# Towards Unification of Hallucination Detection and Fact Verification for Large Language Models

Weihang Su
swh22@mails.tsinghua.edu.cn
DCST, Tsinghua University
Beijing, China

Jianming Long
DCST, Tsinghua University
Beijing, China

Changyue Wang
DCST, Tsinghua University
Beijing, China

Shiyu Lin*
DCST, Tsinghua University
Beijing, China

Jingyan Xu*
DCST, Tsinghua University
Beijing, China

Ziyi Ye
Fudan University
Shanghai, China

Qingyao Ai[†]
aiqy@tsinghua.edu.cn
DCST, Tsinghua University
Beijing, China

Yiqun Liu
DCST, Tsinghua University
Beijing, China

## Abstract

Large Language Models (LLMs) frequently exhibit hallucinations, generating content that appears fluent and coherent but is factually incorrect. Such errors undermine trust and hinder their adoption in real-world applications. To address this challenge, two distinct research paradigms have emerged: model-centric Hallucination Detection (HD) and text-centric Fact Verification (FV). Despite sharing the same goal, these paradigms have evolved in isolation, using distinct assumptions, datasets, and evaluation protocols. This separation has created a research schism that hinders their collective progress. In this work, we take a decisive step toward bridging this divide. We introduce UniFact[1], a unified evaluation framework that enables direct, instance-level comparison between FV and HD by dynamically generating model outputs and corresponding factuality labels. Through large-scale experiments across multiple LLM families and detection methods, we reveal three key findings: (1) No paradigm is universally superior; (2) HD and FV capture complementary facets of factual errors; and (3) hybrid approaches that integrate both methods consistently achieve state-of-the-art performance. Beyond benchmarking, we provide the first in-depth analysis of why FV and HD diverged, as well as empirical evidence supporting the need for their unification. The comprehensive experimental results call for a new, integrated research agenda toward unifying Hallucination Detection and Fact Verification in LLMs.

## 1 Introduction

Large Language Models (LLMs) have fundamentally reshaped the digital landscape, powering next-generation search engines and automating complex tasks that previously seemed impossible [6, 16, 17, 62]. However, their deployment in high-stakes, real-world applications is critically hampered by a dangerous failure mode: the generation of fluent yet factually incorrect content, a phenomenon

known as hallucination [22, 37, 57, 60]. In domains like healthcare [43, 64, 76], legal [55, 63, 73], and journalism, where factual reliability is a critical requirement, hallucinations can erode user trust, spread misinformation, and lead to harmful consequences. Therefore, developing robust mechanisms to detect these factual errors is not just a technical challenge but a foundational requirement for building trustworthy and reliable intelligent systems.

To address this critical challenge, the research community has developed two distinct paradigms: Fact Verification (FV) [20, 26, 30, 65] and Hallucination Detection (HD) [2, 8, 18, 22, 33, 35, 40, 61]. Fact Verification is a well-established, text-centric research area rooted in information retrieval [9, 11, 21] and fake news detection [42, 47, 83]. Typically, FV evaluates a statement's factuality based on external knowledge sources, such as web documents or databases. A defining characteristic of FV is its origin-agnostic nature: it evaluates a statement's factuality regardless of whether it is human-authored or machine-generated. In this paradigm, the text is treated as an isolated, static artifact, where factuality is assessed based on the consistency between the claim and retrieved evidence. In contrast, Hallucination Detection is a more recent paradigm that emerged alongside the rise of LLMs. The fundamental distinction lies in the object of evaluation: HD specifically targets content generated by LLMs. Unlike FV, which evaluates text in isolation, HD introduces a new dimension of introspective evaluation by leveraging the generative process itself. While HD methods can also employ external knowledge retrieval similar to FV, the mainstream research focuses on model-centric signals, such as token-level uncertainty, activation dynamics, and consistency across sampled generations [8, 37, 61]. Consequently, HD is inherently model-centric, allowing it to identify hallucinations even in the absence of external evidence.

Despite sharing the same ultimate goal of identifying factual errors, these two research directions have evolved in surprising isolation, which is evidenced by their use of distinct benchmarks, evaluation metrics, and publication venues. The root of this divergence lies in a fundamental incompatibility between evaluation paradigms. FV is inherently text-centric: its benchmarks, such as FEVER [65], consist of claims treated as standalone texts, where

---

their origin is not considered. HD methods, however, depend on internal model states that FV settings cannot provide. Attempts to bridge this gap by constructing static benchmarks of LLM outputs (e.g., HELM [60]) remain limited, as such datasets cannot evaluate new or out-of-distribution models, nor capture dynamic generative signals. This incompatibility has led to the development of two parallel research tracks, resulting in redundant efforts and leaving fundamental questions unanswered. It remains unclear whether one paradigm is inherently superior, how their strengths differ, or whether they can be effectively combined. Moreover, their isolation prevents leveraging their complementary nature: HD's internal model signals could augment FV's evidence-based reasoning, while FV's retrieved context could ground HD's uncertainty signals.

To this end, we take a decisive step toward bridging this long-standing divide. We present the first large-scale, systematic empirical study to directly compare Hallucination Detection and Fact Verification under a single, unified setting. Rather than a mere benchmarking effort, our work provides a comprehensive investigation into the fundamental relationship between these paradigms. This unified approach allows us, for the first time, to address three key research questions (RQs) that have remained unexplored due to the fields' historical separation:

(1) **Overall Performance (RQ1)**: How do Hallucination Detection and Fact Verification methods differ in performance when detecting factual errors on the same LLM outputs?
(2) **Complementary Strengths (RQ2)**: Do HD and FV capture distinct aspects of factual inaccuracy, showing complementary strengths, or do they overlap and fail on similar errors?
(3) **Hybrid Performance (RQ3)**: Can the signals from both paradigms be effectively combined? Can such a hybrid detection system surpass either approach in isolation?

This comparison was previously considered infeasible due to the paradigms' fundamental incompatibility. Our work makes this comparison possible by introducing **UniFact**, a unified evaluation framework designed to overcome the fundamental incompatibility that has historically split the field. Unlike static benchmarks that rely on pre-generated, fixed LLM's outputs, UniFact adopts a dynamic evaluation process. It prompts any LLM to generate answers to factual questions in real time. The framework then employs an automated labeling system to determine correctness by comparing the generated output against ground-truth answers and their corresponding pre-labeled reference knowledge. This dynamic process is the key: it simultaneously yields both the textual outputs required for FV methods and the internal model signals needed for HD methods. This design enables the direct, instance-level, head-to-head comparison of both paradigms on the same generated content.

Our large-scale experiments, conducted via UniFact across diverse datasets, model families (e.g., LLaMA, Qwen), and detection methods, yield several crucial insights that directly answer our research questions. First, addressing RQ1, we find that no single paradigm is universally superior; performance varies significantly with the underlying LLM and the nature of the task. Second, and most critically (addressing RQ2), our analysis provides strong empirical evidence of complementarity. HD and FV methods consistently succeed on different subsets of factual errors, confirming they capture

distinct yet synergistic dimensions of factuality. Finally (addressing RQ3), building on this discovery, we demonstrate that simple hybrid methods integrating both paradigms consistently and significantly outperform all individual approaches, establishing a new state-of-the-art in factual error detection.

The implications of these findings are twofold. For the research community, they motivate a paradigm shift from isolated studies toward a unified research agenda that embraces this synergy. For practitioners, our results offer a clear takeaway: robust factuality assessment requires combining internal, model-based signals with external, evidence-based reasoning. In summary, this paper makes the following key contributions:

- **Analytical**: We identify and analyze the conceptual and methodological schism between Hallucination Detection and Fact Verification, highlighting its consequences for progress in AI factuality.
- **Methodological**: We introduce UniFact, a unified evaluation framework. It overcomes this divide and enables direct, fair, and scalable comparison between HD and FV.
- **Empirical**: Through the first large-scale comparative study, we provide strong evidence that HD and FV exhibit distinct yet complementary strengths.
- **Practical**: We demonstrate that simple hybrid strategies combining both paradigms achieve new state-of-the-art performance.

## 2 Problem Formulation

This section establishes formal definitions for Fact Verification (FV) and Hallucination Detection (HD). While these terms are often used interchangeably in recent literature, distinguishing them is crucial for understanding the evaluation landscape. As outlined in the introduction and illustrated in Figure 1, the primary distinction lies in the object of evaluation: FV targets the semantic truthfulness of text in an origin-agnostic manner, whereas HD targets the factual correctness of the text generated by an LLM. We formally define HD and FV in § 2.1 and § 2.2, respectively, and unify their evaluation metrics in § 2.3.

### 2.1 Fact Verification

Fact Verification (FV) is defined as the task of assessing the factual correctness of a textual claim by grounding it in authoritative external evidence. A defining characteristic of FV is that it is origin-agnostic and content-centric: it evaluates the text as a static artifact, independent of whether it was authored by a human or generated by an LLM. Formally, let $y$ denote a textual claim (e.g., a sentence or an atomic proposition), and $\mathcal{E}$ represent a set of trusted external evidence (e.g., retrieved documents, knowledge graph triples). FV treats the verification process as a classification problem between the claim and the evidence, disregarding the source that produced $y$. The task is formulated as estimating the probability that $y$ is non-factual given $\mathcal{E}$:

$$p_{\mathrm{FV}} = D_{\mathrm{FV}}(y, \mathcal{E}) \in [0, 1], \tag{1}$$

where $D_{\mathrm{FV}}$ is a scoring function quantifying the factual error, such that $p_{\mathrm{FV}} \approx 1$ indicates that $y$ conflicts with $\mathcal{E}$ (non-factual), while $p_{\mathrm{FV}} \approx 0$ implies consistency. In this paradigm, the focus is exclusively on the alignment between the claim content and the world knowledge provided by $\mathcal{E}$ [39, 65].
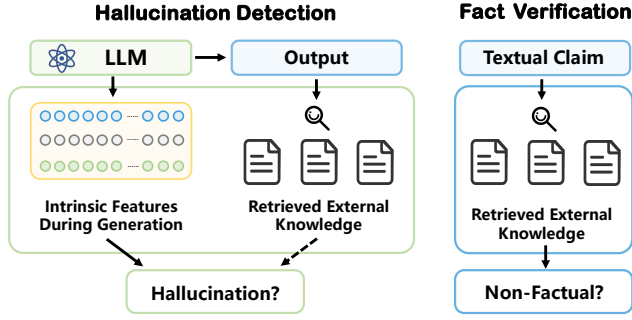
**Figure 1: An illustration of comparison between Hallucination Detection (HD) and Fact Verification (FV). HD targets LLM outputs leveraging both intrinsic features and external knowledge (dashed arrow), whereas FV targets generic textual claims using exclusively retrieved evidence.**

## 2.2 Hallucination Detection

In contrast to FV, Hallucination Detection (HD) is model-centric and specifically targets the factual accuracy of the content generated by LLMs. HD aims to identify instances where the model generates content that is unfaithful to the source input or factually incorrect relative to world knowledge. Unlike FV, which evaluates text independently of its origin, HD is inherently model-centric and targets factual inconsistencies specifically in text generated by an LLM. Formally, let $M$ be a language model that generates an output sequence $y$ given an input context $x$. HD estimates the probability that $y$ contains a hallucination. This estimation can be derived from two distinct sources of signals:

- **Intrinsic Signals ($\mathcal{F}_M$):** Features derived directly from the LLM's generative process, including predictive entropy, activation patterns, cross-sample consistency, etc.
- **Extrinsic Signals ($\mathcal{E}$):** External evidence retrieved to verify the generated content, comparable to FV, but utilized to evaluate the factual integrity of the LLM's response.

