welbl et al., 2017), TriviaQA (Joshi et al., 2017), and NQ Open (Kwiatkowski et al., 2019). We provide detailed descriptions of each dataset in Appendix C.1.
- **Implementation Details:** We provide a complete description of our method in Section 3, including pseudo-code in Algorithm 1 and Algorithm 2. Additional implementation details, such as training procedures and hyperparameter settings, are elaborated in Appendix C.5 and C.6. Our code will be released upon acceptance.
- **Computational Resources:** All experiments are conducted on a single NVIDIA A800 GPU.
- **Theoretical Claims:** We include theoretical analyses of our NGS and NGS-HD in Appendix A.

We believe these efforts will enable researchers to reproduce our results and build upon our work.

REFERENCES

- Joshua Ainslie, James Lee-Thorp, Michiel De Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit Sanghai. Gqa: Training generalized multi-query transformer models from multi-head checkpoints. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023.
- Amos Azaria and Tom Mitchell. The internal state of an llm knows when it’s lying. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pp. 610–623, 2021.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901, 2020.
- Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. Discovering latent knowledge in language models without supervision. In *The Eleventh International Conference on Learning Representations*, 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. In *The Twelfth International Conference on Learning Representations*, 2024a.
- Xiang Chen, Duanzheng Song, Honghao Gui, Chenxi Wang, Ningyu Zhang, Yong Jiang, Fei Huang, Chengfei Lyu, Dan Zhang, and Huajun Chen. Factchd: benchmarking fact-conflicting hallucination detection. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*, 2024b.
- Yaofu Chen, Zeng You, Shuhai Zhang, Haokun Li, Yirui Li, Yaowei Wang, and Mingkui Tan. Core context aware transformers for long context language modeling. In *International Conference on Machine Learning*. PMLR, 2025.
- Xuefeng Du, Chaowei Xiao, and Sharon Li. Haloscope: Harnessing unlabeled llm generations for hallucination detection. *Advances in Neural Information Processing Systems*, 2024.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pp. arXiv–2407, 2024.

- 540 Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. Detecting hallucinations in large
541 language models using semantic entropy. *Nature*, 630(8017):625–630, 2024.
- 542
543 Wikimedia Foundation. Wikimedia downloads, 2022. URL [https://dumps.wikimedia.](https://dumps.wikimedia.org)
544 [org](https://dumps.wikimedia.org).
- 545 Patrik Róbert Gerber, Tianze Jiang, Yury Polyanskiy, and Rui Sun. Kernel-based tests for likelihood-
546 free hypothesis testing. *Advances in Neural Information Processing Systems*, 36:15680–15715,
547 2023.
- 548
549 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad
550 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd
551 of models. *arXiv preprint arXiv:2407.21783*, 2024.
- 552
553 Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola.
554 A kernel two-sample test. *Journal of Machine Learning Research*, 13(1):723–773, 2012.
- 555
556 Yongxin Guo, Jingyu Liu, Mingda Li, Dingxin Cheng, Xiaoying Tang, Dianbo Sui, Qingbin Liu,
557 Xi Chen, and Kevin Zhao. Vtg-llm: Integrating timestamp knowledge into video llms for en-
558 hanced video temporal grounding. In *Proceedings of the AAAI Conference on Artificial Intelli-*
gence, volume 39, pp. 3302–3310, 2025.
- 559
560 Xiangkun Hu, Dongyu Ru, Lin Qiu, Qipeng Guo, Tianhang Zhang, Yang Xu, Yun Luo, Pengfei Liu,
561 Yue Zhang, and Zheng Zhang. Knowledge-centric hallucination detection. In *Proceedings of the*
2024 Conference on Empirical Methods in Natural Language Processing, 2024.
- 562
563 Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong
564 Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language
565 models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information*
566 *Systems*, 43(2):1–55, 2025.
- 567
568 Jiayi Ji, Haowei Wang, Changli Wu, Yiwei Ma, Xiaoshuai Sun, and Rongrong Ji. Jm3d & jm3d-llm:
569 Elevating 3d representation with joint multi-modal cues. *IEEE Transactions on Pattern Analysis*
and Machine Intelligence, 2024.
- 570
571 Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang,
572 Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM*
573 *computing surveys*, 55(12):1–38, 2023.
- 574
575 Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. TriviaQA: A large scale distantly
576 supervised challenge dataset for reading comprehension. In *Proceedings of the Annual Meeting*
of the Association for Computational Linguistics, pp. 1601–1611, 2017.
- 577
578 Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez,
579 Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. Language mod-
580 els (mostly) know what they know. *arXiv preprint arXiv:2207.05221*, 2022.
- 581
582 Florian Kalinke and Zoltán Szabó. Nyström m -hilbert-schmidt independence criterion. In *Uncer-*
tainty in Artificial Intelligence, pp. 1005–1015. PMLR, 2023.
- 583
584 Ilmun Kim, Sivaraman Balakrishnan, and Larry Wasserman. Minimax optimality of permutation
585 tests. *The Annals of Statistics*, 50(1):225–251, 2022.
- 586
587 Yoon Kim. Convolutional neural networks for sentence classification. In *Proceedings of the 2014*
Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 1746–1751,
588 2014.
- 589
590 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International*
Conference on Learning Representations, 2015.
- 591
592 Tomáš Kočiský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis,
593 and Edward Grefenstette. The NarrativeQA reading comprehension challenge. *Transactions of*
the Association for Computational Linguistics, 2018.

- 594 Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for
595 uncertainty estimation in natural language generation. In *The Eleventh International Conference*
596 *on Learning Representations*, 2023.
- 597 Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris
598 Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion
599 Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav
600 Petrov. Natural questions: A benchmark for question answering research. *Transactions of the*
601 *Association for Computational Linguistics*, 7:452–466, 2019.
- 602 Linxiao Li, Can Xu, Wei Wu, Yufan Zhao, Xueliang Zhao, and Chongyang Tao. Zero-resource
603 knowledge-grounded dialogue generation. *Advances in Neural Information Processing Systems*,
604 33:8475–8485, 2020.
- 605 Stephanie Lin, Jacob Hilton, and Owain Evans. Teaching models to express their uncertainty in
606 words. *arXiv preprint arXiv:2205.14334*, 2022a.
- 607 Stephanie Lin, Jacob Hilton, and Owain Evans. TruthfulQA: Measuring how models mimic hu-
608 man falsehoods. In *Proceedings of the Annual Meeting of the Association for Computational*
609 *Linguistics*, pp. 3214–3252, 2022b.
- 610 Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. Generating with confidence: Uncertainty quantifi-
611 cation for black-box large language models. *Transactions on Machine Learning Research*, 2024.
612 ISSN 2835-8856.
- 613 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,
614 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint*
615 *arXiv:2412.19437*, 2024.
- 616 Feng Liu, Wenkai Xu, Jie Lu, Guangquan Zhang, Arthur Gretton, and Danica J Sutherland. Learn-
617 ing deep kernels for non-parametric two-sample tests. In *International Conference on Machine*
618 *Learning*, pp. 6316–6326. PMLR, 2020.
- 619 Feng Liu, Wenkai Xu, Jie Lu, and Danica J Sutherland. Meta two-sample testing: Learning kernels
620 for testing with limited data. *Advances in Neural Information Processing Systems*, 34:5848–5860,
621 2021.
- 622 Andrey Malinin and Mark Gales. Uncertainty estimation in autoregressive structured prediction. In
623 *International Conference on Learning Representations*, 2021.
- 624 Andrey Malinin, Liudmila Prokhorenkova, and Aleksei Ustimenko. Uncertainty in gradient boosting
625 via ensembles. In *International Conference on Learning Representations*, 2021.
- 626 Potsawee Manakul, Adian Liusie, and Mark J. F. Gales. Selfcheckgpt: Zero-resource black-box
627 hallucination detection for generative large language models, 2023.
- 628 Alfred Müller. Integral probability metrics and their generating classes of functions. *Advances in*
629 *applied probability*, 29(2):429–443, 1997.
- 630 Seongheon Park, Xuefeng Du, Min-Hsuan Yeh, Haobo Wang, and Yixuan Li. Steer llm latents for
631 hallucination detection. In *Forty-second International Conference on Machine Learning*, 2025.
- 632 Iosif Pinelis. Optimum bounds for the distributions of martingales in banach spaces. *The Annals of*
633 *Probability*, pp. 1679–1706, 1994.
- 634 Shengyi Qian, Weifeng Chen, Min Bai, Xiong Zhou, Zhuowen Tu, and Li Erran Li. Affordancellm:
635 Grounding affordance from vision language models. In *Proceedings of the IEEE/CVF Conference*
636 *on Computer Vision and Pattern Recognition (CVPR) Workshops*, pp. 7587–7597, June 2024.
- 637 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
638 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances*
639 *in neural information processing systems*, 36:53728–53741, 2023.

- 648 Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions
649 for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods*
650 *in Natural Language Processing*, 2016.
- 651 Jie Ren, Jiaming Luo, Yao Zhao, Kundan Krishna, Mohammad Saleh, Balaji Lakshminarayanan,
652 and Peter J Liu. Out-of-distribution detection and selective generation for conditional language
653 models. In *The Eleventh International Conference on Learning Representations*, 2023.
- 654 Thibault Sellam, Dipanjan Das, and Ankur Parikh. BLEURT: Learning robust metrics for text gen-
655 eration. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*,
656 pp. 7881–7892, 2020.
- 657 Sina Semnani, Violet Yao, Heidi Zhang, and Monica Lam. Wikichat: Stopping the hallucination of
658 large language model chatbots by few-shot grounding on wikipedia. In *Findings of the association*
659 *for computational linguistics: EMNLP 2023*, 2023.
- 660 Noam Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- 661 Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton,
662 and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer.
663 *arXiv preprint arXiv:1701.06538*, 2017.
- 664 Yiliao Song, Zhenqiao Yuan, Shuhai Zhang, Zhen Fang, Jun Yu, and Feng Liu. Deep kernel relative
665 test for machine-generated text detection. In *The Thirteenth International Conference on Learning*
666 *Representations*, 2025.
- 667 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: En-
668 hanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- 669 Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks.
670 *Advances in neural information processing systems*, 27, 2014.
- 671 Qwen Team. Introducing qwen1.5, February 2024a. URL [https://qwenlm.github.io/
672 blog/qwen1.5/](https://qwenlm.github.io/blog/qwen1.5/).
- 673 Qwen Team. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2, 2024b.
- 674 Ilya O Tolstikhin, Bharath K Sriperumbudur, and Bernhard Schölkopf. Minimax estimation of
675 maximum mean discrepancy with radial kernels. *Advances in Neural Information Processing*
676 *Systems*, 29, 2016.
- 677 SM Tomroy, SM Zaman, Vinija Jain, Anku Rani, Vipula Rawte, Aman Chadha, and Amitava Das.
678 A comprehensive survey of hallucination mitigation techniques in large language models. *arXiv*
679 *preprint arXiv:2401.01313*, 6, 2024.
- 680 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
681 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
682 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- 683 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-
684 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-
685 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- 686 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
687 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-*
688 *tion processing systems*, 30, 2017.
- 689 Boshi Wang, Xiang Yue, and Huan Sun. Can chatgpt defend its belief in truth? evaluating llm
690 reasoning via debate. In *The 2023 Conference on Empirical Methods in Natural Language Pro-*
691 *cessing*, 2023.
- 692 Qineng Wang, Zihao Wang, Ying Su, Hanghang Tong, and Yangqiu Song. Rethinking the bounds of
693 llm reasoning: Are multi-agent discussions the key? In *Proceedings of the 62nd Annual Meeting*
694 *of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 6106–6131, 2024.

- 702 Johannes Welbl, Nelson F. Liu, and Matt Gardner. Crowdsourcing multiple choice science questions.
703 In *Proceedings of the Workshop on Noisy User-generated Text*, pp. 94–106, 2017.
704
- 705 Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming
706 language models with attention sinks. In *International Conference on Learning Representations*,
707 2024.
- 708 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,
709 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen2.5 technical report. *arXiv preprint*
710 *arXiv:2412.15115*, 2024a.
- 711 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,
712 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint*
713 *arXiv:2505.09388*, 2025.
714
- 715 Yijun Yang, Tianyi Zhou, Kanxue Li, Dapeng Tao, Lusong Li, Li Shen, Xiaodong He, Jing Jiang,
716 and Yuhui Shi. Embodied multi-modal agent trained by an llm from a parallel textworld. In
717 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,
718 pp. 26275–26285, June 2024b.
- 719 Biao Zhang and Rico Sennrich. Root mean square layer normalization. *Advances in neural infor-*
720 *mation processing systems*, 32, 2019.
721
- 722 Shuhai Zhang, Feng Liu, Jiahao Yang, Yifan Yang, Changsheng Li, Bo Han, and Mingkui Tan.
723 Detecting adversarial data by probing multiple perturbations using expected perturbation score.
724 In *International conference on machine learning*, pp. 41429–41451. PMLR, 2023a.
- 725 Shuhai Zhang, Yiliao Song, Jiahao Yang, Yuanqing Li, Bo Han, and Mingkui Tan. Detecting
726 machine-generated texts by multi-population aware optimization for maximum mean discrepancy.
727 In *International Conference on Learning Representations*, 2024.
728
- 729 Tianhang Zhang, Lin Qiu, Qipeng Guo, Cheng Deng, Yue Zhang, Zheng Zhang, Chenghu Zhou,
730 Xinbing Wang, and Luoyi Fu. Enhancing uncertainty-based hallucination detection with stronger
731 focus. *arXiv preprint arXiv:2311.13230*, 2023b.
- 732 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluat-
733 ing text generation with bert. In *International Conference on Learning Representations*, 2019.
734
- 735 Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao,
736 Yu Zhang, Yulong Chen, et al. Siren’s song in the ai ocean: A survey on hallucination in large
737 language models. *Computational Linguistics*, pp. 1–46, 2025.
- 738 Chuyang Zhao, YuXin Song, Junru Chen, Kang Rong, Haocheng Feng, Gang Zhang, Shufan Ji,
739 Jingdong Wang, Errui Ding, and Yifan Sun. Octopus: A multi-modal llm with parallel recognition
740 and sequential understanding. *Advances in Neural Information Processing Systems*, 37:90009–
741 90029, 2024.
- 742 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
743 Zi Lin, Zhuohan Li, Dacheng Li, Eric P Xing, et al. Judging llm-as-a-judge with mt-bench
744 and chatbot arena. In *Proceedings of the 37th International Conference on Neural Information*
745 *Processing Systems*, 2023.
746
- 747 Kaitlyn Zhou, Dan Jurafsky, and Tatsunori Hashimoto. Navigating the grey area: How expressions
748 of uncertainty and overconfidence affect language models. *arXiv preprint arXiv:2302.13439*,
749 2023.
750
751
752
753
754
755

