

HalluClean: A Unified Framework to Combat Hallucinations in LLMs

Yaxin Zhao¹, Yu Zhang¹

¹Harbin Institute of Technology, Harbin, China
{yxzhao, zhangyu}@ir.hit.edu.cn

Abstract

Large language models (LLMs) have achieved impressive performance across a wide range of natural language processing tasks, yet they often produce hallucinated content that undermines factual reliability. To address this challenge, we introduce HalluClean, a lightweight and task-agnostic framework for detecting and correcting hallucinations in LLM-generated text. HalluClean adopts a reasoning-enhanced paradigm, explicitly decomposing the process into planning, execution, and revision stages to identify and refine unsupported claims. It employs minimal task-routing prompts to enable zero-shot generalization across diverse domains, without relying on external knowledge sources or supervised detectors. We conduct extensive evaluations on five representative tasks—question answering, dialogue, summarization, math word problems, and contradiction detection. Experimental results show that HalluClean significantly improves factual consistency and outperforms competitive baselines, demonstrating its potential to enhance the trustworthiness of LLM outputs in real-world applications.

Code — <https://github.com/tingmuor/HalluClean>

Extended version — <https://arxiv.org/abs/2511.08916>

Introduction

Large language models (LLMs) have revolutionized natural language processing (NLP), powering applications such as conversational agents, content generation, and decision support systems (Chowdhery et al. 2023; Touvron et al. 2023; Bang et al. 2023; Qin et al. 2023). These models leverage large-scale pretraining and are further enhanced via instruction tuning (Chung et al. 2024; Wang et al. 2022b,a) and alignment techniques that optimize for human preferences (Ouyang et al. 2022; Achiam et al. 2023). However, despite their remarkable fluency and versatility, LLMs frequently produce hallucinated or factually incorrect content (Huang et al. 2023; Ji et al. 2023), undermining their reliability in safety-critical contexts.

Existing hallucination mitigation approaches typically fall into two categories. Retrieval-augmented generation methods (Varshney et al. 2023; Cao et al. 2023; Kang, Ni, and Yao 2023; Rawte et al. 2023) query external knowledge

sources to validate or correct model outputs. Meanwhile, supervised detection approaches (Razumovskaia et al. 2024; Zhang et al. 2023; Qiu et al. 2023) rely on human-labeled data to train classifiers that identify hallucinations. While both strategies have shown promise, they suffer from key limitations: retrieval-based methods depend on the availability and accuracy of external sources, and annotation-based methods are costly and poorly generalize to novel hallucination types. Furthermore, hallucination behaviors vary widely across tasks—such as question answering (Zheng, Huang, and Chang 2023), summarization (Cao, Dong, and Cheung 2022), and dialogue (Das, Saha, and Srihari 2022)—yet most prior work focuses narrowly on specific settings (Mündler et al. 2023), limiting scalability and robustness.

To address these challenges, we propose a lightweight, task-agnostic framework for hallucination detection and correction that operates without external knowledge or task-specific supervision. Our approach leverages minimal task descriptions to instantiate a task-adaptive interface, guiding LLMs through a reasoning-enhanced, zero-shot process. Inspired by the plan-and-solve paradigm, we decompose hallucination mitigation into explicit planning and execution phases, enabling models to locate unsupported claims and revise outputs with improved factuality.

We introduce **HalluClean**, a unified framework that detects and corrects hallucinations in LLM-generated outputs via structured reasoning. HalluClean employs compact prompts to elicit multi-step reasoning traces, which serve both to identify hallucinated segments and to guide targeted revision. This plug-and-play, prompt-based design ensures broad applicability across model architectures and NLP tasks, while remaining compatible with open-source LLMs for privacy-preserving deployments.

Our contributions are summarized as follows:

- **HalluClean Framework:** We present a zero-shot hallucination detection and correction framework based on structured reasoning. HalluClean is modular and supports flexible integration with diverse LLMs, including open-source models.
- **Task-Agnostic Generalization:** HalluClean achieves strong performance across a variety of NLP tasks—question answering, summarization, dia-