Therefore, the general formulation for HD is:

$$p_{\text{HD}} = D_{\text{HD}}(y, x, \mathcal{F}_M, \mathcal{E}) \in [0, 1], \qquad (2)$$

Crucially, while external evidence $\mathcal{E}$ is a valid information source for HD, most existing Hallucination Detection approaches rarely incorporate $\mathcal{E}$, reflecting the paradigm's emphasis on leveraging intrinsic signals $\mathcal{F}_M$ as its core methodological advantage. Consequently, HD enables introspective evaluation based solely on model-internal signals, avoiding the need for external retrieval.

## 2.3 Evaluation

Despite their distinct methodologies, HD and FV can be unified under a single evaluation framework. Both tasks can be formulated as a binary classification problem. Specially, let $s_i \in [0, 1]$ denote the predicted score for instance $i$, representing the probability that the text is non-factual (FV) or hallucinated (HD), and let $l_i^* \in \{0, 1\}$ denote the ground-truth label, where $l_i^* = 1$ signifies a negative sample. Standard classification metrics can be adopted to assess performance, specifically Accuracy and the Area Under

the ROC Curve (AUC). AUC is particularly emphasized to evaluate the model's ranking capability independent of specific decision thresholds, addressing the potential class imbalance in evaluation datasets.

## 3 Unified Evaluation Framework

In this section, we introduce our unified evaluation framework, designed to bridge the long-standing divide between Hallucination Detection and Fact Verification. We begin by examining why these two closely related paradigms have historically evolved in isolation, and we analyze the incompatibilities in existing benchmarks that have impeded cross-paradigm evaluation (§3.1). We then present the core design principles underlying our framework (§3.2), followed by a detailed description of its architecture (§3.3).

### 3.1 The Benchmark Incompatibility Challenge

Before detailing our framework design, it is crucial to analyze the technical root causes that have historically segregated these two highly related paradigms. The divergence stems from a fundamental incompatibility between the *text-centric* design of traditional FV datasets and the *model-centric* requirements of HD methods. Traditional FV benchmarks, such as FEVER [65], are inherently text-centric. They consist of static claim-label pairs that are agnostic to the origin of the claim. This structure renders them unsuitable for Hallucination Detection methods (e.g., SelfCheckGPT [37], MIND [60]), which fundamentally rely on analyzing the generation process, such as internal states or cross-sample consistency. These generative signals are absent in static FV datasets, creating an impassable barrier for unified evaluation.

A seemingly intuitive solution to bridge this gap is to construct static benchmarks comprised of pre-generated LLM outputs (e.g., HELM [60]). While valuable, we argue that reliance on static collections of model outputs suffers from two prohibitive limitations in the rapidly evolving landscape of AI:

- **Model Obsolescence:** The pace of LLM development is extraordinary. A dataset constructed with outputs from today's state-of-the-art models (e.g., LLaMA-3) risks becoming an archival snapshot rather than a durable evaluation tool as new architectures emerge.
- **Inability to Evaluate New Models:** Crucially, static benchmarks cannot evaluate HD methods on new or out-of-distribution (OOD) models. Since HD performance is often tightly coupled to specific model characteristics, a framework that cannot test methods on arbitrary new LLMs fundamentally restricts progress and fails to measure true generalizability.

To overcome these structural limitations, a paradigm shift is required from analyzing static artifacts to evaluating dynamic processes. An effective unified framework must not rely on fixed outputs but instead on a core set of instructions (questions) that can trigger generation from *any* target LLM on the fly. This dynamic approach ensures that (1) FV methods receive the textual content they require, (2) HD methods gain access to the real-time internal states of the generating model, and (3) the evaluation remains perpetually relevant regardless of model evolution. This insight directly motivates the design of UniFact.
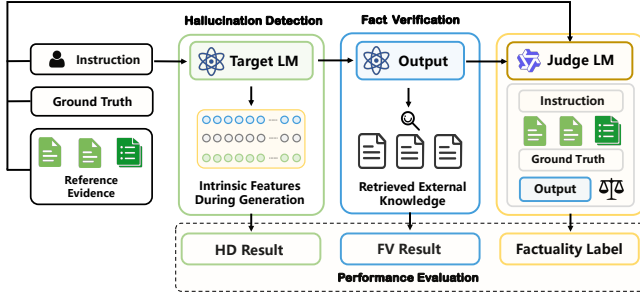
**Figure 2: An illustration of the UniFact evaluation framework. The pipeline generates answers dynamically and evaluates both hallucination detection and fact verification against a shared factuality label.**

## 3.2 Design Principles

To address the limitations of static benchmarks, our unified framework is built upon three foundational design principles that ensure a robust, equitable, and future-proof evaluation.

(1) **Model-Agnosticism via On-the-Fly Generation.** To combat model obsolescence, the framework abandons static datasets of outputs. Instead, it uses a core set of factual questions as the instruction to generate outputs dynamically from *any* target LLM. This ensures the framework's perpetual relevance for evaluating methods on current and future LLMs.

(2) **Equitable and Principled Comparison.** The framework is feasible for both hallucination detection and fact verification methods. It provides HD methods with access to the model's internal states during on-the-fly generation, while supplying FV methods with the generated text and an external corpus. Despite these different mechanisms, both are evaluated on the identical, ultimate task: predicting the factual correctness of a textual content (the output of that LLM).

(3) **Scalability through Automated Labeling.** To move beyond the slow and subjective process of manual annotation, the framework employs an automatic labeling system. It automatically determines the factuality of any generated response by comparing it against the ground-truth answer and relevant knowledge, enabling large-scale, objective, and reproducible evaluations across any number of LLMs.

## 3.3 Framework Architecture and Workflow

As illustrated in Figure 2, UniFact operates as a modular pipeline composed of three distinct stages: dynamic generation, automated annotation, and unified evaluation. This workflow transforms a static instruction set into a dynamic, reliably labeled benchmark that can evaluate both HD and FV paradigms.

*3.3.1 Stage 1: Dynamic Instance Generation.* The process initiates with an input triplet $(q, A^*, E^*)$, where $q$ is a factual question, $A^*$ represents the set of validated answers, and $E^*$ denotes authoritative evidence. Unlike static benchmarks, we treat $q$ as a prompt to trigger the target LLM $M_{target}$ to produce a response $y_{gen}$ in real-time.

To accommodate Hallucination Detection (HD) methods that require access to internal model states (white-box) or multiple

stochastic samples (black-box), this stage facilitates the real-time extraction of all relevant generative signals during the decoding process. Consequently, it allows for the construction of a model-specific instance that encapsulates both the textual output $y_{gen}$ and the requisite features, ensuring that downstream HD methods can obtain the necessary raw data for analysis.

*3.3.2 Stage 2: Reference-Based Automated Annotation.* To establish a ground-truth factuality label for the dynamically generated answer $y_{gen}$, we employ a reference-based verification procedure using an external evaluator model $M_{eval}$. The evaluator receives the tuple $(q, y_{gen}, A^*, E^*)$ and is tasked with a discriminative judgment: determining whether $y_{gen}$ is semantically consistent with the provided ground truth $A^*$ and evidence $E^*$. This process yields a rigorously annotated evaluation instance, formally denoted as $(q, y_{gen}, l^*)$, where $l^*$ represents the binary factuality label (Consistent/Inconsistent).

It is important to emphasize that the role of $M_{eval}$ is fundamentally distinct from that of $M_{target}$. While the target model engages in open-ended generation, $M_{eval}$ operates under a fixed, rubric-governed judgment protocol. The rubric specifies precise criteria for factual agreement, contradiction, partial support, and unsupported assertions, thereby enforcing a controlled and reproducible decision-making process. Crucially, $M_{eval}$ is granted privileged access to the ground-truth answer and evidence, allowing it to anchor its judgment in the authoritative reference materials rather than relying on its own parametric knowledge. As a result, the task reduces to a constrained consistency-checking problem rather than a knowledge-generation problem, substantially narrowing the space for evaluator hallucination and minimizing model subjectivity.

This evidence-rich setup also provides strong empirical reliability. As detailed in Section 3.4.3, human annotators exhibit 97.42% agreement with $M_{eval}$ on hallucination labels and 99.02% agreement on non-hallucination labels. These results confirm that, under a well-specified rubric and with explicit access to $(A^*, E^*)$, the evaluator model produces factuality labels $l^*$ that closely track human judgments and are sufficiently dependable for large-scale automated annotation.

*3.3.3 Stage 3: Unified Evaluation Interface.* The final stage bridges the paradigm gap by evaluating HD and FV methods on the identical prediction target $l^*$, despite their differing input requirements.

- **Hallucination Detection (HD):** These methods receive the generated text $y_{gen}$ along with the intrinsic model features captured in Stage 1. They are evaluated on their ability to detect non-factual content using only self-contained signals.
- **Fact Verification (FV):** These methods treat $y_{gen}$ as a claim to be verified. Crucially, to simulate a realistic verification scenario, FV methods are *not* given access to the ground-truth evidence $E^*$. Instead, they must retrieve supporting documents from an external corpus (e.g., Wikipedia) to predict the label.

This design ensures a fair yet distinct comparison: both paradigms aim to predict the same ground truth $l^*$, but HD is tested on its utilization of internal model uncertainty, while FV is tested on its ability to retrieve and reason over external knowledge.

## 3.4 Implementation Details

This section outlines key implementation details of our framework, covering dataset construction, judge model configuration, and the human verification procedure.

