

The Digital Dunning-Kruger Effect: Decoupling Hallucinations via Geometric Hidden-state Observation for Semantic Truthfulness

Anonymous ACL submission

Abstract

Large Language Models (LLMs) often generate overconfident yet factually incorrect hallucinations. Current detection paradigms suffer from a trade-off between the high accuracy of computationally expensive black-box methods and the inability of white-box methods to detect stubborn hallucinations. To bridge this gap, we propose **GHOST** (Geometric Hidden-state Observation for Semantic Truthfulness), an efficient white-box framework for hallucination detection in LLMs. We distinguish between two hallucination mechanisms: *confused hallucinations*, marked by internal reasoning instability, and *stubborn hallucinations*, characterized by premature layer-wise convergence. By integrating internal geometric dynamics with output probability distributions, GHOST constructs a high-dimensional feature space for non-linear truthfulness classification. Extensive evaluations on FinanceBench, RAGTruth, HaluEval, and PopQA show that GHOST outperforms white-box baselines and achieves competitive black-box performance while reducing computational overhead by over 90%, offering a robust solution for real-time detection.

1 Introduction

Large Language Models (LLMs) have achieved remarkable breakthroughs in the field of Natural Language Processing (NLP) (Zhao et al., 2025; Farquhar et al., 2024a). However, underlying these capabilities lies a critical flaw: the models occasionally generate factually incorrect, logically fallacious, or unsubstantiated statements with high confidence (Ji et al., 2023). This phenomenon, termed *hallucination* as illustrated, severely impedes the practical deployment of LLMs in high-stakes domains such as healthcare, law, and finance, and remains a pivotal challenge awaiting resolution in the field (Zhang et al., 2025b).

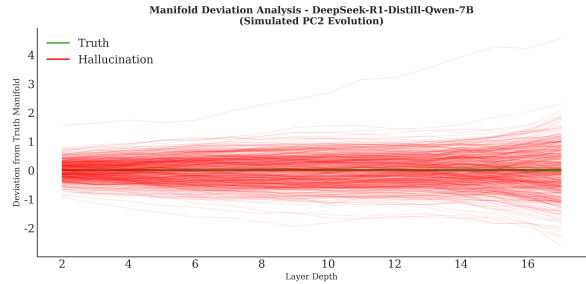


Figure 1: Manifold Deviation Analysis illustrating the Digital Dunning-Kruger Effect.

The academic community categorizes hallucination detection into black-box and white-box paradigms. Black-box methods (Manakul et al., 2023; Cohen et al., 2023; Min et al., 2023) utilize multi-sample consistency or post-hoc verification, yet prohibitive computational overhead and inference latency impede real-time deployment. Conversely, white-box approaches (Farquhar et al., 2024a; Chuang et al., 2024) leverage output logits or singular internal metrics. These coarse-grained indicators often fail to exploit rich internal dynamics, yielding limited discriminative power when models exhibit overconfidence in erroneous knowledge.

Recent investigations into internal mechanisms provide granular perspectives. While Xu et al. (Xu et al., 2020) employed V-usable information to assess model faithfulness, their analysis remains confined to final output layers and neglects dynamic transitions. Similarly, Kim et al. (Kim et al., 2025) introduced Layer-wise Information Deficiency (LI), yet this framework attributes hallucinations solely to information loss. Consequently, LI fails to identify stubborn hallucinations arising from pre-training biases or erroneous memorization.

We posit that hallucination reflects a cognitive dissonance analogous to the Dunning-Kruger Effect (Kruger and Dunning, 1999), where models

071 overestimate their competence. Specifically, we
072 identify two distinct mechanisms: *Stubborn Hallu-*
073 *cinations*, characterized by premature convergence
074 and high epistemic overconfidence despite insuffi-
075 cient factual grounding, and *Confused Hallucina-*
076 *tions*, manifested as internal reasoning instability
077 when resolving conflicting semantic signals. Ex-
078 periments have found that the reasoning process of
079 the truth is smooth, while illusions are fluctuating
080 and part of them highly overlap with the truth, as
081 shown in the figure 1, visualizes the inter-layer hid-
082 den state evolution via second principal component
083 (PC2) (Chen et al., 2024) deviation.

084 To capture these subtle cognitive signatures,
085 we propose **GHOST** (Geometric Hidden-state
086 Observation for Semantic Truthfulness). Diverg-
087 ing from previous approaches dependent on static
088 representations, GHOST conceptualizes the infer-
089 ence process as a dynamic latent trajectory within a
090 high-dimensional semantic manifold. Our primary
091 contributions are summarized as follows:

- 092 • **Cognitive-Inspired Taxonomy:** We formal-
093 ize a novel hallucination taxonomy grounded
094 in the Dunning-Kruger Effect. By distinguish-
095 ing *Confused Hallucinations* from *Stubborn*
096 *Hallucinations*, we provide a theoretical founda-
097 tion for why conventional uncertainty met-
098 rics fail during high-confidence erroneous gen-
099 eration.
- 100 • **The GHOST Framework:** We introduce a
101 white-box framework leveraging internal geo-
102 metric dynamics. By quantifying *Mental Tur-*
103 *turbulence* and *Stubbornness* across hidden lay-
104 ers, GHOST constructs a multi-dimensional
105 feature space to disentangle truthful reasoning
106 from deceptive generation without structural
107 modifications.
- 108 • **Efficiency and Generalizability:** Evalua-
109 tions on FinanceBench, RAGTruth, HaluEval,
110 and PopQA demonstrate that GHOST
111 significantly outperforms existing white-box
112 baselines. Remarkably, GHOST achieves per-
113 formance competitive with black-box meth-
114 ods while reducing computational overhead
115 by over 90%, facilitating robust real-time de-
116 ployment.