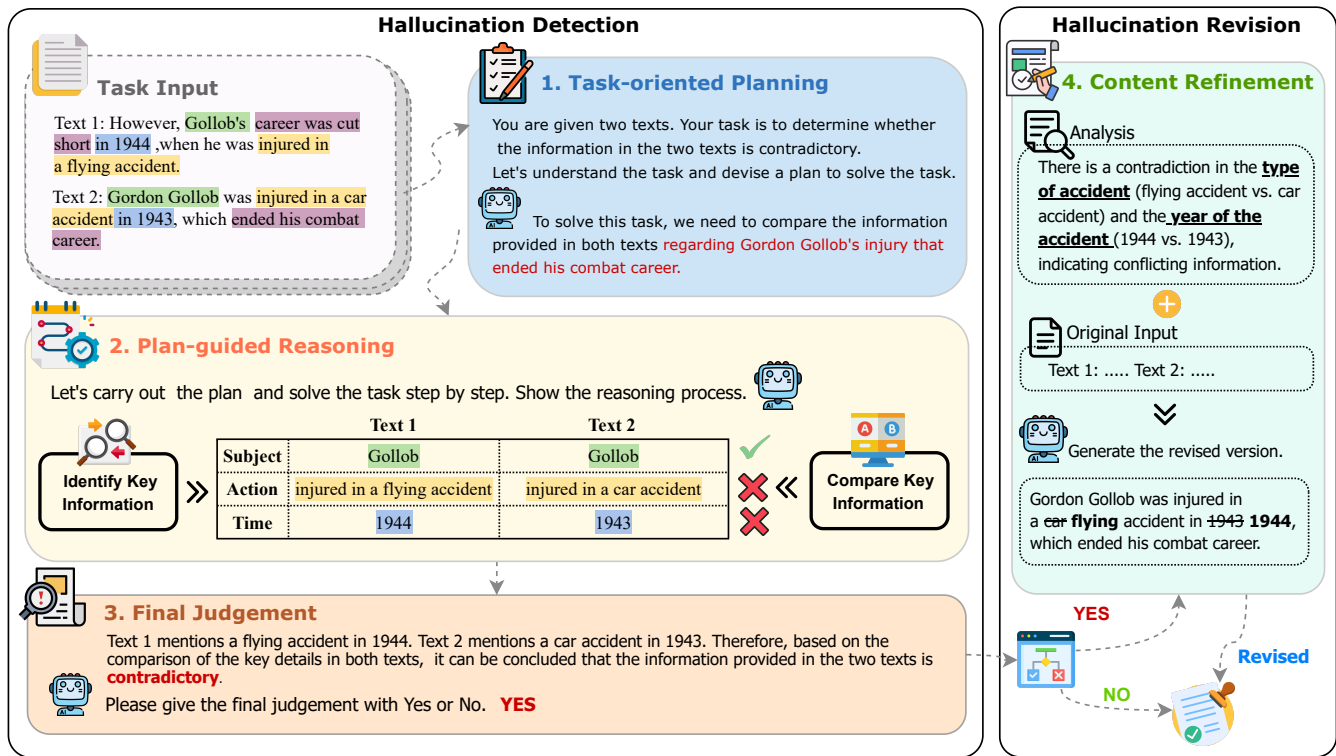


Figure 1: Overview of the HalluClean framework. It consists of two modules: hallucination detection and revision. The detection module generates a task-specific plan, performs step-by-step reasoning, and makes a final judgment. If a hallucination is detected, the revision module revises the content based on the identified reasoning to eliminate hallucinated information.

logue, math word problems, and contradiction detection—without requiring task-specific fine-tuning.

- **Domain-Level Robustness:** We demonstrate HalluClean’s effectiveness in domain-sensitive settings such as medicine and finance, highlighting its potential for deployment in real-world, high-stakes applications.

Related Work

Hallucinations in LLMs

Hallucinations in large language models (LLMs) have been extensively studied, focusing on their causes (Pan et al. 2023; Chen and Shu 2024; Kasai et al. 2024; Wang et al. 2023a; Lee et al. 2022; Yao et al. 2023), evaluation methodologies (Lin, Hilton, and Evans 2022; Lee et al. 2022; Min et al. 2023; Li et al. 2023a), and behavioral analysis (Zhao et al. 2023; Dong et al. 2024; Li et al. 2023b). Many studies have investigated ways to mitigate hallucinations through retrieval-augmented generation (Peng et al. 2023; Varshney et al. 2023; Kang, Ni, and Yao 2023) and supervised fine-tuning (Elaraby et al. 2023; Razumovskaia et al. 2024; Zhang et al. 2023). To improve LLM reliability, researchers have explored prompting-based solutions. For example, Si et al. (2022) proposed simple yet effective prompts that enhance GPT-3’s factual accuracy, while Mitchell et al. (2022) introduced a two-model framework in which one model generates responses and another evaluates

their logical coherence. More recently, Mündler et al. (2023) found that 17.7% of ChatGPT-generated sentences contain self-contradictions and proposed a three-step pipeline to detect and mitigate them without relying on external knowledge. Despite these advancements, hallucination detection and mitigation in LLMs remain challenging. Our method uses task-adaptive prompts in a zero-shot setting to guide LLMs in detecting and revising hallucinated content by eliciting and leveraging model reasoning.

Advancements in Prompting Techniques

Prompting strategies have played a crucial role in enhancing the reasoning abilities of LLMs. Chain-of-Thought (CoT) prompting (Wei et al. 2022) explicitly structures intermediate reasoning steps, significantly improving model performance on complex reasoning tasks. Building upon this, various enhancements have been proposed, including prompt ensembling (Wang et al. 2022a; Li et al. 2022; Fu et al. 2022), problem decomposition (Zhou et al. 2022; Khot et al. 2022; Dua et al. 2022), and structured planning methods (Yao et al. 2022; Huang et al. 2022; Wang et al. 2023b; Liu et al. 2023). To reduce manual effort and computational overhead, zero-shot CoT prompting (Kojima et al. 2022) was introduced, allowing LLMs to autonomously generate reasoning steps without the need for labeled exemplars. However, existing methods primarily focus on general reasoning tasks, and limited work has explored their applicabil-

| Task Type | Task Routing Prompt |
|--------------------|--|
| Question Answering | You are provided with a question and its corresponding answer. Your task is to determine whether the answer contains hallucinated content. |
| Dialogue Systems | You are provided with a dialogue history and its corresponding response. Your task is to determine whether the response contains hallucinated content. |
| Summarization | You are provided with a document and its corresponding summary. Your task is to determine whether the summary contains hallucinated content. |
| Math Word Problems | You are provided with a math word problem. Your task is to determine whether the problem is unanswerable. |
| Self-contradiction | You are given two texts. Your task is to determine whether the information in the two texts is contradictory. |

Table 1: Task-oriented routing prompts for different NLP applications. These concise instructions guide the model to understand the specific hallucination detection objective for each task type.

ity in addressing hallucinations within LLM-generated text.

In this work, we propose a new prompt-based mechanism to enhance the effectiveness of hallucination detection and ensuring more reliable correction.

Method

We introduce **HalluClean**, a task-agnostic framework for hallucination detection and correction in LLMs. It leverages structured reasoning guided by minimal task prompts and operates in zero-shot settings.

Task-Based Categorization of LLM Hallucinations

LLMs support diverse applications—summarization, QA, dialogue, and problem solving—but frequently generate factually inconsistent or logically contradictory content, known as *hallucinations*. These typically arise when generated outputs lack grounding in verifiable knowledge or logic. To address this, we adopt a task-based categorization of hallucinations across five representative NLP scenarios:

Question Answering Hallucinations manifest as unsupported claims, misinterpreted context, or factual errors that deviate from the input or common knowledge.