*3.4.1 Dataset Construction.* The foundation of our framework is a static set of diverse instructions, which we constructed by sourcing questions from five factual QA datasets. As our paradigm requires an automated system to verify the factual correctness of generated outputs, we deliberately selected datasets where each instance provides the three essential components for this process: a question (serving as the instruction), a ground-truth answer set, and relevant knowledge in the form of supporting relevant documents (enabling reliable factuality labeling). The five chosen datasets fulfill these criteria and were selected to cover a spectrum of hallucination challenges. TriviaQA (TQA) [23], NQ-Open (NQ) [29], and PopQA (PQA) [36] are open-domain benchmarks that assess factual recall across diverse generations. Meanwhile, 2WikiMultihopQA (2WQA) [19] evaluates multi-hop reasoning, for which we use the Bridge and Comparison question types (denoted as Bridge and Comp). HotpotQA (HQA) [77] also targets multi-hop reasoning, from which we select Comparison-type questions (referred to as HComp). For each dataset, we randomly sample 500 questions as test instances, ensuring a balanced and statistically meaningful evaluation. For the external knowledge base, we utilize the Wikipedia dumps processed by DPR [25][2] as the external corpus.

*3.4.2 Judge Model Configuration.* For the LLM-as-Judge component ($M_{\text{eval}}$), we employ the Qwen-2.5-32B model [75], a high-capacity LLM demonstrating strong instruction-following capabilities. Specifically, we utilize a carefully constructed prompt (see Appendix A.1) that explicitly instructs $M_{\text{eval}}$ to identify factual discrepancies based strictly on the provided ground-truth answer set $A^*$ and evidence $E^*$. This configuration ensures that the judge performs a reference-based verification rather than utilizing its internal parametric memory, thereby aligning the automated evaluation with the rubric-defined objective.

*3.4.3 Human Verification.* To validate the reliability of the Uni-Fact framework, we conducted a manual verification of the labels produced by the automatic labeling system illustrated in Figure 2. The human-annotation process is conducted as follows. First, we filtered out instances where the target LLM refused to answer (e.g., responses containing "I'm unable to xxx" or "I am not sure"). This resulted in a set of 1,602 generated answers, of which the LLM-as-Judge identified 890 as containing hallucinations and 712 as factually correct. We then tasked four undergraduate students to verify these labels manually. The human annotators agreed with the judge's "hallucination" label in 97.42% of cases and with the "non-hallucination" label in 99.02% of cases. This high level of agreement confirms the reliability of our auto-labeling approach. It is crucial to note that this high accuracy is achieved because the automatic labeling system is provided with the ground truth answer ($A^*$) as well as the evidence ($E^*$) from the dataset. This contrasts with the real-world scenario for FV methods, which must retrieve their own, often imperfect, evidence from an external corpus.

[2]https://github.com/facebookresearch/DPR/tree/main

## 4 Unified Empirical Study

In this section, we conduct a systematic empirical study to bridge the gap between Hallucination Detection and Fact Verification. We first introduce our study design and experimental setup. We then present a detailed analysis of three proposed research questions, exploring the comparative performance, complementary strengths, and integration potential of these two paradigms.

### 4.1 Study Design

To bridge the historical gap between Hallucination Detection and Fact Verification, we conduct a large-scale empirical study designed to systematically compare these two paradigms within a unified framework. Our investigation is structured to provide clear and direct answers to the following three research questions (RQs):

- **RQ1: A Head-to-Head Performance Comparison.** How do SOTA Hallucination Detection and Fact Verification methods perform when evaluated on the same LLM-generated content? Are there consistent winners?
- **RQ2: Uncovering Complementary Strengths.** Do these two paradigms possess unique, complementary strengths in identifying different types of factual errors, or do their capabilities largely overlap?
- **RQ3: Performance of Hybrid Methods.** Can we create a more robust and accurate system by combining signals from both paradigms, outperforming any single approach?

To answer these questions, we organize our empirical study into three logical and progressive steps. First, to answer **RQ1**, we implement a wide array of classic and SOTA baselines from both paradigms and benchmark their performance head-to-head on a unified testbed. This provides the first direct comparison of their relative effectiveness. Second, building on these initial results, we quantitatively analyze the synergy between the two paradigms to address **RQ2**. We demonstrate that the mutual enhancement achieved by integrating HD and FV significantly surpasses the gains from combining methods within a single paradigm. Furthermore, we conduct a detailed case study in Section 4.5 to add qualitative depth to our analysis of RQ2. Finally, motivated by the strong evidence of their synergy, we investigate **RQ3** by designing and evaluating two simple yet effective hybrid methods. We test fusion at both the final probability score level and the system pipeline level, showing that even straightforward integration can significantly boost overall performance.

### 4.2 Experimental Setup

*4.2.1 Backbone Models and Implementation Details.* All experiments are conducted using open-source pre-trained large language models (LLMs). We select LLMs of varying sizes and architectural series, including **Qwen2.5-14B-Instruct** [75] and **LLaMA-3.1-8B-Instruct** [38]. All models are deployed via their official Hugging Face implementations. We adopt the default hyperparameters and chat templates provided in the official repositories, with the only modification being the use of greedy decoding to ensure reproducibility of the results. All experiments are conducted on NVIDIA A100 GPUs with 40GB VRAM.

**Table 1: Performance of Hallucination Detection and Fact Verification methods under our proposed unified evaluation framework. The best results in each column are highlighted in bold, and the second-best results are underlined. The reported metrics in the table are Area Under the Curve (AUC) scores.**

| | | Meta-Llama-3.1-8B-Instruct | | | | | | Qwen2.5-14B-Instruct | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | Bridge | Comp | HComp | NQ | PQA | TQA | Bridge | Comp | HComp | NQ | PQA | TQA |
| HD | LNPP | 0.6557 | 0.6950 | 0.7259 | 0.7361 | 0.6103 | 0.7476 | 0.5278 | 0.7189 | 0.7062 | 0.6466 | 0.6129 | 0.6912 |
| | LNPE | 0.6991 | 0.7648 | 0.7540 | 0.7370 | 0.5853 | 0.7524 | 0.4953 | <u>0.7389</u> | 0.7211 | 0.6820 | 0.6242 | 0.7022 |
| | SCG-MQA | 0.4785 | 0.6799 | 0.6841 | 0.7082 | 0.6675 | 0.7290 | 0.5351 | 0.6294 | 0.6386 | 0.6129 | 0.6026 | 0.7291 |
| | SCG-NG | 0.3598 | 0.3291 | 0.4152 | 0.6553 | 0.6649 | 0.7109 | 0.4266 | 0.5559 | 0.5424 | 0.6012 | 0.6701 | 0.6524 |
| | SCG-BS | 0.4097 | 0.4194 | 0.6229 | 0.7631 | 0.5750 | 0.7672 | 0.5333 | 0.5566 | 0.6674 | 0.6844 | 0.6541 | 0.6942 |
| | SCG-NLI | 0.5082 | 0.5870 | 0.7010 | 0.7149 | 0.5513 | 0.7537 | 0.5701 | 0.6209 | 0.6722 | 0.6401 | 0.4790 | 0.7031 |
| | SAPLMA | 0.4335 | 0.4714 | 0.6451 | 0.5758 | 0.7356 | 0.6013 | 0.5365 | 0.4368 | 0.5969 | 0.5357 | 0.5505 | 0.6164 |
| | EUBHD | 0.5904 | 0.6442 | 0.6591 | 0.6714 | 0.7012 | 0.6613 | 0.4546 | 0.5607 | 0.7154 | 0.6836 | **0.8130** | 0.6639 |
| | MIND | 0.5675 | 0.6898 | 0.7296 | 0.6098 | 0.8043 | 0.7472 | <u>0.7214</u> | 0.4790 | 0.6126 | 0.6232 | 0.5560 | 0.6242 |
| | PTrue | 0.5185 | 0.5027 | 0.4403 | 0.5593 | 0.2880 | 0.3944 | 0.5400 | **0.7446** | 0.6965 | 0.7048 | 0.6239 | 0.7325 |
| | SE | 0.4446 | 0.4851 | 0.4260 | 0.5436 | 0.3141 | 0.5280 | 0.5317 | 0.5519 | 0.5340 | 0.5195 | 0.4339 | 0.5351 |
| | SEU | 0.5136 | 0.6669 | 0.7085 | 0.7212 | 0.5623 | 0.7508 | 0.5291 | 0.5914 | 0.7284 | 0.7019 | 0.5533 | 0.7449 |
| | SIndex | 0.3991 | 0.5299 | 0.6004 | 0.6938 | 0.3156 | 0.6594 | 0.4959 | 0.5939 | 0.6501 | 0.6716 | 0.4311 | 0.7128 |
| FV | LLM-Q | 0.5034 | 0.5162 | 0.6456 | 0.7405 | 0.5080 | 0.6954 | 0.5901 | 0.5840 | 0.7392 | 0.7086 | 0.6963 | 0.7907 |
| | BERT-Q | 0.6485 | 0.7607 | 0.7409 | 0.6154 | <u>0.8256</u> | 0.7215 | 0.4946 | 0.6399 | 0.7099 | 0.5956 | 0.7617 | 0.7019 |
| | LLM-QA | 0.5354 | 0.5707 | 0.6597 | 0.7537 | 0.4874 | 0.7019 | 0.5054 | 0.5881 | <u>0.7691</u> | 0.7474 | 0.6831 | 0.8120 |
| | BERT-QA | 0.6362 | 0.7813 | 0.7393 | 0.6165 | 0.8141 | 0.7364 | 0.5086 | 0.6367 | 0.6986 | 0.6290 | 0.7518 | 0.7338 |
| Unify | BERT-Q+EUBHD | 0.6414 | 0.8126 | 0.7612 | 0.6846 | **0.8552** | 0.7485 | 0.4688 | 0.6391 | 0.7614 | 0.6609 | 0.8074 | 0.7394 |
| | BERT-Q+LNPE | **0.7268** | <u>0.8128</u> | <u>0.7902</u> | 0.7062 | 0.7902 | <u>0.7904</u> | 0.4845 | 0.7285 | 0.7668 | 0.6901 | 0.7830 | 0.7563 |
| | BERT-QA+MIND | 0.5738 | 0.7669 | 0.7572 | 0.6300 | 0.8063 | 0.7692 | **0.7239** | 0.5726 | 0.6801 | 0.6500 | 0.6839 | 0.6904 |
| | LLM-QA+LNPE | 0.6369 | 0.6903 | 0.7540 | **0.8222** | 0.5392 | 0.7902 | 0.5004 | 0.6687 | **0.8230** | **0.8028** | 0.7062 | **0.8597** |
| | Pipeline | <u>0.7026</u> | **0.8176** | **0.7989** | <u>0.7877</u> | 0.8089 | **0.8252** | 0.4958 | 0.7239 | 0.7491 | <u>0.7567</u> | <u>0.8098</u> | <u>0.8175</u> |

*4.2.2 Hallucination Detection Baselines.* We evaluate the following set of hallucination detection baselines. **Semantic Entropy (SE)** [27] computes token-level entropy and applies clustering to refine uncertainty estimates. **Semantic Embedding Uncertainty (SEU)** [18] captures response variability through pairwise embedding similarity analysis. **SINdex** [1] extends SE by modeling both intra-group and inter-group inconsistencies. **Length-Normalized Predictive Entropy (LNPE)** [35] estimates predictive entropy and normalizes it by sequence length. **PTrue** [24] estimates the probability that a generation is truthful using token-level likelihoods. **Length-Normalized Predictive Probability (LNPP)** [37] evaluates hallucination likelihood based on length-normalized token probabilities. **SAPLMA** [4] trains a classifier that detects hallucination based on the LLM's activation values. **MIND** [60] detects an LLM's hallucination based on the final layer's token embeddings. **EUBHD** [80] represents a state-of-the-art uncertainty-based method with enhanced predictive distribution modeling. Finally, we evaluate the **SelfCheckGPT (SCG)** [37] family, which detects hallucinations based on sampling consistency. We consider four specific variants: SCG-BS, SCG-MQA, SCG-NLI, and SCG-NG.