2 Related Works 117

2.1 Hallucination Detection Paradigms 118

Black-box Methods primarily assess the veracity
119 of responses through output-level consistency or
120 external verification. **SelfCheckGPT** (Manakul
121 et al., 2023) and **LM-Polygraph** (Fadeeva et al.,
122 2023) employ stochastic sampling to quantify se-
123 mantic consistency, while **FactScore** (Min et al.,
124 2023) introduces retrieval-augmented verification
125 for fine-grained factual checking. Despite their
126 high precision, these paradigms are often hindered
127 by prohibitive sampling latency and substantial
128 computational overhead, limiting their utility in
129 real-time applications. 130

White-box Methods aim to circumvent these
131 costs by leveraging model-internal signals. **Se-**
132 **matic Entropy** (Farquhar et al., 2024b) formal-
133 izes uncertainty estimation at the semantic level,
134 and **Lookback Lens** (Chuang et al., 2024) uti-
135 lizes attention maps to identify contextual hallu-
136 cinations. However, these paradigms frequently
137 rely on coarse-grained uncertainty metrics, which
138 may fail to detect "stubborn" hallucinations where
139 the model generates erroneous content with high
140 confidence. 141

2.2 Probing Internal Representations 142

A growing body of research suggests that the hid-
143 den states of LLMs inherently encode latent truth-
144 fulness. **SAPLMA** (Azaria and Mitchell, 2023)
145 and **ITI** (Li et al., 2023) demonstrate that linear
146 probes or directions within intermediate represen-
147 tations can effectively disentangle truthful behaviors
148 from deceptive ones, enabling targeted inference-
149 time interventions. **INSIDE** (Chen et al., 2024)
150 further extends this by introducing EigenScore to
151 measure consistency within the spectral domain of
152 internal states, while **PRISM** (Zhang et al., 2025a)
153 leverages prompting to align hidden states with
154 more salient truth-related manifolds. Nevertheless,
155 these approaches predominantly focus on static
156 feature distributions at specific layers, potentially
157 oversimplifying the complex dynamic evolution of
158 hidden states across the entire Transformer stack. 159

2.3 Dynamic Geometric and Topological Analysis 160

Our work aligns with the emerging trend of ana-
162 lyzing the dynamic trajectory of LLMs. **LI** (Kim
163 et al., 2025) analyzes cross-layer information dy-
164 namics and is most directly motivated by ambigu-
165

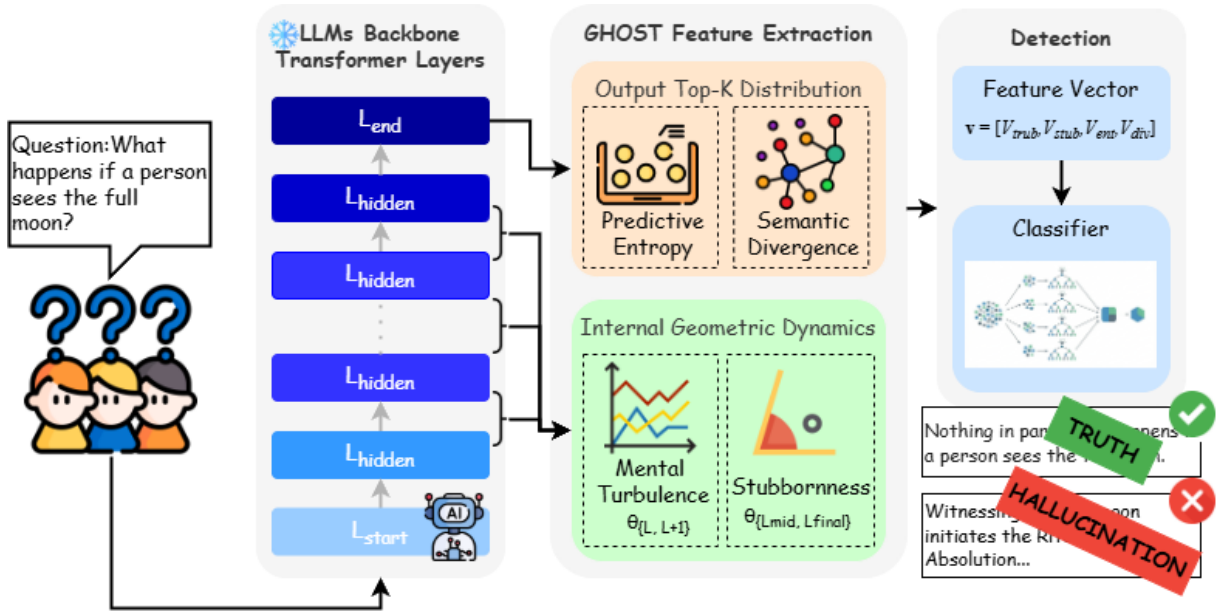


Figure 2: The schematic architecture of the proposed GHOST framework. The system integrates internal geometric trajectory features derived from hidden states with external semantic distribution features from output logits. These multi-dimensional indicators are processed by a non-linear classifier to identify confused and stubborn hallucinations.

ous prompts / unanswerable questions. **END** (Wu et al., 2025) utilizes cross-layer entropy signals to adjust decoding for factuality, yet it remains primarily focused on decoding-time mitigation. **TOHA** (Bazarova et al., 2025) explores the topological divergence of attention graphs and is restricted to RAG scenarios. **GHOST** distinguishes itself by providing a unified parametric framework that captures both transient turbulence and persistent rigidity in geometric trajectories, offering a more robust detection for both confused and stubborn hallucinations.

3 Method

Methodology. This section elucidates the high-dimensional features constituting the GHOST framework. Unlike black-box approaches that rely on external retrieval or multi-model sampling consistency (Goel et al., 2025), our method exploits internal state dynamics during generation. Accordingly, we construct a high-dimensional feature space integrating geometric dynamics and semantic distributions. This approach utilizes non-linear classifiers to capture the complex topological structures of truthfulness, as illustrated in Figure 2.

3.1 Problem Formulation

Given a Large Language Model \mathcal{M} and an input prompt x , the model generates a response sequence

$y = \{y_1, y_2, \dots, y_N\}$. Our objective is to construct a detection function $f(x, y; \mathcal{M}) \rightarrow \{0, 1\}$ that identifies whether the response y contains factual errors. Unlike previous methods that rely on the hidden state of a specific "key token" (e.g., the last token), we leverage the geometric dynamics of the entire generation process. We denote the hidden state of the l -th layer for the i -th token in the sequence as $h_l^{(i)} \in \mathbb{R}^d$, where $l \in [L_{start}, L_{end}]$ and d represents the hidden dimension.

Truthful responses consistently occupy a narrow manifold with smooth trajectories, suggesting stable semantic propagation. Hallucinations, however, demonstrate a clear bifurcation. One subset of red trajectories remains tightly clustered with the truthful paths, reflecting the "hyper-stability" of stubborn hallucinations where the model prematurely commits to incorrect priors. Another subset exhibits significant divergence starting around layer 15, representing confused hallucinations characterized by late-stage reasoning disarray. This empirical observation provides a solid geometric foundation for our proposed taxonomy of hallucination mechanisms.

3.2 Internal Geometric Dynamics

Large Language Model (LLM) inference constitutes a hidden state evolution trajectory within the inter-layer semantic space. To capture global re-

221 sponse characteristics, we employ a *Calculate-then-*
 222 *Average* aggregation strategy. Specifically, we first
 223 compute geometric metrics for each token in the
 224 sequence y , and subsequently average these values
 225 across the total sequence length N to derive the
 226 final feature vector. This approach ensures that the
 227 detector captures both transient reasoning uncer-
 228 tainties and persistent stubbornness throughout the
 229 entire generation chain.