Dialogue Systems Errors often stem from entity mismatches—substituting similar, dissimilar, or cross-type entities—leading to factual incoherence with the dialogue history.

Summarization Generated summaries may include unverifiable details or fabricate facts not grounded in the source text, often misrepresenting entities or relations.

Math Word Problems Under-specified or ill-posed problems cause hallucinations when essential constraints are missing, variables are vague, or assumptions violate logic (e.g., negative quantities where not allowed).

Self-contradiction Contradictions within the same response (e.g., mutually exclusive statements) signal hallucination. Such contradictions appear in 17.7% of ChatGPT-generated sentences (Mündler et al. 2023).

HalluClean: A Unified Framework

Based on our analysis of hallucination patterns, we propose HalluClean, the framework is composed of two main modules: reasoning-enhanced hallucination detection module and targeted revision module. HalluClean is designed to adapt to various task settings and requires no task-specific fine-tuning. Figure 1 provides an overview of HalluClean’s architecture. The framework first detects hallucination through structured reasoning, and then modifies the hallucinated parts based on the rationale. The framework uses modular prompt templates, enabling easy adaptation across tasks and LLM architectures. The following section describes each component of our framework.

Hallucination Detection We employ the structural-reasoning-enhanced detection module to assess the factual consistency of the generated output. This component constitutes the core technical innovation of our approach. Rather than directly prompting for a binary classification, we guide the model through a structured three-step reasoning process. This process yields a reliable binary judgment indicating whether hallucination is present, along with a detailed reasoning trace explaining the rationale behind this determination.

Hallucination Revision If any hallucination is detected, the framework activates the revision module. Rather than modifying the content directly, the model performs revision based on the reasoning trace generated during detection. This ensures that the correction is guided by explicit analysis, improving the reliability and quality of the revision. By preserving accurate content and focusing only on identified issues, this targeted strategy makes the correction process more precise and controllable.

The detection and revision modules together form a unified pipeline for identifying and correcting hallucinations. This design ensures factual consistency, improves interpretability, and enhances control over the generation process.

Task-Oriented Routing For each supported task type, we design a concise task-specific prompt that provides the

| Method | QA | | DA | | SUM | | MWP | | SC | |
|---|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | R | Q | R | Q | R | Q | R | Q | R | Q |
| LLM-Direct Ask (Detection → Revision) | | | | | | | | | | |
| GPT-3.5-turbo | 20.5% | 12.0% | 57.5% | 53.0% | 14.5% | 7.0% | 42.0% | 13.0% | 30.7% | 30.7% |
| GPT-4o-mini | 39.5% | 22.5% | 84.5% | 79.5% | 30.0% | 30.0% | 47.5% | 14.0% | 72.7% | 72.7% |
| Llama-3-70B | 30.5% | 18.5% | 74.5% | 68.5% | 19.5% | 19.5% | 83.0% | 32.5% | 49.3% | 49.3% |
| DeepSeek-V3 | 49.0% | 32.0% | 86.5% | 78.0% | 41.0% | 41.0% | 40.0% | 17.0% | 35.3% | 35.3% |
| DeepSeek-R1 | 62.0% | 45.5% | 74.0% | 67.5% | 36.5% | 36.0% | 42.0% | 28.0% | 25.3% | 25.3% |
| Existing Baselines (Detection→Revision(with rationale); GPT-3.5-turbo) | | | | | | | | | | |
| Step-by-Step | 13.0% | 10.0% | 54.5% | 51.5% | 13.0% | 13.0% | 44.0% | 37.5% | 53.3% | 53.3% |
| Plan-and-Solve | 20.5% | 12.5% | 11.5% | 10.5% | 3.5% | 3.5% | 53.0% | 44.3% | 53.3% | 53.3% |
| ChatProtect | 38.0% | 24.0% | 79.5% | 74.0% | 23.0% | 22.5% | 80.5% | 37.9% | 79.3% | 79.3% |
| Ours-GPT-3.5-turbo | 72.5% | 25.5% | 89.0% | 83.0% | 59.5% | 59.0% | 75.5% | 45.0% | 87.3% | 79.3% |
| Ours-Deepseek-V3 | 74.0% | 37.5% | 92.5% | 86.0% | 54.5% | 55.0% | 75.5% | 41.0% | 87.3% | 62.7% |

Table 2: The effectiveness of the Framework HalluClean: R denotes the hallucination reduction rate after applying the revision module, and Q denotes the revision success rate, which reflects the quality of corrections.

model with minimal yet sufficient context to understand its objective. These prompts serve as task adapters, allowing HalluClean to flexibly operate across diverse applications without requiring fine-tuning or additional training data. Table 1 presents examples of task routing prompts for different NLP tasks.

Structural Reasoning Mechanism The core innovation of HalluClean lies in its structural reasoning mechanism for hallucination detection. While direct classification can identify obvious hallucinations, we find that step-by-step reasoning significantly improves detection accuracy, especially for subtle or complex hallucination cases. Our approach draws inspiration from cognitive science literature on human reasoning, which emphasizes the role of structured thinking in error detection and verification (Kahneman 2011).

We implement this insight through a four-step prompt-based inference mechanism, as illustrated in Figure 1 and detailed below:

Step 1: Task-oriented Planning The first step guides the model to develop a systematic approach tailored to the specific task and input. The planning prompt follows this template:

[INPUT] Task Input
[TASK] Task Description
Let’s understand the task and devise a plan to solve the task.

This planning step serves multiple important functions: it encourages metacognitive reflection before analysis, creates task-specific verification strategies, and breaks complex detection tasks into manageable sub-components. For example in 1, when analyzing a potential contradiction between two statements, the plan might involve identifying key entities, extracting their relationships, and systematically comparing these elements.

Step 2: Plan-guided Reasoning In the second step, the model implements the verification plan developed in Step 1:

[INPUT] Task Input
[PLAN] Result from Step-1
Let’s carry out the plan and solve the task step by step.
Show the reasoning process.