*4.2.3 Fact Verification Baselines.* To provide a comprehensive comparison, we implement baselines representing the two dominant paradigms in automated fact-checking: LLM-based Verification [39,

74]) and NLI-based Verification [49, 65]). To ensure a fair comparison within our unified framework, we standardize the retrieval module across all methods, using BM25 [44] for evidence retrieval from the DPR Wikipedia corpus [25]. Implementation details for all baselines are provided in Appendix B.2.

- **LLM-based Verification.** Recent works such as FActScore [39] and SAFE [74] have demonstrated the effectiveness of using LLMs to verify claims against retrieved evidence. Following this paradigm, we implement two variants based on their query formulation strategies:
  - **LLM-Q:** This baseline adopts the standard "Retrieve-then-Verify" workflow. Given a question, we retrieve evidence passages using the question as the query. An evaluator LLM is then prompted to verify whether the target LLM's generated response is supported by the retrieved evidence.
  - **LLM-QA:** This variant retrieves evidence using both the question and the generated answer. This allows the verifier to capture evidence relevant to claims made in the response. The same evaluator LLM is used for the final verification.
- **NLI Verification.** We also evaluate the classical "Retriever-Reader" architecture, which remains a strong baseline in fact-checking benchmarks like FEVER [65]. Adapting the architecture from Soleimani et al. [49], we implement two variants that utilize a fine-tuned BERT classifier as the verifier, distinguished by their retrieval strategies:

- **BERT-Q:** This baseline retrieves evidence passages using only the question as the query. The retrieved passages are then concatenated with the question and answer to form the input for a BERT-based classifier, which predicts whether the answer contains factual errors.
- **BERT-QA:** This variant retrieves evidence using both the question and the generated answer. Similar to LLM-QA, this strategy aims to capture evidence that is semantically closer to the specific claims in the response. The same BERT-based classifier is employed to predict the verification label based on the retrieved context.

*Implementation Note on Standardization.* While these baselines are conceptually grounded in SOTA fact verification frameworks like FActScore [39] and BERT-based verification systems [49], we adapt their implementations to fit a unified evaluation setting. Specifically, we standardize the retrieval backbone (BM25) and the knowledge corpus (Wikipedia) across all experiments. This control eliminates potential bias introduced by different retrieval systems or proprietary databases used in original implementations. Consequently, to distinguish our standardized adaptations from the specific system configurations in prior work, we denote them using descriptive names (e.g., LLM-QA, BERT-Q) rather than their original names (e.g., FActScore).

### 4.3 RQ1: A Head-to-Head Performance Analysis

The results presented in Table 1 reveal a distinct performance divergence between Hallucination Detection (HD) and Fact Verification (FV), indicating that neither paradigm offers a universally superior solution. We observe that the performance of HD methods is unstable and depends largely on the specific LLM being evaluated. For instance, methods that achieve high scores on the LLaMA family often see a significant drop when applied to Qwen models. This instability suggests that reliance on internal signals (e.g., predictive entropy) for hallucination detection may not transfer reliably across different model families. In contrast, FV methods demonstrate greater stability. Notably, incorporating the generated answer into retrieval queries (the QA-based strategy) consistently outperforms question-only retrieval across most datasets. This confirms that for FV, the primary bottleneck is finding the right evidence to match the claim, rather than the model architecture itself.

Beyond the stability trends, a fine-grained comparison of individual datasets further confirms that neither paradigm maintains a strict dominance. As detailed in Table 1, the "winner" fluctuates significantly depending on the dataset and model combination. For instance, on Meta-Llama-3.1-8B-Instruct, HD methods secure the top performance on four out of six datasets, whereas FV methods are superior on two out of six datasets (e.g., PQA and Comp). However, this advantage disappears on Qwen2.5-14B-Instruct, where the two paradigms split the victories evenly. This lack of a consistent winner strongly suggests that the two paradigms are sensitive to different types of factual errors. Rather than being redundant, their performance divergence implies a potential for orthogonality, which directly motivates our quantitative analysis of their complementarity in the following section (RQ2). Note that the performance of hybrid methods (the "Unify" block) is distinct from this head-to-head comparison and will be discussed in detail under RQ3.

### 4.4 Paradigm Synergy Analysis (RQ2)

To quantitatively address **RQ2**, we conduct a comparative analysis of method synergy within and across HD and FV paradigms. Specifically, we investigate whether combining methods from different paradigms (i.e., HD+FV) yields greater complementarity than combining methods within the same paradigm (i.e., HD+HD or FV+FV). To this end, we use three complementary metrics to evaluate the synergy between an arbitrary pair of methods, where the methods are denoted as $M_A$ and $M_B$. Let $C_A$ and $C_B$ be the sets of samples correctly classified by $M_A$ and $M_B$, respectively (regardless of whether they belong to the same or different paradigms), and let $W_A$ and $W_B$ be the corresponding sets of incorrectly classified samples. The total set of samples is $S$. We formalize the synergy between any two methods using the following three metrics:

- **Average Complementarity Score (ACS) [28].** This metric measures the proportion of samples correctly classified by exactly one of the two methods, quantifying their degree of non-overlapping correctness. Equivalently, ACS corresponds to the classical disagreement rate between two classifiers.

$$\text{ACS}(M_A, M_B) = \frac{|(C_A \cap W_B) \cup (W_A \cap C_B)|}{|S|}. \quad (3)$$

- **Average Synergy Gain (ASG) [46].** This metric captures the performance improvement of an ideal oracle ensemble over the best-performing individual method. It reveals the potential accuracy gain achievable in principle by combining $M_A$ and $M_B$.

$$\text{ASG}(M_A, M_B) = \frac{|C_A \cup C_B|}{|S|} - \max\left(\frac{|C_A|}{|S|}, \frac{|C_B|}{|S|}\right). \quad (4)$$

In other words, ASG is the fraction of samples that are correctly classified by the oracle ensemble but misclassified by the better of the two individual methods.

- **Average Error Correction Rate (AECR) [28].** This metric quantifies the capacity of one method to compensate for another's errors. We define the directional error correction rate of method $M_A$ for $M_B$, denoted as $R(A \text{ for } B)$, as the conditional probability that $M_A$ is correct given $M_B$ has failed:

$$R(A \text{ for } B) = P(C_A \mid W_B) = \frac{|C_A \cap W_B|}{|W_B|}, \quad (5)$$

$$R(B \text{ for } A) = P(C_B \mid W_A) = \frac{|C_B \cap W_A|}{|W_A|}. \quad (6)$$

The AECR is the symmetrical average of these two directional rates[3]:

$$\text{AECR}(M_A, M_B) = \frac{1}{2}\left(R(A \text{ for } B) + R(B \text{ for } A).\right) \quad (7)$$

For each of the three metrics, we compute the value for all unordered pairs of methods within the HD paradigm, within the FV paradigm, and across both paradigms. We then report the average over method pairs in each group. The resulting scores are presented in Table 2, providing quantitative evidence that the HD and FV paradigms possess complementary strengths. The notably higher ACS for cross-paradigm pairings indicates that their successes are more mutually exclusive than those of methods within the same

---

[3]AECR is applied only to pairs where both methods make at least one error (i.e., $|W_A|, |W_B| > 0$), which holds for all methods in our experiments.

**Table 2: Synergy and complementarity analysis between and within paradigms. Cross-paradigm pairings consistently outperform intra-paradigm comparisons across all metrics, demonstrating strong complementary strengths.**

| Method Pairing | ACS | ASG | AECR |
|---|---|---|---|
| **Intra-HD Average** | 0.315 | 0.118 | 0.503 |
| **Intra-FV Average** | 0.379 | 0.102 | 0.496 |
| **Cross-Paradigm Average** | **0.428** | **0.144** | **0.634** |

paradigm, suggesting that they excel on different subsets of the problem space. Similarly, the superior ASG score reveals that combining HD and FV methods offers the greatest potential for performance improvement, exceeding any intra-paradigm combination. Most importantly, the AECR exhibits a substantial increase for cross-paradigm pairings. This implies that the failure modes of the two paradigms are largely distinct: when one paradigm makes an error, the other is considerably more likely to be correct. Overall, this analysis supports the view that HD and FV are not merely different but genuinely complementary, providing an empirical basis for the development of hybrid models that integrate both paradigms.

## 4.5 Case Study

While the previous section established the statistical complementarity between HD and FV, this section investigates the underlying causes of their divergent performance. Through a detailed case study, we examine specific scenarios where each paradigm tends to fail. We highlight two distinct failure patterns: Fact Verification is fundamentally limited by its reliance on retrieval quality, often leading to errors when relevant evidence is missing. Conversely, Hallucination Detection is prone to false alarms caused by semantic flexibility, where the model's diverse phrasing is misinterpreted as hallucination. Analyzing these patterns helps clarify the nature of their synergy and motivates the design of the hybrid integration methods proposed in the subsequent section.