230 3.2.1 Mental Turbulence

231 Drawing from the psychological theory of **Cogni-**
 232 **tive Dissonance** (Festinger, 1957), which describes
 233 the mental discomfort experienced when holding
 234 conflicting beliefs, we hypothesize that *Confused*
 235 *Hallucinations* arise when a model struggles to
 236 resolve contradictory internal information. This
 237 internal conflict manifests geometrically as drastic
 238 shifts in the hidden state trajectory between ad-
 239 jacent layers. We quantify this phenomenon as
 240 **Mental Turbulence**.

241 For the i -th token in the generated sequence,
 242 its turbulence score $v_{turb}^{(i)}$ measures the average
 243 cosine deviation across the selected layer range
 244 $[L_{start}, L_{end}]$:

$$245 v_{turb}^{(i)} = \frac{1}{L_{end} - L_{start}} \sum_{l=L_{start}}^{L_{end}-1} \left(1 - \frac{h_l^{(i)} \cdot h_{l+1}^{(i)}}{\|h_l^{(i)}\| \|h_{l+1}^{(i)}\|} \right) \quad (1)$$

246 The final turbulence feature for the entire response
 247 is defined as $V_{turb} = \frac{1}{N} \sum_{i=1}^N v_{turb}^{(i)}$. A significant
 248 increase in V_{turb} serves as a computational proxy
 249 for cognitive dissonance, indicating high internal
 250 conflict and instability in the model’s reasoning
 251 path.

252 3.2.2 Stubbornness

253 Conversely, to capture the "illusory superiority"
 254 akin to the **Dunning-Kruger Effect** (Kruger and
 255 Dunning, 1999), we introduce the **Stubbornness**
 256 metric. This cognitive bias in LLMs manifests
 257 as *Stubborn Hallucinations*, where the model con-
 258 verges to a conclusion prematurely, despite a lack
 259 of factual grounding. This reflects a digital counter-
 260 part to the "peak of inflated expectations," where
 261 high confidence masks low informational compe-
 262 tence.

263 We quantify Stubbornness by measuring the simi-
 264 larity between intermediate layer states and the

final layer representation $h_{final}^{(i)}$ for each token:

$$265 v_{stub}^{(i)} = \frac{1}{L_{end} - L_{start} + 1} \sum_{l=L_{start}}^{L_{end}} \frac{h_l^{(i)} \cdot h_{final}^{(i)}}{\|h_l^{(i)}\| \|h_{final}^{(i)}\|} \quad (2)$$

266 The global stubbornness feature is computed as
 267 $V_{stub} = \frac{1}{N} \sum_{i=1}^N v_{stub}^{(i)}$. Within the GHOST frame-
 268 work, the coupling of high V_{stub} and low V_{turb}
 269 provides a distinct geometric fingerprint of stub-
 270 born hallucinations. This profile stands in sharp
 271 contrast to authentic reasoning, which typically ex-
 272 hibits moderate turbulence as the model iteratively
 273 refines its semantic output.
 274

275 3.3 Output Distribution Analysis

276 Beyond internal geometric features, the output
 277 layer probability landscape encapsulates critical
 278 semantic uncertainty. We introduce two metrics
 279 based on Top- K predictions to complement the
 280 internal hidden state analysis.

281 3.3.1 Predictive Entropy

282 Entropy is a classical uncertainty metric. Given
 283 context x , the model yields a next-token distribu-
 284 tion $P(w|x)$. We select the top- K candidates (set-
 285 ting $K = 10$ for all experiments) and calculate
 286 their Shannon entropy over the top- K probability
 287 distribution. This metric depends entirely on out-
 288 put probabilities without accessing the embedding
 289 space:

$$290 V_{ent} = - \sum_{k=1}^K p_k \log p_k \quad (3)$$

291 where p_k represents the normalized probability of
 292 the k -th candidate token.

293 3.3.2 Semantic Divergence

294 To distinguish lexical synonyms from factual con-
 295 fusion, we measure the geometric dispersion of
 296 candidates. We utilize the static *Input Embedding*
 297 *Layer* to retrieve vectors $\{e_1, \dots, e_K\}$ for the top- K
 298 candidates ($K = 10$), capturing intrinsic semantic
 299 discrepancies independent of contextual processing.
 300 Their average pairwise cosine distance is defined as:
 301

$$302 V_{div} = \frac{1}{K(K-1)} \sum_{i=1}^K \sum_{j \neq i}^K \left(1 - \frac{e_i \cdot e_j}{\|e_i\| \|e_j\|} \right) \quad (4)$$

303 3.4 Non-linear Hallucination Detector

304 We project responses into a four-dimensional fea-
 305 ture space \mathbf{v} . Since linear classifiers fail to capture

Table 1: Comparison of different classifiers within the GHOST framework using **Qwen2.5-1.5B** features. **Random Forest** consistently outperforms other classifiers across all datasets. While XGBoost achieves competitive performance, Random Forest demonstrates superior robustness and stability in modeling the complex, non-linear geometric manifolds of hallucinations.

Classifier	PopQA		FinanceBench		RAGTruth		HaluEval	
	AUPRC↑	F1↑	AUPRC↑	F1↑	AUPRC↑	F1↑	AUPRC↑	F1↑
Logistic Regression	0.8215	0.7842	0.8510	0.6534	0.7721	0.7105	0.8115	0.6890
SVM (RBF)	0.9054	0.8612	0.8745	0.7420	0.9123	0.8045	0.8321	0.7856
XGBoost	0.9382	0.9105	0.9216	0.7645	0.9410	0.8656	0.8805	0.8212
Random Forest	0.9801	0.9529	0.9552	0.8214	0.9596	0.8993	0.9175	0.8698

these non-linear decision boundaries, we benchmarked SVM, XGBoost, and Random Forest using rigorous randomized search cross-validation. Random Forest consistently yielded superior performance by robustly modeling high-dimensional feature interactions within the semantic manifold. Results are summarized in Table 1.

4 Experiments

In this chapter, we evaluate the effectiveness of our proposed GHOST method through extensive experiments on four mainstream hallucination evaluation benchmarks: FinanceBench, RAGTruth, PopQA, and HaluEval. We demonstrate empirically that GHOST is a rapid and effective approach for detecting whether the responses generated by LLMs contain hallucinations.

4.1 Experimental Setup

Baselines. We benchmark GHOST against six competitive baselines across three categories. Logit-based metrics include **Predictive Entropy** (Malinin and Gales, 2021), derived from output probability distributions. Internal-state methods encompass **INSIDE** (Chen et al., 2024), utilizing hidden state covariance eigenvalues, **LI** (Kim et al., 2025), quantifying layer-wise information deficiency, and **LapEigvals** (Binkowski et al., 2025), a spectral approach based on attention map graph Laplacians. For black-box consistency, we evaluate **SelfCheckGPT** (Manakul et al., 2023) using five sampled responses.