During reasoning, the model systematically applies each verification step defined in the plan to validate the input content. The structured nature of this process ensures comprehensive and consistent examination, minimizing the risk of overlooking subtle inconsistencies. Furthermore, by following an explicit plan, the model generates transparent and interpretable reasoning traces that not only support the final judgment but also facilitate human verification and analysis.

Step 3: Final Judgment This step synthesizes the detailed analysis into a conclusive judgment:

[INPUT] Task Input
[ANALYSIS] Result from Step-2
Please conclude whether the [INPUT] contains hallucinated content with Yes or No.

This step produces a binary judgment indicating whether hallucinations are present. If hallucinations are detected, the input proceeds to the subsequent revision phase for correction.

Step 4: Content Refinement The final step corrects the response based on the hallucination identification analysis:

[INPUT] Task Input
[ANALYSIS] Result from Step-2
Given the analysis explaining why [INPUT] contains hallucinated content. Generate a revised version without hallucinations.

This step refines the original response by leveraging the reasoning behind hallucination identification, aiming to produce a factually consistent revision.

Our structural reasoning approach offers several key advantages. It reduces the risk of oversight through step-by-step analysis, provides transparent reasoning traces for interpretability, adapts verification strategies to specific tasks

| Method | QA | | DA | | SUM | | MWP | | SC | |
|---|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. |
| LLM-Direct Ask | | | | | | | | | | |
| GPT-3.5-turbo | 33.5% | 59.3% | 62.8% | 66.0% | 24.7% | 55.8% | 50.9% | 59.5% | 46.0% | 64.0% |
| GPT-4o-mini | 52.7% | 64.5% | 76.5% | 74.0% | 45.5% | 64.0% | 61.7% | 70.5% | 84.2% | 72.7% |
| Llama-3-70B | 44.7% | 62.3% | 66.4% | 62.3% | 32.4% | 59.3% | 83.4% | 83.5% | 65.8% | 73.6% |
| DeepSeek-V3 | 62.2% | 70.3% | 75.4% | 71.8% | 55.0% | 66.5% | 55.6% | 68.0% | 52.0% | 67.3% |
| DeepSeek-R1 | 67.6% | 70.3% | 71.0% | 69.8% | 49.3% | 62.5% | 65.2% | 72.0% | 40.0% | 62.0% |
| Existing Baselines (GPT-3.5-turbo) | | | | | | | | | | |
| Step-by-Step | 22.0% | 54.0% | 61.4% | 65.8% | 22.1% | 54.3% | 55.7% | 65.0% | 68.1% | 75.0% |
| SelfCheckGPT | 43.3% | 43.8% | 19.9% | 27.8% | 53.1% | 37.3% | 25.8% | 54.0% | 5.7% | 12.0% |
| Plan-and-Solve | 32.0% | 56.5% | 19.3% | 52.0% | 6.7% | 51.3% | 66.9% | 73.8% | 66.4% | 73.0% |
| ChatProtect | 51.4% | 64.0% | 72.0% | 69.3% | 36.7% | 60.3% | 74.0% | 71.8% | 83.8% | 84.7% |
| Ours-GPT-3.5-turbo | 67.8% | 66.5% | 74.3% | 69.3% | 65.9% | 69.2% | 80.3% | 81.5% | 87.0% | 87.0% |
| Ours-Deepseek-V3 | 71.5% | 70.5% | 77.1% | 72.5% | 62.9% | 67.5% | 89.1% | 89.5% | 76.1% | 80.3% |
| Ours-Llama-3-70B | 70.6% | 69.0% | 74.0% | 68.8% | 46.5% | 61.5% | 85.6% | 86.0% | 80.8% | 83.3% |

Table 3: Comparison of hallucination detection performance between our method, existing methods under a unified GPT-3.5-turbo backbone, and direct classification baselines across various LLMs. Best results are highlighted in **bold**.

| Model | QA | | DA | | SUM | | MWP | | SC | |
|------------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. |
| Direct Ask | 33.5% | 59.3% | 62.8% | 66.0% | 24.7% | 55.8% | 50.9% | 59.5% | 46.0% | 64.0% |
| +Task-oriented Routing | 39.3% | 59.0% | 69.8% | 67.8% | 60.0% | 65.3% | 50.3% | 59.5% | 82.5% | 83.3% |
| +Structural Reasoning | 67.8% | 66.5% | 74.3% | 69.3% | 65.9% | 69.2% | 80.3% | 81.5% | 87.0% | 87.0% |

Table 4: Ablation study of the HalluClean framework. We evaluate the impact of removing the task-oriented routing and structural-reasoning mechanism.

via task-oriented planning, and guides targeted revisions by identifying what to fix and why. All of this is achieved in a single execution.

Experiments

Dataset We collect evaluation data from four established hallucination detection benchmarks: 1. **HaluEval (Li et al. 2023a)**: Covers hallucinated samples across three task types—question answering, knowledge-grounded dialogue, and text summarization. 2. **UMWP (Sun et al. 2024)**: Evaluates hallucination in math word problems (MWPs) by identifying questions with no or non-unique solutions. Such unanswerable questions are known to induce hallucinations in LLMs and are often used to test whether models can recognize ill-posed or unsolvable problems—similar to how educators gauge student understanding with trick questions. 3. **ChatProtect (Mündler et al. 2023)**: Focuses on self-contradictory hallucinations, where a language model produces logically inconsistent statements within the same context. 4. **HaluBench (Ravi et al. 2024)**: A domain-specific benchmark composed of hallucinated QA examples in the medical and financial domains, sourced from CovidQA, PubMedQA, and FinanceBench.

We demonstrate the effectiveness of HalluClean by eval-

uating it across multiple hallucination-prone NLP tasks in zero-shot setting, including: Question Answering (QA), Dialogue (DA), Summarization (SUM), Math Word Problems (MWPs) and Self-contradictory Hallucinations (SC).