*4.5.1 Fact Verification.* Our analysis suggests that the effectiveness of FV methods is fundamentally limited by the performance of the retrieval module and the availability of relevant external knowledge. As demonstrated in Table 3, this reliance on retrieval creates a critical vulnerability when the system fails to provide necessary evidence. Specifically, when the retriever cannot locate relevant passages (for instance, due to keyword mismatch in the LLM-Q setting), the verification module is compelled to make judgments without sufficient context. This absence of evidence results in unstable model behavior: the model may either show unfounded skepticism (leading to false alarms, as in Case 1) or accept false claims without verification (leading to missed errors, as in Case 2). In contrast, when the retrieval process is augmented with the generated answer (LLM-QA), the system is more likely to obtain relevant supporting documents, thereby allowing the FV model to make more robust predictions.

*4.5.2 Hallucination Detection.* Conversely, HD methods are prone to false alarms triggered by the inherent semantic flexibility of language generation. Table 4 reveals a clear dichotomy: HD methods

**Table 3: Failure Analysis of Fact Verification (FV). FV judgments are unstable and heavily dependent on retrieval quality. When retrieval fails to surface relevant evidence, FV models act erratically, leading to both False Alarms (Case 1) and Missed Detections (Case 2).**

| Impact of Retrieval Quality on Verification Verdicts | |
|---|---|
| *Case 1: Retrieval Failure → False Alarm (Unjustified Skepticism)* | |
| **Question** | Who plays Cindy Lou Who on the Grinch? |
| **LLM Output** | The role of Cindy Lou Who in the 2000 live-action film "How the Grinch Stole Christmas" is played by **Taylor Momsen**. |
| **Ground Truth** | Taylor Michel Momsen *(Label: Factual)* |
| **State: Retrieval Failure** | **Evidence:** *[No relevant passage regarding Taylor Momsen]* <br> **Verdict: Hallucination** *(Error: False Alarm)* |
| **State: Retrieval Success** | **Evidence:** *Cindy Lou Who is a generous young girl who was introduced in the book "How the Grinch Stole Christmas!" In the 2000 live action film, "How the Grinch Stole Christmas!" she is played by actress Taylor Momsen.* <br> **Verdict: Factual** *(Correct)* |
| *Case 2: Retrieval Failure → Missed Detection (Blind Confirmation)* | |
| **Question** | Where did the theme song from SWAT come from? |
| **LLM Output** | The theme song from the popular 1970s-80s TV show "S.W.A.T." was composed by **Michel Legrand and Al Burton**. |
| **Ground Truth** | Written by Barry De Vorzon *(Label: Hallucination)* |
| **State: Retrieval Failure** | **Evidence:** *[No relevant passage regarding composer]* <br> **Verdict: Factual** *(Error: Missed Detection)* |
| **State: Retrieval Success** | **Evidence:** *Theme from "S.W.A.T." is an instrumental song written by Barry De Vorzon and performed by American funk group Rhythm Heritage, released on their debut album "Disco-Fied".* <br> **Verdict: Hallucination** *(Correct)* |

perform robustly on deterministic answers (Group A) but degrade significantly on expressive responses (Group B). In the latter case, a model may be factually confident about an event (e.g., the formation year of a band) yet lexically uncertain about how to articulate it. Current HD methods struggle to disentangle this valid semantic diversity from actual falsehoods. As a result, the high entropy caused by diverse phrasing is misinterpreted as hallucination, causing the detector to flag correct generations solely because they are phrased with variety.

*4.5.3 Implications for Hybrid Design.* Our qualitative analysis reveals that HD and FV operate with orthogonal error boundaries, providing the theoretical grounding for their integration. FV functions as a precision-oriented anchor: when it successfully retrieves evidence, its judgment is highly reliable, but it performs poorly when the retrieval module fails to fetch relevant evidence. Conversely, HD functions as a recall-oriented safety net: it can capture potential errors even without external knowledge, but suffers from noise due to semantic flexibility. Therefore, an optimal hybrid strategy should not assign equal importance to both. Instead, it should prioritize FV's judgment when retrieved evidence is relevant, while relying on HD's intrinsic signals when the retrieval result is not relevant. This conditional complementarity directly motivates the hierarchical design of our proposed hybrid framework in RQ3.

**Table 4: Failure Analysis of Hallucination Detection (HD). HD methods struggle to distinguish *Factual Uncertainty* from *Lexical Uncertainty*. While they excel at verifying deterministic facts (Group A), they frequently generate False Alarms on semantically flexible responses (Group B).**

| Impact of Semantic Flexibility on HD Verdicts |
| --- |
| **Group A: Deterministic Facts (Low Lexical Entropy)** |
| *Characteristic: Single, concise answers. Low token-level uncertainty.* |
| **Q1** How many states are in the United States? <br> **LLM** There are 50 states in the United States. <br> **Result** **Correct Verdict** (Low Uncertainty → Factual) |
| **Q2** What is the water cycle also known as? <br> **LLM** The water cycle is also known as the hydrologic cycle. <br> **Result** **Correct Verdict** (Low Uncertainty → Factual) |
| **Group B: Semantically Flexible Responses (High Lexical Entropy)** |
| *Characteristic: Valid facts expressed with high variation (word choice/order), causing false spikes in uncertainty metrics.* |
| **Q3** Which band formed first, Awolnation or Foo Fighters? <br> **LLM** **Foo Fighters formed first.** Foo Fighters were formed in 1994 by Nirvana's drummer Dave Grohl, following the band's dissolution... <br> **Truth** Foo Fighters formed first. <br> **Result** **False Alarm** (High Uncertainty due to phrasing → Hallucination) |
| **Q4** Who directed the classic 30s western Stagecoach? <br> **LLM** **The classic 1939 western 'Stagecoach' was directed by John Ford.** It is considered one of the greatest Westerns ever made... <br> **Truth** John Ford. <br> **Result** **False Alarm** (High Uncertainty due to phrasing → Hallucination) |

## 4.6 RQ3: Investigation of Hybrid Strategies

Building on the findings in RQ2, which demonstrated that Hallucination Detection (HD) and Fact Verification (FV) exhibit statistical complementarity and mechanistically distinct failure modes, we now address **RQ3**: Can we operationalize this synergy into high-performance hybrid systems? To this end, we propose and evaluate two integration paradigms: a **Score-Level Fusion** approach and an **Evidence-Aware Pipeline** designed to address the specific structural limitations of FV identified in our case study.

### 4.6.1 Hybrid Method Design.

*Strategy I: Score-Level Fusion.* Since HD and FV capture different aspects of hallucination, merging their signals can smooth out individual errors. We propose a straightforward linear combination of their normalized scores. Let $S_{\text{HD}}$ and $S_{\text{FV}}$ be the probability scores from a Hallucination Detection method (e.g., LNPE) and a Fact Verification method (e.g., LLM-QA), respectively. To ensure consistency, we map both scores to represent the likelihood of a hallucination (i.e., higher scores indicate lower quality). The hybrid score is calculated as:

$$S_{\text{Hybrid}} = \lambda S_{\text{HD}} + (1 - \lambda)S_{\text{FV}}. \tag{8}$$

In our experiments, we set $\lambda = 0.5$, effectively taking the average of the two signals. We adopt this unweighted setting to demonstrate that the hybrid method is robust and effective even without dataset-specific hyperparameter tuning. We test pairings such as **LLM-QA + LNPE** to verify these gains.

*Strategy II: Evidence-Aware Pipeline.* While Strategy I enhances robustness through score-level fusion, it does not explicitly address the failure modes identified in Section 4.5. Specifically, FV systems are reliable when evidence is explicitly retrieved (yielding "Supported" or "Contradicted" verdicts) but become unreliable when relevant evidence is absent (i.e., "Not Enough Information" or NEI). To address this limitation, we design an *Evidence-Aware Pipeline* that prioritizes external grounding, falling back to internal signals only when retrieval fails:

(1) **Step 1: Retrieval & Verification.** An FV module (specifically LLM-QA) retrieves evidence and classifies the claim–evidence relationship as *Supported*, *Contradicted*, or *NEI*.
(2) **Step 2: Conditional Branching.** We employ a cascading decision mechanism. If the FV yields a definitive prediction (*Supported* or *Contradicted*), the system directly accepts this result. Conversely, if the output is *NEI*, the system acknowledges the lack of external support. Instead of forcing a potentially erroneous verification, the pipeline bypasses the FV and falls back to an HD method (e.g., **LNPE**), deriving the final score solely from the model's internal uncertainty signals.

This conditional design effectively mitigates the limitations of FV when retrieval results are not relevant. Crucially, decoupling retrieval dependence from internal signals prevents the propagation of errors caused by low-quality evidence retrieval, ensuring robust detection even when external knowledge is unavailable.

### 4.6.2 Performance Analysis.
The experimental results in Table 1 validate the efficacy of integrating model-centric and data-centric paradigms. A comprehensive examination of the "Unify" block against individual baselines reveals that hybrid strategies consistently establish a new upper bound for performance across diverse datasets and model architectures. This consistent superiority confirms that internal uncertainty signals and external evidence retrieval operate in orthogonal feature spaces. Regarding the specific strategies, the Score-Level Fusion approach demonstrates that a simple linear combination can effectively mitigate the high variance inherent in single-paradigm methods. By aggregating divergent signals, this strategy smoothens the noise from individual prediction errors, yielding a more stable evaluation that outperforms standalone HD or FV methods in most scenarios. Furthermore, the Evidence-Aware Pipeline exhibits an even more distinct advantage, frequently achieving the optimal performance among all compared methods. This performance leap can be attributed to the pipeline's conditional mechanism, which addresses the fundamental "blind spot" of retrieval-based verification. In instances where external knowledge is insufficient (the NEI cases), forcing a verification verdict often introduces noise; the pipeline avoids this by dynamically retreating to internal uncertainty signals, which remain a reliable proxy for truthfulness in the absence of evidence.

Notably, hybrid methods mitigate the performance instability observed in RQ1. Unlike standalone HD methods, which exhibit fluctuations between LLaMA and Qwen families, hybrid approaches maintain robustness across diverse backbones. This suggests that combined signals effectively compensate for the miscalibration of internal confidence inherent in standalone HD. Consequently, the proposed frameworks not only enhance accuracy but also offer a generalized solution independent of model-specific variations.

# 5 Related Work

This section reviews existing work in Fact Verification (FV) and Hallucination Detection (HD), highlighting the lack of unified benchmarks that motivates our proposed **UniFact** framework.