Models. We evaluate GHOST across four state-of-the-art LLMs strategically selected for their architectural and functional diversity. **Qwen2.5-1.5B-Instruct** (Yang et al., 2024) acts as a representative for *Small Language Models*, allowing us to probe capacity-induced hallucinations in resource-constrained scenarios. To ensure architectural generalization beyond the LLaMA lineage,

we include **Gemma-3-4B-IT** (Gemma Team and Google DeepMind, 2025), which features distinct configurations like GeGLU activations. **Mistral-7B-Instruct-v0.3** (AI, 2024) serves as the industry-standard baseline to verify the practical relevance of our metrics in widely deployed 7B-scale models. Finally, we incorporate **DeepSeek-R1-Distill-Qwen-7B** (Guo et al., 2025) to examine the complex internal trajectories of *Reasoning Models*. This allows us to verify whether GHOST can effectively decode the "thinking processes" and logic loops inherent in Reinforcement Learning-optimized models, which are particularly susceptible to confused and stubborn hallucinations derived from distillation.

Datasets. We evaluate our method on four benchmarks targeting distinct hallucination facets. **FinanceBench** (Islam et al., 2023) provides complex financial questions to test the model’s accuracy in professional domains. For retrieval-augmented scenarios, we utilize **RAGTruth** (Wu et al., 2023) to assess hallucinations occurring within external knowledge integration. Regarding general knowledge assessment, we employ the **HaluEval** (Li et al., 2023) QA subset containing 10,000 synthesized pairs. Additionally, **PopQA** (Mallen et al., 2023) enables investigation into long-tail entity hallucinations. We specifically analyze the subset where models fail to recall correct entities to determine if intrinsic metrics effectively distinguish genuine knowledge from the hallucination of obscure facts.

Ground Truth Annotation. For generative tasks, obtaining reliable binary labels is critical. We employed gpt-oss-120b (OpenAI, 2025), a large-scale open-source language model, as an automated annotator. The judge is instructed to classify the response as Hallucination only if it contradicts the reference or contains fabricated information, while treating semantic paraphrases as Truth.

Table 2: **Main Results across Multiple LLMs.** We evaluate hallucination detection performance (**AUPRC** and **F1-Score**) on four different base models across diverse tasks: PopQA (Long-tail QA), FinanceBench (Domain-specific), RAGTruth (RAG context), and HaluEval (General). Our method GHOST achieves consistent SOTA performance, particularly demonstrating absolute advantages on the reasoning model DeepSeek-R1. The best results are highlighted in **bold**.

Base Model	Method	PopQA		FinanceBench		RAGTruth		HaluEval		Average	
		AUPRC \uparrow	F1 \uparrow	AUPRC \uparrow	F1 \uparrow	AUPRC \uparrow	F1 \uparrow	AUPRC \uparrow	F1 \uparrow	AUPRC \uparrow	F1 \uparrow
Qwen2.5-1.5B	Predictive Entropy	0.6390	0.6813	0.5123	0.5707	0.6980	0.7041	0.7538	0.7296	0.6508	0.6714
	INSIDE	0.7842	0.7215	0.7456	0.6320	0.8123	0.7544	0.7912	0.7456	0.7833	0.7134
	LI	0.7215	0.7045	0.6845	0.5912	0.7564	0.7123	0.7623	0.7105	0.7312	0.6796
	LapEigvals	0.8528	0.8512	0.8546	0.8190	0.8374	0.8600	0.8846	0.8524	0.8574	0.8457
	<i>SelfCheckGPT</i>	0.8845	0.8410	0.8612	0.7432	0.9120	0.8655	0.9245	0.8512	0.8956	0.8252
	GHOST (Ours)	0.9801	0.9529	0.9552	0.8214	0.9596	0.8993	0.9175	0.8698	0.9531	0.8859
DeepSeek-R1-7B	Predictive Entropy	0.7412	0.6945	0.7256	0.7180	0.7510	0.7234	0.7842	0.7456	0.7505	0.7204
	INSIDE	0.8245	0.7612	0.8123	0.7845	0.8034	0.7567	0.8321	0.8112	0.8181	0.7784
	LI	0.9123	0.8545	0.9012	0.8845	0.9234	0.8912	0.9056	0.8745	0.9106	0.8762
	LapEigvals	0.9312	0.8645	0.9256	0.8912	0.9412	0.9045	0.9123	0.8845	0.9276	0.8862
	<i>SelfCheckGPT</i>	0.9902	0.9045	0.9654	0.9412	0.9884	0.9423	0.9712	0.9485	0.9788	0.9341
	GHOST (Ours)	0.9876	0.9111	0.9696	0.9583	0.9925	0.9536	0.9777	0.9539	0.9819	0.9442
gemma-3-4b	Predictive Entropy	0.7215	0.6510	0.7056	0.5842	0.5312	0.4756	0.6432	0.6180	0.6504	0.5822
	INSIDE	0.9142	0.8567	0.8954	0.7412	0.6723	0.6012	0.8145	0.7824	0.8241	0.7454
	LI	0.8654	0.8105	0.8423	0.7045	0.6356	0.5704	0.7712	0.7412	0.7786	0.7067
	LapEigvals	0.9180	0.8589	0.9012	0.7456	0.6685	0.5984	0.8204	0.7795	0.8270	0.7456
	<i>SelfCheckGPT</i>	0.9654	0.8980	0.9387	0.7912	0.7023	0.6285	0.8512	0.8195	0.8644	0.7843
	GHOST (Ours)	0.9623	0.9006	0.9423	0.7826	0.7070	0.6337	0.8592	0.8238	0.8677	0.7852
Mistral-7B-Instruct	Predictive Entropy	0.6874	0.6012	0.7092	0.5442	0.6215	0.5658	0.5842	0.5312	0.6506	0.5606
	INSIDE	0.8712	0.7612	0.8984	0.6892	0.7845	0.7164	0.7384	0.6756	0.8231	0.7106
	LI	0.8254	0.7212	0.8512	0.6531	0.7432	0.6789	0.6995	0.6402	0.7798	0.6734
	LapEigvals	0.8756	0.7654	0.9012	0.6912	0.7892	0.7201	0.7423	0.6795	0.8271	0.7141
	<i>SelfCheckGPT</i>	0.9123	0.7985	0.9412	0.7214	0.8312	0.7495	0.7712	0.7045	0.8640	0.7435
	GHOST (Ours)	0.9171	0.8012	0.9458	0.7255	0.8263	0.7544	0.7774	0.7115	0.8667	0.7482

We conducted a rigorous human evaluation on a stratified subset of 500 randomly selected samples. We calculated Cohen’s Kappa coefficient (κ) to measure the inter-rater reliability between the human experts and the automated judge. The resulting score of $\kappa = 0.82$ indicates a strong alignment between human judgment and the automated annotator, validating the quality of our ground truth labels for large-scale evaluation.