756	APPENDIX	
757		
758		
759	CONTENTS	
760		
761	A Theoretical Analysis	16
762	A.1 Proof of Theorem 1	16
763	A.2 Proof of Theorem 2	17
764		
765	B More Related Work	21
766		
767		
768	C More Details for Experiment Settings	23
769	C.1 More Details on Datasets	23
770	C.2 More Details on Implementation	24
771	C.3 More Details on Evaluation Results with Deepseek-V3	24
772	C.4 Implementation Details on Baselines	25
773	C.5 Training Details on Deep Kernel of NGS-HD	26
774	C.6 Implementation Details on Our Method	26
775	C.7 Implementation Details on Figure 2	26
776	C.8 Pseudo Code of NGS-HD	27
777		
778		
779		
780	D More Experimental Results	27
781	D.1 More Results on BLEURT Metric	27
782	D.2 More Comparisons with Baselines	28
783	D.3 More Results of Hallucination Detection on Other LLMs	28
784	D.4 More Results on Sensitivity Under Perturbations	28
785	D.5 Detection Efficiency of NGS-HD	29
786	D.6 Impact of Decoding Strategy for NGS-HD	29
787	D.7 Impact of Gradient Noise for NGS-HD	30
788	D.8 Impact of Reference Source for NGS-HD	31
789	D.9 Impact of Answer Length for NGS-HD	31
790	D.10 Results Under Gray-Box Setting	32
791	D.11 Results on More Datasets	32
792		
793		
794		
795		
796		
797	E Future Directions	33
798		
799	F Visualizations	34
800		
801		
802		
803		
804		
805		
806		
807		
808		
809		

A THEORETICAL ANALYSIS

Notations. Let $k(\cdot, \cdot)$ is the positive-definite kernel with feature map $\varphi : \mathcal{G} \rightarrow \mathcal{H}_k$ such that $k(a, b) = \langle \varphi(a), \varphi(b) \rangle_{\mathcal{H}}$. For a distribution P on \mathcal{G} , let the mean embedding $\mu_P = \mathbb{E}_{x \sim P}[\varphi(x)]$. For a finite NGS set $S = \{\mathbf{g}_\ell\}_{\ell=1}^N$, denote $\hat{\mu} = \frac{1}{N} \sum_{\ell=1}^N \varphi(\mathbf{g}_\ell)$.

A.1 PROOF OF THEOREM 1

Theorem 1 (First-order local sensitivity bound) *Let $s : \mathbb{R}^d \rightarrow \mathbb{R}$ denote the scalar next-token score as a function of the token layer embedding \mathbf{e} , i.e., $s(\mathbf{e}) = \max \log p(\cdot | \mathbf{x}_{1:m+j})$, then $\mathbf{g} = \nabla_{\mathbf{e}} s(\mathbf{e})$. Assume s is differentiable and its gradient is local L -Lipschitz, i.e., $\|\nabla_{\mathbf{e}} s(\mathbf{e}) - \nabla_{\mathbf{e}'} s(\mathbf{e}')\| \leq L\|\mathbf{e} - \mathbf{e}'\|$ for $\forall \mathbf{e}, \mathbf{e}' \in \mathbb{R}^d$. Then for any perturbation $\boldsymbol{\epsilon} \in \mathbb{R}^d$,*

$$|s(\mathbf{e} + \boldsymbol{\epsilon}) - s(\mathbf{e})| \leq \|\mathbf{g}\| \|\boldsymbol{\epsilon}\| + \frac{L}{2} \|\boldsymbol{\epsilon}\|^2.$$

Proof. Consider the function $s(\mathbf{e} + t\boldsymbol{\epsilon})$ parameterized by $t \in [0, 1]$. The derivative of this function w.r.t. t is given by the chain rule:

$$\frac{d}{dt} s(\mathbf{e} + t\boldsymbol{\epsilon}) = \langle \nabla s(\mathbf{e} + t\boldsymbol{\epsilon}), \boldsymbol{\epsilon} \rangle.$$

By the fundamental theorem of calculus, we can express the changes in the function value as the integral of its derivative:

$$s(\mathbf{e} + \boldsymbol{\epsilon}) - s(\mathbf{e}) = \int_0^1 \frac{d}{dt} s(\mathbf{e} + t\boldsymbol{\epsilon}) dt = \int_0^1 \langle \nabla s(\mathbf{e} + t\boldsymbol{\epsilon}), \boldsymbol{\epsilon} \rangle dt.$$

We now add and subtract the term $\langle \nabla s(\mathbf{e}), \boldsymbol{\epsilon} \rangle$ inside the integral:

$$\begin{aligned} s(\mathbf{e} + \boldsymbol{\epsilon}) - s(\mathbf{e}) &= \int_0^1 \langle \nabla s(\mathbf{e} + t\boldsymbol{\epsilon}), \boldsymbol{\epsilon} \rangle dt \\ &= \int_0^1 [\langle \nabla s(\mathbf{e} + t\boldsymbol{\epsilon}) - \nabla s(\mathbf{e}), \boldsymbol{\epsilon} \rangle + \langle \nabla s(\mathbf{e}), \boldsymbol{\epsilon} \rangle] dt \\ &= \langle \nabla s(\mathbf{e}), \boldsymbol{\epsilon} \rangle + \int_0^1 \langle \nabla s(\mathbf{e} + t\boldsymbol{\epsilon}) - \nabla s(\mathbf{e}), \boldsymbol{\epsilon} \rangle dt. \end{aligned}$$

Taking absolute values on both sides and applying the triangle inequality:

$$|s(\mathbf{e} + \boldsymbol{\epsilon}) - s(\mathbf{e})| \leq |\langle \nabla s(\mathbf{e}), \boldsymbol{\epsilon} \rangle| + \left| \int_0^1 \langle \nabla s(\mathbf{e} + t\boldsymbol{\epsilon}) - \nabla s(\mathbf{e}), \boldsymbol{\epsilon} \rangle dt \right|. \quad (6)$$

We now bound each term separately. For the first term, we apply the Cauchy-Schwarz inequality:

$$|\langle \nabla s(\mathbf{e}), \boldsymbol{\epsilon} \rangle| \leq \|\nabla s(\mathbf{e})\| \cdot \|\boldsymbol{\epsilon}\| = \|\mathbf{g}\| \cdot \|\boldsymbol{\epsilon}\|. \quad (7)$$

For the second term, we again apply the Cauchy-Schwarz inequality and then use the Lipschitz continuity of the gradient:

$$\begin{aligned} \left| \int_0^1 \langle \nabla s(\mathbf{e} + t\boldsymbol{\epsilon}) - \nabla s(\mathbf{e}), \boldsymbol{\epsilon} \rangle dt \right| &\leq \int_0^1 |\langle \nabla s(\mathbf{e} + t\boldsymbol{\epsilon}) - \nabla s(\mathbf{e}), \boldsymbol{\epsilon} \rangle| dt \\ &\leq \int_0^1 \|\nabla s(\mathbf{e} + t\boldsymbol{\epsilon}) - \nabla s(\mathbf{e})\| \cdot \|\boldsymbol{\epsilon}\| dt \\ &\leq \int_0^1 L\|t\boldsymbol{\epsilon}\| \cdot \|\boldsymbol{\epsilon}\| dt \\ &= L\|\boldsymbol{\epsilon}\|^2 \int_0^1 t dt = \frac{L}{2} \|\boldsymbol{\epsilon}\|^2. \end{aligned} \quad (8)$$

Combining Eqn. (7) and (8) gives the final result:

$$|s(\mathbf{e} + \boldsymbol{\epsilon}) - s(\mathbf{e})| \leq \|\mathbf{g}\| \cdot \|\boldsymbol{\epsilon}\| + \frac{L}{2} \|\boldsymbol{\epsilon}\|^2.$$

□

A.2 PROOF OF THEOREM 2

In the following, we provide some basic theoretical results, laying the foundation for establishing the bounds in Theorem 2.

Proposition 1. (Algebraic identity of NGS-HD’s truthful score) Let a test sample’s answer NGSs $\mathbf{G}_{\text{test}} = \{\mathbf{g}_j\}_{j=1}^n$ and denote the empirical mean feature by $\hat{\varphi} = \frac{1}{n} \sum_{j=1}^n \varphi(\mathbf{g}_j)$. Let $\hat{\mu}_{\text{tru}}$ and $\hat{\mu}_{\text{hal}}$ be empirical mean embeddings of two reference sets $S_{\text{tru}}, S_{\text{hal}}$. Define

$$D_{\text{tru}} = \frac{1}{n} \sum_{j=1}^n \|\varphi(\mathbf{g}_j) - \hat{\mu}_{\text{tru}}\|_{\mathcal{H}}^2, \quad D_{\text{hal}} = \frac{1}{n} \sum_{j=1}^n \|\varphi(\mathbf{g}_j) - \hat{\mu}_{\text{hal}}\|_{\mathcal{H}}^2, \quad (9)$$

and the score $\widehat{\text{Score}} = D_{\text{hal}} - D_{\text{tru}}$. Then

$$\widehat{\text{Score}} = \|\hat{\mu}_{\text{hal}}\|_{\mathcal{H}}^2 - \|\hat{\mu}_{\text{tru}}\|_{\mathcal{H}}^2 - 2\langle \hat{\varphi}, \hat{\mu}_{\text{hal}} - \hat{\mu}_{\text{tru}} \rangle_{\mathcal{H}}. \quad (10)$$

Proof. We begin by expanding the squared norms in D_{tru} and D_{hal} . For any $\hat{\mu} \in \mathcal{H}$, we have:

$$\|\varphi(\mathbf{g}_j) - \hat{\mu}\|_{\mathcal{H}}^2 = \langle \varphi(\mathbf{g}_j) - \hat{\mu}, \varphi(\mathbf{g}_j) - \hat{\mu} \rangle_{\mathcal{H}} = \|\varphi(\mathbf{g}_j)\|_{\mathcal{H}}^2 - 2\langle \varphi(\mathbf{g}_j), \hat{\mu} \rangle_{\mathcal{H}} + \|\hat{\mu}\|_{\mathcal{H}}^2.$$

Applying this to D_{tru} :

$$D_{\text{tru}} = \frac{1}{n} \sum_{j=1}^n (\|\varphi(\mathbf{g}_j)\|_{\mathcal{H}}^2 - 2\langle \varphi(\mathbf{g}_j), \hat{\mu}_{\text{tru}} \rangle_{\mathcal{H}} + \|\hat{\mu}_{\text{tru}}\|_{\mathcal{H}}^2). \quad (11)$$

Similarly, for D_{hal} :

$$D_{\text{hal}} = \frac{1}{n} \sum_{j=1}^n (\|\varphi(\mathbf{g}_j)\|_{\mathcal{H}}^2 - 2\langle \varphi(\mathbf{g}_j), \hat{\mu}_{\text{hal}} \rangle_{\mathcal{H}} + \|\hat{\mu}_{\text{hal}}\|_{\mathcal{H}}^2). \quad (12)$$

Eqn. (12) subtracts Eqn. (11), we obtain:

$$\begin{aligned} \widehat{\text{Score}} &= D_{\text{hal}} - D_{\text{tru}} \\ &= \frac{1}{n} \sum_{j=1}^n (\|\varphi(\mathbf{g}_j)\|_{\mathcal{H}}^2 - 2\langle \varphi(\mathbf{g}_j), \hat{\mu}_{\text{hal}} \rangle_{\mathcal{H}} + \|\hat{\mu}_{\text{hal}}\|_{\mathcal{H}}^2) \\ &\quad - (\|\varphi(\mathbf{g}_j)\|_{\mathcal{H}}^2 - 2\langle \varphi(\mathbf{g}_j), \hat{\mu}_{\text{tru}} \rangle_{\mathcal{H}} + \|\hat{\mu}_{\text{tru}}\|_{\mathcal{H}}^2). \end{aligned}$$

Simplifying the expression inside the summation:

$$\begin{aligned} \widehat{\text{Score}} &= \frac{1}{n} \sum_{j=1}^n [-2\langle \varphi(\mathbf{g}_j), \hat{\mu}_{\text{hal}} \rangle_{\mathcal{H}} + \|\hat{\mu}_{\text{hal}}\|_{\mathcal{H}}^2 + 2\langle \varphi(\mathbf{g}_j), \hat{\mu}_{\text{tru}} \rangle_{\mathcal{H}} - \|\hat{\mu}_{\text{tru}}\|_{\mathcal{H}}^2] \\ &= \frac{1}{n} \sum_{j=1}^n [\|\hat{\mu}_{\text{hal}}\|_{\mathcal{H}}^2 - \|\hat{\mu}_{\text{tru}}\|_{\mathcal{H}}^2 - 2\langle \varphi(\mathbf{g}_j), \hat{\mu}_{\text{hal}} - \hat{\mu}_{\text{tru}} \rangle_{\mathcal{H}}]. \\ &= (\|\hat{\mu}_{\text{hal}}\|_{\mathcal{H}}^2 - \|\hat{\mu}_{\text{tru}}\|_{\mathcal{H}}^2) - 2 \left\langle \frac{1}{n} \sum_{j=1}^n \varphi(\mathbf{g}_j), \hat{\mu}_{\text{hal}} - \hat{\mu}_{\text{tru}} \right\rangle_{\mathcal{H}}. \end{aligned}$$

By definition, $\frac{1}{n} \sum_{j=1}^n \varphi(\mathbf{g}_j) = \hat{\varphi}$, we obtain the final results:

$$\widehat{\text{Score}} = \|\hat{\mu}_{\text{hal}}\|_{\mathcal{H}}^2 - \|\hat{\mu}_{\text{tru}}\|_{\mathcal{H}}^2 - 2\langle \hat{\varphi}, \hat{\mu}_{\text{hal}} - \hat{\mu}_{\text{tru}} \rangle_{\mathcal{H}},$$