## 5.1 Fact Verification

*5.1.1 Traditional Pipeline-based Methods.* The establishment of the FEVER benchmark [65] formalized FV as a multi-stage pipeline consisting of document retrieval, evidence selection, and textual entailment (NLI). Early approaches, such as the baseline provided by Thorne et al. [65], relied on TF-IDF retrieval and simple classifiers, suffering from error propagation between stages. Subsequent works addressed these limitations through end-to-end neural architectures. For instance, NSMN [41] proposed a unified neural semantic matching network for joint retrieval and verification, while GEAR [82] introduced graph-based reasoning to model dependencies across multiple evidence sentences. With the advent of pre-trained language models, BERT-based approaches [49] further improved performance by enhancing contextual representations for both retrieval and claim verification. Extensions like FEVEROUS [3] and SciFact [69] expanded this paradigm to structured data (tables) and scientific domains, respectively.

*5.1.2 Fact Verification in the Era of LLMs.* With the rapid advancement of LLMs in natural language processing, an increasing number of studies explore their application to fact verification. Techniques such as retrieval-augmented generation (RAG [13, 53, 54, 58, 59, 66, 71]) and few-shot in-context learning (ICL [7, 14, 72]) enable LLMs to perform accurate truth assessments and generate supporting evidence by incorporating relevant contextual knowledge. These methods typically use a retrieval module (either a simple retriever [15, 34, 44, 50–52, 79] or a more sophisticated retrieval system [10, 45, 56, 67, 68]) to identify evidence related to the statement from knowledge sources and then provide both the statement and retrieved evidence as input to the LLM, guiding it to make judgments through ICL. For instance, Singal et al. [48] proposes a method that combines evidence extraction via retrieval with truth prediction using ICL with LLMs. Deng et al. [12] further investigates how to improve the reliability of RAG systems by reordering retrieved documents based on credibility, thereby mitigating the impact of erroneous information on LLM predictions. Furthermore, Zhang and Gao [81] designs a hierarchical prompting method that guides LLMs to decompose complex statements into more manageable sub-statements and verify each progressively, aiming to improve the robustness of fact-checking complex claims.

## 5.2 Hallucination Detection

*5.2.1 White-Box Methods: Leveraging Internal States.* White-box methods require access to the model's parameters and activations [60, 70]. Early works focused on token-level probability distributions, using metrics like perplexity or entropy to quantify uncertainty [35]. For instance, EUBHD [80] propose an uncertainty-based hallucination detection method based on predictive distribution modeling. More recent approaches leverage deeper internal representations. SAPLMA [5] and MIND [61] train supervised classifiers based on the model's hidden states to detect hallucinations. IN-SIDE [8] analyzes the covariance matrix of internal states across multiple sampled outputs to detect inconsistencies. From another angle, HD-NDEs [2] employs neural differential equations to model the dynamic evolution of hidden states.

*5.2.2 Black-Box Methods: Behavioral Consistency.* For closed-source models where internal states are inaccessible, black-box methods rely on sampling-based consistency checks. The prevailing hypothesis is that if an LLM "knows" a fact, it will generate consistent answers across multiple stochastic samples; if it is hallucinating, the answers will diverge. SelfCheckGPT [37] is a representative approach that samples multiple responses and measures their consistency using NLI or LLM-based prompting. Other methods, such as SEU [18], quantify this consistency via semantic embeddings, while InterrogateLLM [78] reconstructs the query from the answer to check for semantic drift. While HD methods excel at capturing the model's intrinsic uncertainty, they often lack external grounding. A model can be "confidently wrong" (low uncertainty, high consistency) due to mimicked falsehoods in pre-training data, a failure mode that internal signals alone often fail to detect.

## 5.3 The Evaluation Gap

A fundamental barrier to unifying these two paradigms lies in the incompatibility of existing benchmarks. Traditional FV benchmarks like FEVER [65] consist of static claim-evidence pairs. While suitable for text-based verification, they are unusable for HD methods. Conversely, benchmarks designed for LLM evaluation, such as TruthfulQA [32] or HaluEval [31], often rely on pre-generated text or fixed multiple-choice formats. These static benchmarks suffer from rapid obsolescence; they cannot evaluate HD methods on new architectures (e.g., LLaMA-3 vs. LLaMA-2) nor facilitate the extraction of real-time internal signals required for white-box detection. This dichotomy has forced researchers to evaluate HD and FV in isolation. To address this, our framework, **UniFact**, introduces a dynamic evaluation paradigm. By generating responses on-the-fly and automatically verifying them against external references, UniFact simultaneously provides the internal signals required by HD and the textual evidence required by FV, enabling the first rigorous head-to-head comparison of these complementary approaches.

# 6 Conclusion

In this work, we bridge the historical gap between Hallucination Detection (HD) and Fact Verification (FV) by introducing UniFact, a unified dynamic factuality evaluation framework for LLMs. Through the first direct comparison of these two paradigms, we reveal that they represent orthogonal dimensions of factuality assessment rather than redundant approaches. Specifically, while FV provides robust grounding when external evidence is available, HD leverages internal uncertainty to serve as a crucial safeguard when retrieval fails. Capitalizing on this synergy, we demonstrate that integrating these distinct signals establishes a new state-of-the-art in factual error detection. Ultimately, this paper advocates for a paradigm shift from isolated research tracks toward unified factuality assessment. We hope UniFact serves as a foundation for future research, paving the way for trustworthy AI systems that seamlessly integrate internal introspection with external knowledge.

# References

[1] Samir Abdaljalil, Hasan Kurban, Parichit Sharma, Erchin Serpedin, and Rachad Atat. 2025. SINdex: Semantic INconsistency Index for Hallucination Detection in LLMs. *CoRR* abs/2503.05980 (2025). arXiv:2503.05980 doi:10.48550/ARXIV.2503.05980

[2] ACL ARR 2024 December Submission581 Authors. 2024. HD-NDEs: Neural Differential Equations for Hallucination Detection in LLMs. ACL Rolling Review.

[3] Rami Aly, Zhijiang Guo, Michael Schlichtkrull, James Thorne, Andreas Vlachos, Christos Christodoulopoulos, Oana Cocarascu, and Arpit Mittal. 2021. Feverous: Fact extraction and verification over unstructured and structured information. *arXiv preprint arXiv:2106.05707* (2021).

[4] Amos Azaria and Tom Mitchell. 2023. The internal state of an llm knows when its lying. *arXiv preprint arXiv:2304.13734* (2023).

[5] Amos Azaria and Tom M. Mitchell. 2023. The Internal State of an LLM Knows When It's Lying. In *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, 967–976. doi:10.18653/V1/2023.FINDINGS-EMNLP.68

[6] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.

[7] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, et al. 2020. Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.).

[8] 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, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net. https://openreview.net/forum?id=Zj12nzlQbz

[9] Jiangui Chen, Ruqing Zhang, Jiafeng Guo, Yixing Fan, and Xueqi Cheng. 2022. GERE: Generative evidence retrieval for fact verification. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2184–2189.

[10] Xuesong Chen, Ziyi Ye, Xiaohui Xie, Yiqun Liu, Xiaorong Gao, Weihang Su, Shuqi Zhu, Yike Sun, Min Zhang, and Shaoping Ma. 2022. Web search via an efficient and effective brain-machine interface. In *Proceedings of the fifteenth ACM international conference on web search and data mining*. 1569–1572.

[11] Zhendong Chen, Siu Cheung Hui, Fuzhen Zhuang, Lejian Liao, Fei Li, Meihuizi Jia, and Jiaqi Li. 2022. EvidenceNet: Evidence fusion network for fact verification. In *Proceedings of the ACM web conference 2022*. 2636–2645.

[12] Boyi Deng, Wenjie Wang, Fengbin Zhu, Qifan Wang, and Fuli Feng. 2025. CrAM: Credibility-Aware Attention Modification in LLMs for Combating Misinformation in RAG. *Proceedings of the AAAI Conference on Artificial Intelligence* 39, 22 (Apr. 2025), 23760–23768. doi:10.1609/aaai.v39i22.34547

[13] Qian Dong, Qingyao Ai, Hongning Wang, Yiding Liu, Haitao Li, Weihang Su, Yiqun Liu, Tat-Seng Chua, and Shaoping Ma. 2025. Decoupling Knowledge and Context: An Efficient and Effective Retrieval Augmented Generation Framework via Cross Attention. In *Proceedings of the ACM on Web Conference 2025*.

[14] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Baobao Chang, et al. 2024. A survey on in-context learning. In *Proceedings of the 2024 conference on empirical methods in natural language processing*. 1107–1128.

[15] Yan Fang, Jingtao Zhan, Qingyao Ai, Jiaxin Mao, Weihang Su, Jia Chen, and Yiqun Liu. 2024. Scaling laws for dense retrieval. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1339–1349.

[16] Xueyang Feng, Zhi-Yuan Chen, Yujia Qin, Yankai Lin, Xu Chen, Zhiyuan Liu, and Ji-Rong Wen. 2024. Large language model-based human-agent collaboration for complex task solving. *arXiv preprint arXiv:2402.12914* (2024).

[17] Adam Fourney, Gagan Bansal, Hussein Mozannar, Cheng Tan, Eduardo Salinas, Friederike Niedtner, Grace Proebsting, Griffin Bassman, Jack Gerrits, Jacob Alber, et al. 2024. Magentic-one: A generalist multi-agent system for solving complex tasks. *arXiv preprint arXiv:2411.04468* (2024).

[18] Yashvir S. Grewal, Edwin V. Bonilla, and Thang D. Bui. 2024. Improving Uncertainty Quantification in Large Language Models via Semantic Embeddings. *CoRR* abs/2410.22685 (2024). arXiv:2410.22685 doi:10.48550/ARXIV.2410.22685

[19] Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing a multi-hop QA dataset for comprehensive evaluation of reasoning steps. *arXiv preprint arXiv:2011.01060* (2020).

[20] Nan Hu, Zirui Wu, Yuxuan Lai, Xiao Liu, and Yansong Feng. 2022. Dual-Channel Evidence Fusion for Fact Verification over Texts and Tables. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, 5232–5242. doi:10.18653/v1/2022.naacl-main.384

[21] Xuming Hu, Zhaochen Hong, Zhijiang Guo, Lijie Wen, and Philip Yu. 2023. Read it twice: Towards faithfully interpretable fact verification by revisiting evidence. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2319–2323.

[22] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2023. A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions. *CoRR* abs/2311.05232 (2023). arXiv:2311.05232 doi:10.48550/ARXIV.2311.05232

[23] Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. 2017. TriviaQA: A Large Scale Distantly Supervised Challenge Dataset for Reading Comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*. Association for Computational Linguistics, 1601–1611. doi:10.18653/V1/P17-1147

[24] Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, et al. 2022. Language Models (Mostly) Know What They Know. *CoRR* abs/2207.05221 (2022). arXiv:2207.05221 doi:10.48550/ARXIV.2207.05221

[25] Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906* (2020).