To validate the reliability of the automatic judge, we additionally conduct a double-blind human evaluation with three domain experts. Annotators are blind to both the model method and the judge predictions; each item is independently labeled by all three experts, and disagreements are resolved by majority vote.

Evaluation Metrics. We treat hallucination de-

tection as a binary classification task, employing **AUPRC** and **F1-score** as the primary evaluation metrics. **AUPRC** serves as our core threshold-independent metric, as it provides a more robust assessment of global discriminative performance under the class imbalances often found in hallucination benchmarks.

4.2 Main Results

Table 2 presents a detailed performance comparison of different methods across the FinanceBench, RAGTruth, HaluEval, and PopQA datasets. To provide a comprehensive evaluation, we report both the AUROC and F1-score for each approach.

Consistent Superiority Across Architectures and Tasks. GHOST achieves state-of-the-art performance across all evaluated LLM architectures and benchmarks. Results in Table 2 show that

Table 3: Ablation study on feature groups using **Qwen-2.5-1.5B** as the base model. We report AUPRC and F1 scores on PopQA, FinanceBench, RAGTruth, and HaluEval datasets.

Method (Features)	PopQA		FinanceBench		RAGTruth		HaluEval		Average	
	AUPRC↑	F1↑	AUPRC↑	F1↑	AUPRC↑	F1↑	AUPRC↑	F1↑	AUPRC↑	F1↑
All Features (GHOST)	0.9801	0.9529	0.9552	0.8214	0.9596	0.8993	0.9175	0.8698	0.9531	0.8859
w/o Entropy	0.9412	0.9105	0.9234	0.7845	0.9312	0.8645	0.8945	0.8312	0.9226	0.8477
w/o Divergence	0.9654	0.9387	0.9412	0.8012	0.9456	0.8812	0.9012	0.8542	0.9384	0.8688
w/o Turbulence	0.9123	0.8845	0.8912	0.7456	0.8845	0.8234	0.8567	0.7912	0.8862	0.8112
w/o Stubbornness	0.9785	0.9510	0.9512	0.8195	0.9570	0.8954	0.9134	0.8655	0.9500	0.8829

Table 4: Main performance comparison of GHOST across various datasets and models. The top section reports In-Distribution (ID) performance. The bottom section shows the mean Out-of-Distribution (OOD) performance when training on one dataset and testing on others.

Base Model	PopQA		FinanceBench		RAGTruth		HaluEval		
	AUPRC	F1	AUPRC	F1	AUPRC	F1	AUPRC	F1	
Mistral-7B-Instruct	0.9171	0.8012	0.9458	0.7255	0.8263	0.7544	0.7774	0.7115	
DeepSeek-R1-7B	0.9876	0.9111	0.9696	0.9583	0.9925	0.9536	0.9777	0.9539	
gemma-3-4b	0.9623	0.9006	0.9423	0.7826	0.7070	0.6337	0.8592	0.8238	
Qwen2.5-1.5B	0.9801	0.9529	0.9552	0.8214	0.9596	0.8993	0.9175	0.8698	
Average (ID)	0.9618	0.8915	0.9532	0.8220	0.8714	0.8103	0.8830	0.8398	
<i>Cross-Dataset OOD Generalization (Mean Performance across Models)</i>									
OOD (PopQA Train)	–	–	0.8845	0.7104	0.8123	0.7245	0.7956	0.7412	
OOD (FinanceBench Train)	0.7214	0.6532	–	–	0.6145	0.5234	0.6523	0.5845	
OOD (RAGTruth Train)	0.8934	0.8215	0.8412	0.7012	–	–	0.8412	0.7956	
OOD (HaluEval Train)	0.8545	0.7912	0.8234	0.6845	0.8312	0.7645	–	–	

GHOST outperforms established white-box baselines and the high-cost SelfCheckGPT in Average AUPRC and F1-score. On the Qwen2.5-1.5B model, our method reaches a remarkable 0.9801 AUPRC, demonstrating that capturing geometric trajectories provides a more precise signal for hallucination than traditional uncertainty or static probing.

Exceptional Efficacy in Reasoning-Intensive Models. GHOST exhibits a distinct advantage when applied to reasoning-enhanced models like DeepSeek-R1-7B. Our method reaches an average AUPRC of 0.9819 on this backbone, surpassing the robust SelfCheckGPT baseline. While traditional predictive entropy fails to capture subtle deceptive patterns in complex logical deductions, the Mental Turbulence metric effectively quantifies the geometric instability inherent in flawed reasoning chains to enable near-perfect detection.

Robust Generalization in Domain-Specific and RAG Scenarios. GHOST effectively addresses hallucination detection in professional contexts and retrieval-augmented environments. In FinanceBench, GHOST maintains an AUPRC above 0.94 across various models, significantly exceed-

ing consistency-based baselines. This suggests that GHOST identifies false confidence when factual grounding is absent through the Stubbornness metric, proving its utility in distinguishing genuine knowledge from the hallucination of obscure facts.

4.3 Ablation Study: Dissecting the GHOST

To verify the effectiveness of each feature component in GHOST, we conducted an ablation study by systematically removing specific feature groups, as shown in Table 3.

The ablation results in Table 3 indicate that the Turbulence feature group is the most critical component of GHOST. Removing Turbulence leads to the largest degradation across all datasets, reducing the average AUPRC from 0.9531 to 0.8862 (0.0669) and the average F1 from 0.8859 to 0.8112 (0.0747). This supports our hypothesis that prediction instability provides a strong signal for hallucination detection.

Among the remaining groups, Entropy contributes the next most substantial improvements: excluding Entropy decreases the average AUPRC to 0.9226 (0.0305) and the average F1 to 0.8477 (0.0382). In contrast, removing Divergence

causes only a modest drop (average AUPRC 0.9384, 0.0147; average F1 0.8688, 0.0171), and removing Stubbornness has a minimal effect (average AUPRC 0.9500, 0.0031; average F1 0.8829, 0.0030). Overall, Turbulence provides the most distinct and indispensable signal, while Entropy offers complementary gains and Divergence/Stubbornness appear partially redundant given the current classifier.

4.4 Generalization and OOD Robustness

To evaluate the transferability of GHOST’s internal geometric features, we conduct Out-of-Distribution experiments across two dimensions: **Cross-Dataset OOD** and **Cross-Model OOD**. The specific performance is shown in the table 4.