Evaluation Metrics The effectiveness of the proposed framework is evaluated along three dimensions: 1. **Hallucination Reduction Rate**: To measure the effectiveness of the revision step, we compute the hallucination reduction rate by comparing the number of hallucinations detected before and after revision. Specifically, we first identify hallucinations in the original outputs, then re-evaluate the revised outputs using GPT-4o-mini. 2. **Revision Success Rate**: To assess revision quality, we compute BERTScore between each revised output and its gold reference. A revision is considered *acceptable* if the BERTScore (Zhang et al. 2019) exceeds 0.85 (chosen to balance strict semantic fidelity and flexibility in surface expression). The revision success rate is the proportion of acceptable revisions among all hallucinations identified before revision. For multi-word problems (MWPs), revision quality evaluation is performed based on unanswerable reason categories, using exact match between predicted and gold labels. 3. **Hallucination Detection**: Since hallucination detection is formulated as a binary classification task, we evaluate its effectiveness using two

| Model | CovidQA | | PubmedQA | | FinanceBench | | Overall | |
|--------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. |
| GPT-3.5-turbo | 9.5% | 52.5% | 7.7% | 52.0% | 11.0% | 51.5% | 9.4% | 52.0% |
| GPT-4o-mini | 53.2% | 67.5% | 68.7% | 74.5% | 19.4% | 50.0% | 47.1% | 64.0% |
| Llama-3-70B | 7.7% | 52.0% | 16.5% | 54.5% | 7.4% | 50.0% | 10.5% | 52.2% |
| DeepSeek-V3 | 81.6% | 84.0% | 71.0% | 77.5% | 21.1% | 55.0% | 57.9% | 72.2% |
| DeepSeek-R1 | 68.4% | 75.5% | 76.4% | 79.0% | 47.9% | 63.0% | 64.2% | 72.5% |
| Ours-GPT-3.5-turbo | 91.7% | 92.0% | 81.7% | 81.0% | 73.4% | 76.5% | 82.3% | 83.2% |

Table 5: The effectiveness of the detection in real world specific-domain.

standard metrics: F1 score and accuracy, computed against human-annotated gold labels from the original benchmarks. Accuracy reflects overall correctness, while F1 provides additional insight into the model’s balance between precision and recall, ensuring that both metrics are measured with respect to verified ground truth.

Results and Analysis

Hallucination Revision Performance We evaluate the effectiveness of our HalluClean framework in mitigating hallucinations by comparing model performance before and after applying our framework. As baselines, we consider several mainstream LLMs that (i) directly detect hallucinations and generate revised outputs when necessary, and (ii) variants that incorporate intermediate rationales during the detection stage.

Table 2 shows the performance of HalluClean in reducing hallucinations across five tasks. Our method achieves the highest reduction rate (R) and revision quality (Q) on most tasks. Overall, HalluClean delivers more consistent and effective correction across all evaluated scenarios.

We perform an ablation study of the HalluClean framework to evaluate the impact of HalluClean’s task-oriented routing and structural-reasoning mechanism in HalluClean. As shown in Table 4, each module individually contributes to performance improvement, confirming their complementary roles in the framework.

Hallucination Detection Performance Table 3 presents a comprehensive comparison of hallucination detection performance across five tasks: QA, DA, SUM, MWPs, and SC. Our method consistently outperforms both direct LLM judgment and existing baselines under a unified GPT-3.5-turbo backbone.

When using a stronger backbone such as DeepSeek-V3, our method achieves further performance gains. Specifically, Ours-DeepSeek-V3 attains the highest F1 scores on QA, DA, and MWPs, and achieves the highest accuracy across all five tasks. Moreover, the competitive performance of Ours-Llama-3-70B highlights the practicality of our method when deployed with open-source backbones, offering a compelling solution for resource-constrained or privacy-sensitive applications.

Ablation Study Table 4 presents the ablation results of the HalluClean framework, illustrating the impact of its two

| Method | Vanilla | | Retrieval Aug. | |
|---------------|--------------|--------------|----------------|--------------|
| | F1 | Acc. | F1 | Acc. |
| GPT-3.5-turbo | 33.5% | 59.3% | 56.2% | 65.3% |
| +Ours | 67.8% | 66.5% | 80.4% | 82.3% |

Table 6: Evaluation of Detection Performance with Retrieval-Augmented Strategy

core components: Task-Oriented Routing and Structural-reasoning mechanism. Starting from a Direct Ask baseline, we incrementally add these modules and observe consistent performance improvements across all five tasks. Overall, both modules contribute complementary benefits.

| Method | HalluQA | | CMHE-HD | |
|---------------|--------------|--------------|--------------|--------------|
| | F1 | Acc. | F1 | Acc. |
| GPT-3.5-turbo | 7.0% | 46.5% | 21.9% | 50.0% |
| +Ours | 41.6% | 55.0% | 57.3% | 51.5% |

Table 7: Evaluation of Detection Performance under Cross-Lingual Settings

Domain-Specific Evaluation in Real-World Applications

We further evaluate the effectiveness of our method in domain-specific hallucination detection, focusing on the medical and financial fields. As shown in Table 5 Across all three domain-specific datasets, our method consistently achieves the highest F1 score and accuracy, demonstrating its robustness and effectiveness in hallucination detection within specialized fields.

Integration with Retrieval-Augmented Generation To evaluate whether our method complements external knowledge, we test hallucination detection on the QA task from HalEval (Li et al. 2023a) under two settings: *vanilla* (no external knowledge) and *retrieval-augmented* (with background knowledge).