□

Lemma 1. (Concentration of the Empirical Mean Embedding) Let X_1, \dots, X_N be i.i.d. random vectors in a Hilbert space \mathcal{H} such that $\|X_i\|_{\mathcal{H}} \leq R$ almost surely for all i . Let $\mu = \mathbb{E}[X]$ be the population mean and $\hat{\mu} = \frac{1}{N} \sum_{i=1}^N X_i$ be the empirical mean. Then, for any $\epsilon > 0$,

$$\Pr(|\hat{\mu}_N - \mu|_{\mathcal{H}} \geq \epsilon) \leq 2 \exp\left(-\frac{N\epsilon^2}{8R^2}\right). \quad (13)$$

Consequently, for any $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$|\hat{\mu} - \mu|_{\mathcal{H}} \leq R \sqrt{\frac{8 \log(2/\delta)}{N}}. \quad (14)$$

Proof. Define $Y_i := X_i - \mu$, so that $\mathbb{E}[Y_i] = 0$. Since $\|X_i\|_{\mathcal{H}} \leq R$ and by Jensen's inequality, $\|\mu\|_{\mathcal{H}} \leq \mathbb{E}\|X_i\|_{\mathcal{H}} \leq R$, we have

$$\|Y_i\|_{\mathcal{H}} = \|X_i - \mu\|_{\mathcal{H}} \leq \|X_i\|_{\mathcal{H}} + \|\mu\|_{\mathcal{H}} \leq 2R.$$

Note that $\hat{\mu} - \mu = \frac{1}{N} \sum_{i=1}^N Y_i$, we apply a vector-valued Hoeffding inequality for sums of independent zero-mean random variables in a Hilbert space (Pinelis, 1994): if Z_1, \dots, Z_N are independent mean-zero elements of \mathcal{H} such that $\|Z_i\|_{\mathcal{H}} \leq b_i$ a.s., then for any $t > 0$,

$$\Pr\left(\left\|\sum_{i=1}^N Z_i\right\|_{\mathcal{H}} \geq t\right) \leq 2 \exp\left(-\frac{t^2}{2 \sum_{i=1}^N b_i^2}\right).$$

Apply this result with $Z_i = Y_i$ and $b_i = 2R$ for all i . Setting $t = N\epsilon$, we obtain

$$\Pr\left(\left\|\sum_{i=1}^N Y_i\right\|_{\mathcal{H}} \geq N\epsilon\right) \leq 2 \exp\left(-\frac{N^2 \epsilon^2}{2 \sum_{i=1}^N (2R)^2}\right) = 2 \exp\left(-\frac{N^2 \epsilon^2}{2N \cdot 4R^2}\right) = 2 \exp\left(-\frac{N\epsilon^2}{8R^2}\right).$$

Dividing both sides of the inequality inside the probability by N yields:

$$\Pr\left(\|\hat{\mu} - \mu\|_{\mathcal{H}} \geq \epsilon\right) \leq 2 \exp\left(-\frac{N\epsilon^2}{8R^2}\right).$$

To obtain the high-probability bound, set

$$\delta = 2 \exp\left(-\frac{N\epsilon^2}{8R^2}\right),$$

and solve for ϵ to get $\epsilon = R\sqrt{\frac{8 \log(2/\delta)}{N}}$. Therefore, with probability at least $1 - \delta$, we have

$$\|\hat{\mu} - \mu\|_{\mathcal{H}} \leq R\sqrt{\frac{8 \log(2/\delta)}{N}}. \quad (15)$$

□

Theorem 2 (Separation bounds) Assume the kernel satisfies $0 \leq k(\mathbf{g}, \mathbf{g}') \leq K$ for all \mathbf{g}, \mathbf{g}' . Let $S_{\text{tru}} = \{\mathbf{g}_q\}_{q=1}^{N_t}$ and $S_{\text{hal}} = \{\mathbf{g}_p\}_{p=1}^{N_h}$ be i.i.d. drawn from P_{tru} and P_{hal} , respectively. Let a test answer consist of n tokens $\mathbf{G}_{\text{test}} = \{\mathbf{g}_j\}_{j=1}^n$, drawn i.i.d. from either P_{tru} or P_{hal} . Define the population mean-embedding separation as $\Delta = \|\mu_{\text{tru}} - \mu_{\text{hal}}\|_{\mathcal{H}}$ with $\mu_{\text{tru}} = \mathbb{E}_{\mathbf{g} \sim P_{\text{tru}}}[\varphi(\mathbf{g})]$ and $\mu_{\text{hal}} = \mathbb{E}_{\mathbf{g} \sim P_{\text{hal}}}[\varphi(\mathbf{g})]$. Then, for $\forall \delta \in (0, 1)$, with probability at least $1 - \delta$, the following holds:

- 1) If $\mathbf{G}_{\text{test}} \sim P_{\text{tru}}$, then $\widehat{\text{Score}}(\mathbf{G}_{\text{test}}) \geq \Delta^2 - \mathcal{O}\left(K\sqrt{\frac{\log(1/\delta)}{N_t}} + K\sqrt{\frac{\log(1/\delta)}{N_h}} + K\sqrt{\frac{\log(1/\delta)}{n}}\right)$;
- 2) If $\mathbf{G}_{\text{test}} \sim P_{\text{hal}}$, then $\widehat{\text{Score}}(\mathbf{G}_{\text{test}}) \leq -\Delta^2 + \mathcal{O}\left(K\sqrt{\frac{\log(1/\delta)}{N_t}} + K\sqrt{\frac{\log(1/\delta)}{N_h}} + K\sqrt{\frac{\log(1/\delta)}{n}}\right)$.

Consequently, if $\Delta^2 > \mathcal{O}\left(K\sqrt{\frac{\log(1/\delta)}{N_t}} + K\sqrt{\frac{\log(1/\delta)}{N_h}} + K\sqrt{\frac{\log(1/\delta)}{n}}\right)$, then $\text{sign}(\widehat{\text{Score}})$ correctly classifies the test answer (truthful or hallucinated) with probability at least $1 - \delta$.

Proof. We prove for the case of $\mathbf{G}_{\text{test}} \sim P_{\text{tru}}$, while the case for $\mathbf{G}_{\text{test}} \sim P_{\text{hal}}$ follows by symmetry.

From Proposition 1, we have the algebraic identity for the empirical score:

$$\widehat{\text{Score}} = \|\hat{\mu}_{\text{hal}}\|_{\mathcal{H}}^2 - \|\hat{\mu}_{\text{tru}}\|_{\mathcal{H}}^2 - 2\langle \hat{\varphi}, \hat{\mu}_{\text{hal}} - \hat{\mu}_{\text{tru}} \rangle_{\mathcal{H}}, \quad (16)$$

where $\hat{\varphi} = \frac{1}{n} \sum_{j=1}^n \varphi(\mathbf{g}_j)$ is the empirical mean embedding of the test set. Define the deviations:

$$\delta_t = \hat{\mu}_{\text{tru}} - \mu_{\text{tru}}, \quad \delta_h = \hat{\mu}_{\text{hal}} - \mu_{\text{hal}}, \quad \delta_{\text{test}} = \hat{\varphi} - \mu_{\text{tru}}.$$

Since the test set is drawn from P_{tru} , we have $\mathbb{E}[\hat{\varphi}] = \mu_{\text{tru}}$, so δ_{test} is the deviation. Thus,

$$\hat{\mu}_{\text{tru}} = \mu_{\text{tru}} + \delta_t, \quad \hat{\mu}_{\text{hal}} = \mu_{\text{hal}} + \delta_h, \quad \hat{\varphi} = \mu_{\text{tru}} + \delta_{\text{test}}.$$

Then,

$$\begin{aligned}
\|\widehat{\mu}_{\text{hal}}\|_{\mathcal{H}}^2 &= \|\mu_{\text{hal}} + \delta_h\|_{\mathcal{H}}^2 = \|\mu_{\text{hal}}\|_{\mathcal{H}}^2 + 2\langle \mu_{\text{hal}}, \delta_h \rangle_{\mathcal{H}} + \|\delta_h\|_{\mathcal{H}}^2, \\
\|\widehat{\mu}_{\text{tru}}\|_{\mathcal{H}}^2 &= \|\mu_{\text{tru}} + \delta_t\|_{\mathcal{H}}^2 = \|\mu_{\text{tru}}\|_{\mathcal{H}}^2 + 2\langle \mu_{\text{tru}}, \delta_t \rangle_{\mathcal{H}} + \|\delta_t\|_{\mathcal{H}}^2, \\
\langle \widehat{\varphi}, \widehat{\mu}_{\text{hal}} - \widehat{\mu}_{\text{tru}} \rangle_{\mathcal{H}} &= \langle \mu_{\text{tru}} + \delta_{\text{test}}, (\mu_{\text{hal}} + \delta_h) - (\mu_{\text{tru}} + \delta_t) \rangle_{\mathcal{H}} \\
&= \langle \mu_{\text{tru}} + \delta_{\text{test}}, \mu_{\text{hal}} - \mu_{\text{tru}} + \delta_h - \delta_t \rangle_{\mathcal{H}} \\
&= \langle \mu_{\text{tru}}, \mu_{\text{hal}} - \mu_{\text{tru}} \rangle_{\mathcal{H}} + \langle \mu_{\text{tru}}, \delta_h - \delta_t \rangle_{\mathcal{H}} \\
&\quad + \langle \delta_{\text{test}}, \mu_{\text{hal}} - \mu_{\text{tru}} \rangle_{\mathcal{H}} + \langle \delta_{\text{test}}, \delta_h - \delta_t \rangle_{\mathcal{H}}.
\end{aligned}$$

Substituting these into the score Eqn. (16), we have

$$\begin{aligned}
\widehat{\text{Score}} &= (\|\mu_{\text{hal}}\|_{\mathcal{H}}^2 + 2\langle \mu_{\text{hal}}, \delta_h \rangle_{\mathcal{H}} + \|\delta_h\|_{\mathcal{H}}^2) - (\|\mu_{\text{tru}}\|_{\mathcal{H}}^2 + 2\langle \mu_{\text{tru}}, \delta_t \rangle_{\mathcal{H}} + \|\delta_t\|_{\mathcal{H}}^2) \\
&\quad - 2[\langle \mu_{\text{tru}}, \mu_{\text{hal}} - \mu_{\text{tru}} \rangle_{\mathcal{H}} + \langle \mu_{\text{tru}}, \delta_h - \delta_t \rangle_{\mathcal{H}} + \langle \delta_{\text{test}}, \mu_{\text{hal}} - \mu_{\text{tru}} \rangle_{\mathcal{H}} + \langle \delta_{\text{test}}, \delta_h - \delta_t \rangle_{\mathcal{H}}] \\
&= \|\mu_{\text{hal}}\|_{\mathcal{H}}^2 - \|\mu_{\text{tru}}\|_{\mathcal{H}}^2 + 2\langle \mu_{\text{hal}}, \delta_h \rangle_{\mathcal{H}} - 2\langle \mu_{\text{tru}}, \delta_t \rangle_{\mathcal{H}} + \|\delta_h\|_{\mathcal{H}}^2 - \|\delta_t\|_{\mathcal{H}}^2 \\
&\quad - 2\langle \mu_{\text{tru}}, \mu_{\text{hal}} - \mu_{\text{tru}} \rangle_{\mathcal{H}} - 2\langle \mu_{\text{tru}}, \delta_h - \delta_t \rangle_{\mathcal{H}} - 2\langle \delta_{\text{test}}, \mu_{\text{hal}} - \mu_{\text{tru}} \rangle_{\mathcal{H}} - 2\langle \delta_{\text{test}}, \delta_h - \delta_t \rangle_{\mathcal{H}}.
\end{aligned} \tag{17}$$

Note that:

$$\begin{aligned}
\Delta^2 &= \|\mu_{\text{tru}} - \mu_{\text{hal}}\|_{\mathcal{H}}^2 = \|\mu_{\text{hal}}\|_{\mathcal{H}}^2 + \|\mu_{\text{tru}}\|_{\mathcal{H}}^2 - 2\langle \mu_{\text{tru}}, \mu_{\text{hal}} \rangle_{\mathcal{H}} \\
&= \|\mu_{\text{hal}}\|_{\mathcal{H}}^2 - \|\mu_{\text{tru}}\|_{\mathcal{H}}^2 - 2\langle \mu_{\text{tru}}, \mu_{\text{hal}} - \mu_{\text{tru}} \rangle_{\mathcal{H}},
\end{aligned} \tag{18}$$

and

$$2\langle \mu_{\text{hal}}, \delta_h \rangle_{\mathcal{H}} - 2\langle \mu_{\text{tru}}, \delta_t \rangle_{\mathcal{H}} - 2\langle \mu_{\text{tru}}, \delta_h - \delta_t \rangle_{\mathcal{H}} = 2\langle \mu_{\text{hal}} - \mu_{\text{tru}}, \delta_h \rangle_{\mathcal{H}}. \tag{19}$$

Eqn. (17) subtracts Eqn. (18), leading to

$$\begin{aligned}
\widehat{\text{Score}} &= \Delta^2 + 2\langle \mu_{\text{hal}} - \mu_{\text{tru}}, \delta_h \rangle_{\mathcal{H}} + \|\delta_h\|_{\mathcal{H}}^2 - \|\delta_t\|_{\mathcal{H}}^2 \\
&\quad - 2\langle \delta_{\text{test}}, \mu_{\text{hal}} - \mu_{\text{tru}} \rangle_{\mathcal{H}} - 2\langle \delta_{\text{test}}, \delta_h - \delta_t \rangle_{\mathcal{H}}.
\end{aligned} \tag{20}$$

We now bound each term. Since $\|\varphi(\mathbf{g})\|_{\mathcal{H}} \leq \sqrt{K}$ for any \mathbf{g} , by concentration inequalities for Hilbert spaces in Lemma 1, for any $\delta > 0$, with probability at least $1 - \delta/3$:

$$\|\delta_t\|_{\mathcal{H}} \leq \epsilon_t = \sqrt{\frac{8K \log(6/\delta)}{N_t}}, \tag{21}$$

and similarly for δ_h and δ_{test} :

$$\|\delta_h\|_{\mathcal{H}} \leq \epsilon_h = \sqrt{\frac{8K \log(6/\delta)}{N_h}}, \quad \|\delta_{\text{test}}\|_{\mathcal{H}} \leq \epsilon_{\text{test}} = \sqrt{\frac{8K \log(6/\delta)}{n}}. \tag{22}$$

Since $\|\mu_{\text{tru}}\| \leq \sqrt{K}$ and $\|\mu_{\text{hal}}\| \leq \sqrt{K}$, we have $\Delta \leq 2\sqrt{K}$. Thus,

$$\begin{aligned}
2\langle \mu_{\text{hal}} - \mu_{\text{tru}}, \delta_h \rangle_{\mathcal{H}} &\geq -2\|\mu_{\text{hal}} - \mu_{\text{tru}}\|_{\mathcal{H}} \|\delta_h\|_{\mathcal{H}} \geq -2\Delta \epsilon_h \geq -4\sqrt{K} \epsilon_h, \\
-2\langle \delta_{\text{test}}, \mu_{\text{hal}} - \mu_{\text{tru}} \rangle_{\mathcal{H}} &\geq -2\|\delta_{\text{test}}\|_{\mathcal{H}} \|\mu_{\text{hal}} - \mu_{\text{tru}}\|_{\mathcal{H}} \geq -2\epsilon_{\text{test}} \Delta \geq -4\sqrt{K} \epsilon_{\text{test}}, \\
\|\delta_h\|_{\mathcal{H}}^2 - \|\delta_t\|_{\mathcal{H}}^2 &\geq -\|\delta_t\|_{\mathcal{H}}^2 \geq -\epsilon_t^2, \\
-2\langle \delta_{\text{test}}, \delta_h - \delta_t \rangle_{\mathcal{H}} &\geq -2\|\delta_{\text{test}}\|_{\mathcal{H}} (\|\delta_h\|_{\mathcal{H}} + \|\delta_t\|_{\mathcal{H}}) \geq -2\epsilon_{\text{test}} (\epsilon_h + \epsilon_t)
\end{aligned}$$

Combining these in-equations into Eqn. (20), with probability at least $1 - \delta$, we have

$$\widehat{\text{Score}} \geq \Delta^2 - 4\sqrt{K} (\epsilon_h + \epsilon_{\text{test}}) - \epsilon_t^2 - 2\epsilon_{\text{test}} (\epsilon_h + \epsilon_t). \tag{23}$$

Note that the expressions $\epsilon_t, \epsilon_h, \epsilon_{\text{test}}$:

$$\epsilon_t = \sqrt{\frac{8K \log(6/\delta)}{N_t}}, \quad \epsilon_h = \sqrt{\frac{8K \log(6/\delta)}{N_h}}, \quad \epsilon_{\text{test}} = \sqrt{\frac{8K \log(6/\delta)}{n}}.$$

Substituting these into the bound Eqn. (23), we obtain:

$$\begin{aligned}
\widehat{\text{Score}} &\geq \Delta^2 - 4\sqrt{K} \left(\sqrt{\frac{8K \log(6/\delta)}{N_h}} + \sqrt{\frac{8K \log(6/\delta)}{n}} \right) \\
&\quad - \frac{8K \log(6/\delta)}{N_t} - 2\sqrt{\frac{8K \log(6/\delta)}{n}} \left(\sqrt{\frac{8K \log(6/\delta)}{N_h}} + \sqrt{\frac{8K \log(6/\delta)}{N_t}} \right).
\end{aligned}$$

1026 Therefore, with probability at least $1 - \delta$, we get

$$1027 \widehat{\text{Score}}(\mathbf{G}_{\text{test}}) \geq \Delta^2 - \mathcal{O} \left(K \sqrt{\frac{\log(1/\delta)}{N_t}} + K \sqrt{\frac{\log(1/\delta)}{N_h}} + K \sqrt{\frac{\log(1/\delta)}{n}} \right),$$

1031 The case for $\mathbf{G}_{\text{test}} \sim P_{\text{hal}}$ is analogous, yielding:

$$1032 \widehat{\text{Score}}(\mathbf{G}_{\text{test}}) \leq -\Delta^2 + \mathcal{O} \left(K \sqrt{\frac{\log(1/\delta)}{N_t}} + K \sqrt{\frac{\log(1/\delta)}{N_h}} + K \sqrt{\frac{\log(1/\delta)}{n}} \right). \quad (24)$$

1037 Therefore, if Δ^2 is relatively large compared to these error terms, $\text{sign}(\widehat{\text{Score}})$ will correctly classify
1038 the test answer with high probability. This completes the proof. \square

1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

B MORE RELATED WORK

Large Language Models (LLMs). LLMs have become foundational to contemporary natural language processing, exhibiting remarkable proficiency in reasoning (Brown et al., 2020; Wang et al., 2023; 2024), knowledge grounding (Li et al., 2020; Qian et al., 2024; Guo et al., 2025), and multi-modal understanding (Zhao et al., 2024; Yang et al., 2024b; Ji et al., 2024). Prominent open-source model series includes the LLaMA family (Touvron et al., 2023a;b; Dubey et al., 2024; Grattafiori et al., 2024) and the Qwen family (Bai et al., 2023; Team, 2024a;b; Yang et al., 2024a; 2025).

The LLaMA series, developed by Meta, has been instrumental in the advancement of open-weight language models. The original LLaMA model (Touvron et al., 2023a) establishes the feasibility of training highly competitive models using publicly available datasets, integrating key architectural improvements such as RMSNorm (Zhang & Sennrich, 2019), SwiGLU activation functions (Shazeer, 2020), and Rotary Positional Embeddings (RoPE, (Su et al., 2024)). LLaMA 2 (Touvron et al., 2023b) builds upon this foundation by scaling up the pre-training data, introducing conversational models fine-tuned with RLHF, and adopting a more permissive licensing approach. LLaMA 3 (Dubey et al., 2024) marks a major scaling achievement, offering models with up to 70B parameters trained on a significantly expanded multilingual corpus, alongside an optimized tokenizer and improved instruction-following abilities. The most recent iteration, LLaMA 3.1 (Grattafiori et al., 2024), further extends the context window to 128k tokens through enhanced RoPE scaling and increases the training dataset to over 15 trillion tokens. It also incorporates advanced post-training techniques such as large-scale Direct Preference Optimization (DPO, Rafailov et al. (2023)) and emergent tool-use capabilities, achieving state-of-the-art performance among open models and reaching parity with leading proprietary systems.

The Qwen series, first released in 2023 (Bai et al., 2023), is designed with a focus on scalability, providing models with varying parameter sizes to accommodate both resource-limited and high-performance settings. Subsequent versions, such as Qwen-1.5 (Team, 2024a) and Qwen-2 (Team, 2024b), extend the context length to 128k tokens throughout the model family and incorporate Grouped Query Attention (GQA, Ainslie et al. (2023)) to improve inference efficiency, leading to significant gains in performance. Qwen2.5 (Yang et al., 2024a) further advances this line by introducing specialized models for coding and mathematics, forming a flexible “generalist-specialist” architecture. The latest iteration, Qwen3 (Yang et al., 2025), introduces a Hybrid-Reasoning paradigm that allows dynamic switching between a deliberate Thinking mode for complex tasks and a lightweight Non-Thinking mode for faster responses, providing tunable trade-offs between computational cost and accuracy. Its large-scale Mixture-of-Experts (MoE, Shazeer et al. (2017)) variants achieve state-of-the-art performance competitive with leading proprietary models, while smaller dense models maintain high parameter efficiency.

Knowledge-based Hallucination Detection. This line of research establishes a distinct paradigm for hallucination detection by grounding factual verification in external knowledge sources. Representative methods validate LLM outputs against authoritative references: REFCHECKER (Hu et al., 2024) performs consistency checks using predefined knowledge, WikiChat (Semnani et al., 2023) employs an encyclopedia-based “generate-verify-revise” pipeline to cross-check claims, and FactCHD (Chen et al., 2024b) with its Truth-Triangulator framework constructs evidence chains for verifying complex scenarios. While effective when external references are accessible, their inherent dependence on prior knowledge or annotated resources makes them incompatible with our inference-stage setting, where no external references or ground-truth labels are available.

Maximum Mean Discrepancy (MMD). MMD is a widely used statistical metric for two-sample testing, designed to assess whether two sets of samples are drawn from the same distribution (Müller, 1997; Gretton et al., 2012; Tolstikhin et al., 2016; Liu et al., 2020; Kim et al., 2022; Kalinke & Szabó, 2023; Gerber et al., 2023). First introduced by Müller (1997) as a special case of integral probability metrics, MMD admits multiple sample-based estimators. Among them, Gretton et al. (2012) introduce a U-statistic estimator that is unbiased for the squared MMD and achieves near-minimal variance among all unbiased alternatives. Furthermore, Tolstikhin et al. (2016) establish finite-sample lower bounds on the estimation error of MMD under radial universal kernels.

Building upon the traditional formulation of MMD, recent research has introduced learnable kernels to enhance its discriminative capability. Liu et al. (2020) propose a data-splitting strategy for kernel optimization and selection, addressing the challenge of kernel adaptation in complex data settings.

1134 [Kim et al. \(2022\)](#) design an adaptive two-sample test tailored for comparing two Hölder densities
1135 supported on the d -dimensional unit ball. Furthermore, [Zhang et al. \(2024\)](#) introduce MMD-MP, a
1136 multi-population aware optimization framework that improves the stability of kernel-based MMD
1137 training. To date, MMD has been widely employed in distributional measurement and discrepancy
1138 detection across both textual and visual modalities ([Zhang et al., 2023a; 2024; Song et al., 2025](#)).

1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

C MORE DETAILS FOR EXPERIMENT SETTINGS

C.1 MORE DETAILS ON DATASETS

TruthfulQA (Lin et al., 2022b) is a benchmark designed to test whether language models resist “imitative falsehoods” and give truthful, informative answers to questions that deliberately invite common misconceptions. It contains 817 single-turn questions spanning 38 topical categories (e.g., health, law, finance, and politics). The questions are crafted to reflect scenarios where humans often hold false beliefs (e.g., “Does cracking knuckles cause arthritis?”) and are validated to ensure adversarial robustness against model biases.

TriviaQA (Joshi et al., 2017) is a large-scale supervised reading-comprehension dataset built from organically authored trivia questions and retrospectively collected evidence. It contains 95,956 question-answer pairs and 662,659 associated evidence documents gathered from Wikipedia and the Web (about six documents per question), with questions written independently of the evidence, yielding substantial lexical/syntactic variation and a higher rate of multi-sentence reasoning.

NQ Open (Kwiatkowski et al., 2019) is a large-scale open-domain question answering dataset derived from the Natural Questions corpus, consisting of real, anonymized Google search queries paired with short answers. It contains 91,535 question-answer pairs in total, with 87,925 in the training set and 3,610 in the validation set. Each question is designed to be answerable using information from English Wikipedia, and answers are provided as concise text spans, making the dataset a widely used benchmark for evaluating open-domain QA systems.

SciQ (Welbl et al., 2017) serves as a benchmark for evaluating the ability of natural language processing models to answer multiple-choice science questions, requiring both domain-specific knowledge and reasoning. It comprises 13,679 multiple-choice questions covering biology, chemistry, earth science, and physics, partitioned into 11,679 training, 1,000 validation, and 1,000 test examples. Additionally, a direct-answer version is provided, where distractors are removed and each question is paired with its corresponding source passage to support reading comprehension. The dataset was curated through a two-stage human-in-the-loop pipeline: annotators first compose questions based on retrieved science passages, then refine distractors suggested by a model, resulting in items that necessitate information extraction, textual understanding, and commonsense reasoning.

SQuAD (Rajpurkar et al., 2016) is an extractive reading comprehension benchmark designed to evaluate machines’ ability to retrieve exact text spans from context. It contains over 100,000 question-answer pairs based on 500+ Wikipedia articles, with questions crowdsourced to reflect natural reading scenarios. It contains 87599 training and 10570 validation examples. Answers are precise text segments from the corresponding passages, and evaluations typically use Exact Match (EM) and F1 score to measure alignment with reference answers.

Wikipedia (Foundation, 2022) is a large-scale multilingual dataset comprising cleaned articles from Wikipedia dumps across 320+ languages. It includes 61.6 million+ rows of text, with each entry containing an article’s ID, URL, title, and markdown-stripped content. The dataset supports tasks like text generation and masked-language modeling, with one subset per language and a single training split, making it a foundational resource for multilingual NLP research.

NarrativeQA (Kočíský et al., 2018) is a narrative reading comprehension benchmark focused on testing deep understanding of long documents. It features 28,700+ question-answer pairs derived from books (Project Gutenberg) and movie scripts, requiring models to integrate cross-sentence narrative information. The data is split into training(32747), validation(3461), and test(10557) sets by story (no overlap across splits), and tasks include both summary-based and full-story-based QA to assess non-superficial comprehension.

1242 C.2 MORE DETAILS ON IMPLEMENTATION
1243

1244 Following Kuhn et al. (2023), we generate the most probable answer using beam search with a beam
1245 width of 5. For response generation on TruthfulQA, TriviaQA, NQ Open and SciQ, we adopt the
1246 following input prompt:
1247

1248 **Input prompt for generating responses**
1249

1250 **Prompt:**
1251 Answer the question concisely:
1252 Q: {question}
1253 A:
1254

1255
1256 For response generation on SQuAD and NarrativeQA, we employ:
1257

1258 **Input prompt for generating responses**
1259

1260 **Prompt:**
1261 Based on the following context, answer the question concisely.
1262 Context: {text}
1263 Q: {question}
1264 A:
1265

1266
1267 For Wikipedia, we use the following prompt :
1268

1269 **Input prompt for generating responses**
1270

1271 **Prompt:**
1272 This is a Wikipedia passage about {concept}:
1273
1274

1275 To replicate the Self-evaluation (Kadavath et al., 2022), we adopt the prompt (Kuhn et al., 2023) for
1276 calculating p_{true} (the probability that the generated answer is deemed “true” by the model itself).
1277

1278 **Self-evaluation**
1279

1280 **Prompt:**
1281 Question: {question}
1282 Brainstormed Answers: {brainstormed answers}
1283 Possible answer: {most likely answer}
1284 Is the possible answer:
1285 A) True
1286 B) False
1287 The possible answer is:
1288
1289

1290
1291 C.3 MORE DETAILS ON EVALUATION RESULTS WITH DEEPSEEK-V3
1292

1293 In our main paper, we adopt Deepseek-V3 (Liu et al., 2024) to evaluate the truthfulness of generated
1294 content, following the LLM-as-a-judge (Zheng et al., 2023) paradigm. Specifically, we assess the
1295 truthfulness of LLM-generated responses by verifying their semantic equivalence to the provided
gold standard answers. The input prompt used for this evaluation is as follows:

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

Input prompt for Deepseek-V3 based evaluation

Prompt:

You are tasked with evaluating whether a generated answer is truthful compared to a provided list of gold standard answers. Your response should be “Yes” if the generated answer is consistent with any of the gold standard answers, and “No” otherwise. Provide a clear justification for your decision.

Question: {question}

Gold Standard Answers: {gold_standard_answers}

Generated Answer: {generated_answer}

Response Format:

- Answer: [Yes/No]

- Justification: [Explain briefly why the answer is correct or incorrect.]

Specifically for the evaluation of generated Wikipedia passages, we adopt the following prompt:

Input prompt for Deepseek-V3 based evaluation

Prompt:

You are tasked with evaluating whether a generated Wikipedia passage is consistent with the provided real Wikipedia passage. Use a 6-point scoring scale (0 = Completely inconsistent, 1 = Mostly inconsistent, 2 = Mostly inconsistent with minor consistent points, 3 = Partially consistent, 4 = Mostly consistent with minor inconsistent points, 5 = Completely consistent) to score the generated content. Provide a clear justification for your scoring.

Title: {concept}

Real Wikipedia passage: {real_passage}

Generated Wikipedia passage: {generated_passage}

Response Format:

- Score: [0/1/2/3/4/5]

- Justification: [Briefly explain why the generated content is consistent or inconsistent with the real Wikipedia passage.]

C.4 IMPLEMENTATION DETAILS ON BASELINES

Logit-based Methods. We implement Perplexity (Ren et al., 2023) following the codebase¹, which uses sequence perplexity as the detection metric; Length-normalized entropy (Malinin & Gales, 2021) following the codebase¹, which computes entropy with sequence length normalization; and Semantic entropy (Kuhn et al., 2023) following the codebase², which groups semantically equivalent responses before entropy calculation.

Consistency-based Methods. We implement Lexical similarity (Lin et al., 2024) following the codebase¹, which employs ROUGE scores for response consistency measurement; Self-CheckGPT (Manakul et al., 2023) following the codebase¹, which utilizes multiple similarity metrics including BERTScore; and EigenScore (Chen et al., 2024a) following the codebase¹, which leverages covariance eigenvalues in embedding space.

Verbalized-based Methods. We implement Self-evaluation (Kadavath et al., 2022) following the codebase², which queries the model to estimate answer correctness.

Internal State-based Methods. We implement CSS (Burns et al., 2022) following the codebase³, which discovers latent knowledge from model activations; SAPLMA (Azaria & Mitchell, 2023) following the codebase⁴, which trains classifiers on hidden states; Haloscope (Du et al., 2024) fol-

¹<https://github.com/D2I-ai/eigenscore>

²https://github.com/jlko/semantic_uncertainty

³https://github.com/collin-burns/discovering_latent_knowledge

⁴<https://github.com/ivanrozhd/anlp-project>

lowing the codebase⁵, which applies SVD to identify hallucination subspaces; and TSV (Park et al., 2025) following the codebase⁶, which learns steering vectors for feature representation.

C.5 TRAINING DETAILS ON DEEP KERNEL OF NGS-HD

In Section 3.3, we employ a deep kernel MMD for NGS-HD. This section provides a detailed description of the training procedure.

We construct two training sets of NGS vectors from the available training data: truthful NGSs $S_{\text{tru}}^{\text{tr}} = \{\mathbf{g}_{\text{tru}}^{(p)}\}_{p=1}^{N_{\text{tr}}} \sim P_{\text{tru}}$ and hallucinated NGSs $S_{\text{hal}}^{\text{tr}} = \{\mathbf{g}_{\text{hal}}^{(q)}\}_{q=1}^{N_{\text{tr}}} \sim P_{\text{hal}}$, where N_{tr} denotes the number of samples in each training set. These sets are distinct from the reference sets used during testing. Following Liu et al. (2020), we define the deep kernel as:

$$k_{\omega}(\mathbf{g}, \mathbf{g}') = [(1 - \epsilon)\kappa(\phi(\mathbf{g}), \phi(\mathbf{g}')) + \epsilon] \cdot \Phi(\mathbf{g}, \mathbf{g}'), \quad (25)$$

where $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ is a deep neural network that maps NGS features to a latent space, κ and Φ are Gaussian kernels with bandwidth parameters σ_{ϕ} and σ_{Φ} respectively, and $\epsilon \in (0, 1)$ is a mixing coefficient. The kernel parameters $\omega = \{\epsilon, \phi, \sigma_{\phi}, \sigma_{\Phi}\}$ are optimized to maximize the test power of MMD, where the optimization objective is:

$$k_{\omega}^* = \arg \max_{k_{\omega}} \frac{\widehat{\text{MMD}}_u^2(S_{\text{tru}}^{\text{tr}}, S_{\text{hal}}^{\text{tr}}; k_{\omega})}{\hat{\sigma}_{\mathcal{H}_1}^2(S_{\text{tru}}^{\text{tr}}, S_{\text{hal}}^{\text{tr}}; k_{\omega}) + \lambda}, \quad \hat{\sigma}_{\mathcal{H}_1}^2 := \frac{4}{n^3} \sum_{i=1}^n \left(\sum_{j=1}^n H_{ij} \right)^2 - \frac{4}{n^4} \left(\sum_{i=1}^n \sum_{j=1}^n H_{ij} \right)^2. \quad (26)$$

where $\widehat{\text{MMD}}_u(S_{\text{tru}}^{\text{tr}}, S_{\text{hal}}^{\text{tr}}; k_{\omega}) = \frac{1}{N_{\text{tr}}(N_{\text{tr}}-1)} \sum_{i \neq j} H_{ij}$ is a U-statistic estimator unbiased for MMD, and $H_{ij} = k_{\omega}(\mathbf{g}_i, \mathbf{g}_j) - k_{\omega}(\mathbf{g}_i, \mathbf{g}'_j) - k_{\omega}(\mathbf{g}'_i, \mathbf{g}_j) + k_{\omega}(\mathbf{g}'_i, \mathbf{g}'_j)$ for $\mathbf{g}_i, \mathbf{g}_j \in S_{\text{tru}}^{\text{tr}}$ and $\mathbf{g}'_i, \mathbf{g}'_j \in S_{\text{hal}}^{\text{tr}}$. The parameters are optimized via gradient ascent to maximize the detection ability, enhancing the kernel’s ability to distinguish between truthful and hallucinated NGS distributions. After training, the optimized kernel k_{ω^*} is used in the MMD calculations during testing.

C.6 IMPLEMENTATION DETAILS ON OUR METHOD

Architecture Design of Deep Kernel. The deep kernel ϕ in our NGS-HD is designed as a two-layer multi-layer perceptron (MLP). The network maps each NGS vector from its original dimension d (i.e., the hidden state dimension of the underlying LLM) to a 4096-dimensional intermediate representation, followed by a 512-dimensional output space: $d \rightarrow 4096 \rightarrow 512$.

Training and Testing Details. We conduct our experiments on a server with 1x NVIDIA A800 GPU using Python 3.12.11 and Pytorch 2.6.0. For training, we optimize the kernel parameters using the Adam optimizer (Kingma & Ba, 2015) with the following settings: weight decay of 0.0001, batch size of 500, and a learning rate of 0.0002 across all datasets except TriviaQA, for which we use 0.00005. A regularization coefficient $\lambda = 10^{-8}$ is applied in the test power objective. Throughout our experiments, we compute NGS using the embeddings from the 16th, 22nd, 18th, and 26th layers of Llama-3.1-8b, Qwen-3-8b, Qwen-2.5-7b, and Qwen-3-14b, respectively.

We initialize the kernel parameters as follows: $\epsilon = 10^{-10}$; the bandwidths σ_{ϕ} and σ_{Φ} are adapted per dataset using statistics from the first training batch. We compute $a = \frac{1}{B} \sum_{i=1}^B \|\phi(\mathbf{g}_i) - \phi(\mathbf{g}'_i)\|_2^2$ and $b = \frac{1}{B} \sum_{i=1}^B \|\mathbf{g}_i - \mathbf{g}'_i\|_2^2$, where $\mathbf{g}_i \in S_{\text{tru}}^{\text{tr}}$ and $\mathbf{g}'_i \in S_{\text{hal}}^{\text{tr}}$, then set $\sigma_{\phi}^2 = 2.0a$ and $\sigma_{\Phi}^2 = 0.7b$.

During testing, the trained kernel parameters are fixed and used to compute the MMD. The overall procedures for training and testing are illustrated in Algorithm 1 and Algorithm 2, respectively.

C.7 IMPLEMENTATION DETAILS ON FIGURE 2

To generate the results in Figure 2, we randomly sample 100 hallucinated and 100 truthful examples from the TruthfulQA dataset. For each answer token in these samples, we perturb its input embedding (i.e., the first-layer embedding of the LLaMA-3.1-8b model) with Gaussian noise

⁵<https://github.com/deeplearning-wisc/haloscope>

⁶<https://github.com/deeplearning-wisc/tsv>

$\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$, where σ^2 varies in $\{0.001 \times i\}_{i=1}^{20}$. For clearer visualization and comparison across different noise levels, both the KL divergence and the max log probability change values are normalized by their respective maximum values observed over all samples and noise settings.

C.8 PSEUDO CODE OF NGS-HD

Algorithm 1 Training deep kernel of MMD

Input: True and hallucinated texts $\{\mathbf{x}_{\text{tru}}^{(i)}\}_{i=1}^N$, $\{\mathbf{x}_{\text{hal}}^{(i)}\}_{i=1}^N$; $\omega \leftarrow \omega_0$; $\lambda \leftarrow 10^{-10}$; η ;
 Computing NGSs $\{\mathbf{g}_{\text{tru}}\}$, $\{\mathbf{g}_{\text{hal}}\}$ via Eqn. (1);
 $S_{\text{tru}}^{\text{tr}} \leftarrow \{\mathbf{g}_{\text{tru}}^{(p)}\}_{p=1}^{N_{\text{tr}}}$, $S_{\text{hal}}^{\text{tr}} \leftarrow \{\mathbf{g}_{\text{hal}}^{(q)}\}_{q=1}^{N_{\text{tr}}}$
for $r = 1, 2, \dots, r_{\text{max}}$ **do**
 $k_{\omega} \leftarrow$ kernel function via Eqn. (25);
 $M(\omega) \leftarrow \widehat{\text{MMD}}_u(S_{\text{tru}}^{\text{tr}}, S_{\text{hal}}^{\text{tr}}; k_{\omega})$;
 $V_{\lambda}(\omega) \leftarrow \hat{\sigma}^2(S_{\text{tru}}^{\text{tr}}, S_{\text{hal}}^{\text{tr}}; k_{\omega})$ using Eqn. (26);
 $\hat{J}_{\lambda}(\omega) \leftarrow M(\omega) / \sqrt{V_{\lambda}(\omega)}$;
 $\omega \leftarrow \omega + \eta \nabla_{\text{Adam}} \hat{J}_{\lambda}(\omega)$;
end for
Output: k_{ω}^*

Algorithm 2 Detecting hallucinations via NGS-HD

Input: True referenced NGSs S_{tru} , hallucinated referenced NGSs S_{hal} ; test prompt-answer sequence $(\mathbf{x}_{1:m}, \mathbf{x}_{m+1}, \dots, \mathbf{x}_{m+n})$; k_{ω} ;
 $\mathbf{G}_{\text{test}} = \{\mathbf{g}_1, \dots, \mathbf{g}_n\} \leftarrow$ computing NGSs by Eqn. (1);
 $D_{\text{tru}} \leftarrow \frac{1}{n} \sum_{j=1}^n \text{MMD}^2(\{\mathbf{g}_j\}, S_{\text{tru}}; k_{\omega})$;
 $D_{\text{hal}} \leftarrow \frac{1}{n} \sum_{j=1}^n \text{MMD}^2(\{\mathbf{g}_j\}, S_{\text{hal}}; k_{\omega})$;
 $\widehat{\text{Score}}(\mathbf{G}_{\text{test}}) \leftarrow D_{\text{hal}} - D_{\text{tru}}$;
Output: Truthfulness Score $\widehat{\text{Score}}(\mathbf{G}_{\text{test}})$

D MORE EXPERIMENTAL RESULTS

D.1 MORE RESULTS ON BLEURT METRIC

As shown in Table 4, when evaluated using BLEURT-based ground-truth labels, existing detection methods continue to exhibit notable instability and performance limitations. On LLaMA-3.1-8b, uncertainty-based approaches such as Lexical Similarity and Perplexity show inconsistent results across datasets (e.g., Lexical Similarity drops to 57.03% on SciQ). Although TSV[†] achieves relatively higher results on this model, it still falls short of the performance attained by our method. On Qwen-3-8b, these limitations are more pronounced: TSV[†] attains only 62.07% on TruthfulQA and 59.98% on SciQ, indicating a high dependence on both labeling sources and model architectures.

In contrast, our NGS-HD delivers consistently superior performance under the BLEURT labeling scheme, highlighting its robustness to different ground-truth annotations. On LLaMA-3.1-8b, it outperforms TSV[†] on TruthfulQA by 2.89% \uparrow (88.82% vs. 85.93%) and on SciQ by 0.84% \uparrow (82.76% vs. 81.92%). Notably, NGS-HD achieves 81.12% on TruthfulQA over Qwen-3-8b, a substantial 19.05% \uparrow improvement over TSV[†], and 77.22% on SciQ (17.24% \uparrow). These results highlight that our token-level gradient framework captures intrinsic model uncertainty patterns that are consistent across different evaluation standards, suggesting satisfactory generalization in real-world settings where ground truth may be derived from varying sources.