4.5 Efficiency Analysis

We evaluate the computational efficiency of the GHOST framework on a high-performance server equipped with a 4 NVIDIA GeForce RTX 3090 GPU and dual Intel Xeon Gold 6326 CPUs. As presented in Table 5, the standard inference process for the Mistral-7B-Instruct model requires an average of 3.05 seconds per query.

Table 5: Efficiency comparison between SelfCheckGPT and GHOST. *Gen. Time* represents the average time for generating a single response. *Add. Latency* denotes the additional wall-clock time required for hallucination detection.

Method	Mechanism	Gen. Time	Add. Latency	Overhead
SelfCheckGPT	Sampling ($N = 5$)	3.05s	12.20s	400.0%
GHOST (Ours)	Vectorized Extr.	3.05s	0.14s	4.6%

GHOST utilizes a fully vectorized extraction mechanism integrated into the model’s single forward pass. By employing optimized broadcasting to compute internal metrics, GHOST incurs a negligible latency of only 0.14s—a marginal 4.6% overhead. This efficiency highlights GHOST’s suitability for real-time deployment, significantly outperforming existing baselines.

In contrast, SelfCheckGPT imposes a prohibitive computational burden due to its reliance on stochastic consistency. Requiring N additional sampling sequences, it introduces a generation-bound latency increase of approximately 400%. Such overhead renders SelfCheckGPT impractical for latency-sensitive applications despite its detection efficacy.

5 Conclusion

First, we provide empirical evidence of the **Digital Dunning-Kruger Effect** in Large Language Models, revealing that internal confidence often functions as a deceptive proxy for semantic veracity. Building on this observation, we introduce **GHOST**, a training-free framework for hallucination detection. Diverging from prior approaches reliant on coarse-grained signals, GHOST leverages hierarchical geometric trajectories coupled with token-level dynamics, enabling a high-resolution characterization of model behaviors intrinsically linked to factual errors.

Second, we formalize a novel hallucination taxonomy grounded in observable geometric manifestations, distinguishing *confused hallucinations* from *stubborn hallucinations*. This categorization provides a rigorous operational lens for diagnosing specific error modes, facilitating more targeted evaluation and mitigation strategies in complex reasoning tasks.

Finally, we demonstrate the empirical efficacy of GHOST across four diverse benchmarks. Our results show that GHOST consistently outperforms established baselines in AUPRC and F1-score while maintaining superior computational efficiency, thereby offering a practical solution for real-time hallucination monitoring.

Limitations

Despite its effectiveness, several limitations of GHOST warrant further study. First, this work primarily evaluates decoder-only Transformer architectures; the applicability of geometric dynamics to encoder-decoder or non-Transformer structures remains unexplored. Second, as a white-box method, GHOST requires access to internal hidden states, limiting its use in closed-source, API-only scenarios. Finally, our current taxonomy may not fully capture complex reasoning errors, such as multi-step logical fallacies. Future research will focus on developing generalized, zero-shot geometric indicators to reduce reliance on labeled datasets.

Ethics Statement

This research adheres to the ethical guidelines prescribed by the academic community and focuses on improving the reliability of Large Language Models. Our methodology utilizes publicly available, open-source datasets and models, ensuring that no

private or personally identifiable information is processed during the experiments.

We acknowledge that while GHOST is designed to detect and mitigate hallucinations, it is not a definitive oracle for truth. There is a potential risk of false negatives, where incorrect model outputs may remain undetected, and false positives, which could lead to the suppression of creative or subjective content. Users should exercise caution and not rely solely on GHOST for high-stakes decision-making in critical domains such as medical or legal sectors without human oversight.

Furthermore, we are committed to the democratization of AI safety tools. By providing a computationally efficient, white-box detection framework, we aim to reduce the energy consumption associated with resource-intensive black-box verification methods. We do not foresee any direct negative social impacts arising from this work, provided it is used as a transparency-enhancing tool rather than a mechanism for automated censorship.

References

Mistral AI. 2024. Mistral 7b v0.3 model card. <https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.3>. Accessed: 2025-12-25.

Amos Azaria and Tom Mitchell. 2023. [The internal state of an LLM knows when it’s lying](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 967–976, Singapore. Association for Computational Linguistics.

Alexandra Bazarova, Aleksandr Yugay, Andrey Shulga, Alina Ermilova, Andrei Volodichev, Konstantin Polev, Julia Belikova, Rauf Parchiev, Dmitry Simakov, Maxim Savchenko, Andrey Savchenko, Serguei Barannikov, and Alexey Zaytsev. 2025. [Hallucination detection in llms with topological divergence on attention graphs](#). *Preprint*, arXiv:2504.10063.

Jakub Binkowski, Denis Janiak, Albert Sawczyn, Bogdan Gabrys, and Tomasz Jan Kajdanowicz. 2025. [Hallucination detection in LLMs using spectral features of attention maps](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 24354–24385, Suzhou, China. Association for Computational Linguistics.

Chao Chen, Kai Liu, Ze Chen, Yi Gu, Yue Wu, Mingyuan Tao, Zhihang Fu, and Jieping Ye. 2024. [INSIDE: LLMs’ internal states retain the power of hallucination detection](#). In *The Twelfth International Conference on Learning Representations*.

Yung-Sung Chuang, Linlu Qiu, Cheng-Yu Hsieh, Ranjay Krishna, Yoon Kim, and James R. Glass. 2024.

[Lookback lens: Detecting and mitigating contextual hallucinations in large language models using only attention maps](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 1419–1436, Miami, Florida, USA. Association for Computational Linguistics.

Roi Cohen, May Hamri, Mor Geva, and Amir Globerson. 2023. [LM vs LM: Detecting factual errors via cross examination](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12621–12640, Singapore. Association for Computational Linguistics.

Ekaterina Fadeeva, Roman Vashurin, Akim Tsvigun, Artem Vazhentsev, Sergey Petrakov, Kirill Fedyanin, Daniil Vasilev, Elizaveta Goncharova, Alexander Panchenko, Maxim Panov, Timothy Baldwin, and Artem Shelmanov. 2023. [LM-polygraph: Uncertainty estimation for language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 446–461, Singapore. Association for Computational Linguistics.

Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, Yarin Chen, and Yarin Gal. 2024a. [Detecting hallucinations in large language models using semantic entropy](#). *Nature*, 630(8017):625–630.

Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. 2024b. [Detecting hallucinations in large language models using semantic entropy](#). *Nature*, 630(8017):625–630.

Leon Festinger. 1957. *A theory of cognitive dissonance*. Stanford university press.

Gemma Team and Google DeepMind. 2025. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*.