As shown in Table 6, our method significantly outperforms direct judgment by GPT-3.5-turbo in both the vanilla and retrieval-augmented settings. Notably, when augmented with retrieval, our approach achieves substantial gains in both F1 (80.4%) and accuracy (82.3%), demonstrating its

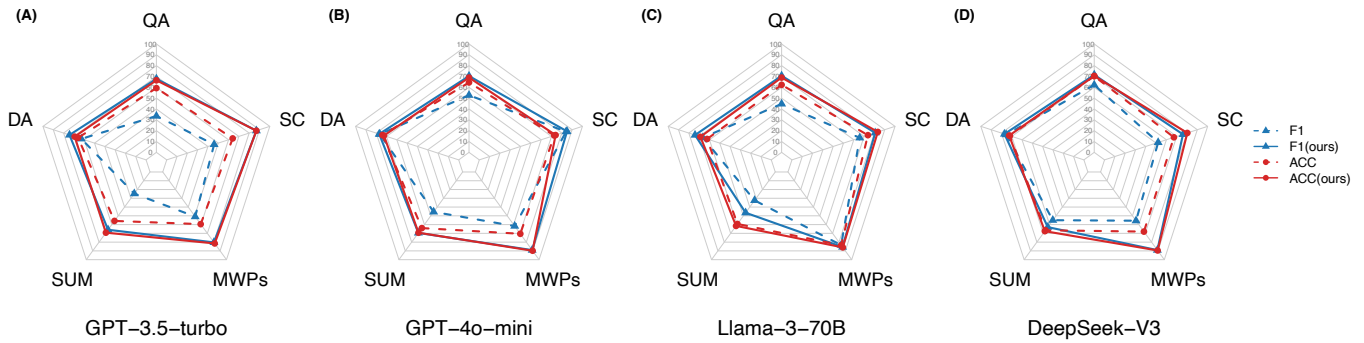


Figure 2: Detection adaptability across backbone LLMs. F1 and Acc denote the F1 score and accuracy of hallucination detection, respectively.

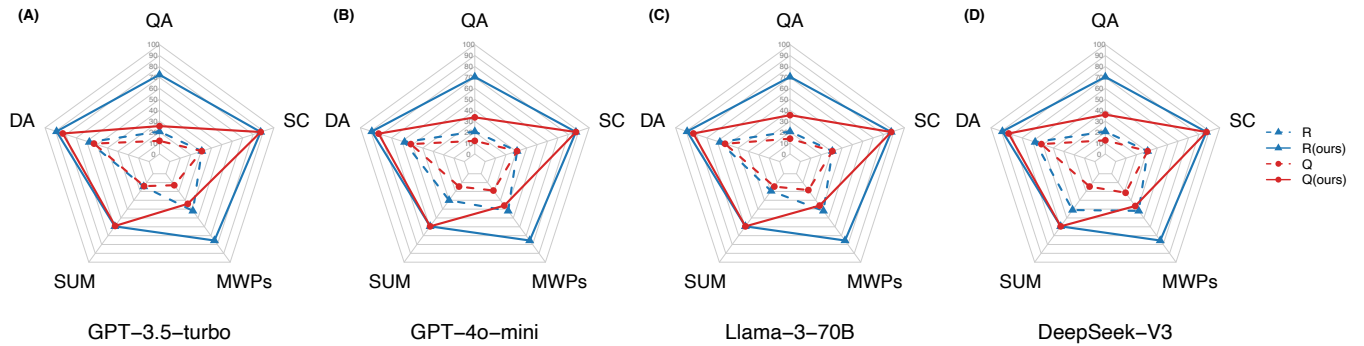


Figure 3: Revision module adaptability across backbone LLMs. R represents hallucination reduction rate, and Q represents revision success rate.

strong ability to leverage external information for more accurate hallucination detection.

Cross-Lingual Transferability To evaluate cross-lingual generalization, we test our method on two Chinese hallucination detection benchmarks: **HalluQA** and **CMHE-HD**, each with 200 samples (100 hallucinated, 100 faithful), using GPT-3.5-turbo as the backbone. As shown in Table 7, our method outperforms the GPT-3.5-turbo baseline on both datasets, demonstrating strong cross-lingual adaptability.

Module-Level Adaptability Across Backbone LLMs To evaluate the generalization ability of our framework, we assess the adaptability of **Detection Module** and **Revision Module** across five tasks using different backbone LLMs. Figures 2 and 3 evaluate the adaptability of our hallucination detection and revision modules across multiple backbone LLMs and task types.

Our detection module consistently improves F1 and accuracy across all evaluated tasks and backbone models. The performance gain is especially notable in QA and SC tasks, demonstrating strong adaptability to diverse reasoning types. Notably, GPT-3.5-turbo exhibit the most significant gains (e.g., over 41% F1 improvement on summarization and self-contradiction detection, and over 30% on QA and MWP). The revision module mitigates performance disparities across tasks, enabling each model to achieve more

balanced and consistent results.

These findings underscore the robustness and cross-task generalization of our detection and revision modules in hallucination identification, consistently performing well across different task types and model architectures.

Conclusion

We presented HalluClean, a lightweight and generalizable framework for detecting and correcting hallucinations in language model outputs. In contrast to prior approaches that rely on extensive fine-tuning or external knowledge sources, HalluClean operates in a zero-shot setting through structured reasoning. It achieves strong performance across a wide range of hallucination-prone tasks—including question answering, summarization, dialogue, mathematical reasoning, and self-contradiction detection—without task-specific supervision. Moreover, HalluClean supports local deployment with open-source LLMs, making it particularly suitable for privacy-sensitive or resource-constrained scenarios. These attributes establish HalluClean as a practical, interpretable, and broadly applicable solution for enhancing the factual consistency and trustworthiness of LLM-generated content.

Acknowledgments

We thank the anonymous reviewers for their valuable feedback. We also thank Anqi Zhang, Changyu Xu, and Ruiheng Liu for their help. This work was supported by the National Natural Science Foundation of China (No. 62476066).