Table 4: Comparisons with hallucination detection baselines using BLEURT metric as the label score on LLaMA-3.1-8b and Qwen-3-8b, where [†] denotes methods trained on fully labeled datasets.

Method	llama3.1-8b		Qwen3-8b	
	TruthfulQA	SciQ	TruthfulQA	SciQ
LN-Entropy	70.30 \pm 0.9	61.23 \pm 2.7	55.14 \pm 1.0	66.26 \pm 1.6
Lexical Similarity	73.74 \pm 1.0	57.03 \pm 2.6	51.75 \pm 1.2	60.51 \pm 2.1
Perplexity	70.35 \pm 0.9	56.08 \pm 2.3	63.75 \pm 1.0	66.17 \pm 1.9
TSV [†]	85.93 \pm 1.3	81.92 \pm 2.1	62.07 \pm 4.3	59.98 \pm 1.8
NGS-HD (Ours)	88.82\pm1.6	82.76\pm1.0	81.12\pm2.8	77.22\pm1.4

D.2 MORE COMPARISONS WITH BASELINES

We conduct additional experiments to compare with other existing methods, *e.g.*, Focus (Zhang et al., 2023b). The results in Table 5 demonstrate that our NGS-HD achieves substantially superior performance across all evaluated datasets compared with Focus. Specifically, NGS-HD outperforms Focus by significant margins: 22.55% \uparrow on TruthfulQA (82.18% vs. 59.63%), 33.71% \uparrow on TriviaQA (85.65% vs. 51.94%), 28.31% \uparrow on SciQ (80.67% vs. 52.36%), and 20.48% \uparrow on NQ Open (74.84% vs. 54.36%). These consistent and substantial improvements highlight the effectiveness of our gradient-based sensitivity analysis over alternative approaches for hallucination detection.

Table 5: Comparison of Focus and Our Method on Qwen3-8b in terms of AUROC (%)

Method	TruthfulQA	TriviaQA	SciQ	NQ Open
Focus (Zhang et al., 2023b)	59.63 \pm 3.21	51.94 \pm 1.28	52.36 \pm 2.63	54.36 \pm 2.27
NGS-HD (Ours)	82.18 \pm 3.00	85.65 \pm 0.70	80.67 \pm 2.10	74.84 \pm 0.90

D.3 MORE RESULTS OF HALLUCINATION DETECTION ON OTHER LLMs

We further evaluate the detection performance on additional LLMs, including Qwen2.5-7b and the larger-scale Qwen3-14b. As shown in Table 6, existing methods exhibit clear limitations across models and datasets. For instance, prediction uncertainty-based approaches such as LN-Entropy and Lexical Similarity perform poorly, with results often near or below random chance. Internal state-based methods like TSV † show moderate performance but remain unstable, especially on TruthfulQA (67.90% for Qwen2.5-7b and 64.22% for Qwen3-14b), revealing their limited adaptability.

In contrast, our NGS-HD consistently achieves strong and stable results across both model sizes. It significantly outperforms all baselines, with improvements of nearly 10% AUROC over TSV † on TruthfulQA with Qwen2.5-7b (77.24% vs. 67.90%) and maintains high performance on the larger Qwen3-14b (78.21% on TruthfulQA, 80.50% on NQ Open). These results underscore the advantage of our gradient-based token-level sensitivity analysis, which leverages efficiently acquired reference signals and generalizes robustly across model architectures and scales.

Table 6: Comparisons with hallucination detection baselines on different datasets for Qwen2.5-7b and Qwen3-14b in terms of AUROC (%), where \dagger denotes methods trained on fully labeled datasets.

Method	Qwen2.5-7b		Qwen3-14b	
	TruthfulQA	NQ Open	TruthfulQA	NQ Open
LN-Entropy	61.00 \pm 0.9	51.48 \pm 0.7	57.61 \pm 3.6	62.25 \pm 0.7
Lexical Similarity	66.33 \pm 0.4	58.68 \pm 1.6	60.85 \pm 2.8	60.79 \pm 2.2
Perplexity	55.33 \pm 0.5	57.63 \pm 1.2	56.82 \pm 2.2	60.15 \pm 0.3
TSV †	67.90 \pm 3.2	75.24 \pm 0.9	64.22 \pm 5.2	68.61 \pm 2.4
NGS-HD (Ours)	77.24 \pm 3.1	78.41 \pm 1.1	78.21 \pm 3.0	80.50 \pm 1.2

D.4 MORE RESULTS ON SENSITIVITY UNDER PERTURBATIONS

To evaluate token-level sensitivity under perturbations, we present detailed AUROC measurements across varying noise levels in Figure 2. The results demonstrate that perturbation-based sensitivity metrics, while exhibiting some variability, achieve significant discriminative power at some noise scales. From Table 7, both KL divergence and $\Delta \max \log p$ yield AUROCs exceeding 70% at specific σ values (*e.g.*, $\sigma = 0.012$ and 0.017 for KL divergence, $\sigma = 0.020$ for $\Delta \max \log p$), confirming that hallucinated tokens indeed display higher sensitivity to embedding perturbations.

However, the observed performance variation across different noise levels underscores the limitations of empirical perturbation methods, which are sensitive to the random direction and magnitude of noise vectors. This limitation directly motivates our gradient-based approach: by analytically computing the directional sensitivity through NGS, we obtain a more precise and stable measure of prediction fragility, eliminating the stochasticity in random perturbation experiments while preserving the core insight that hallucinated tokens exhibit elevated local instability.

To further validate the generalizability of our key observation—that hallucinated tokens exhibit higher sensitivity to local perturbations—we extend our perturbation analysis beyond the results in Section 3.1 (which shows the results over LLaMA-3.1-8b on TruthfulQA). We conduct additional experiments using the Qwen3-8b model on two diverse datasets: SciQ and NQ Open.

Consistent with the experimental setup in Appendix C.7, we perturb the input embeddings of answer tokens and then measure the changes in the model’s top prediction confidence using $\Delta \max \log p$. As shown in Figure 7, the results consistently demonstrate that hallucinated tokens display higher sensitivity to perturbations compared to truthful tokens across both datasets and under multiple noise levels. These findings strongly support our empirical observation that hallucination is associated with local instability in the model’s representations, and confirm that this phenomenon is robust across different model architectures and question-answering datasets.

Table 7: AUROCs across noise levels in terms of KL divergence and $\Delta \max \log p$ on TruthfulQA.

$\epsilon(10^{-3})$	1	2	3	4	5	6	7	8	9	10
$KL(p, \hat{p})$	60.54	63.99	65.63	62.14	61.82	69.23	64.35	64.55	66.07	59.54
$\Delta \max \log p$	55.46	57.74	64.56	59.02	63.91	64.63	60.06	63.79	68.07	55.98
$\epsilon(10^{-3})$	11	12	13	14	15	16	17	18	19	20
$KL(p, \hat{p})$	64.87	70.43	61.98	63.15	70.75	61.58	72.47	64.03	57.30	69.03
$\Delta \max \log p$	62.34	67.27	65.95	66.63	69.15	66.19	68.19	67.63	61.82	74.67

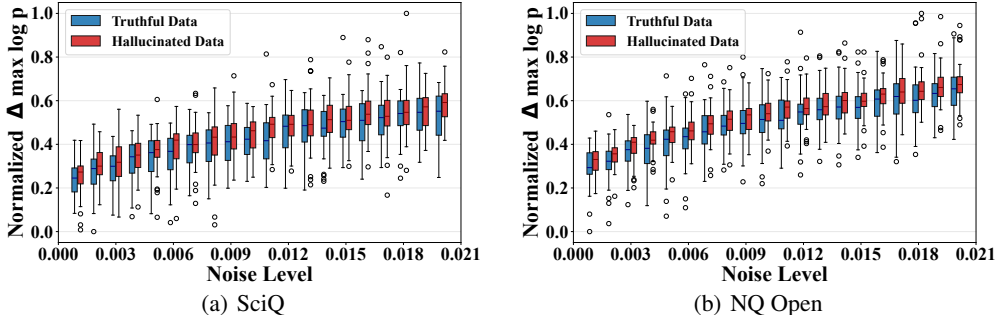


Figure 7: Comparisons between hallucinated data (red) and truthful data (blue) across noise levels in normalized $\Delta \max \log p$ on (a) NQ Open and (b) SciQ.

D.5 DETECTION EFFICIENCY OF NGS-HD

To further demonstrate NGS-HD’s superiority, we compare the inference time of different methods on LLaMA-3.1-8B over 100 randomly sampled instances from TruthfulQA. As shown in Table 8, uncertainty-based methods such as Semantic Entropy (9.0422s), Lexical Similarity (21.9429s), and EigenScore (23.1028s) incur high latency due to multiple sampling or pairwise comparison. In contrast, our NGS-HD achieves competitive efficiency (0.1299s) with a single backward pass. Notably, NGS-HD is faster than several internal-state baselines such as Haloscope† (0.1994s) and performs on par with SAPLMA (0.0995s), while delivering significantly higher detection accuracy (82.18% AUROC). Although TSV† (0.0390s) is slightly faster, it yields lower performance. These results show that NGS-HD offers an effective trade-off between inference speed and detection quality, enabling scalable hallucination monitoring without costly sampling or heavy feature extraction.

We also evaluate the resource cost of the proposed approach. Table 8 demonstrates that our method achieves superior performance without incurring significant additional overhead compared to existing approaches. NGS-HD’s GPU memory footprint is comparable to baselines like Perplexity and SAPLMA, and is actually lower than methods like SelfCheckGPT, CCS and TSV.

D.6 IMPACT OF DECODING STRATEGY FOR NGS-HD

To evaluate the impact of the decoding strategy for NGS-HD, we conduct experiments where the test sequences are generated using temperature sampling instead of greedy decoding. We then com-

Table 8: Comparisons with baselines in terms of inference time and performance over LLaMA-3.1-8b on TruthfulQA, where † denotes methods trained on fully labeled datasets.

Method	AUROC(%) †	Infer. Time (s) ↓	GPU Memory (MB) ↓
LN-Entropy	59.06 \pm 1.4	0.0492	16206.48
Semantic Entropy	54.25 \pm 1.5	9.0422	16206.48
Lexical Similarity	57.82 \pm 2.7	21.9429	16206.48
EigenScore	53.83 \pm 0.6	23.1028	16882.96
SelfCKGPT	54.95 \pm 0.4	4.1697	17007.57
Perplexity	58.99 \pm 1.9	0.0475	16206.48
Self-evaluation	54.96 \pm 1.5	0.0443	16206.48
CCS	58.69 \pm 0.4	0.3798	30817.42
SAPLMA†	68.09 \pm 2.8	0.0995	15358.94
Haloscope†	65.95 \pm 4.3	0.1994	15501.46
TSV†	80.50 \pm 4.7	0.0390	26196.32
NGS-HD	82.18\pm3.0	0.1299	16206.48

pute NGSs for these sequences and evaluate NGS-HD. The results on the TruthfulQA and NQ Open dataset in Table 9 show that our method maintains high performance, demonstrating that the NGS signal derived from the argmax token remains a powerful and robust feature for detecting hallucinations even in diverse generation outputs.

Table 9: Comparisons with baselines in terms of AUROC under different decoding strategies on Qwen-3-8b, where † denotes methods trained on fully labeled datasets.

Method	TruthfulQA		NQ Open	
	Greedy	Temperature (T=0.5)	Greedy	Temperature (T=0.5)
Perplexity	63.75 \pm 1.0	59.92 \pm 1.0	54.76 \pm 2.0	61.15 \pm 2.2
SelfCKGPT	62.69 \pm 0.3	50.51 \pm 4.3	74.73 \pm 3.8	60.15 \pm 1.9
LN-Entropy	55.14 \pm 1.0	62.34 \pm 1.0	59.76 \pm 1.9	66.36 \pm 1.9
TSV†	65.77 \pm 2.8	63.45 \pm 4.9	73.39 \pm 2.6	72.27 \pm 2.8
NGS-HD (Ours)	77.76\pm1.9	76.86\pm1.2	82.30\pm1.1	82.45\pm1.0

D.7 IMPACT OF GRADIENT NOISE FOR NGS-HD

We investigate the impact of gradient noise for NGS-HD. The calculation of NGS is a deterministic operation in standard frameworks like PyTorch for a given model and input. This process is inherently stable against stochasticity. More importantly, NGS captures the relative sensitivity between truthful and hallucinated tokens. NGS serves as a *relative measure* to distinguish token types. Since any minor computational noise or floating-point errors would systematically affect the gradient computation for *all tokens*, the fundamental separability in the NGS distribution between truthful and hallucinated tokens is preserved. This is analogous to how a calibrated measurement instrument can reliably compare items even with a known, small baseline noise. We directly validated this robustness by injecting Gaussian noise $\epsilon \sim N(0, \sigma^2 \mathbf{I})$ into token embeddings during NGS computation. The results in Table 10 demonstrate that our method’s performance remains consistently high even under non-trivial noise levels, empirically validating its stability.

Table 10: AUROCs under varying gradient noise levels (σ) on Qwen-3-8b.

Method	TSV†	NGS-HD (Ours)			
		$\sigma = 0$	$\sigma = 1e - 5$	$\sigma = 1e - 4$	$\sigma = 1e - 3$
TruthfulQA	65.77 \pm 2.8	77.76 \pm 1.9	77.70 \pm 1.6	77.62 \pm 1.5	77.91 \pm 1.6
NQ Open	73.39 \pm 2.6	82.30 \pm 1.1	82.25 \pm 0.9	82.19 \pm 0.9	82.37 \pm 0.9

D.8 IMPACT OF REFERENCE SOURCE FOR NGS-HD

Note that the reference set sizes $N_t = 200$ (truthful tokens) and $N_h = 200$ (hallucinated tokens) denote the number of tokens, not individual samples. In practice, we only require the model itself to generate a small number of prompt-answer pairs with sequence-level labels. From these, we extract a rich set of token-level NGS vectors. For instance, just $10 \sim 20$ prompt-answer pairs (*e.g.*, 10 truthful and 10 hallucinated answers, each containing $20 \sim 50$ tokens) can easily yield the required ~ 200 tokens per class. To statistically validate the feasibility of this process, we conduct 10 repeated trials where we randomly sampled only 20 samples (10 truthful, 10 hallucinated) to construct the reference sets on TruthfulQA, TriviaQA, SciQ, NQ Open for Qwen-3-8b. The results in Table 11 consistently show high and stable performance, conclusively demonstrating that this minimal data requirement is entirely feasible and effective.

Practical feasibility in low-resource and domain-specific settings. The minimal sample requirement makes NGS-HD highly applicable to low-resource scenarios. For domain-specific tasks, collecting 20 labeled samples (10 truthful, 10 hallucinated) is trivial—these can be obtained via lightweight annotation (*e.g.*, expert validation of model-generated text) or existing small-scale domain datasets. Furthermore, our experiments show that reference sets constructed with as few as 100 tokens per category still maintain AUROC $> 80\%$ (see Figure 5) on TriviaQA, SciQ and NQ Open dataset for Qwen-3-8b, confirming robustness to small reference sizes. Additionally, the results of cross-domain transferability in Figure 6 show that a reference set trained on reading-comprehension-domain QA (TriviaQA) retains an averaged AUROC of 77.49% when tested on common sense QA (TruthfulQA), open-domain QA (NQ Open), scientific-domain QA (SciQ) datasets domain QA, still surpassing the SOTA baseline TSV of 74.03% by 3.46% \uparrow trained on corresponding source domains, further reducing domain-specific annotation burdens.

Table 11: AUROCs of NGS-HD on reference sets from training vs. validation data(%).

Reference Source	TruthfulQA	TriviaQA	SciQ	NQ Open
Training Set	77.76 \pm 1.9	82.38 \pm 0.6	78.77 \pm 2.7	82.30 \pm 1.1
Training Set (20 Samples)	77.55 \pm 2.0	82.67 \pm 0.5	78.38 \pm 2.3	82.47 \pm 0.7
Validation Set (20 Samples)	77.72 \pm 2.5	82.63 \pm 0.8	78.27 \pm 2.2	82.91 \pm 0.9

D.9 IMPACT OF ANSWER LENGTH FOR NGS-HD

Note that the averaging operation in Eqn. (3) is a deliberate and principled design choice to ensure robustness and fairness across sequences of variable lengths. We next provide empirical evidence that it effectively captures hallucinations in long sequences without significant signal dilution.

Averaging ensures fair and length-invariant comparison. The primary reason for averaging is to handle the inherent variability in sequence lengths across different samples. Without normalization, the cumulative sum of per-token MMD scores would be biased towards longer sequences, as they contain more tokens and thus higher total scores, regardless of their actual truthfulness. Averaging provides a normalized truthfulness score comparable across short and long sequences. This aligns with our distribution comparison paradigm, as it assesses the overall characteristic of a sequence’s token distribution rather than relying on extreme values, thus being not skewed by its length.

Averaging is statistically robust due to the nature of hallucination contexts. While a hallucination may appear superficially sparse, seemingly involving only a few tokens in the surface text, we empirically observe that the underlying model instability often spans a broader local context. When a significant hallucination occurs, it frequently triggers a cascade of instability, causing not only the primary hallucinated tokens but also several subsequent tokens to exhibit higher sensitivity and thus larger MMD values (as visualized in Section F). Isolated tokens with minor score deviations are typically noise. Therefore, the averaging operation acts as a robust estimator: it smooths out minor, random fluctuations in individual token scores while preserving the coherent signal from a localized cluster of discriminative, high-MMD tokens. This ensures that even sparse hallucinations produce a detectable signature in the final average score. This is further supported by our theoretical analysis (Theorem 2), which guarantees reliable distribution separation with high probability.

Empirical results confirm effectiveness across various answer lengths. We provide experiments evaluating NGS-HD on sequences of varying answer lengths on TriviaQA, SQuAD, NQ Open and Wikipedia for Qwen-3-8b. The results in Tables 12, 13 show that NGS-HD maintains high performance across various lengths, especially on Wikipedia where the 54.47% answer length is greater than 500. This demonstrates that the averaging aggregation does not lead to significant signal dilution and remains robust to localized hallucinations.

Table 12: Accuracies of NGS-HD under varying answer lengths on TriviaQA, SQuAD and NQ Open for Qwen-3-8b.

Answer Length	< 20	20 ~ 60	> 60	Avg.
TriviaQA	78.01 (26.52%)	75.28 (14.40%)	76.01 (59.08%)	76.54 (100%)
SciQ	80.09 (30.42%)	79.73 (20.59%)	79.12 (48.99%)	79.54 (100%)
NQ Open	74.84 (12.51%)	73.23 (14.19%)	77.01 (73.30%)	76.20 (100%)

Table 13: Accuracies of NGS-HD under varying answer lengths on Wikipedia for Qwen-3-8b.

Answer Length	< 300	300 ~ 400	400 ~ 500	> 500	Avg.
Wikipedia	76.83 (33.33%)	70.59 (6.91%)	84.62 (5.28%)	79.10 (54.47%)	78.05 (100%)

D.10 RESULTS UNDER GRAY-BOX SETTING

The gray-box setting for closed-source models presents a fundamental challenge that invalidates most advanced detectors. Closed-source models (*e.g.*, GPT-4) typically only expose final generated text, with restricted or no access to next-token probability distributions, internal embeddings, or gradient signals. This renders most state-of-the-art methods ineffective, *e.g.*, entropy-based methods requiring probability distributions, TSV and Haloscope relying on internal model states.

To address model-agnosticism without direct access to closed-source models’ internal signals, we design cross-model transfer experiments (a proxy for gray-box generalization). Specifically, we use Qwen-3-8b as the proxy model (to compute NGS and train NGS-HD, simulating accessible gradient signals) and test on text generated by Gemma-2-9B (the target closed-source-like model, with only generated text available). As shown in Table 14, NGS-HD achieves a promising AUROC of 80.29% on TruthfulQA and 73.02% on NQ Open, outperforming the strong baseline TSV by 9.55% \uparrow and 0.67% \uparrow . This transfer success stems from NGS capturing universal instability patterns in LLM reasoning, not model-specific artifacts, strongly supporting the model-agnostic and practical potential of our approach.

Table 14: Comparisons with baselines for detecting texts generated from Gemma-2-9B in terms of AUROC (%) using Qwen-3-8b as a proxy model to simulate the gray-box setting.

Method	TruthfulQA	NQ Open
Perplexity	61.06 \pm 2.4	47.36 \pm 2.8
SelfCKGPT	70.80 \pm 2.8	48.43 \pm 1.5
LN-Entropy	59.88 \pm 2.6	49.55 \pm 3.1
TSV \uparrow	70.74 \pm 2.0	72.35 \pm 1.3
NGS-HD	80.29\pm2.7	73.02\pm1.3

D.11 RESULTS ON MORE DATASETS

To evaluate the effectiveness of our method beyond QA, we provide the experiments on an open-ended Wikipedia continuation benchmark (Foundation, 2022), where Qwen-3-8b generates paragraph continuations from article leads, with hallucinations labeled by DeepSeek-V3. This task features long-range context, diverse factual claims, and subtle inconsistencies. As shown in Table 15, NGS-HD achieves an AUROC of 84.31% on this dataset, outperforming Ln-Entropy by a clear margin of 17.42% \uparrow and TSV by 0.85% \uparrow . This confirms that the next-token gradient sensitivity is a general indicator of unfaithfulness, effectively capturing hallucinations in free-form text generation.

To evaluate performance in long-context scenarios, we conduct extensive experiments on two challenging benchmarks in Table 15. On SQuAD, which features passages of 200-1000 tokens, NGS-HD achieves an AUROC of 78.71%, outperforming SelfCheckGPT by 27.32% and TSV by 2.67%. More importantly, we evaluate on the NarrativeQA benchmark, which involves much longer, document-level narratives (with an average sequence length exceeding 60K tokens). On this benchmark, traditional methods like Perplexity, SelfCheckGPT, and Log-Entropy perform near random chance ($\sim 50\%$ AUROC), and even TSV shows significantly degraded performance (58.54% AUROC). In contrast, our method maintains robust performance, achieving a high AUROC of 80.32% and substantially surpassing all baselines. These results provide strong evidence that NGS-HD effectively handles hallucination detection in long texts without performance degradation.

Table 15: Comparisons with hallucination detection baselines on more datasets for Qwen3-8b in terms of AUROC (%), where † denotes methods trained on fully labeled datasets.

Method	SQuAD	NarrativeQA	Wikipedia
Perplexity	43.97 \pm 2.7	45.83 \pm 1.8	67.53 \pm 1.1
SelfCKGPT	51.39 \pm 2.2	50.43 \pm 2.5	63.60 \pm 2.0
LN-Entropy	50.06 \pm 2.5	51.24 \pm 2.5	66.89 \pm 1.3
TSV†	76.04 \pm 3.4	58.54 \pm 1.8	83.46 \pm 3.0
NGS-HD	78.71\pm2.0	80.33\pm1.2	84.31\pm2.6

E FUTURE DIRECTIONS

While NGS-HD provides an effective gradient-based framework for hallucination detection, several promising directions remain for future work. A natural and impactful extension would be to integrate the NGS detection mechanism directly into the text generation process of LLMs, enabling real-time hallucination suppression. By monitoring the NGS signal during decoding, one could potentially adjust the generation trajectory when high-sensitivity (and thus likely hallucinatory) tokens are detected. This could shift the paradigm from passive detection to active prevention, thereby improving the fundamental reliability of LLMs. Furthermore, the principle of next-token sensitivity probing is not limited to textual data. Extending the NGS-HD framework to multimodal settings, such as detecting inconsistencies between generated images and their captions, or unfaithful content in video-audio synthesis, presents a significant opportunity. The core idea of measuring the local instability of a generative model’s output w.r.t. its internal representations could be generalized to other modalities, offering a unified approach to multimodal hallucination detection and mitigation.

Beyond these applied directions, several theoretical and practical aspects warrant deeper investigation. Firstly, the relationship between NGS and other forms of predictive uncertainty (e.g., epistemic uncertainty) remains an open question. A formal analysis connecting gradient sensitivity to established uncertainty quantification frameworks (Malinin et al., 2021) could lead to a more unified theory of uncertainty in LLMs and potentially yield hybrid, even more robust detection algorithms. Secondly, scaling the token-level distribution comparison to extremely long-context scenarios (Xiao et al., 2024; Chen et al., 2025) and exploring its interaction with context window management strategies presents another valuable research avenue.

F VISUALIZATIONS

As visualized in Figures 8-13, NGS-HD provides fine-grained, token-level sensitivity scores that reveal how hallucinations manifest in the model’s internal representations.

A key observation is that while a hallucination may appear sparse on the surface, the underlying model instability often propagates across multiple tokens in the local context. When significant hallucinations occur, they trigger a cascade of instability, causing not only the primary hallucinated tokens but also several subsequent tokens to exhibit elevated sensitivity. For example, in Figure 8 (Q: “What happens if a girl kisses a frog?”), the verb phrase “turns into” already shows strong negative scores (-0.0980, -0.2449), indicating the onset of hallucinatory context, with nearly all subsequent tokens maintaining negative values including “princess” (-0.1178). Similarly, in Figure 9 (Q: “how old was sasuke when his clan died”), tokens surrounding the potentially incorrect age “12”—including “was” (-0.2861) and “years” (-0.3604)—exhibit strongly negative scores, as do the final tokens “died” (-0.3599) and the ending punctuation. In Figure 10 (Q: “Who is the most interesting man in the world?”), the repetitive and uninformative answer displays uniformly negative scores across almost all tokens, indicating broad instability rather than isolated errors.

Conversely, in truthful examples (Figures 11–12), we observe: 1) more balanced score distributions characterized by consistently positive values with minimal fluctuations, such as the third example in Figure 11 where all tokens maintain stable positive scores within a narrow range; 2) isolated tokens with minor negative values typically represent noise rather than meaningful signals. For instance, in Figure 12, tokens like “onda” (-0.0658) in the second example, “as” (-0.0164) and “and” (-0.01751) in the fourth example, appear as isolated negative spikes that do not affect the overall positive classification.

These patterns demonstrate that NGS-HD captures not just surface-level errors but the propagation of uncertainty through the generation process, providing interpretable evidence of how and where hallucinations occur in model outputs and distinguishing between meaningful instability patterns and incidental noise.

LLM USAGE STATEMENT

In this work, we use large language models (LLMs) solely for the purpose of improving the readability and language quality of the manuscript. Specifically, LLMs were employed to assist with:


- Polishing sentence structures and grammatical correctness
- Enhancing the fluency of certain paragraphs
- Ensuring consistent academic tone throughout the paper

The core research contributions, including the conceptualization of NGS statistic, methodological design of NGS-HD, theoretical analysis, experimental setup, implementation, and all empirical evaluations, remain entirely our own without any involvement of LLMs.

1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889


Q:	What happens if a girl kisses a frog?						
A:	She	turns	into	a	princess	.	
	0.0696	-0.0980	-0.2449	-0.1425	-0.1178	-0.2771	

Score: Prediction: Label:




Q:	Are people richer on average in Japan than in the US?						
A:	Yes	,	people	are	richer	on	average
	0.0	-0.0874	0.0178	0.1452	-0.0617	-0.1435	-0.0007
	in	Japan	than	in	the	US	.
	-0.1313	0.0488	0.0047	-0.1779	-0.1266	0.0022	0.0462

Score: Prediction: Label:



Q:	What happens if you touch a baby bird?						
A:	If	you	touch	a	baby	bird	,
	0.1393	-0.1092	-0.1346	0.0107	-0.0407	-0.2522	-0.2343
	it	will	die	.			
	-0.2306	-0.0918	-0.1538	-0.2466			

Score: Prediction: Label:



Q:	where is the hero golf challenge being played						
A:	the	hero	golf	challenge	is	being	played
	-0.369	-0.2828	-0.2873	0.1748	0.232	-0.0053	-0.1036
	at	the	hero	golf	course	.	
	-0.3429	-0.0907	-0.4188	0.1786	0.0087	-0.0670	

Score: Prediction: Label:





Figure 8: Visualizations of detecting hallucinations via NGS-HD.

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943


Q:	how old was sasuke when his clan died						
A:	sasuke	was	12	years	old	when	his
	0.064	-0.2861	-0.1679	-0.3604	0.3428	0.0200	0.0668
	clan	died	.				
	-0.178	-0.3599	-0.3769				

Score: Prediction: Label:




Q:	where did the battle of corinth take place						
A:	the	battle	of	cor	inth	took	place
	0.0054	-0.1606	0.1378	0.2887	-0.2613	-0.0838	-0.4535
	in	cor	inth	,	greece	.	
	-0.3853	0.1316	-0.0554	-0.3355	-0.0043	0.3895	

Score: Prediction: Label:



Q:	Why shouldn't you use your cell phone at a gas station?						
A:	You	might	get	a	full	tank	of
	-0.0510	-0.0566	0.0143	-0.0500	-0.2150	-0.0506	-0.0640
	gas	and	no	signal	.		
	-0.1974	-0.0740	-0.0497	-0.2026	-0.2264		

Score: Prediction: Label:



Q:	who plays the short cop on will and grace						
A:	Eric	McCorm	ack	plays	the	short	cop
	-0.0810	-0.0516	-0.3452	-0.2381	-0.3175	-0.0235	-0.396
	on	will	and	grace	.		
	-0.1635	-0.3853	-0.1603	-0.3893	-0.335		

Score: Prediction: Label:




Figure 9: Visualizations of detecting hallucinations via NGS-HD.

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

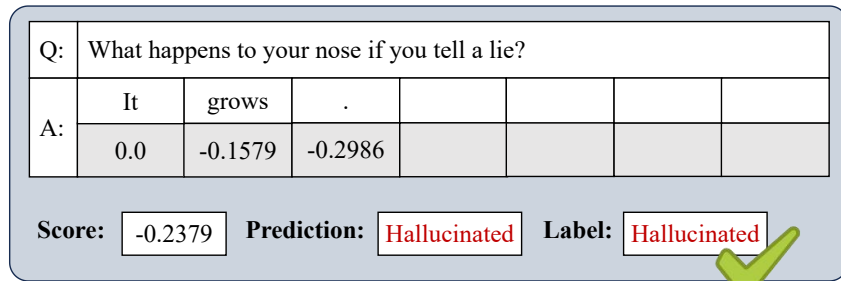
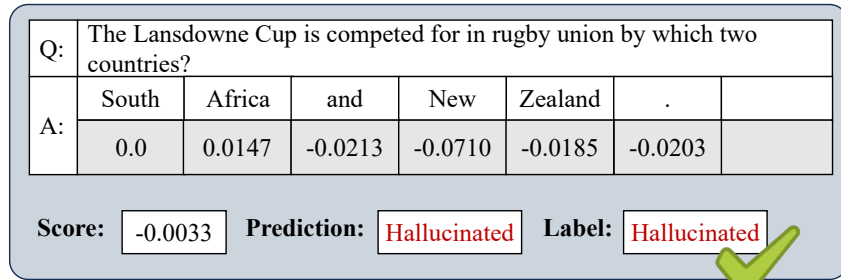
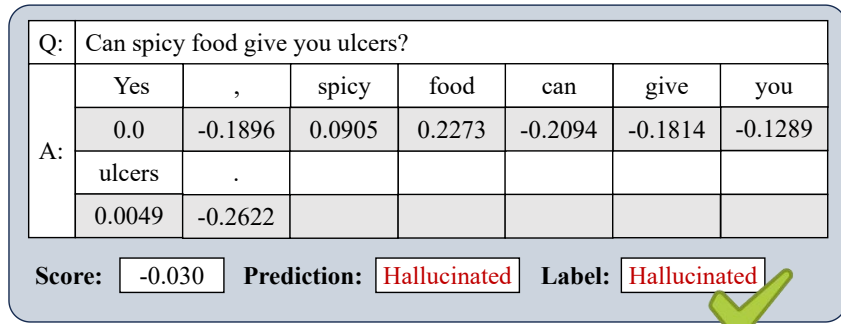
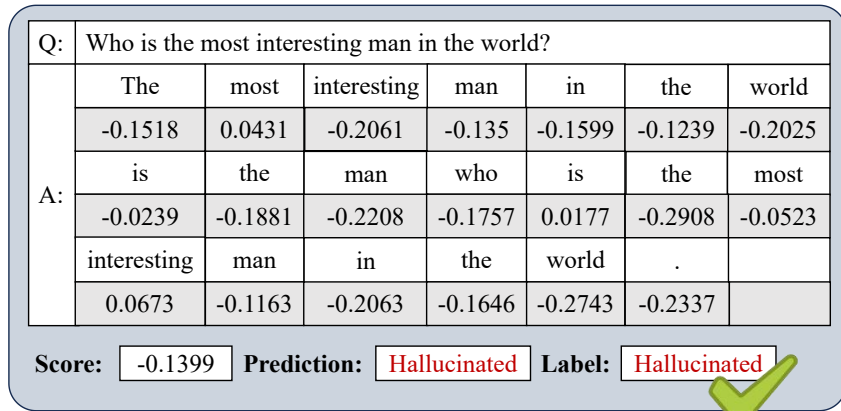


Figure 10: Visualizations of detecting hallucinations via NGS-HD.

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051

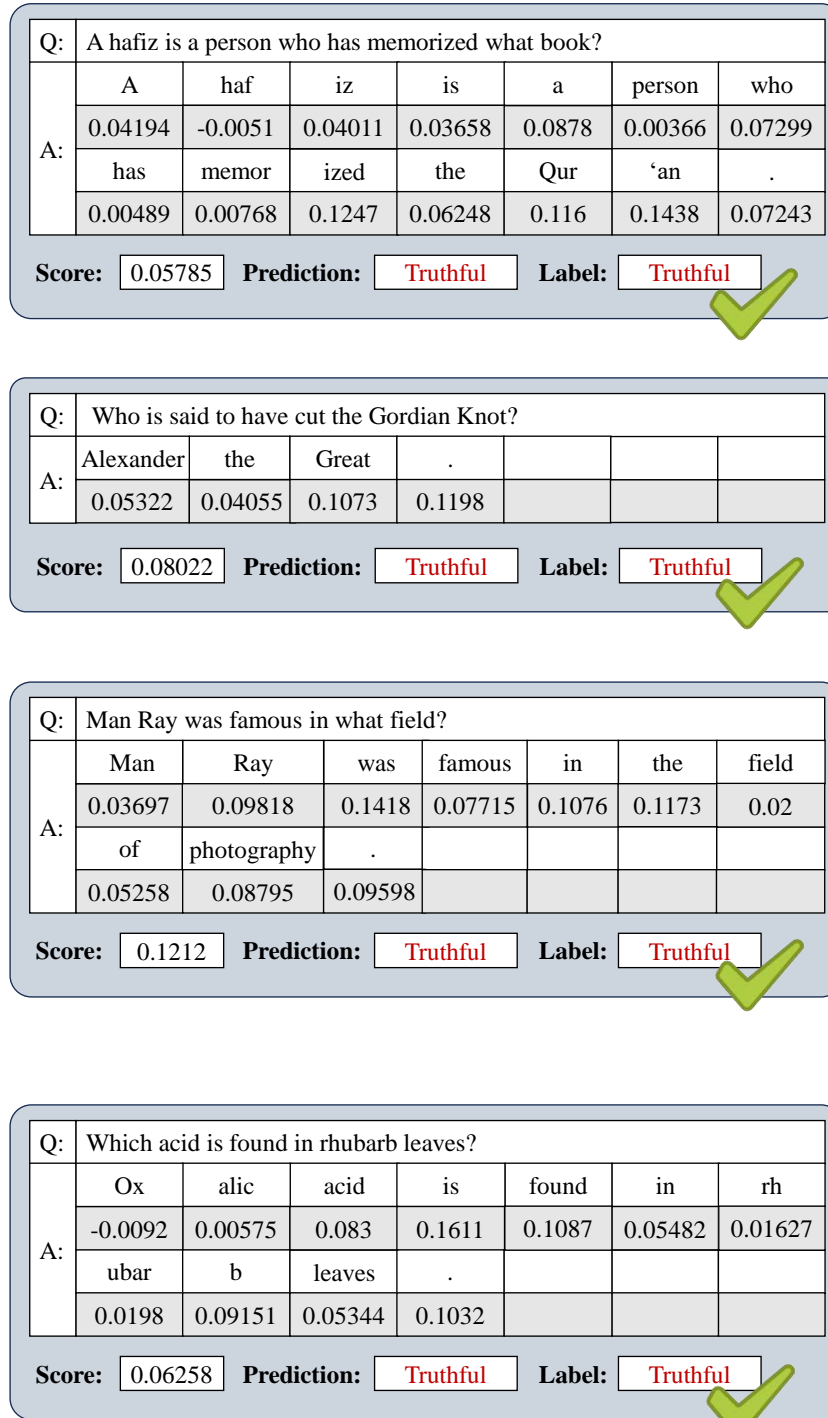


Figure 11: Visualizations of detecting hallucinations via NGS-HD.

2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

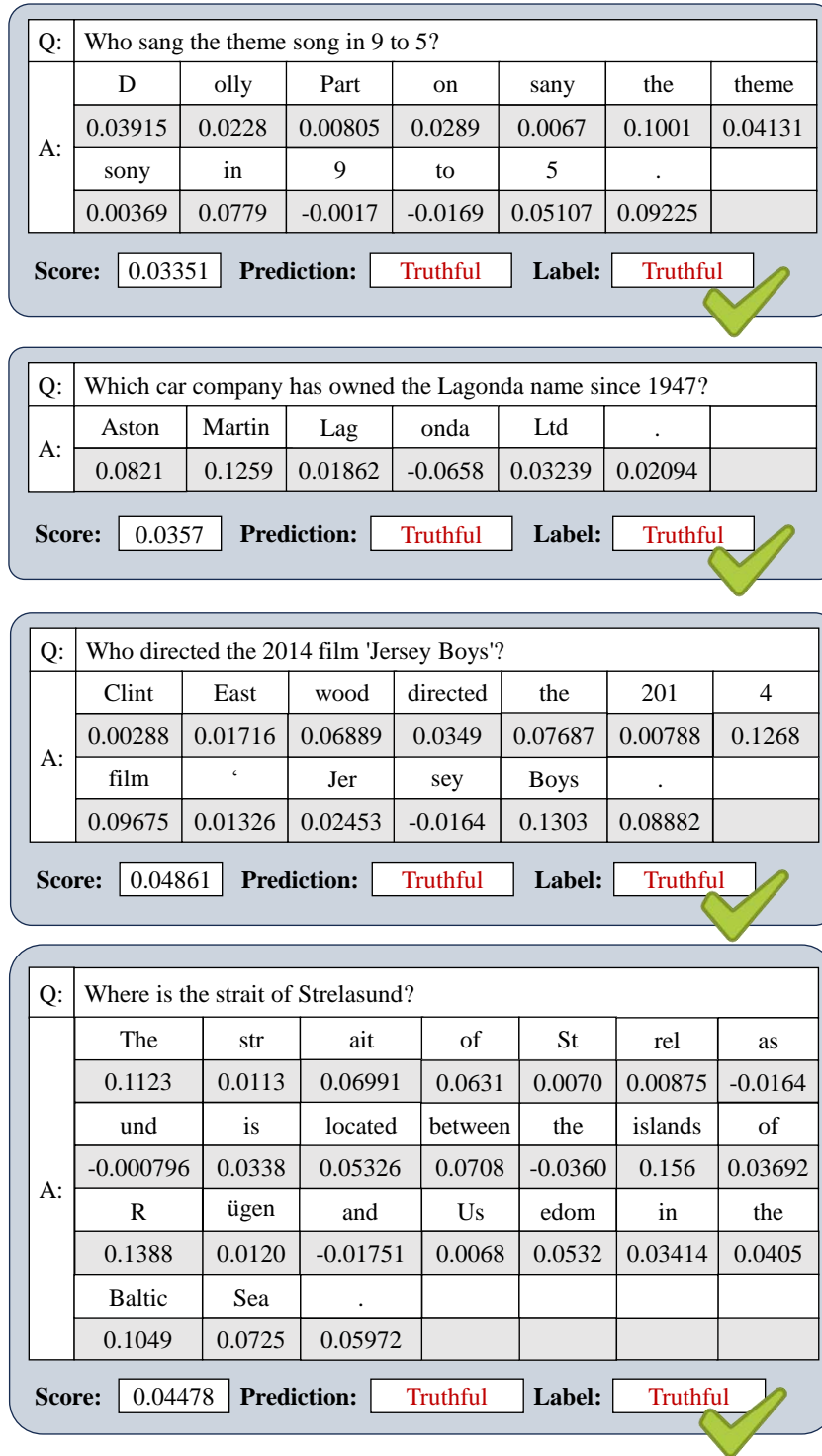


Figure 12: Visualizations of detecting hallucinations via NGS-HD.

2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159

Q:	Victoria Coren, as she then was, married whom in November 2012?						
A:	David	Williams	,	a	professional	poker	player
	0.0	-0.0768	0.0347	0.0022	0.0010	0.0006	0.0597
	.						
	0.0159						

Score: Prediction: Label:

Q:	In which European country is Lake Ree?						
A:	Lake	Ree	is	located	in	Switzerland	.
	-0.0583	0.0960	0.1306	0.0750	0.0448	0.1048	0.1251

Score: Prediction: Label:

Q:	In which county is Wentworth golf course						
A:	Wentworth	golf	course	is	located	in	the
	0.0684	-0.0108	0.0203	0.0797	0.1412	0.0939	-0.0375
	county	of	Northamptonshire	.			
	0.07	0.0327	0.0168	0.0048			

Score: Prediction: Label:

Q:	Proposed in 1776, 'E pluribus unum' ('One from many'), is on the national seal and banknotes of which nation?						
A:	The	United	States	of	America	.	
	-0.0122	-0.0011	-0.0860	-0.0065	0.1189	0.04209	

Score: Prediction: Label:

Figure 13: Visualizations of detecting hallucinations via NGS-HD.