Aman Goel, Daniel Schwartz, and Yanjun Qi. 2025. [Zero-knowledge LLM hallucination detection and mitigation through fine-grained cross-model consistency](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 1982–1999, Suzhou (China). Association for Computational Linguistics.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruipu Xu, Qi Zhu, Shirong Ma, Pei Wang, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.

Pranab Islam, Anand Kannappan, Douwe Kiela, Rebecca Qian, Nino Scherrer, and Bertie Vidgen. 2023. [Financebench: A new benchmark for financial question answering](#). *Preprint*, arXiv:2311.11944.

Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. [Survey of hallucination in natural language generation](#). *ACM Comput. Surv.*, 55(12).

661	Hazel Kim, Tom A. Lamb, Adel Bibi, Philip Torr, and Yarin Gal. 2025. Detecting LLM hallucination through layer-wise information deficiency: Analysis of ambiguous prompts and unanswerable questions . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 32298–32310, Suzhou, China. Association for Computational Linguistics.	Zhang. 2023. Ragtruth: A hallucination corpus for developing trustworthy retrieval-augmented language models . <i>Preprint</i> , arXiv:2401.00396.	718 719 720
662			
663			
664			
665		Yilun Xu, Shengjia Zhao, Jiaming Song, Russell Stewart, and Stefano Ermon. 2020. A theory of usable information under computational constraints . In <i>International Conference on Learning Representations</i> .	721 722 723 724
666			
667			
668			
669	Justin Kruger and David Dunning. 1999. Unskilled and unaware of it: How difficulties in recognizing one’s own incompetence lead to inflated self-assessments . <i>Journal of Personality and Social Psychology</i> , 77(6):1121–1134.	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, and 1 others. 2024. Qwen2.5 technical report. <i>arXiv preprint arXiv:2412.15115</i> .	725 726 727 728 729
670			
671			
672			
673			
674	Junyi Li, Xiaoxue Cheng, Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. 2023. HaluEval: A large-scale hallucination evaluation benchmark for large language models . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 6449–6464, Singapore. Association for Computational Linguistics.	Fujie Zhang, Peiqi Yu, Biao Yi, Baolei Zhang, Tong Li, and Zheli Liu. 2025a. Prompt-guided internal states for hallucination detection of large language models . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 21806–21818, Vienna, Austria. Association for Computational Linguistics.	730 731 732 733 734 735 736
675			
676			
677			
678			
679			
680			
681	Andrey Malinin and Mark Gales. 2021. Uncertainty estimation in autoregressive structured prediction . In <i>Proceedings of the International Conference on Learning Representations (ICLR)</i> .	Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. 2025b. Siren’s song in the ai ocean: A survey on hallucination in large language models . <i>Computational Linguistics</i> , pages 1–46.	737 738 739 740 741 742 743
682			
683			
684			
685	Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 9802–9822, Toronto, Canada. Association for Computational Linguistics.	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, and 3 others. 2025. A survey of large language models . <i>Preprint</i> , arXiv:2303.18223.	744 745 746 747 748 749 750
686			
687			
688			
689			
690			
691			
692			
693	Potsawee Manakul, Adian Liusie, and Mark Gales. 2023. SelfCheckGPT: Zero-resource black-box hallucination detection for generative large language models . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 9004–9017, Singapore. Association for Computational Linguistics.	A Implementation Details	751
694			
695			
696			
697			
698			
699			
700	Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. Factscore: Fine-grained atomic evaluation of factual precision in long form text generation . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 12076–12100.	A.1 Detailed Dataset Analysis and Characteristics	752 753
701			
702			
703			
704			
705			
706			
707	OpenAI. 2025. gpt-oss-120b & gpt-oss-20b model card . <i>Preprint</i> , arXiv:2508.10925.	To evaluate the robustness of GHOST across diverse hallucination scenarios, we curate a benchmark suite encompassing four distinct domains. The statistical distribution and partitioning of these datasets are summarized in Table 6.	754 755 756 757 758
708			
709	Jialiang Wu, Yi Shen, Sijia Liu, Yi Tang, Sen Song, Xiaoyi Wang, and Longjun Cai. 2025. Improve decoding factuality by token-wise cross layer entropy of large language models . In <i>Findings of the Association for Computational Linguistics: NAACL 2025</i> , pages 3912–3921, Albuquerque, New Mexico. Association for Computational Linguistics.	<ul style="list-style-type: none"> • HaluEval (General): This large-scale benchmark provides 10,000 samples covering general knowledge. It is instrumental for training our baseline classifiers and assessing the model’s fundamental factuality alignment in open-domain conversations. • PopQA (Long-tail Knowledge): Focused on low-popularity entities, PopQA challenges the model’s ability to distinguish between internal knowledge and parasitic memory. This dataset is crucial for validating the "Digital 	759 760 761 762 763 764 765 766 767 768 769
710			
711			
712			
713			
714			
715			
716	Yuanhao Wu, Juno Zhu, Siliang Xu, Kashun Shum, Cheng Niu, Randy Zhong, Juntong Song, and Tong		
717			