References

- Achiam, J.; Adler, S.; Agarwal, S.; Ahmad, L.; Akkaya, I.; Aleman, F. L.; Almeida, D.; Altschmidt, J.; Altman, S.; Anadkat, S.; et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Bang, Y.; Cahyawijaya, S.; Lee, N.; Dai, W.; Su, D.; Wilie, B.; Lovenia, H.; Ji, Z.; Yu, T.; Chung, W.; et al. 2023. A Multitask, Multilingual, Multimodal Evaluation of ChatGPT on Reasoning, Hallucination, and Interactivity. 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)*, 675–718.
- Cao, H.; An, Z.; Feng, J.; Xu, K.; Chen, L.; and Zhao, D. 2023. A Step Closer to Comprehensive Answers: Constrained Multi-Stage Question Decomposition with Large Language Models. *CoRR*.
- Cao, M.; Dong, Y.; and Cheung, J. C. K. 2022. Hallucinated but Factual! Inspecting the Factuality of Hallucinations in Abstractive Summarization. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 3340–3354.
- Chen, C.; and Shu, K. 2024. Combating misinformation in the age of llms: Opportunities and challenges. *AI Magazine*, 45(3): 354–368.
- Chowdhery, A.; Narang, S.; Devlin, J.; Bosma, M.; Mishra, G.; Roberts, A.; Barham, P.; Chung, H. W.; Sutton, C.; Gehrmann, S.; et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240): 1–113.
- Chung, H. W.; Hou, L.; Longpre, S.; Zoph, B.; Tay, Y.; Fedus, W.; Li, Y.; Wang, X.; Dehghani, M.; Brahma, S.; et al. 2024. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70): 1–53.
- Das, S.; Saha, S.; and Srihari, R. K. 2022. Diving Deep into Modes of Fact Hallucinations in Dialogue Systems. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, 684–699.
- Dong, Q.; Xu, J.; Kong, L.; Sui, Z.; and Li, L. 2024. Statistical knowledge assessment for large language models. *Advances in Neural Information Processing Systems*, 36.
- Dua, D.; Gupta, S.; Singh, S.; and Gardner, M. 2022. Successive Prompting for Decomposing Complex Questions. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 1251–1265.
- Elaraby, M.; Lu, M.; Dunn, J.; Zhang, X.; Wang, Y.; Liu, S.; Tian, P.; Wang, Y.; and Wang, Y. 2023. Halo: Estimation and reduction of hallucinations in open-source weak large language models. *arXiv preprint arXiv:2308.11764*.
- Fu, Y.; Peng, H.; Sabharwal, A.; Clark, P.; and Khot, T. 2022. Complexity-based prompting for multi-step reasoning. In *The Eleventh International Conference on Learning Representations*.
- Huang, L.; Yu, W.; Ma, W.; Zhong, W.; Feng, Z.; Wang, H.; Chen, Q.; Peng, W.; Feng, X.; Qin, B.; et al. 2023. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*.
- Huang, W.; Abbeel, P.; Pathak, D.; and Mordatch, I. 2022. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *International conference on machine learning*, 9118–9147. PMLR.
- Ji, Z.; Lee, N.; Frieske, R.; Yu, T.; Su, D.; Xu, Y.; Ishii, E.; Bang, Y. J.; Madotto, A.; and Fung, P. 2023. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12): 1–38.
- Kahneman, D. 2011. *Thinking, fast and slow*. macmillan.
- Kang, H.; Ni, J.; and Yao, H. 2023. Ever: Mitigating hallucination in large language models through real-time verification and rectification. *arXiv preprint arXiv:2311.09114*.
- Kasai, J.; Sakaguchi, K.; Le Bras, R.; Asai, A.; Yu, X.; Radev, D.; Smith, N. A.; Choi, Y.; Inui, K.; et al. 2024. RE-ALTIME QA: what’s the answer right now? *Advances in Neural Information Processing Systems*, 36.
- Khot, T.; Trivedi, H.; Finlayson, M.; Fu, Y.; Richardson, K.; Clark, P.; and Sabharwal, A. 2022. Decomposed prompting: A modular approach for solving complex tasks. *arXiv preprint arXiv:2210.02406*.
- Kojima, T.; Gu, S. S.; Reid, M.; Matsuo, Y.; and Iwasawa, Y. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35: 22199–22213.
- Lee, N.; Ping, W.; Xu, P.; Patwary, M.; Fung, P. N.; Shoeybi, M.; and Catanzaro, B. 2022. Factuality enhanced language models for open-ended text generation. *Advances in Neural Information Processing Systems*, 35: 34586–34599.
- Li, J.; Cheng, X.; Zhao, W. X.; Nie, J.-Y.; and Wen, J.-R. 2023a. HaluEval: A Large-Scale Hallucination Evaluation Benchmark for Large Language Models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 6449–6464.
- Li, X. L.; Shrivastava, V.; Li, S.; Hashimoto, T.; and Liang, P. 2023b. Benchmarking and improving generator-validator consistency of language models. *arXiv preprint arXiv:2310.01846*.
- Li, Y.; Lin, Z.; Zhang, S.; Fu, Q.; Chen, B.; Lou, J.-G.; and Chen, W. 2022. On the advance of making language models better reasoners. *arXiv preprint arXiv:2206.02336*.
- Lin, S.; Hilton, J.; and Evans, O. 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.
- Liu, B.; Jiang, Y.; Zhang, X.; Liu, Q.; Zhang, S.; Biswas, J.; and Stone, P. 2023. Llm+ p: Empowering large language models with optimal planning proficiency. *arXiv preprint arXiv:2304.11477*.

- Min, S.; Krishna, K.; Lyu, X.; Lewis, M.; Yih, W.-t.; Koh, P.; Iyyer, M.; Zettlemoyer, L.; and Hajishirzi, H. 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.
- Mitchell, E.; Noh, J. J.; Li, S.; Armstrong, W. S.; Agarwal, A.; Liu, P.; Finn, C.; and Manning, C. D. 2022. Enhancing self-consistency and performance of pre-trained language models through natural language inference. *arXiv preprint arXiv:2211.11875*.
- Mündler, N.; He, J.; Jenko, S.; and Vechev, M. 2023. Self-contradictory hallucinations of large language models: Evaluation, detection and mitigation. *arXiv preprint arXiv:2305.15852*.
- Ouyang, L.; Wu, J.; Jiang, X.; Almeida, D.; Wainwright, C.; Mishkin, P.; Zhang, C.; Agarwal, S.; Slama, K.; Ray, A.; et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744.
- Pan, Y.; Pan, L.; Chen, W.; Nakov, P.; Kan, M.-Y.; and Wang, W. Y. 2023. On the risk of misinformation pollution with large language models. *arXiv preprint arXiv:2305.13661*.
- Peng, B.; Galley, M.; He, P.; Cheng, H.; Xie, Y.; Hu, Y.; Huang, Q.; Liden, L.; Yu, Z.; Chen, W.; et al. 2023. Check your facts and try again: Improving large language models with external knowledge and automated feedback. *arXiv preprint arXiv:2302.12813*.
- Qin, C.; Zhang, A.; Zhang, Z.; Chen, J.; Yasunaga, M.; and Yang, D. 2023. Is ChatGPT a general-purpose natural language processing task solver? *arXiv preprint arXiv:2302.06476*.
- Qiu, Y.; Embar, V.; Cohen, S. B.; and Han, B. 2023. Think while you write: Hypothesis verification promotes faithful knowledge-to-text generation. *arXiv preprint arXiv:2311.09467*.
- Ravi, S. S.; Mielczarek, B.; Kannappan, A.; Kiela, D.; and Qian, R. 2024. Lynx: An open source hallucination evaluation model. *arXiv preprint arXiv:2407.08488*.
- Rawte, V.; Chakraborty, S.; Pathak, A.; Sarkar, A.; Tonmoy, S.; Chadha, A.; Sheth, A. P.; and Das, A. 2023. The troubling emergence of hallucination in large language models—an extensive definition, quantification, and prescriptive remediations. *arXiv preprint arXiv:2310.04988*.
- Razumovskaia, E.; Vulić, I.; Marković, P.; Cichy, T.; Zheng, Q.; Wen, T.-H.; and Budzianowski, P. 2024. Dial beinfo for faithfulness: Improving factuality of information-seeking dialogue via behavioural fine-tuning. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, 17139–17152.
- Si, C.; Gan, Z.; Yang, Z.; Wang, S.; Wang, J.; Boyd-Graber, J.; and Wang, L. 2022. Prompting gpt-3 to be reliable. *arXiv preprint arXiv:2210.09150*.
- Sun, Y.; Yin, Z.; Guo, Q.; Wu, J.; Qiu, X.; and Zhao, H. 2024. Benchmarking Hallucination in Large Language Models Based on Unanswerable Math Word Problem. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, 2178–2188.
- Touvron, H.; Martin, L.; Stone, K.; Albert, P.; Almahairi, A.; Babaei, Y.; Bashlykov, N.; Batra, S.; Bhargava, P.; Bhosale, S.; et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Varshney, N.; Yao, W.; Zhang, H.; Chen, J.; and Yu, D. 2023. A stitch in time saves nine: Detecting and mitigating hallucinations of llms by validating low-confidence generation. *arXiv preprint arXiv:2307.03987*.
- Wang, F.; Mo, W.; Wang, Y.; Zhou, W.; and Chen, M. 2023a. A causal view of entity bias in (large) language models. *arXiv preprint arXiv:2305.14695*.
- Wang, Y.; Kordi, Y.; Mishra, S.; Liu, A.; Smith, N. A.; Khashabi, D.; and Hajishirzi, H. 2022a. Self-instruct: Aligning language models with self-generated instructions. *arXiv preprint arXiv:2212.10560*.
- Wang, Y.; Mishra, S.; Alipoormolabashi, P.; Kordi, Y.; Mirzaei, A.; Arunkumar, A.; Ashok, A.; Dhanasekaran, A. S.; Naik, A.; Stap, D.; et al. 2022b. Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. *arXiv preprint arXiv:2204.07705*.
- Wang, Z.; Cai, S.; Chen, G.; Liu, A.; Ma, X.; and Liang, Y. 2023b. Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents. *arXiv preprint arXiv:2302.01560*.
- Wei, J.; Wang, X.; Schuurmans, D.; Bosma, M.; Xia, F.; Chi, E.; Le, Q. V.; Zhou, D.; et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35: 24824–24837.
- Yao, S.; Zhao, J.; Yu, D.; Du, N.; Shafran, I.; Narasimhan, K.; and Cao, Y. 2022. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*.
- Yao, Y.; Wang, P.; Tian, B.; Cheng, S.; Li, Z.; Deng, S.; Chen, H.; and Zhang, N. 2023. Editing large language models: Problems, methods, and opportunities. *arXiv preprint arXiv:2305.13172*.
- Zhang, H.; Diao, S.; Lin, Y.; Fung, Y. R.; Lian, Q.; Wang, X.; Chen, Y.; Ji, H.; and Zhang, T. 2023. R-tuning: Teaching large language models to refuse unknown questions. *arXiv preprint arXiv:2311.09677*.
- Zhang, T.; Kishore, V.; Wu, F.; Weinberger, K. Q.; and Artzi, Y. 2019. BERTScore: Evaluating Text Generation with BERT. *CoRR*, abs/1904.09675.
- Zhao, Y.; Yan, L.; Sun, W.; Xing, G.; Meng, C.; Wang, S.; Cheng, Z.; Ren, Z.; and Yin, D. 2023. Knowing what llms do not know: A simple yet effective self-detection method. *arXiv preprint arXiv:2310.17918*.
- Zheng, S.; Huang, J.; and Chang, K. C.-C. 2023. Why Does ChatGPT Fall Short in Providing Truthful Answers? *arXiv:2304.10513*.
- Zhou, D.; Schärli, N.; Hou, L.; Wei, J.; Scales, N.; Wang, X.; Schuurmans, D.; Cui, C.; Bousquet, O.; Le, Q.; et al. 2022. Least-to-most prompting enables complex reasoning in large language models. *arXiv preprint arXiv:2205.10625*.