[26] Amrith Krishna, Sebastian Riedel, and Andreas Vlachos. 2022. ProoFVer: Natural Logic Theorem Proving for Fact Verification. *Transactions of the Association for Computational Linguistics* 10 (2022), 1013–1030. doi:10.1162/tacl_a_00503

[27] Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. Semantic Uncertainty: Linguistic Invariances for Uncertainty Estimation in Natural Language Generation. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*.

[28] Ludmila I Kuncheva and Christopher J Whitaker. 2003. Measures of diversity in classifier ensembles and their relationship with the ensemble accuracy. *Machine learning* 51, 2 (2003), 181–207.

[29] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics* 7 (2019), 453–466.

[30] Patrick Lewis, Ethan Perez, , et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems* 33 (2020), 9459–9474.

[31] Junyi Li, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. 2023. Halueval: A large-scale hallucination evaluation benchmark for large language models. *arXiv preprint arXiv:2305.11747* (2023).

[32] Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Truthfulqa: Measuring how models mimic human falsehoods. In *Proceedings of the 60th annual meeting of the association for computational linguistics (volume 1: long papers)*. 3214–3252.

[33] Qiang Liu, Xinlong Chen, Yue Ding, Shizhen Xu, Shu Wu, and Liang Wang. 2025. Attention-guided Self-reflection for Zero-shot Hallucination Detection in Large Language Models. arXiv:2501.09997 [cs.CL] https://arxiv.org/abs/2501.09997

[34] Yixiao Ma, Yueyue Wu, Weihang Su, Qingyao Ai, and Yiqun Liu. 2023. CaseEncoder: A Knowledge-enhanced Pre-trained Model for Legal Case Encoding. *arXiv preprint arXiv:2305.05393* (2023).

[35] Andrey Malinin and Mark Gales. 2020. Uncertainty estimation in autoregressive structured prediction. *arXiv preprint arXiv:2002.07650* (2020).

[36] 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 *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 9802–9822.

[37] Potsawee Manakul, Adian Liusie, and Mark JF Gales. 2023. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models. *arXiv preprint arXiv:2303.08896* (2023).

[38] Meta. 2024. Llama-3.2-1B-Instruct. https://huggingface.co/meta-llama/Llama-3.2-1B-Instruct Accessed: 2024-09.

[39] 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 *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. 12076–12100.

[40] Yixin Nie, Haonan Chen, and Mohit Bansal. 2019. Combining fact extraction and verification with neural semantic matching networks. In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence*. AAAI Press, Article 842, 8 pages. doi:10.1609/aaai.v33i01.33016859

[41] Yixin Nie, Haonan Chen, and Mohit Bansal. 2019. Combining fact extraction and verification with neural semantic matching networks. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 33. 6859–6866.

[42] Ray Oshikawa, Jing Qian, and William Yang Wang. 2020. A survey on natural language processing for fake news detection. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*. 6086–6093.

[43] Cheng Peng, Xi Yang, Aokun Chen, Kaleb E Smith, Nima PourNejatian, Anthony B Costa, Cheryl Martin, Mona G Flores, Ying Zhang, Tanja Magoc, et al. 2023. A study of generative large language model for medical research and healthcare.

*NPJ digital medicine* 6, 1 (2023), 210.

[44] Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: BM25 and beyond. *Foundations and Trends® in Information Retrieval* 3, 4 (2009), 333–389.

[45] Alireza Salemi and Hamed Zamani. 2024. Towards a search engine for machines: Unified ranking for multiple retrieval-augmented large language models. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 741–751.

[46] Catherine A Shipp and Ludmila I Kuncheva. 2002. Relationships between combination methods and measures of diversity in combining classifiers. *Information fusion* 3, 2 (2002), 135–148.

[47] Kai Shu, Suhang Wang, and Huan Liu. 2019. Beyond news contents: The role of social context for fake news detection. In *Proceedings of the twelfth ACM international conference on web search and data mining*. 312–320.

[48] Ronit Singal, Pransh Patwa, Parth Patwa, Aman Chadha, and Amitava Das. 2024. Evidence-backed Fact Checking using RAG and Few-Shot In-Context Learning with LLMs. In *Proceedings of the Seventh Fact Extraction and VERification Workshop (FEVER)*. Association for Computational Linguistics, Miami, Florida, USA, 91–98. doi:10.18653/v1/2024.fever-1.10

[49] Amir Soleimani, Christof Monz, and Marcel Worring. 2020. Bert for evidence retrieval and claim verification. In *European Conference on Information Retrieval*. Springer, 359–366.

[50] Weihang Su, Qingyao Ai, Xiangsheng Li, Jia Chen, Yiqun Liu, Xiaolong Wu, and Shengluan Hou. 2023. Wikiformer: Pre-training with Structured Information of Wikipedia for Ad-hoc Retrieval. *arXiv preprint arXiv:2312.10661* (2023).

[51] Weihang Su, Qingyao Ai, Yueyue Wu, Yixiao Ma, Haitao Li, and Yiqun Liu. 2023. Caseformer: Pre-training for Legal Case Retrieval. *arXiv preprint arXiv:2311.00333* (2023).

[52] Weihang Su, Qingyao Ai, Yueyue Wu, Anzhe Xie, Changyue Wang, Yixiao Ma, Haitao Li, Zhijing Wu, Yiqun Liu, and Min Zhang. 2025. Pre-training for Legal Case Retrieval Based on Inter-Case Distinctions. *ACM Transactions on Information Systems* 43, 5 (2025), 1–27.

[53] Weihang Su, Qingyao Ai, Jingtao Zhan, Qian Dong, and Yiqun Liu. 2025. Dynamic and parametric retrieval-augmented generation. In *Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 4118–4121.

[54] Weihang Su, Qian Dong, Qingyao Ai, and Yiqun Liu. 2025. SIGIR-AP 2025 Tutorial Proposal: Dynamic and Parametric Retrieval-Augmented Generation. *3rd International ACM SIGIR Conference on Information Retrieval in the Asia Pacific* (2025).

[55] Weihang Su, Yiran Hu, Anzhe Xie, Qingyao Ai, Quezi Bing, Ning Zheng, Yun Liu, Weixing Shen, and Yiqun Liu. 2024. STARD: A Chinese Statute Retrieval Dataset Derived from Real-life Queries by Non-professionals. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.). Association for Computational Linguistics, Miami, Florida, USA, 10658–10671. doi:10.18653/v1/2024.findings-emnlp.625

[56] Weihang Su, Xiangsheng Li, Yiqun Liu, Min Zhang, and Shaoping Ma. 2023. Thuir2 at ntcir-16 session search (ss) task. *arXiv preprint arXiv:2307.00250* (2023).

[57] Weihang Su, Yichen Tang, Qingyao Ai, Changyue Wang, Zhijing Wu, and Yiqun Liu. 2024. Mitigating entity-level hallucination in large language models. In *Proceedings of the 2024 Annual International ACM SIGIR Conference on Research and Development in Information Retrieval in the Asia Pacific Region*. 23–31.

[58] Weihang Su, Yichen Tang, Qingyao Ai, Zhijing Wu, and Yiqun Liu. 2024. DRAGIN: Dynamic Retrieval Augmented Generation based on the Real-time Information Needs of Large Language Models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, Bangkok, Thailand, 12991–13013. doi:10.18653/v1/2024.acl-long.702

[59] Weihang Su, Yichen Tang, Qingyao Ai, Junxi Yan, Changyue Wang, Hongning Wang, Ziyi Ye, Yujia Zhou, and Yiqun Liu. 2025. Parametric Retrieval Augmented Generation. In *Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '25), July 13–18, 2025, Padua, Italy*. doi:10.1145/3726302.3729957

[60] Weihang Su, Changyue Wang, Qingyao Ai, Yiran Hu, Zhijing Wu, Yujia Zhou, and Yiqun Liu. 2024. Unsupervised real-time hallucination detection based on the internal states of large language models. *arXiv preprint arXiv:2403.06448* (2024).

[61] Weihang Su, Changyue Wang, Qingyao Ai, Yiran Hu, Zhijing Wu, Yujia Zhou, and Yiqun Liu. 2024. Unsupervised Real-Time Hallucination Detection based on the Internal States of Large Language Models. In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, 14379–14391. doi:10.18653/V1/2024.FINDINGS-ACL.854

[62] Weihang Su, Anzhe Xie, Qingyao Ai, Jianming Long, Jiaxin Mao, Ziyi Ye, and Yiqun Liu. 2025. SurGE: A Benchmark and Evaluation Framework for Scientific Survey Generation. *arXiv preprint arXiv:2508.15658* (2025).

[63] Weihang Su, Baoqing Yue, Qingyao Ai, Yiran Hu, Jiaqi Li, Changyue Wang, Kaiyuan Zhang, Yueyue Wu, and Yiqun Liu. 2025. JuDGE: Benchmarking Judgment Document Generation for Chinese Legal System. In *Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '25), July 13–18, 2025, Padua, Italy*. doi:10.1145/3726302.3730295

[64] Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. 2023. Large language models in medicine. *Nature medicine* 29, 8 (2023), 1930–1940.

[65] James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a Large-scale Dataset for Fact Extraction and VERification. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. Association for Computational Linguistics, New Orleans, Louisiana, 809–819. doi:10.18653/v1/N18-1074

[66] Yiteng Tu, Weihang Su, Yujia Zhou, Yiqun Liu, and Qingyao Ai. 2025. RbFT: Robust Fine-tuning for Retrieval-Augmented Generation against Retrieval Defects. *arXiv preprint arXiv:2501.18365* (2025).

[67] Yiteng Tu, Weihang Su, Yujia Zhou, Yiqun Liu, Fen Lin, Qin Liu, and Qingyao Ai. 2025. Generalized Pseudo-Relevance Feedback. *arXiv preprint arXiv:2510.25488* (2025).

[68] Yiteng Tu, Zhichao Xu, Tao Yang, Weihang Su, Yujia Zhou, Yiqun Liu, Fen Lin, Qin Liu, and Qingyao Ai. 2022. Reinforcement Learning to Rank Using Coarse-grained Rewards. *arXiv e-prints* (2022), arXiv–2208.

[69] David Wadden, Shanchuan Lin, Kyle Lo, Lucy Lu Wang, Madeleine van Zuylen, Arman Cohan, and Hannaneh Hajishirzi. 2020. Fact or fiction: Verifying scientific claims. *arXiv preprint arXiv:2004.14974* (2020).

[70] Changyue Wang, Weihang Su, Qingyao Ai, and Yiqun Liu. 2025. Joint Evaluation of Answer and Reasoning Consistency for Hallucination Detection in Large Reasoning Models. *arXiv preprint arXiv:2506.04832* (2025).

[71] Changyue Wang, Weihang Su, Qingyao Ai, Yichen Tang, and Yiqun Liu. 2025. Knowledge editing through chain-of-thought. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*. 10684–10704.

[72] Changyue Wang, Weihang Su, Qingyao Ai, Yujia Zhou, and Yiqun Liu. 2025. Decoupling Reasoning and Knowledge Injection for In-Context Knowledge Editing. *arXiv preprint arXiv:2506.00536* (2025).

[73] Changyue Wang, Weihang Su, Hu Yiran, Qingyao Ai, Yueyue Wu, Cheng Luo, Yiqun Liu, Min Zhang, and Shaoping Ma. 2024. LeKUBE: A Legal Knowledge Update BEnchmark. *arXiv preprint arXiv:2407.14192* (2024).

[74] Jerry Wei, Chengrun Yang, Xinying Song, Yifeng Lu, Nathan Hu, Jie Huang, Dustin Tran, Daiyi Peng, Ruibo Liu, Da Huang, et al. 2024. Long-form factuality in large language models. *Advances in Neural Information Processing Systems* 37 (2024), 80756–80827.

[75] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 Technical Report. *arXiv preprint arXiv:2412.15115* (2024).

[76] Rui Yang, Ting Fang Tan, Wei Lu, Arun James Thirunavukarasu, Daniel Shu Wei Ting, and Nan Liu. 2023. Large language models in health care: Development, applications, and challenges. *Health Care Science* 2, 4 (2023), 255–263.

[77] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600* (2018).

[78] Yakir Yehuda, Itzik Malkiel, Oren Barkan, Jonathan Weill, Royi Ronen, and Noam Koenigstein. 2024. InterrogateLLM: Zero-Resource Hallucination Detection in LLM-Generated Answers. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, Bangkok, Thailand, 9333–9347. doi:10.18653/v1/2024.acl-long.506

[79] ChengXiang Zhai. 2008. Statistical language models for information retrieval. *Synthesis lectures on human language technologies* 1, 1 (2008), 1–141.

[80] Tianhang Zhang, Lin Qiu, Qipeng Guo, Cheng Deng, Yue Zhang, Zheng Zhang, Chenghu Zhou, Xinbing Wang, and Luoyi Fu. 2023. Enhancing Uncertainty-Based Hallucination Detection with Stronger Focus. *arXiv preprint arXiv:2311.13230* (2023).

[81] Xuan Zhang and Wei Gao. 2023. Towards LLM-based Fact Verification on News Claims with a Hierarchical Step-by-Step Prompting Method. In *Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, Jong C. Park, Yuki Arase, Baotian Hu, Wei Lu, Derry Wijaya, Ayu Purwarianti, and Adila Alfa Krisnadhi (Eds.). Association for Computational Linguistics, Nusa Dua, Bali, 996–1011. doi:10.18653/v1/2023.ijcnlp-main.64

[82] Jie Zhou, Xu Han, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. 2019. GEAR: Graph-based evidence aggregating and reasoning for fact verification. *arXiv preprint arXiv:1908.01843* (2019).

[83] Xinyi Zhou and Reza Zafarani. 2020. A survey of fake news: Fundamental theories, detection methods, and opportunities. *ACM Computing Surveys (CSUR)* 53, 5 (2020), 1–40.

# A    Prompt Templates

This appendix provides the complete prompt templates used in our unified evaluation framework.

## A.1    LLM-as-Judge Prompt

The LLM-as-Judge is a critical component for dynamic ground truth annotation. We use the following prompt template when both ground truth answers and passages are available:

---

**Prompt:**

You are an expert evaluator assessing the factual accuracy of a Generated Answer based *only* on the provided Golden Answer and Golden Passages for a given Question. Determine if the Generated Answer introduces significant factual errors, hallucinations, or contradictions compared to the trusted information.

[Instruction]
Evaluate the Generated Answer based *strictly* on the Golden Answer and Golden Passages. Do not use external knowledge. Identify if there are clear factual inaccuracies or information present in the Generated Answer that is not supported by or contradicts the provided golden information.

[Question]
{question}

[Golden Answer]
{golden_answer}

[Golden Passages]
{golden_passages_formatted}

[Generated Answer]
{generated_answer}

[Evaluation Task]
Choose the option that best describes the Generated Answer's factuality based *only* on the provided Golden Answer and Golden Passages:
(A) Accurate: The Generated Answer is factually consistent with the Golden Answer and Golden Passages. It does not contain significant inaccuracies or unsupported claims.
(B) Inaccurate/Hallucinated: The Generated Answer contains significant factual errors, contradictions, or information not supported by the Golden Answer and Golden Passages.

Provide your final choice by writing only the letter (A or B).
Choice:

---

When ground truth passages are unavailable, we use a simplified version that evaluates the generated answer solely against the golden answer, omitting the [Golden Passages] section.

## A.2    Fact Verification Prompt (LLM-based)

For LLM-based fact verification methods (LLM-Q and LLM-QA), we employ the following prompt to assess the hallucination level of generated answers against retrieved evidence:

---

**Prompt:**

Based on the following information, determine if the generated answer contains hallucinations (inaccurate or fabricated information):

Question: {question}

Generated Answer: {answer}

Reference Documents:
{retrieved_passages}

Please carefully analyze whether the generated answer is consistent with the information in the reference documents. If the answer contains information not mentioned in or contradicting the reference documents, it should be considered as containing hallucinations.

Rate the hallucination level on a scale from 0 to 1:
- 0: Completely reliable (no hallucination)
- 1: Completely unreliable (severe hallucination)

Output only a number between 0 and 1, no explanation needed.

---

# B    Implementation Details and Hyperparameters

## B.1    Answer Generation and Hallucination Detection

The main answer is generated using sampling-based decoding with the following hyperparameters:

- **Max New Tokens:** 30
- **Temperature:** 0.8
- **Top-p:** 0.9

For methods requiring multiple samples (e.g., SelfCheckGPT, Semantic Entropy, SEU, SIndex, PTrue), we generate 5 additional samples using the same hyperparameters. All hallucination detection methods follow their original hyperparameters as reported in their respective papers.

## B.2    Fact Verification Methods

*B.2.1    Evidence Retrieval with BM25.* All fact verification methods rely on BM25 retrieval from a Wikipedia corpus. We use the preprocessed Wikipedia dump provided by DPR [25], which consists of 21 million passages, each containing approximately 100 words. The corpus is indexed using Elasticsearch with default BM25 parameters (k1=1.2, b=0.75). For each query, we retrieve the top-3 most relevant passages based on BM25 scoring.

We implement two retrieval strategies:

- **Question-Only Retrieval:** Uses only the original question as the search query. This approach is query-agnostic and does not depend on the generated answer.
- **Question+Answer Retrieval:** Concatenates the question and the generated answer to form the search query. This often yields more targeted evidence passages that are directly relevant to the specific claims in the generated answer.

Both retrieval strategies use the same Elasticsearch backend with BM25 scoring. Retrieved passages are then used as evidence for both LLM-based and BERT-based fact verification methods.

*B.2.2    LLM-based Fact Verification.* We implement two LLM-based fact verification methods that differ only in their retrieval strategy:

- **LLM-Q:** Retrieves passages using the question only, then uses an LLM (Qwen-2.5-32B) to verify the answer against retrieved evidence. The generation uses greedy decoding.
- **LLM-QA:** Similar to LLM-Q but retrieves passages using both question and answer as the query.

*B.2.3    BERT-based Fact Verification.* We employ a BERT-based classifier initialized with `bert-base-uncased` to predict whether a generated answer constitutes a hallucination. We train two variants that differ only in the retrieval query used to obtain supporting passages: **BERT-Q** retrieves passages using the question alone, while **BERT-QA** retrieves passages using the concatenated question and answer.

*Input Representation.* For a question $q$, a generated answer $a$, and a set of retrieved passages $\{p_1, p_2, \ldots, p_k\}$, we construct the input sequence by first concatenating the retrieved passages, then appending the question and answer as separate components. Formally, let $P = p_1 \oplus p_2 \oplus \cdots \oplus p_k$ denote the concatenated passage text, where $\oplus$ represents string concatenation with space delimiters. The complete input representation is:

$$X(q, a, P) = [\text{CLS}]\, P\, [\text{SEP}]\, q\, [\text{SEP}]\, a\, [\text{SEP}], \tag{9}$$

where [CLS] serves as the classification token and [SEP] tokens separate the passage, question, and answer components. This structure enables the model to jointly encode the retrieved evidence and the question-answer pair for verification.

*Model Architecture.* We feed the tokenized input $X(q, a, P)$ into the BERT encoder to obtain contextualized embeddings. The classifier extracts the final-layer [CLS] token embedding and applies a linear transformation for binary classification:

$$h_{\text{CLS}} = \text{BERT}_{[\text{CLS}]}(X(q, a, P)), \tag{10}$$

$$\mathbf{z} = \text{Dropout}(h_{\text{CLS}}; p = 0.2), \tag{11}$$

$$\text{logits} = \mathbf{W}\mathbf{z} + \mathbf{b}, \tag{12}$$

where $h_{\text{CLS}} \in \mathbb{R}^{768}$ is the contextualized embedding of the [CLS] token, and $\mathbf{W} \in \mathbb{R}^{2 \times 768}$, $\mathbf{b} \in \mathbb{R}^2$ are learned parameters. The final prediction probability for the hallucination class is obtained via softmax normalization.

The model is optimized using the cross-entropy loss:

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^{N} \left[ y_i \log(\hat{p}_i) + (1 - y_i) \log(1 - \hat{p}_i) \right], \tag{13}$$

where $\hat{p}_i = \text{softmax}(\text{logits}_i)_1$ is the predicted probability of hallucination, $y_i \in \{0, 1\}$ denotes the ground truth label (0 for accurate, 1 for hallucinated), and $N$ is the batch size.

*Training Configuration.* Both BERT classifiers are trained for 5 epochs with a learning rate of $2 \times 10^{-5}$ and batch size of 16, using the AdamW optimizer with linear warmup (10% of total steps) and gradient clipping (max norm 1.0). Training data is automatically generated using our unified evaluation framework by collecting LLM-generated answers with corresponding LLM-as-Judge labels, and we use a disjoint slice of questions from those reserved for evaluation to prevent leakage. We perform a stratified 90/10 train-validation split and select the model checkpoint with the highest AUROC on the validation set.