770	Dunning-Kruger Effect," as models often exhibit high confidence in these long-tail facts despite frequent hallucinations.	818
771		819
772		820
773	• FinanceBench (Domain-Specific): Representing high-stakes professional reasoning, FinanceBench requires precise numerical and conceptual accuracy. It serves to test GHOST’s efficacy in detecting "Stubborn Hallucinations" where models generate plausible but incorrect financial deductions.	821
774		822
775		823
776		824
777		
778		
779		
780	• RAGTruth (Retrieval-Augmented): This dataset evaluates hallucinations in the context of external documents. It allows us to analyze how geometric trajectories shift when a model is constrained by provided context versus relying solely on internal parameters, providing insights into "Confused Hallucinations."	
781		
782		
783		
784		
785		
786		
787	As shown in Table 6, we adopt a stratified 80/20 split across all datasets to ensure the Random Forest classifier is trained on a representative distribution of both truthful and hallucinated responses while maintaining a rigorous held-out set for performance reporting.	
788		
789		
790		
791		
792		
793	A.2 Computational Resources and Environment	
794		
795	All experiments were conducted on a server equipped with four NVIDIA GeForce RTX 3090 GPUs with 24GB VRAM. The software environment was built upon Python 3.10 and the HuggingFace Transformers (v4.40.0) library. The inference process utilized greedy decoding with a temperature of zero and a batch size of one to simulate a real-time streaming scenario. For the parallel execution of baseline models, we disabled tokenizer parallelism by setting <code>TOKENIZERS_PARALLELISM</code> to false and pre-loaded all model weights in the main process to prevent resource contention. The total GPU time for extracting geometric features across all four datasets and four base models was approximately 24 hours.	
796		
797		
798		
799		
800		
801		
802		
803		
804		
805		
806		
807		
808		
809		
810	A.3 Model Configuration and Layer Selection	
811	To capture the most representative internal reasoning dynamics while excluding shallow lexical encoding noise and deep output formatting convergence, we apply a dynamic layer selection strategy relative to the total model depth L . Based on empirical observations in our pilot studies, the most salient information regarding semantic truthfulness	
812		
813		
814		
815		
816		
817		
	is located in the middle-to-late stages of the network. We specifically extract hidden states from the relative depth interval of $0.1L$ to $0.9L$. This adaptive strategy ensures that GHOST effectively captures core geometric features across varying model scales. The specific layer indices for each backbone model are detailed in Table 7.	825
	A.4 Training Protocol and Classifier Optimization	826
	To ensure reproducibility, we employ a stratified train test split across all datasets, partitioning each corpus into 80% training and 20% testing sets. Stratification preserves both class proportions and dataset-source distributions across splits. For out-of-distribution experiments, models are trained on the full training set of a single source dataset and evaluated on the complete test sets of unseen datasets or model families. We adopt a Random Forest classifier as the primary non-linear predictor to map the geometric feature vector $\mathbf{v} = [V_{turb}, V_{stub}, V_{ent}, V_{div}]$ to the binary truthfulness label due to its robustness to noisy features.	827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
	• Handling Class Imbalance. To address class imbalance in our datasets, we use the <code>balanced_subsample</code> strategy, which reweights classes within each bootstrap sample to reduce bias toward the majority class.	840
		841
		842
		843
		844
	• Hyperparameter Optimization. All hyperparameters are tuned using <code>RandomizedSearchCV</code> with 5-fold cross-validation on the training set. We sample 50 configurations from a shared hyperparameter sampling distribution for all experiments and fix the random seed. The same tuning protocol and search space are applied across all datasets.	845
		846
		847
		848
		849
		850
		851
		852
		853
	• Selected Configuration. The selected model uses 750 trees with a maximum depth of 40 and <code>max_features</code> set to the square root of the feature dimensionality. Specific sampling distributions and selected values are summarized in Table 8.	854
		855
		856
		857
		858
		859
	A.5 SelfCheckGPT Configuration	860
	SelfCheckGPT operates by sampling multiple stochastic responses from the model and measuring their consistency with the original response. We adopted the SelfCheck-NLI variant utilizing	861
		862
		863
		864

Table 6: Detailed statistics and domain characteristics of the benchmarks used in GHOST evaluation. For all datasets, we maintain a stratified 80/20 split for training the geometric classifier and reporting the final detection performance.

Dataset	Total Samples	Train (80%)	Test (20%)	Knowledge Domain
HaluEval	10,000	8,000	2,000	General / Open-domain
PopQA	1,400	1,120	280	Long-tail / Entity-centric
FinanceBench	1,200	960	240	Financial Reasoning
RAGTruth	2,500	2,000	500	Retrieval-Augmented

Table 7: Configuration of backbone models and layer selection strategy. We exclude the first and last 10% of layers to filter initial embedding noise and final distribution convergence.

Base Model	Parameters	Total Layers (L)	Selection Ratio	Selected Indices
Qwen2.5-1.5B-Instruct	1.5B	28	$0.1L \rightarrow 0.9L$	3 \rightarrow 25
Gemma-3-4B-IT	4B	28	$0.1L \rightarrow 0.9L$	3 \rightarrow 25
Mistral-7B-v0.3-Inst	7B	32	$0.1L \rightarrow 0.9L$	3 \rightarrow 29
DeepSeek-R1-Dist-7B	7B	32	$0.1L \rightarrow 0.9L$	3 \rightarrow 29

Table 8: Hyperparameter sampling distributions and the selected Random Forest configuration obtained via randomized search with 5-fold cross-validation.

Parameter	Sampling Distribution	Selected Value
n_estimators	UniformInt(300, 1000)	750
max_depth	{10, 20, 30, 40, 50, None}	40
min_samples_split	UniformInt(2, 20)	8
min_samples_leaf	UniformInt(1, 10)	2
max_features	{'sqrt', 'log2'}	'sqrt'
class_weight	{'balanced', 'balanced_subsample'}	subsample

a deberta-v3-large-mnli model as the entailment scorer. For each query, we generated $N = 5$ stochastic samples using temperature $T = 0.7$ and Top-p $p = 0.9$ with a limit of 200 new tokens. We filtered out degenerate samples to ensure the validity of the consistency check. The final hallucination score was defined as the maximum contradiction probability across all constituent sentences using Max-Prob aggregation. The total inference time for GHOST adds less than 5% latency compared to a standard forward pass, which stands in sharp contrast to the multi-pass requirement of SelfCheck-GPT.

A.6 Evaluation Metrics Calculation

We employ AUPRC (Area Under the Precision-Recall Curve) and F1-score as primary metrics. AUPRC is calculated using the trapezoidal rule via the scikit-learn implementation, providing a threshold-independent measure that is sensitive to the minority (hallucination) class. The F1-score is reported at the optimal threshold determined by the classifier during the validation phase.

B Prompt Templates

To ensure fair evaluation and alignment with the instruction-tuning stage of each model, we strictly adhere to the official chat templates provided by the respective tokenizers. The prompt construction consists of two stages: *Input Construction* (formatting the task content) and *Chat Formatting* (applying model-specific control tokens).

B.1 Input Construction

Depending on the dataset type, we format the user input content as follows:

- **Standard QA (TruthfulQA, PopQA):** The input consists solely of the question.

{question}

- **Context-Aware QA (HaluEval):** When external knowledge or a passage is provided, we prepend it to the question.

Context:

{passage}\n\nQuestion:

{question}

B.2 Model-Specific Formatting

We utilize the apply_chat_template function from the HuggingFace transformers library to automatically apply the correct control tokens. For models requiring manual formatting (e.g., Gemma), we implement the official prompt structure. The specific templates used in our experiments are detailed in Table 9.

Table 9: Chat templates applied to different model families. {Input_Content} refers to the string constructed in the Input Construction phase.

Model Family	Template Structure
Gemma-3-4B-IT	We use the official turn-based control tokens: <start_of_turn>user\n {Input_Content} <end_of_turn>\n<start_of_turn>model\n
Qwen2.5 / DeepSeek-R1	We utilize the standard ChatML-like format via the tokenizer: < im_start >user\n {Input_Content} < im_end >\n< im_start >assistant\n
Mistral-7B-v0.3	We utilize the standard instruction format via the tokenizer: [INST] {Input_Content} [/INST]
<i>Fallback / Generic</i>	In cases where the tokenizer template is unavailable, we default to: User: {Input_Content}\nAssistant: