# Heaven-Sent or Hell-Bent? Benchmarking the Intelligence and Defectiveness of LLM Hallucinations

Chengxu Yang
Wuhan University of Technology
Hubei Key Laboratory of
Transportation Internet of Things
BreathingCORE
China
311368@whut.edu.cn

Jingling Yuan*
Wuhan University of Technology
Hubei Key Laboratory of
Transportation Internet of Things
China
yjl@whut.edu.cn

Siqi Cai
Wuhan University of Technology
Hubei Key Laboratory of
Transportation Internet of Things
China
csiqi@whut.edu.cn

Jiawei Jiang
Wuhan University
China
jiawei.jiang@whu.edu.cn

Chuang Hu*
Wuhan University
China
chuanghu@um.edu.mo

## Abstract

Hallucinations in large language models (LLMs) are commonly regarded as errors to be minimized. However, recent perspectives suggest that some hallucinations may encode creative or epistemically valuable content, a dimension that remains underquantified in current literature. Existing hallucination detection methods primarily focus on factual consistency, struggling to handle heterogeneous scientific tasks and balance creativity with accuracy. To address these challenges, we propose HIC-Bench, a novel evaluation framework that categorizes hallucinations into Intelligent Hallucinations (IH) and Defective Hallucinations (DH), enabling systematic investigation of their interplay in LLM creativity. HIC-Bench features three core characteristics: (1) Structured IH/DH Assessment. using a multi-dimensional metric matrix integrating Torrance Tests of Creative Thinking (TTCT) metrics (Originality, Feasibility, Value) with hallucination-specific dimensions (scientific plausibility, factual deviation); (2) Cross-Domain Applicability. spanning ten scientific domains with open-ended innovation tasks; and (3) Dynamic Prompt Optimization. leveraging the Dynamic Hallucination Prompt (DHP) to guide models toward creative and reliable outputs. The evaluation process employs multiple LLM judges, averaging scores to mitigate bias, with human annotators verifying IH/DH classifications. Experimental results reveal a nonlinear relationship between IH and DH, demonstrating that creativity and correctness can be jointly optimized. These insights position IH as a catalyst for creativity and reveal the ability of LLM hallucinations to drive scientific innovation.Additionally, the HIC-Bench offers a valuable platform for advancing research into the creative intelligence of LLM hallucinations.[1]

## Keywords

Large Language Models, Natural language generation, Benchmark, Dataset

*Corresponding Authors
[1]The code and dataset are available at https://github.com/chujiguangniao/HIC-bench

## 1 Introduction

In recent years, large language models (LLMs) have achieved remarkable progress in diverse domains, including natural language processing (NLP), complex reasoning, and scientific discovery [19, 31]. Notably, their generative capabilities have been successfully applied to scenarios such as medical diagnosis, financial analysis, and scientific research [17, 43]. However, a critical challenge that has emerged is the phenomenon of hallucinations in their generated output, which has become a significant bottleneck in their deployment, as extensively documented in the literature. Hallucinations are typically characterized as instances where the model's output deviates from factual accuracy or user expectations, posing a substantial risk in applications where reliability is paramount. To address this issue, the research community has developed several evaluation benchmarks aimed at quantifying and analyzing hallucinations, including TruthfulQA [29], UHGEval [28], and HalluDial [33]. These benchmarks primarily focus on assessing the truthfulness and reliability of the content generated by LLMs.

However, the essence of this phenomenon goes far beyond mere errors. Research suggests that LLM hallucinations reflect, to some extent, the core attributes of human creative cognition, specifically the capacity for divergent exploration and recombination beyond the boundaries of established knowledge [18, 23, 31]. This ability to transcend factual constraints bears a striking resemblance to the unconstrained imagination exhibited by humans in artistic innovation and scientific breakthroughs. For example, in the domain of protein design, "hallucinated" proteins generated by LLM structures that, although absent, were subjected to experimental validation have been demonstrated to exhibit stable configurations and functional properties in a range of tested scenarios [3, 4]. Similarly, in the field of robotic navigation, a novel paradigm termed "Learning from Hallucinations" has been proposed [44], in which innovative patterns emerging from hallucinations are leveraged to optimize path planning for robotic systems. These findings collectively indicate that hallucinations, far from being solely a deficiency, may also represent an emergent manifestation of creative reasoning.

Drawing on this emerging perspective, this paper investigates the relationship between LLM hallucinations and creativity. As shown in Figure 1, we propose a distinction between two categories

of hallucinations: Defective Hallucinations (**DH**), which encompass content contradicting established facts or scientific principles [20], and Intelligent Hallucinations (**IH**), which, though misaligned with reality, are grounded in plausible reasoning and exhibit innovative characteristics [37]. It is posited that IH constitutes a distinct dimension of LLMs' generative capabilities and may catalyze scientific advancements [22]. To this end, we explore strategies to mitigate DH prevalence while preserving and enhancing IH proportion, augmenting LLMs' utility in creative tasks. This investigation departs from the conventional paradigm focused on "reducing hallucinations," offering a novel lens to elucidate parallels between LLMs and human cognition.

To substantiate this perspective, we introduce **HIC-Bench** (Hallucination & Innovation Classification Benchmark), a novel evaluation benchmark designed to quantify the hallucination performance of LLMs in open-ended creative tasks. Unlike prior work, which predominantly focuses on eliciting induced hallucinations, HIC-Bench is informed by LLMs creativity assessment tasks [9, 32, 39] and integrates hallucination detection methodologies [28] to construct an open-domain question-answering dataset spanning ten scientific disciplines. This dataset, meticulously crafted to incorporate domain-specific core principles and real-world requirements, aims to authentically emulate the reasoning processes inherent in scientific innovation. By employing varied prompt strategies, we systematically examine their impact on the prevalence of DH and IH. Experimental findings suggest that reducing DH does not necessarily precipitate a decline in innovativeness, indicating that it may be feasible to diminish DH while preserving or even enhancing the proportion of IH.

In summary, our key contributions are as follows:

- We introduce a new perspective to extract the most effective hallucination component in LLMs and define it as a valuable resource for advancing scientific innovation.
- We propose HIC-Bench, a comprehensive evaluation benchmark that integrates the assessment of LLMs' scientific creativity with hallucination analysis distinguished by intelligent and defective, incorporating a dataset spanning ten scientific domains to evaluate hallucination performance in creative tasks.
- We introduce the Dynamic Hallucination Prompt (DHP) pipeline, which facilitates dynamic refinement of model generation by iteratively refining the construction of positive and negative examples, enabling the mitigation of the defective nature of hallucinations while enhancing their intelligent aspect in creative scientific contexts.
- Evaluation results elucidate that the relationship between intelligent hallucinations and defective hallucinations does not exhibit a simple positive correlation but appears nonlinear, modulated by multiple factors, suggesting that intelligent hallucinations can be preserved or enhanced under suitable constraints to foster the development of novel scientific concepts.

## 2 Related Works

Our work is inspired by studies in LLM creativity assessment, hallucination classification, and evaluation methodologies.
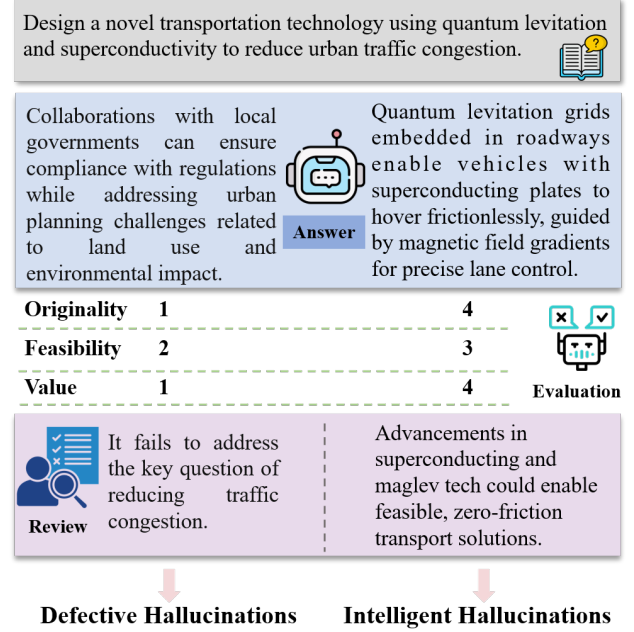


**Figure 1: Example Instances of Hallucination Types, Evaluation, and Review in HIC-Bench**

### 2.1 Hallucinations in LLMs

LLMs frequently generate hallucinations, outputs lacking factual grounding, a phenomenon classified by various frameworks. Huang et al. categorize hallucinations into factual and faithful types [20]. Li et al. separate them into intrinsic and extrinsic forms [6]. Zhang et al. define them as Input-Conflicting, Context-Conflicting, and Fact-Conflicting [46]. Furthermore, benchmarks like HaluEval [27], UHGEval [28], and HalluDial [33] quantitatively analyze these hallucinations but often neglect their creative potential. Building on this foundation, HIC-Bench uniquely distinguishes Intelligent Hallucinations (IH), which promote creativity, from Defective Hallucinations (DH), which indicate inaccuracies, emphasizing their role in scientific innovation through a human creativity perspective.

### 2.2 Assessing Creativity in LLMs

The creative capacity of LLMs in generative artificial intelligence has lately elicited considerable scholarly interest. Prevailing methods for evaluating creativity, grounded in the Torrance Tests of Creative Thinking (TTCT) [41], derive from Guilford's framework of divergent thinking [14]. These methodologies assess four principal dimensions: (1) originality: the aptitude for devising novel concepts; (2) elaboration: the proficiency in refining and expanding notions; (3) fluency: the capability to generate abundant ideas; (4) flexibility: the facility for reasoning across disparate domains. Extending prior investigations on LLMs' creativity [32], we propose an evaluation system specifically calibrated for hallucination-driven creativity. Two critical challenges persist in incorporating hallucinations into this paradigm: (1) redefining metrics of innovation for

LLMs in scientific endeavors to delineate IH from DH; (2) modifying and augmenting established creativity frameworks to precisely measure IH's contributions.

## 2.3 LLMs-as-a-judge

As LLMs increasingly permeate natural language processing tasks, their functionality has expanded beyond content generation to include automated evaluation of text, code, and scientific outputs. Their capacity to assess quality in open-ended questions, creative tasks, and complex reasoning scenarios has garnered attention [26]. Research indicates that, in certain contexts, LLMs' evaluations align closely with human expert judgments [47], underscoring their potential in automated assessment. Existing frameworks for LLMs typically focus on tasks with fixed answers, such as text quality, code accuracy, and factual consistency, where performance depends on matching reference standards. Open-ended tasks, like scientific innovation evaluation or conceptual divergence testing, however, lack singular correct responses. Drawing on TTCT metrics and previous definitions of LLM's creative capacity [22], we propose a refined Prompt design framework for evaluating scientific creativity. This framework establishes metrics across creativity and hallucination dimensions, while enforcing strict variable control to ensure consistency and interpretability in LLMs' assessments.

## 3 HIC-bench: From scientific hallucinations to controllable innovation

In this section, we elaborate on the design philosophy and core components of HIC-Bench, a modular and extensible benchmark designed to evaluate hallucinations of LLMs in scientific creativity tasks. Departing from traditional hallucination studies focused solely on error correction, HIC-Bench leverages cognitive mechanisms of human creative thinking to systematically explore the dual nature of LLMs' outputs, namely DH and IH. Through a meticulously designed cross-domain task set, dynamic generation strategies, and a multi-tiered evaluation system, HIC-Bench offers researchers a structured platform to deeply investigate the role of LLMs' hallucinations in scientific innovation. Its core workflow is illustrated in Figure 2.

## 3.1 Framework Structures

HIC-Bench addresses two pivotal research questions: Do traditional hallucination mitigation strategies inadvertently curtail LLMs' creative capacity? Is it feasible to devise a scientific approach that diminishes DH while amplifying IH, thereby optimizing both innovation and reliability? Beyond serving as an assessment tool, HIC-Bench embodies a paradigm shift, moving from suppressing hallucinations to cultivating their beneficial forms. We combine task-driven goals with the intrinsic properties of LLM into a three-layer modular system that supports cross-model evaluation, scalable dataset integration, and highly repeatable automated workflows:

**Cross-disciplinary Task Set**. HIC-Bench constructs an open-ended task set spanning ten pivotal scientific domains, replicating the generative challenges LLMs face in real-world scientific innovation, with questions designed to embody theoretical depth, interdisciplinary complexity, and cutting-edge relevance. Unlike prior datasets assessing only knowledge breadth or factual accuracy, this

approach deliberately incorporates "**fuzzy factual boundaries**", eliciting outputs that blend hypothetical, predictive, and analogical elements diverging from reality. Such a design unlocks LLMs' latent creativity and hallucination expression. The dataset's question formulation examines not merely knowledge retrieval proficiency, but emphasizes LLMs' capacity for generating "useful hallucinations" within uncharted problem spaces, fostering insights that transcend conventional boundaries.

**Multi-strategy Generation Control**. A suite of Prompt strategies, including Strict Constraint Prompt (SCP), Relaxed Constraint Prompt (RCP), RAG, and CoT among others, is harnessed to explore their influence on hallucination profiles in LLMs systematically. This approach investigates the interplay between innovation and factual fidelity, adapting output characteristics to varying scientific contexts.

**Multi-dimensional Evaluation System**. Conventional hallucination detection hinges on factual consistency or semantic similarity, yet we argue that in scientific creativity tasks, certain hallucinations constitute "high-dimensional cognitive outputs", resisting simplistic classification as errors. HIC-Bench establishes a multi-dimensional evaluation system rooted in creativity assessment frameworks, precisely identifying and quantifying IH. This system merges key TTCT metrics, such as originality, feasibility, and value, with dimensions like scientific plausibility and factual deviation from hallucination analysis, forming a specialized metric matrix tailored to adjudicate LLMs' outputs. We introduce a composite metric, Intelligent-Fidelity Score (**IFS**), which gauges the dynamic equilibrium between creativity and factual integrity in generated results. The evaluation employs multiple LLM judges, averaging their scores to mitigate bias, with human annotators verifying intelligent versus defective hallucination (IH vs. DH) classifications, ensuring robust and interpretable assessments. Beyond enabling quantitative recognition of IH, this framework lays an actionable foundation for constructing and modulating "valuable hallucinations" in future research. For further details, please refer to Appendix A.2

## 3.2 Evaluation Metrics

We integrate metrics from human creativity assessment tasks to devise an evaluation system tailored for LLMs. Through multi-dimensional metrics and a composite scoring formula, this system quantifies the scientific merit and innovative potential of LLMs' outputs.

*3.2.1 LLM Creativity Assessment.* Creativity hinges on novelty, often termed originality or uniqueness, and value, encompassing utility, effectiveness, or relevance [8, 22, 40]. Grounded in this conceptualization, alongside the TTCT framework and LLMs' generative traits, we establish three creativity metrics suited for scientific tasks:

**Originality(Or)**. This metric gauges the novelty and breakthroughs of generated content within academic research, evaluating whether it introduces theoretically valuable perspectives, methodological innovations, or discoveries. Originality demands that LLMs transcend existing paradigms, yielding pioneering outputs at the frontiers of knowledge rather than mere recombinations of prior work.
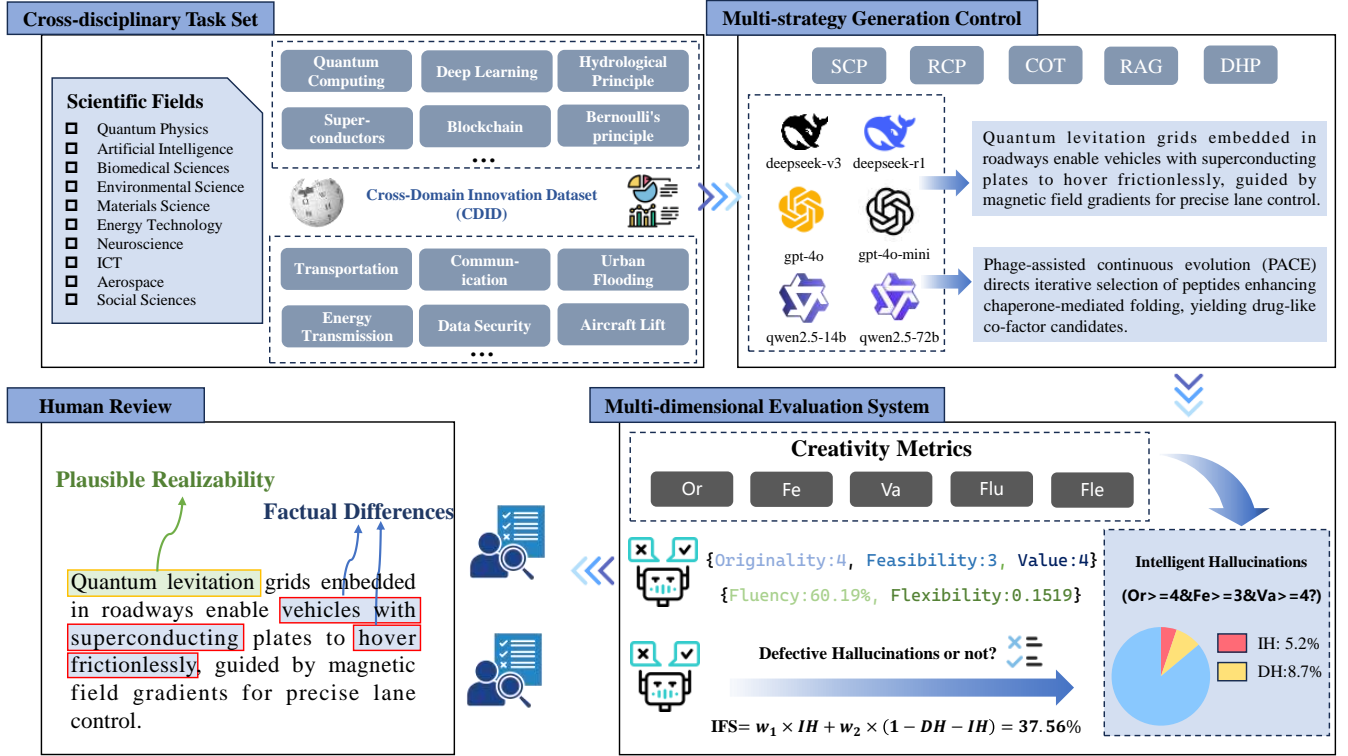
**Figure 2: Design of the HIC-Bench Framework: We first construct a cross-disciplinary dataset by integrating Wikipedia principles with real-world problems. We then evaluate the model's responses across various generation strategies, analyzing its creativity and hallucinations through automated evaluation. Finally, human reviewers perform token-level analysis on the model's evaluation outputs to enhance their reliability."Or" for Originality, "Fe" for Feasibility, "Va" for Value, "Flu" for Fluency, and "Fle" for Flexibility.**

**Feasibility(Fe)**. This assesses the realizability of generated content under current scientific conditions and technical frameworks, encompassing the rigor of theoretical derivations, the practicality of experimental designs, and the viability of technical pathways. For LLMs' outputs, it examines whether proposed solutions adhere to discipline-specific standards of scientific validation.

**Value(Va)**. Value or contribution evaluates the potential impact of generated content on disciplinary advancement spanning theoretical construction methodological innovation and practical utility. High-quality academic outputs should advance critical scientific questions, forge new research paradigms, or offer valuable interdisciplinary insights.

**Fluency&Flexibility(Flu&Fle)**. Given that LLMs excel at generating abundant ideas and reasoning across knowledge domains [12, 21], Fluency and Flexibility are not prioritized in this paper. However, to maintain compliance with the TTCT standards, results for Fluency and Flexibility are still presented.

Each of the first three metrics is scored on a 5-point Likert scale, with explicit criteria defined to ensure transparency and consistency, thereby mitigating ambiguity in LLMs' evaluation. For Fluency, similarity assessment of model outputs is conducted using

SimCSE [13]. Meanwhile, flexibility is evaluated by calculating the variance in scores across different questions.

*3.2.2 IH and DH.* We design the evaluation of IH and DH to classify hallucinations based on the aforementioned metrics, leveraging them to distinguish the Intelligent potential and limitations of LLM outputs in a structured manner.

**IH**. Hallucinations in LLMs share parallels with human innovative thinking, producing reasoning outcomes that diverge from the status quo yet potentially harbor groundbreaking ideas for scientific advancement [23]. This intelligence is described as results or novel idea combinations unlikely to emerge from most individuals [37]. To isolate IH from standard knowledge-based responses, we define generated content as valuable IH if it satisfies three creativity metrics, originality and value each scoring $\geq 4$, and feasibility $\geq 3$ on a 5-point scale, reflecting high innovation and worth despite partial factual divergence, while retaining plausible realizability. The IH ratio is formulated in Equation 1, where $N_{IH}$ denotes the count of IH instances, and $N_{total}$ represents the total number of generated outputs.

$$IH_{ratio} = \frac{N_{IH}}{N_{total}} \quad (1)$$

**DH**. Generated content classified as DH includes outputs marred by factual errors, logical inconsistencies, or severe violations of scientific principles, lacking innovation or practical utility. We employ LLMs for automated scientific validation of responses, supplemented by human evaluation to confirm fidelity. Consistency is further assessed using kwPrec [28], a keyword-segmentation metric. The DH ratio is expressed in Equation 2, where $N_{DH}$ indicates the count of DH instances.

$$DH_{ratio} = \frac{N_{DH}}{N_{total}} \tag{2}$$

Notably, $IH_{ratio} + DH_{ratio} \neq 1$, as outputs neither innovative nor factually errant are classified as neutral noise responses, distinct from both IH and conventional DH traits.

*3.2.3 IFS.* To evaluate the balanced performance of LLMs across creativity and hallucination tendencies, we propose the IFS, a unified metric that integrates results from creativity assessment and hallucination classification. This score offers a comprehensive measure of the quality of generated content. The computation of IFS is detailed in Equation 3, where $w_1$ and $w_2$ are weight parameters summing to 1, enabling adaptability to diverse task requirements. For instance, creative writing tasks may prioritize imaginative hallucinations with a higher $w_1$, whereas medical applications demand greater accuracy with an elevated $w_2$. In our evaluation, $w_1$ and $w_2$ are set to 0.6 and 0.4 respectively, prioritizing intelligent hallucinations while ensuring output fidelity.

$$IFS = w_1 \times IH_{ratio} + w_2 \times (1 - DH_{ratio} - IH_{ratio}) \tag{3}$$

## 3.3 Dataset

We construct a cross-domain dataset supporting HIC-Bench's evaluation capabilities, spanning ten scientific fields: Quantum Physics, Artificial Intelligence, Biomedical Sciences, Environmental Science, Materials Science, Energy Technology, Neuroscience, Information and Communication Technology(ICT), Aerospace, and Social Sciences. These domains encompass a wide range from foundational science to applied technology, ensuring the assessment captures diversity and representativeness.

**Design of Open-ended Questions and Hallucination Mechanisms**. Open-ended questions are crafted with the intent of unleashing LLMs' creative potential while grounding their outputs in reasoned speculation and innovative reasoning drawn from established knowledge. Task descriptions explicitly require LLMs to harness domain-specific insights, such as accounting for quantum entanglement constraints in quantum communication protocols, or incorporating recent deep learning advances into disease diagnostic methods. This structure promotes novel solutions, eliciting outputs that transcend simple recombinations of existing knowledge and offer plausible conjectures not fully bound by reality.

**Cross-Domain Innovation Dataset (CDID)**. Drawing on core principles and techniques extracted from Wikipedia, we build a knowledge foundation, enriched with real-world frontier challenges like transportation and communication, shaping open-ended, cross-disciplinary innovation tasks. Assessing LLMs' innovative diversity comprehensively, we elicit ten responses per task, analyzing 6000 responses across six models and ten domains, covering 100 distinct innovation tasks. Response generation employs diverse strategies

under meticulous control to maintain variability. Additionally, we establish a specialized knowledge base dataset (CDKB) for these strategies, blending Wikipedia data with LLMs' in-depth analyses of task-specific principles, thereby enhancing knowledge precision.

The CDID dataset contains 100 open-domain innovation tasks. This scale is consistent with contemporary practices in the evaluation of generative language models, where task quality and design often outweigh dataset size in determining benchmarking effectiveness. Empirical evidence suggests that compact yet diverse benchmarks can still offer meaningful perspectives on model behavior. For instance, HumanEval [10] and Bamboogle [36], containing 164 and 125 items respectively, have become standard references for code generation and question answering tasks. Furthermore, to ensure statistical rigor, we also conduct significance tests across all compared models. For further details, please refer to Appendix B
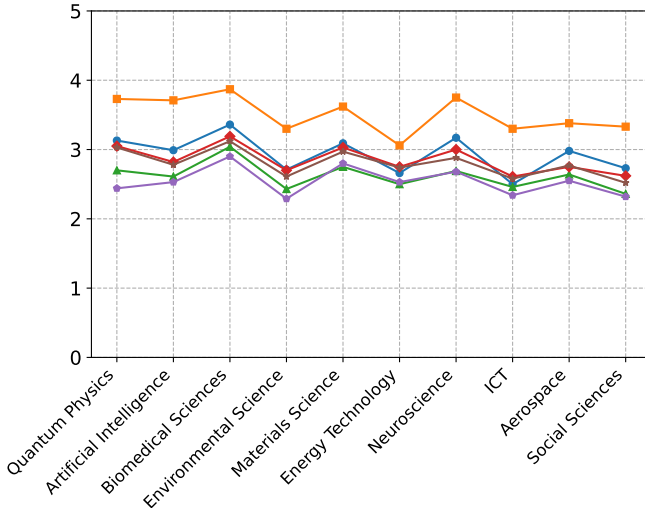
## 4 EXPERIMENTS & DISCUSSION

This chapter employs HIC-Bench to evaluate LLMs' creativity and hallucination performance on the CDID dataset while concurrently exploring the influence of diverse prompt strategies on IH and DH. Experiments utilize the CDID dataset and CDKB knowledge base, encompassing two parallel investigations: the impact of hallucination mitigation techniques on preserving IH, and methods to reduce DH while maintaining IH. These studies advance the understanding of balancing factual fidelity with innovative output in scientific contexts, identifying dynamic adjustments to enhance model performance.

**Selected LLM Models**. Six technically diverse LLMs are selected to evaluate HIC-Bench's applicability across varied architectures and optimization approaches: gpt-4o-2024-11-20 [21] serves as a benchmark for general-purpose performance; gpt-4o-mini [2] represents lightweight models optimized via knowledge distillation; qwen2.5-14b-instruct [5] and qwen2.5-72b-instruct [45] enable analysis of parameter scale effects; deepseek-v3 [30], refined through reinforcement learning, acts as a foundation model; and deepseek-r1 [16] excels in specialized CoT reasoning capabilities. This selection establishes a systematic evaluation framework, facilitating comparative analysis of model compression techniques, parameter scaling, reinforcement learning strategies, and dedicated reasoning architectures.
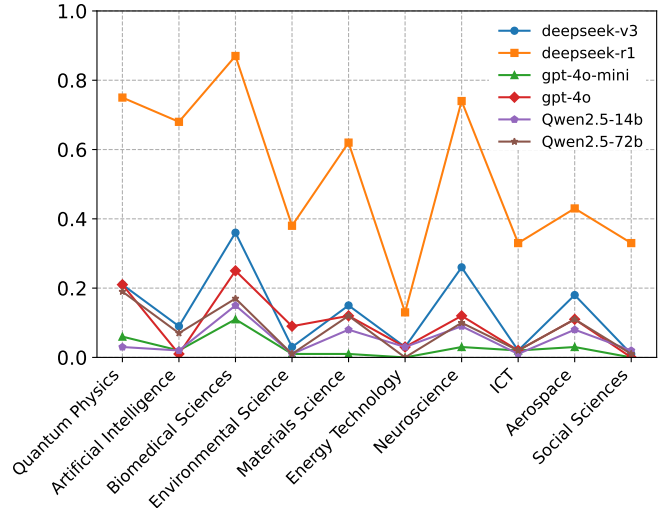
**Experimental Hyperparameters**. Generation employs two foundational strategies: SCP (Strict Constraint Prompt) and RCP (Relaxed Constraint Prompt). SCP enforces stringent control through prompts demanding factual adherence, while RCP fosters diversity by encouraging innovative thinking. Tuning hyperparameters across varied LLMs presents complex challenges, as models differ in architecture and respond distinctly to identical configurations. The principle of harmonizing output stability with creative latitude while preserving cross-model consistency shapes the settings. Generation temperature, set at 1.0, balances creativity with scientific rigor, producing nuanced responses for open-ended tasks. Evaluation temperature, fixed at 0, ensures precision in distinguishing IH and DH during the assessment. The maximum token length is set to 70. These configurations enable a thorough evaluation of LLMs' hallucination and creativity performance, supporting consistent comparisons across diverse prompt strategies.

**Table 1: Creativity and Hallucination Performance of LLMs Across Temperature Settings (T). "Or" for Originality, "Fe" for Feasibility, "Va" for Value, "Flu" for Fluency, and "Fle" for Flexibility, Where ↑ Denotes Higher Values Are Better and ↓ Denotes Lower Values Are Better**

| LLM | T | Creativity | | | | | Hallucination | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Or↑ | Fe↑ | Va↑ | Flu↓ | Fle↓ | IH↑ | DH↓ | IFS↑ |
| deepseek-v3 | 1.0 | 2.93 | 3.82 | 3.22 | 63.00% | 0.15 | 13.40% | 3.20% | 41.40% |
| | 0.4 | 2.90 | 3.79 | 3.18 | 62.46% | 0.14 | 12.40% | 1.40% | 41.92% |
| gpt-4o-mini | 1.0 | 2.62 | 3.73 | 2.98 | 63.54% | 0.09 | 2.90% | 5.00% | 38.58% |
| | 0.4 | 2.57 | 3.74 | 2.97 | 63.74% | 0.08 | 2.20% | 4.50% | 38.64% |
| gpt-4o | 1.0 | 2.85 | 3.86 | 3.14 | 62.39% | 0.14 | 9.60% | 1.30% | 41.40% |
| | 0.4 | 2.75 | 3.82 | 3.06 | 63.00% | 0.14 | 7.10% | 0.60% | 41.18% |
| qwen2.5-14b | 1.0 | 2.54 | 3.65 | 2.96 | 59.35% | 0.12 | 5.20% | 8.70% | 37.56% |
| | 0.4 | 2.47 | 3.67 | 2.93 | 59.17% | 0.13 | 4.70% | 7.90% | 37.40% |
| qwen2.5-72b | 1.0 | 2.80 | 3.76 | 3.17 | 65.51% | 0.11 | 8.00% | 1.70% | 40.92% |
| | 0.4 | 2.69 | 3.74 | 3.13 | 66.13% | 0.11 | 6.60% | 0.90% | 40.85% |



(a) Or across scientific domains



(b) IH across scientific domains

**Figure 3: Comparative Analysis of Originality and IH Across Scientific Domains**

## 4.1 Evaluating Creativity and Hallucination Dynamics in LLMs

**Temperature Effects on Creativity and Hallucination**. To investigate how temperature adjustments influence LLMs' creativity and hallucination profiles, we compare settings of 1.0 as the baseline and 0.4 for conservative outputs, using the standard SCP prompt to assess HIC-Bench under controlled variations. Deepseek-r1 [16] is excluded due to its lack of temperature support. Table 1 indicates that the decrease in temperature to 0.4 reduces DH in all models, with gpt-4o-mini [2] declining from 5.00% to 4.50% and showing the most significant change, though IH also decreases, with deepseek-v3 dropping from 13.40% to 12.40%. This behavior reflects temperature's regulation of the model's generation distribution, where

lower temperatures lead models to favor high-probability conservative outputs, enhancing precision at the expense of creative exploration [38]. IFS scores remain largely stable across temperatures, suggesting that temperature adjustments primarily redistribute IH and DH without altering the overall innovation-factuality balance. These findings, consistent with prior studies, affirm HIC-Bench's capability to capture nuanced behavioral shifts in LLMs under temperature variations.

**Comparative Analysis of Originality and IH Across Scientific Domains**. We examine the generative creativity of large language models across diverse disciplinary domains by evaluating two key dimensions: Originality (Or) and Intelligent Hallucinations (IH), under SCP prompts. Figure 3 (a) presents the distribution of originality scores, where most models exhibit heightened creativity

in Biomedical Sciences and Aerospace, while comparatively lower scores are observed in Environmental Science, Social Sciences, and Energy Technology. This variation suggests that generative originality is not uniform but instead shaped by the epistemic structure of each domain. Figure 3 (b) displays the corresponding IH proportions, which follow a broadly similar trend. Among all evaluated systems, deepseek-r1 consistently outperforms its counterparts on both metrics, reflecting its advanced capacity for creative text generation. These findings collectively highlight the domain-sensitive dynamics of LLM creativity, demonstrating systematic shifts in performance across disciplinary boundaries.

## 4.2 Prompt Strategy Effects on Model Performance

Reducing hallucinations in LLMs without compromising efficacy is challenging [25], as they include DH and a creative facet as IH. As such, reducing DH while preserving IH is a key research goal. However, overly strict constraints may suppress creative capabilities [1]. Table 2 compares the SCP, CoT, RAG, and RCP strategies' impact on IH, DH and IFS. More details are provided in the AppendixF.2.

**SCP**. As the base strategy, SCP enforces logical coherence and creativity through structured prompts designed to integrate existing research insights. It aims to balance innovation with factuality, fostering IH generation while maintaining factual integrity in scientific contexts. This approach provides a foundational framework for evaluating hallucination profiles across different models. Results guide HIC Bench's hallucination assessment and support further strategy comparisons.

**CoT**. CoT uses "Let's think step by step" to boost accuracy [24] via incremental reasoning. It lowers DH and raises IFS but may reduce IH in some models like gpt-4o. Model response varies with this approach.

**RAG**. RAG integrates the CDKB dataset for factual grounding, effectively reducing DH. However, this approach also constrains IH, revealing a trade-off between reliability and creative output. highlighting a trade-off between reliability and creativity in open-ended scientific tasks.

**RCP**. RCP relaxes SCP constraints to prioritize innovation and value, boosting IH across models. Unexpectedly, it also reduces DH, suggesting moderated constraints enhance creativity while lowering factual deviations.

## 4.3 Sensitivity Analysis of IFS

We present a horizontal comparison of model performance under different IFS scoring scenarios in the HIC-Bench evaluation. Figure 4 visualizes the results for Intelligent IFS (IIFS, $w_1 = 0.9$) applied in scenarios emphasizing intelligent hallucinations, Balanced IFS (BIFS, $w_1 = 0.6$) used in balanced creativity and fidelity contexts, and Robust IFS (RIFS, $w_1 = 0.1$) utilized in scenarios requiring high accuracy. In the first two metrics IIFS and BIFS deepseek r1 achieves higher scores demonstrating strong performance in creativity driven scenarios. However in accuracy focused environments as measured by RIFS this model is less suitable due to challenges in maintaining precision.

These findings highlight the trade offs in model performance across different scenarios. Models excelling in creativity driven contexts may struggle in precision focused tasks emphasizing the need to select appropriate models based on the specific requirements of each application.
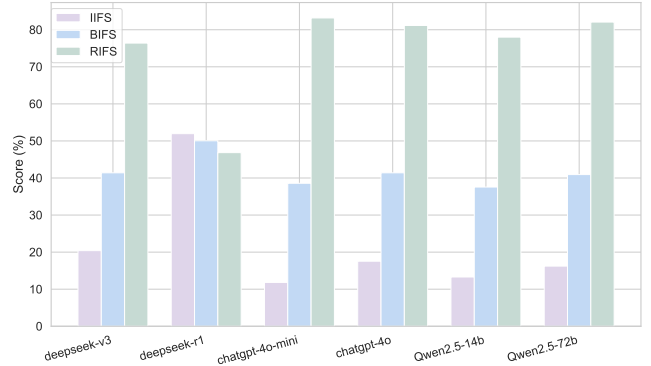


**Figure 4: Comparative analysis of model indicators across scientific fields under SCP**

## 4.4 Dynamic Prompt Effects on Hallucination Profiles

To address the challenge of mitigating DH while enhancing IH, we propose a novel pipeline, Dynamic Hallucination Prompt (DHP). The DHP pipeline optimizes response generation and evaluation in HIC-Bench by dynamically adapting prompts based on real-time feedback. Initially, human experts input positive examples—responses scoring high on Torrance Tests of Creative Thinking (TTCT) metrics and negative examples—responses with high DH rates to establish a baseline. During machine evaluation, responses with IFS scores exceeding the initial positive examples are set as new positive examples, while negative examples are continuously updated with the latest DH responses. This observation motivated the design of the Dynamic Hallucination Prompt (DHP), which provides examples to guide model outputs toward more distinct and reliable responses. Detailed procedures are provided in Appendix D.

Ablation experiments on gpt-4o and gpt-4o-mini use SCP as the baseline, with results shown in Table 3. With DHP positive prompting, DH drops notably to 0.90% for gpt-4o-mini and 0.10% for gpt-4o, while IH increases moderately. When DHP is combined with RCP constraints, IH further rises to 11.10% and 21.10%, with DH reduced to 1.70% and 0.30%, respectively. This setting also yields the highest IFS scores at 41.54% for gpt-4o-mini and 44.10% for gpt-4o, exceeding their SCP baselines. These outcomes underscore the efficacy of DHP in balancing innovation and factuality, outperforming prior strategies in creative generation tasks.

## 5 Limitations and Future Work

HIC-Bench has advanced the understanding of IH and DH interplay in LLMs and the evaluation of mitigation strategies, yet limitations persist. The framework primarily focuses on structured tasks like question-answering, neglecting areas such as literary writing.

**Table 2: Creativity and Hallucination Performance of Six LLMs Under Diverse Prompt Strategies, with CoT and RAG Constructed on the SCP Baseline. Bold values indicate the best-performing model–strategy pairs. "P" denotes the applied prompt strategy.**

| LLM | P | Creativity | | | | | Hallucination | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Or↑ | Fe↑ | Va↑ | Flu↓ | Fle↓ | IH↑ | DH↓ | IFS↑ |
| deepseek-v3 | SCP | 2.93 | **3.82** | 3.22 | **63.00**% | 0.15 | 13.40% | 3.20% | 41.40% |
| | COT | 2.98 | 3.78 | 3.23 | 63.06% | 0.16 | 17.00% | **1.90**% | 42.64% |
| | RAG | 2.84 | 3.81 | 3.15 | 64.31% | 0.14 | 9.80% | 2.20% | 41.08% |
| | RCP | **3.13** | 3.75 | **3.33** | 63.89% | **0.12** | **22.20**% | 2.10% | **43.60**% |
| deepseek-r1 | SCP | 3.50 | 3.75 | 3.63 | **60.19**% | 0.15 | 52.60% | 1.20% | 50.04% |
| | COT | 3.50 | 3.77 | 3.63 | 60.35% | 0.17 | 52.60% | 0.60% | 50.28% |
| | RAG | 3.46 | **3.79** | 3.59 | 60.99% | 0.16 | 49.00% | **0.20**% | 49.72% |
| | RCP | **3.70** | 3.57 | **3.76** | 60.46% | **0.09** | **71.10**% | 0.30% | **54.10**% |
| gpt-4o-mini | SCP | 2.62 | 3.73 | 2.98 | 63.54% | 0.09 | 2.90% | 5.00% | 38.58% |
| | COT | 2.59 | **3.75** | 2.97 | **62.89**% | 0.09 | 2.80% | 4.60% | 38.72% |
| | RAG | 2.58 | 3.73 | 2.95 | 64.65% | 0.08 | 1.70% | 3.50% | 38.94% |
| | RCP | **2.88** | 3.71 | **3.11** | 64.00% | **0.07** | **8.30**% | 2.60% | **40.62**% |
| gpt-4o | SCP | 2.85 | **3.86** | 3.14 | 62.39% | 0.14 | 9.60% | 1.30% | 41.40% |
| | COT | 2.81 | 3.84 | 3.13 | **62.16**% | 0.15 | 9.30% | 0.40% | 41.70% |
| | RAG | 2.81 | 3.83 | 3.11 | 64.50% | 0.14 | 7.90% | **0.10**% | 41.54% |
| | RCP | **3.08** | 3.73 | **3.22** | 62.58% | **0.10** | **18.60**% | 0.50% | **43.52**% |
| qwen2.5-14b | SCP | 2.54 | 3.65 | 2.96 | 59.35% | 0.12 | 5.20% | 8.70% | 37.56% |
| | COT | 2.50 | 3.65 | 2.93 | **58.88**% | 0.12 | 4.30% | 5.20% | 38.78% |
| | RAG | 2.52 | **3.68** | 2.93 | 62.79% | **0.09** | 4.70% | 5.80% | 38.62% |
| | RCP | **2.67** | 3.56 | **2.98** | 59.52% | 0.12 | **6.30**% | 4.50% | **39.46**% |
| qwen2.5-72b | SCP | 2.80 | 3.76 | 3.17 | 65.51% | 0.11 | 8.00% | 1.70% | 40.92% |
| | COT | 2.77 | **3.80** | 3.14 | 65.55% | 0.12 | 7.00% | 1.40% | 40.84% |
| | RAG | 2.69 | 3.79 | 3.13 | 66.71% | 0.11 | 4.80% | 1.30% | 40.44% |
| | RCP | **2.87** | 3.65 | **3.19** | **63.27**% | **0.10** | **10.40**% | **1.20**% | **41.60**% |

**Table 3: Ablation Study on Dynamic Hallucination Prompt (DHP) and RCP with SCP as Baseline.**

| LLM | RCP | DHP | IH↑ | DH↓ | IFS↑ |
|---|---|---|---|---|---|
| gpt-4o-mini | | | 2.90% | 5.00% | 38.58% |
| | ✓ | | 8.30% | 2.60% | 40.62% |
| | | ✓ | 5.50% | 0.90% | 40.74% |
| | ✓ | ✓ | 11.10% | 1.70% | 41.54% |
| gpt-4o | | | 9.60% | 1.30% | 41.40% |
| | ✓ | | 18.60% | 0.50% | 43.52% |
| | | ✓ | 12.40% | 0.10% | 42.44% |
| | ✓ | ✓ | 21.10% | 0.30% | 44.10% |

Future work will extend to broader creativity tasks for a comprehensive assessment. Additionally, balancing IH and DH remains challenging due to potential biases from inductive prompts and the need for refined IH metrics. Future efforts will develop domain-specific prompts and integrate human-involved multi-dimensional evaluations to enhance IH assessment accuracy and manage IH-DH dynamics. Lastly, the framework will be applied to multimodal and cross-lingual benchmarks to validate its generalizability, while examining the ethical implications of promoting IH.

## 6 Conclusion

This paper introduces HIC-Bench, a benchmark tailored to assess the interplay between Intelligent Hallucinations (IH) and Defective Hallucinations (DH) in large language models (LLMs). The framework's reliability is validated through temperature parameter analysis, with results elucidating the hallucination profiles of mainstream LLMs in creative tasks, followed by a comparative evaluation of four strategies to mitigate hallucinations while minimizing their adverse impact on IH. Notably, our findings reveal that the relationship between IH and DH is not simply positively correlated, suggesting that it may be feasible to enhance creative potential while reducing DH. To this end, the Dynamic Hallucination Prompt (DHP) pipeline is introduced, substantially augmenting the intelligent aspects of model hallucinations. Ultimately, this research quantifies the potential of hallucinations as a scientific dream machine, thereby paving new avenues for future hallucination studies.

## 7 AI Ethics

Although our framework distinguishes between Intelligent Hallucinations (IH) and Defective Hallucinations (DH), users must select an appropriate IFS based on the specific context, as shown in Equation 3. For instance, in high accuracy scenarios, increasing the weight of $w_2$ is recommended, while in high innovation scenarios, elevating $w_1$ is more suitable. Furthermore, while IH may contribute to scientific advancement, it remains a form of hallucination. Users should avoid conflating IH outputs with factual content, as this could lead to misinterpreting IH as reality.

## References

[1] Oguz A Acar, Murat Tarakci, and Daan Van Knippenberg. 2019. Creativity and innovation under constraints: A cross-disciplinary integrative review. *Journal of management* 45, 1 (2019), 96–121.

[2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).

[3] Linna An, Derrick R Hicks, Dmitri Zorine, Justas Dauparas, Basile IM Wicky, Lukas F Milles, Alexis Courbet, Asim K Bera, Hannah Nguyen, Alex Kang, et al. 2023. Hallucination of closed repeat proteins containing central pockets. *Nature Structural & Molecular Biology* 30, 11 (2023), 1755–1760.

[4] Ivan Anishchenko, Samuel J Pellock, Tamuka M Chidyausiku, Theresa A Ramelot, Sergey Ovchinnikov, Jingzhou Hao, Khushboo Bafna, Christoffer Norn, Alex Kang, Asim K Bera, et al. 2021. De novo protein design by deep network hallucination. *Nature* 600, 7889 (2021), 547–552.

[5] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609* (2023).

[6] Yejin Bang, Ziwei Ji, Alan Schelten, Anthony Hartshorn, Tara Fowler, Cheng Zhang, Nicola Cancedda, and Pascale Fung. 2025. HalluLens: LLM Hallucination Benchmark. *arXiv preprint arXiv:2504.17550* (2025).

[7] Forrest Bao, Miaoran Li, Rogger Luo, and Ofer Mendelevitch. 2024. HHEM-2.1-Open. doi:10.57967/hf/3240

[8] Frank Barron. 1955. The disposition toward originality. *The Journal of Abnormal and Social Psychology* 51, 3 (1955), 478.

[9] Tuhin Chakrabarty, Philippe Laban, Divyansh Agarwal, Smaranda Muresan, and Chien-Sheng Wu. 2024. Art or artifice? large language models and the false promise of creativity. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 1–34.

[10] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374* (2021).

[11] Qinyuan Cheng, Tianxiang Sun, Wenwei Zhang, Siyin Wang, Xiangyang Liu, Mozhi Zhang, Junliang He, Mianqiu Huang, Zhangyue Yin, Kai Chen, et al. 2023. Evaluating hallucinations in chinese large language models. *arXiv preprint arXiv:2310.03368* (2023).

[12] Jia Chi. 2024. The evolutionary impact of artificial intelligence on contemporary artistic practices. *Communications in Humanities Research* 35 (2024), 52–57.

[13] Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. Simcse: Simple contrastive learning of sentence embeddings. *arXiv preprint arXiv:2104.08821* (2021).

[14] Joy Paul Guilford. 1967. The nature of human intelligence. (1967).

[15] Anisha Gunjal, Jihan Yin, and Erhan Bas. 2024. Detecting and preventing hallucinations in large vision language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 18135–18143.

[16] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948* (2025).

[17] Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xiangliang Zhang. 2024. Large language model based multi-agents: A survey of progress and challenges. *arXiv preprint arXiv:2402.01680* (2024).

[18] Brett A Halperin and Stephanie M Lukin. 2024. Artificial Dreams: Surreal Visual Storytelling as Inquiry Into AI 'Hallucination'. In *Proceedings of the 2024 ACM Designing Interactive Systems Conference*. 619–637.

[19] Yufang Hou, Alessandra Pascale, Javier Carnerero-Cano, Tigran Tchrakian, Radu Marinescu, Elizabeth Daly, Inkit Padhi, and Prasanna Sattigeri. 2024. Wikicontradict: A benchmark for evaluating llms on real-world knowledge conflicts from wikipedia. *Advances in Neural Information Processing Systems* 37 (2024), 109701–109747.

[20] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2025. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems* 43, 2 (2025), 1–55.

[21] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276* (2024).

[22] Mete Ismayilzada, Debjit Paul, Antoine Bosselut, and Lonneke van der Plas. 2024. Creativity in AI: Progresses and Challenges. *arXiv preprint arXiv:2410.17218* (2024).

[23] Xuhui Jiang, Yuxing Tian, Fengrui Hua, Chengjin Xu, Yuanzhuo Wang, and Jian Guo. 2024. A survey on large language model hallucination via a creativity perspective. *arXiv preprint arXiv:2402.06647* (2024).

[24] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems* 35 (2022), 22199–22213.

[25] Minhyeok Lee. 2023. A mathematical investigation of hallucination and creativity in GPT models. *Mathematics* 11, 10 (2023), 2320.

[26] Haitao Li, Qian Dong, Junjie Chen, Huixue Su, Yujia Zhou, Qingyao Ai, Ziyi Ye, and Yiqun Liu. 2024. Llms-as-judges: a comprehensive survey on llm-based evaluation methods. *arXiv preprint arXiv:2412.05579* (2024).

[27] 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).

[28] Xun Liang, Shichao Song, Simin Niu, Zhiyu Li, Feiyu Xiong, Bo Tang, Yezhaohui Wang, Dawei He, Peng Cheng, Zhonghao Wang, et al. 2023. Uhgeval: Benchmarking the hallucination of chinese large language models via unconstrained generation. *arXiv preprint arXiv:2311.15296* (2023).

[29] Stephanie Lin, Jacob Hilton, and Owain Evans. 2021. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958* (2021).

[30] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437* (2024).

[31] Xuan Liu, Jie ZHANG, HaoYang Shang, Song Guo, Yang Chengxu, and Quanyan Zhu. 2025. Exploring Prosocial Irrationality for LLM Agents: A Social Cognition View. In *The Thirteenth International Conference on Learning Representations*. https://openreview.net/forum?id=u8VOQVzduP

[32] Li-Chun Lu, Shou-Jen Chen, Tsung-Min Pai, Chan-Hung Yu, Hung-yi Lee, and Shao-Hua Sun. 2024. Llm discussion: Enhancing the creativity of large language models via discussion framework and role-play. *arXiv preprint arXiv:2405.06373* (2024).

[33] Wen Luo, Tianshu Shen, Wei Li, Guangyue Peng, Richeng Xuan, Houfeng Wang, and Xi Yang. 2024. Halludial: A large-scale benchmark for automatic dialogue-level hallucination evaluation. *arXiv preprint arXiv:2406.07070* (2024).

[34] Potsawee Manakul, Adian Liusie, and Mark Gales. [n. d.]. SelfCheckGPT: Zero-Resource Black-Box Hallucination Detection for Generative Large Language Models. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.

[35] 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.

[36] Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. 2022. Measuring and narrowing the compositionality gap in language models. *arXiv preprint arXiv:2210.03350* (2022).

[37] Vipula Rawte, Amit Sheth, and Amitava Das. 2023. A survey of hallucination in large foundation models. *arXiv preprint arXiv:2309.05922* (2023).

[38] Matthew Renze. 2024. The effect of sampling temperature on problem solving in large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*. 7346–7356.

[39] Kai Ruan, Xuan Wang, Jixiang Hong, Peng Wang, Yang Liu, and Hao Sun. 2024. LiveIdeaBench: Evaluating LLMs' Scientific Creativity and Idea Generation with Minimal Context. *arXiv preprint arXiv:2412.17596* (2024).

[40] Morris I Stein. 1953. Creativity and culture. *The journal of psychology* 36, 2 (1953), 311–322.

[41] E Paul Torrance. 1966. Torrance tests of creative thinking. *Educational and psychological measurement* (1966).

[42] Yuxia Wang, Revanth Gangi Reddy, Zain Muhammad Mujahid, Arnav Arora, Aleksandr Rubashevskii, Jiahui Geng, Osama Mohammed Afzal, Liangming Pan, Nadav Borenstein, Aditya Pillai, et al. 2024. Factcheck-Bench: Fine-Grained Evaluation Benchmark for Automatic Fact-checkers. In *EMNLP (Findings)*.

[43] Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, et al. 2025. The rise and potential of large language model based agents: A survey. *Science China Information Sciences* 68, 2 (2025), 121101.

[44] Xuesu Xiao, Bo Liu, Garrett Warnell, and Peter Stone. 2021. Toward agile maneuvers in highly constrained spaces: Learning from hallucination. *IEEE Robotics and Automation Letters* 6, 2 (2021), 1503–1510.

[45] 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).

[46] Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, et al. 2023. Siren's song in the AI ocean: a survey on hallucination in large language models. *arXiv preprint arXiv:2309.01219* (2023).

[47] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P Xing, et al. 2023. Judging LLM-as-a-judge with MT-bench and Chatbot Arena. In *Proceedings of the 37th International Conference on Neural Information Processing Systems.* 46595–46623.

## A  Prompt Details

### A.1  Multi-strategy Generation Control

**SCP (Strict Constraint Prompt)**
Assume you are an expert in {field}.
Please provide an answer to the following questions, not exceeding 70 tokens. Requirements:
1: Ensure feasibility by grounding in current scientific principles and technological trends;
2: Propose novel concepts or methods, avoiding unsupported speculation;
3: Ensure that the programme has value for the target area;
4: Maintain logical rigor without contradictions or vague statements.
Format: Plain text, no numbering or Markdown.
Question: {question}

**RCP (Role-Constrained Prompt)**
Assume you are an expert in {field}.
Please provide an answer to the following questions, not exceeding 70 tokens. Requirements:
1: Propose novel concepts or methods, avoiding unsupported speculation;
2: Ensure that the programme has value for the target area;
Format: Plain text, no numbering or Markdown.
Question: {question}

**COT (Chain-of-Thought Prompt)**
Assume you are an expert in {field}.
Please provide an answer to the following questions, not exceeding 70 tokens. Requirements:
1: Ensure feasibility by grounding in current scientific principles and technological trends;
2: Propose novel concepts or methods, avoiding unsupported speculation;
3: Ensure that the programme has value for the target area;
4: Maintain logical rigor without contradictions or vague statements.
5: Please think step by step before answering the question.
Format: Plain text, no numbering or Markdown.
Question: {question}

**RAG (Retrieval-Augmented Generation Prompt)**
Assume you are an expert in {field}.
Please provide an answer to the following questions, not exceeding 70 tokens. Requirements:
1: Ensure feasibility by grounding in current scientific principles and technological trends;
2: Propose novel concepts or methods, avoiding unsupported speculation;
3: Ensure that the programme has value for the target area;
4: Maintain logical rigor without contradictions or vague statements.
5: Answers should primarily rely on the provided Wikipedia principles; if information is insufficient, general scientific knowledge may be used, but speculation and contradictions must be avoided:{principle}.
Format: Plain text, no numbering or Markdown.
Question: {question}

### A.2  Evaluation Prompt

For our automated evaluation, we utilized **gpt-4o** and **deepseek-v3 models**, adopting a separate assessment approach. This method reduces the bias toward self-generated responses within the same conversation, ensuring impartial evaluation. Additionally, using two distinct models mitigates the risk of single model bias. Below is our evaluation prompt, which provides detailed explanations for each indicator.

**System Prompt**
You are a rigorous evaluator specialized in assessing the scientific innovation quality of LLM-generated answers. Be highly critical and avoid giving high scores to generic or vague answers.
1. Evaluation Criteria (1-5 Scale):
- Originality: 1 = Common knowledge or widely known concepts with no novelty. 2 = Slight variation of known methods, minimal innovation. 3 = Moderate originality, combining existing ideas in a non-trivial way. 4 = Highly novel approach, limited prior research but plausible. 5 = Breakthrough idea with no direct precedent." If the answer only combines existing concepts without a novel mechanism, the originality score must not exceed 3.
- Feasibility: 1 = Violates fundamental scientific principles. 2 = Theoretical possibility but lacks experimental validation. 3 = Feasible under specific assumptions but challenging to implement. 4 = Technically viable with existing technology. 5 = Readily implementable with minor adaptations of existing techniques.
- Value: 1 = No practical use or improvement over existing methods. 2 = Minor improvement, limited impact." 3 = Meaningful contribution with clear advantages. 4 = Significant improvement over current methods. 5 = Potential revolutionary impact on the field. If the answer lacks technical details or fails to demonstrate security improvements, the value score must not exceed 3.
2. Hallucination Detection:" - If any of the following conditions are met, mark 'Hallucination: Yes': - The answer does not align with the core requirements of the question. - The answer deviates from reality. - The answer contradicts established scientific principles. - The answer provides irrelevant or tangential information without addressing the problem. - The answer contains false information or made-up claims.
3. Scoring Rules: Generic or vague responses must receive lower scores: 'Originality <= 3' & 'Value <= 3'. Ensure a clear distinction between general answers and true innovations—avoid inflated scores.
4. Output format: 'Originality: [1-5] Feasibility: [1-5] Value: [1-5] Hallucination: Yes/No'.

**User Prompt**
[User Questions]:{question}
[Answers to be evaluated]:{answer}

# B  HIC-bench Dataset

The Cross Domain Innovation Dataset (CDID) uses open ended tasks for creativity evaluation as datasets with definitive answers are not suitable. Assessing creativity involves examining the novelty feasibility and value of ideas which are often subjective and context dependent. Open ended tasks allow models to generate diverse responses enabling a more comprehensive evaluation of their ability

to produce innovative and interdisciplinary solutions in scientific contexts.

In each experimental setting, we evaluate 1,000 generated responses per task. This design ensures that, despite the relatively compact number of tasks, the dataset retains strong statistical coverage and response diversity. At the same time, it supports a reliable assessment of model creativity and hallucination patterns, capturing nuanced variations within and across domains. The compact yet dense structure of CDID thus balances breadth with analytical depth, offering a practical and rigorous basis for evaluation.

## B.1  Cross-Domain Innovation Dataset (CDID)

This section introduces the CDID used in HIC-Bench. Table 4 provides a subset of examples from CDID, covering domains such as Quantum Physics and Social Sciences. The dataset is an open-ended question-answering collection constructed by integrating Wikipedia's knowledge organization principles with real-world complex problems. Core concepts from Wikipedia across interdisciplinary domains are combined with practical challenges to form innovative, open-ended questions.

## B.2  Cross-Domain Knowledge Base Dataset (CDKB)

This section introduces the Cross-Domain Knowledge Base Dataset (CDKB) developed for HIC-Bench. CDKB is a specialized knowledge base dataset designed to provide concise yet precise domain knowledge for diverse innovation tasks. For CDKB, we first used an LLM to condense Wikipedia entries, keeping only the core definitions of domain-specific terms and stripping away extraneous text. This ensures the RAG context focuses solely on the target concept without distracting or overly verbose material that could dilute the model's attention or introduce noise. By integrating CDKB into Retrieval-Augmented Generation (RAG) prompts, we enable LLMs to access structured, focused, and principle-aware knowledge, thereby gaining a clearer and more accurate understanding of the underlying concepts relevant to each task. This integration ensures that the generated responses are not only creative but also well-grounded in reliable domain-specific knowledge. Table 5 showcases examples from CDKB.

## B.3  Significance Testing

Given that our tasks are open-ended rather than fixed-response, they must elicit genuinely creative answers; otherwise they would be meaningless for this study. Following recent practice (e.g., HumanEval), we prioritize quality over quantity. We first harvest trending scientific keywords and then expand each into a question plus an associated knowledge base. Every task is designed to solicit highly creative responses, and the multiple answers generated satisfy the fluency requirement of the TTCT. To rule out randomness, we perform significance tests on the IH ratio across models and strategies. Specifically, we apply a two-sided paired permutation test: we take the vector of task-level IH-proportion differences d between two models, randomly flip the sign of each element, and recompute the mean difference 10 000 times under the null hypothesis that the expected difference is zero, thereby constructing the null distribution. The two-tailed p-value was derived from the extremity

**Table 4: Examples from the CDID Dataset: Open-Ended Questions Across Multiple Domains**

| Domain | Principle&Challenge | Question |
|---|---|---|
| Quantum Physics | quantum levitation, traffic | Design a novel transportation technology using quantum levitation and superconductivity to reduce urban traffic congestion. |
| Artificial Intelligence | Graph Neural Networks, social network | Propose a novel social network analysis tool based on Graph Neural Networks (GNN). |
| Biomedical Sciences | tissue engineering, artificial skin fabrication | Propose a novel artificial skin fabrication method based on tissue engineering. |
| Environmental Science | hydrological principles, urban flooding | Propose an urban stormwater management system based on hydrological principles to mitigate urban flooding. |
| Materials Science | biomimetic materials, bulletproof vest | Propose a novel bulletproof vest that enhances protection and comfort, inspired by the structural characteristics of biomimetic materials. |
| Energy Technology | ocean energy utilization technology, electricity | Design a tidal power generation device based on ocean energy utilization technology to supply electricity to coastal areas. |
| Neuroscience | principles of synaptic plasticity, enhance learning abilities | Propose an educational method to enhance learning abilities based on the principles of synaptic plasticity. |
| Information and Communication Technology(ICT) | optical fiber communication principles, cloud computing | Propose a novel ultra-high-speed data transmission system based on optical fiber communication principles to optimize cloud computing center networks. |
| Aerospace | Bernoulli's principle, aircraft lift | Design a novel wing to enhance aircraft lift based on Bernoulli's principle. |
| Social Sciences | social network analysis, prevent cybercrimes | Propose a strategy to prevent cybercrimes and ensure cybersecurity based on social network analysis. |

of the observed mean difference within this distribution. We report the standard error, and the p-value to quantify both magnitude and statistical significance. Most results pass the 1% significance level, and the remainder pass the 5% level. Table 6 summarizes the pairwise significance tests between models under the SCP setting.

## C Human Review

We integrate human review into HIC-Bench to enhance evaluation rigor. Specifically, we conduct a human review on the model's Intelligent Hallucinations (IH) and Defective Hallucinations (DH) outputs, assessing their factuality and feasibility from multiple perspectives. This process ensures a more reliable distinction between intelligent and defective hallucinations in scientific contexts, strengthening the overall assessment of LLMs' creative potential.

Due to the low efficiency and high labor cost of human review as well as the difficulty of evaluating open ended tasks automated evaluation metrics are becoming mainstream. However to ensure the reliability of automated evaluations we conducted a human review of the model classified IH and DH. This section presents examples from the human review process for model outputs classified as IH and DH. Human reviewers perform token level analysis to determine the extent to which the outputs align with reality focusing on the accuracy of concepts and terminology. For IH reviewers identify the hallucinated portions that are not fully realistic while assessing the reasonable and feasible aspects that contribute to their creative value. For DH reviewers evaluate whether the response addresses the given problem and check for obvious factual errors that render the output unrealistic or irrelevant. Additionally reviewers assess

**Table 5: Examples from the CDKB Dataset: Knowledge Base for Cross-Disciplinary Innovation Tasks**

| Domain | Knowledge |
|---|---|
| Quantum Physics | Quantum levitation, also known as quantum locking, occurs when a superconductor is cooled below its critical temperature and expels magnetic fields from its interior (Meissner effect), allowing it to lock in space above a magnet due to flux pinning. |
| Artificial Intelligence | A graph neural network (GNN) is a class of artificial neural networks for processing data that can be represented as graphs. |
| Biomedical Sciences | Gene editing is a type of genetic engineering in which DNA is inserted, deleted, modified, or replaced in the genome of a living organism. |
| Environmental Science | Hydrological principles study the distribution, movement, and properties of water on Earth, encompassing the water cycle processes such as precipitation, evaporation, infiltration, and runoff. |
| Materials Science | Biomimetic materials are synthetic materials designed to imitate the properties and functions of natural biological materials, often resulting in enhanced performance. |
| Energy Technology | Ocean energy utilization encompasses technologies that capture energy from oceanic sources, such as tidal movements, converting it into electricity. |
| Neuroscience | Synaptic plasticity refers to the ability of synapses, the connections between neurons, to strengthen or weaken over time in response to increases or decreases in their activity, playing a crucial role in learning and memory. |
| Information and Communication Technology(ICT) | Optical fiber communication uses light signals transmitted through fiber-optic cables to achieve high-speed data transmission over long distances, essential for optimizing networks in cloud computing centers. |
| Aerospace | Bernoulli's principle states that an increase in the speed of a fluid occurs simultaneously with a decrease in pressure, which is fundamental in understanding how airfoil shapes generate lift in aircraft. |
| Social Sciences | Social network analysis examines the relationships and interactions within a network, aiming to understand how these connections influence individual and group behaviors. |

**Table 6: Significance Tests Between Models under SCP. Values indicate standard errors. *** : $p < 0.01$, ** : $p < 0.05$, * : $p < 0.1$**

| | gpt-4o | gpt-4o-mini | qwen2.5-14b | qwen2.5-72b | deepseek-v3 | deepseek-r1 |
|---|---|---|---|---|---|---|
| gpt-4o | — | 0.016*** | 0.014*** | 0.015** | 0.016** | 0.030*** |
| gpt-4o-mini | | — | 0.009** | 0.011*** | 0.018*** | 0.033*** |
| qwen2.5-14b | | | — | 0.012** | 0.017*** | 0.033*** |
| qwen2.5-72b | | | | — | 0.017*** | 0.032*** |
| deepseek-v3 | | | | | — | 0.030*** |
| deepseek-r1 | | | | | | — |

the feasibility of the proposed ideas to ensure reliable categorization. Through token level analysis we evaluated the rationality of the classifications. Table 7 provides examples of the human review process.

In the Artificial Intelligence domain, we conducted a human evaluation of 600 model responses generated under the SCP strategy, where expert reviewers manually scored and classified the outputs into Intelligent Hallucinations (IH) and Defective Hallucinations (DH) based on criteria such as factual accuracy, creative value, and feasibility. This process involved a detailed analysis of each response, comparing the human assessments with the scores automatically generated by the models to evaluate the reliability and consistency of the automated evaluation system. The results are presented in Table 8.We calculated the precision and recall of

**Table 7: Illustrative Examples of Hallucination with Review Insights**

| Answer | Review |
|---|---|
| Quantum levitation grids embedded in roadways enable vehicles with superconducting plates to hover frictionlessly, guided by magnetic field gradients for precise lane control. | This response exhibits Intelligent Hallucination as it proposes a creative solution with some feasibility—quantum levitation and magnetic field gradients are grounded in physics. However, frictionless hovering and widespread superconducting vehicles are not fully realistic due to current technological and energy constraints, yet the idea holds innovative value for future transportation systems. |
| Stretchable graphene-MXene heterostructure supercapacitors woven into textiles provide high-energy-density storage with rapid charge cycles for self-powered wearables. | This response shows Intelligent Hallucination with an innovative idea. Graphene MXene heterostructures and textile integration are feasible. However high energy density storage and rapid charge cycles in stretchable textiles for self powered wearables are not fully realistic yet. The concept holds value for future wearable energy solutions. |
| A dual acting compound modulating GABA A positive allosteric modulation and NMDA receptor antagonism to balance inhibitory excitatory signaling during sleep wake transitions. | This response shows Intelligent Hallucination with a novel concept. Modulating GABA A and NMDA receptors to balance signaling is theoretically feasible. However achieving precise control during sleep wake transitions is not fully realistic with current pharmacology. The idea offers innovative value for future neurological treatments. |
| Adopt a hemispherical resonator gyroscope array with machine learning based noise filtering enabling sub arcsecond attitude determination for closed loop control stability. | This response shows Intelligent Hallucination with an innovative approach. Hemispherical resonator gyroscopes and machine learning noise filtering are feasible. However achieving sub arcsecond attitude determination for closed loop stability is not fully realistic with current technology. The concept holds value for future precision navigation systems. |
| Collaborations with local governments can ensure compliance with regulations while addressing urban planning challenges related to land use and environmental impact. | This response is classified as Defective Hallucination. It completely fails to address the problem of designing a novel transportation technology using quantum levitation and superconductivity. Instead it focuses on government collaboration and urban planning which are unrelated to the task and lack technological feasibility in this context. |
| Integrating emotion recognition systems could facilitate more nuanced interactions between humans and robots. Adapting behavior based on detected emotional states may improve cooperation safety and overall user experience during collaborative tasks. | This response is classified as Defective Hallucination. It does not address the problem of designing a novel robot control system based on reinforcement learning for autonomous decision making. Instead, it focuses on emotion recognition, which is unrelated to the task and lacks relevance to the specified technical approach. |

LLMs in classifying hallucinations across the 600 responses. For IH, LLMs achieved a precision of 85.4% and a recall of 87.4%, while for DH, the precision was 95.7% and the recall was 88.0%. Despite implementing strict prompts to enforce factual accuracy, LLMs show some leniency in assessing certain problems, leading to some bias in their scoring. Nevertheless, these evaluations generally reflect the variation trends of IH and DH across different prompts, providing valuable insights into hallucination dynamics.

## D  Dynamic Hallucination Prompt (DHP)

In our Fluency evaluation, we calculated the similarity of model responses within the same question. Here, we compute the similarity of responses across different questions as well as the similarity among responses identified as IH and those identified as DH in the SCP. The results are presented in Table 9. The similarity of IH and DH outputs generally exceeds the similarity of answers across questions. This observation motivated the design of the Dynamic Hallucination Prompt (DHP), which provides examples to guide model outputs toward more distinct and reliable responses.

The DHP pipeline is designed to enhance the generation and evaluation of model responses in the HIC-Bench framework by dynamically adapting prompts based on real-time feedback. DHP incorporates iterative prompt refinement by identifying positive examples (high-scoring responses based on Originality, Feasibility, and Value) and negative examples (responses with hallucinations) from each batch of generated answers. This process ensures that subsequent responses are guided toward higher creativity and lower defect rates, while evaluations are saved for further analysis, balancing innovation with reliability across diverse fields.

> **DHP (Dynamic Hallucination Prompt)**
>
> Assume you are an expert in field.
>
> Please provide an answer to the following questions, not exceeding 70 tokens. Requirements:
>
> 1: Ensure feasibility by grounding in current scientific principles and technological trends;
>
> 2: Propose novel concepts or methods, avoiding unsupported speculation;
>
> 3: Ensure that the programme has value for the target area;
>
> 4: Maintain logical rigor without contradictions or vague statements.
>
> Format: Plain text, no numbering or Markdown.
>
> {positive_prompt_examples}
>
> {negative_prompt_examples}
>
> Question: {question}

## E  Model Selection for Evaluation

LLMs demonstrate formidable proficiency in generating scientific text, yet their outputs frequently exhibit hallucinations, deviations from factual or logical coherence. Conventionally, such hallucinations are deemed errors to be eradicated. This paper, however, posits that hallucinations are not uniformly detrimental; certain instances, though diverging from reality, manifest remarkable innovation and foresight, harboring potential value in scientific contexts. To explore this duality, we introduce HIC-Bench, a systematic evaluation framework tailored to analyze LLMs' generative behavior in scientific innovation tasks. For the first time, it models hallucinations' dual attributes: DH, marked by factual inaccuracies, and IH, characterized by novel, scientifically plausible insights.

Six technically diverse LLMs are selected to evaluate HIC-Bench's applicability across varied architectures and optimization approaches: gpt-4o-2024-11-20 [21] serves as a benchmark for general-purpose performance; gpt-4o-mini [2] represents lightweight models optimized via knowledge distillation; qwen2.5-14b-instruct [5] and qwen2.5-72b-instruct [45] enable analysis of parameter scale effects; deepseek-v3 [30], refined through reinforcement learning, acts as a foundation model; and deepseek-r1 [16] excels in specialized CoT reasoning capabilities. This selection establishes a systematic evaluation framework, facilitating comparative analysis of model compression techniques, parameter scaling, reinforcement learning strategies, and dedicated reasoning architectures.

## F  Additional Experimental Results and Analysis

### F.1  Comparative Analysis of Model Indicators Across Scientific Fields

This subsection analyzes the cross domain performance of other key indicators in the HIC-Bench evaluation by examining additional metrics across various scientific fields. Figure 5 provides a detailed comparative analysis focusing on four indicators: Feasibility (Fe) in Subfigure (a) Value (Va) in Subfigure (b) Mean of Originality Feasibility and Value (Mean[Or, Fe, Va]) in Subfigure (c) and DH rates in Subfigure (d). Each Subfigure visualizes performance variations across domains. Feasibility remains relatively consistent across all models indicating a general tendency to produce feasible solutions.

In terms of Value and contribution the deepseek-r1 outperforms others showing stronger performance in Biomedical Science and Neuroscience but lower performance in Environmental Science and Energy Technology. The mean indicator aligns with this trend. For Defective Hallucination (DH) rates Qwen2.5-14b exhibits a notably higher occurrence compared to other models. This is particularly evident in Quantum Physics and Information and Communication Technology (ICT) where the model generates responses with significant factual inaccuracies or irrelevance to the task reflecting a higher propensity for defective outputs in these complex domains.

The statistical analysis reveals substantial differences in model responses across domains. When leveraging IH for innovative generation domain applicability must be carefully considered. For instance fields like Biomedical Science require a low hallucination rate to ensure reliability. In such cases the deepseek-r1 with a higher IH rate may not be suitable despite its strong creativity as the risk of inaccurate outputs could undermine practical utility.

---

**Algorithm 1** Dynamic Hallucination Prompt (DHP)

---

**Require:** Question file, evaluation file, output Excel
**Ensure:** Refined prompts, answer generations, evaluations
 1: Load all questions and principles from Excel
 2: Initialize `positive_example` and `negative_example` to empty
 3: **for all** field ∈ fields **do**
 4:   **for all** question ∈ field **do**
 5:     Build prompt using current `positive_example` and `negative_example`
 6:     Generate answers using the LLM
 7:     Save generated answers
 8:     **for all** answer ∈ answers **do**
 9:       Evaluate answer using LLM to get Originality, Feasibility, Value, Hallucination
10:       **if** Originality ≥ 4 **and** Feasibility ≥ 3 **and** Value ≥ 4 **then**
11:         **if** total score > current best **then**
12:           Update `positive_example`
13:         **end if**
14:       **end if**
15:       **if** Hallucination = Yes **then**
16:         Update `negative_example`
17:       **end if**
18:     **end for**
19:   **end for**
20: **end for**

---

### F.2  Prompt Strategy and performance

Reducing hallucinations in LLMs without compromising their efficacy remains a formidable challenge, given that prior investigations [25] have elucidated, through rigorous mathematical analysis, the intrinsic interplay wherein hallucinations not only manifest as defective deviations from factual accuracy (DH) but also constitute an integral facet of model creativity, specifically IH. Consequently, identifying strategies that mitigate DH while preserving or enhancing IH emerges as a pivotal research endeavor. Although numerous

**Table 8: Comparison of Human and AI Evaluations in the Artificial Intelligence Domain**

| Evaluator | Originality | Feasibility | Value | IH | DH |
|---|---|---|---|---|---|
| Human | 2.87 | 3.85 | 3.04 | 13.0% | 4.2% |
| LLMs | 2.91 | 3.83 | 3.11 | 14.8% | 3.8% |

**Table 9: Similarity Metrics for IH, DH, and All Answers in SCP. All Answer represents the similarity of responses across different questions.**

| LLM | All Answers | IH | DH |
|---|---|---|---|
| deepseek-v3 | 35.99% | 42.10% | 40.26% |
| deepseek-r1 | 38.97% | 41.38% | 37.35% |
| gpt-4o-mini | 37.91% | 48.50% | 43.90% |
| gpt-4o | 32.67% | 40.83% | 35.47% |
| qwen2.5-14b | 34.55% | 39.82% | 36.76% |
| qwen2.5-72b | 35.67% | 42.03% | 39.91% |



(a) Fe across scientific domains

(b) Va across scientific domains

(c) Mean(Or,Fe,Va) across scientific domains
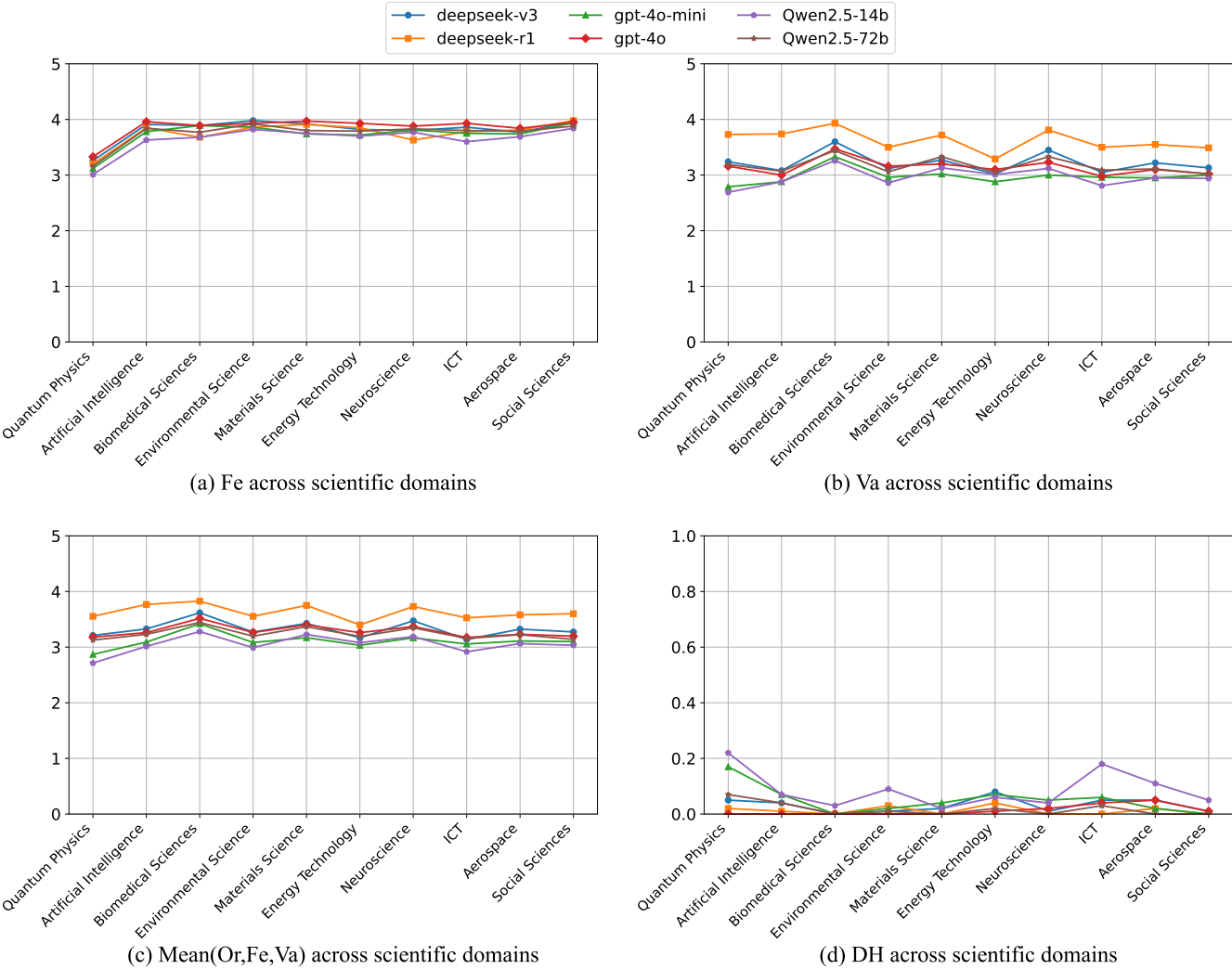
(d) DH across scientific domains

**Figure 5: Comparative analysis of model indicators across scientific fields under SCP**

methodologies have been developed to enhance the reliability of generated content, it has been posited [1] that excessively stringent constraints may suppress model performance in tasks necessitating creativity. Table 2 illustrates the outcomes across varying prompt strategies, through which this section examines the influence exerted by four distinct hallucination mitigation strategies, namely SCP, CoT, RAG, and RCP, on model performance, with particular emphasis on the proportions of IH and DH and the IFS.

**SCP**: Serving as the foundational strategy, SCP enforces rigorous constraints through prompts that mandate the generation of logically coherent outputs devoid of fabricated content, while simultaneously integrating insights from existing research to foster creative responses. Within this experimental framework, prompts are meticulously crafted to enhance the generation of creative outputs, thereby facilitating a comprehensive exploration of their characteristics. Results reported in Table 2, evaluated on the Vectara HHEM benchmark [7], reveal that deepseek-r1 exhibits a notably elevated overall hallucination rate, achieving an IH of 52.60% and an IFS of 50.04% in our benchmark, which surpasses other models such as gpt-4o-mini with an IH of 2.90% and an IFS of 38.58%. This outcome suggests that SCP effectively balances innovation with factuality to a considerable extent. In contrast, gpt-4o-mini, representing distilled architectures, demonstrates a markedly higher DH of 5.00% under identical conditions, indicating a pronounced propensity for generating factually inaccurate content when subjected to stringent constraints.

**CoT**: By incorporating the instruction "Let's think step by step," the CoT strategy substantially enhances response accuracy [24]. This zero-shot approach, which guides models through incremental reasoning processes, optimizes final outputs by generating intermediate inferential steps. Experimental findings indicate that CoT significantly elevates IFS scores while concurrently reducing DH across all models. For instance, deepseek-v3 exhibits a decline in DH from 3.20% under SCP to 1.90% under CoT, with its IFS rising from 41.40% to 42.64%. However, its impact on IH appears to diverge: while the deepseek series, such as deepseek-v3, shows an improved IH from 13.40% to 17.00%, suggesting that structured reasoning augments performance in creative tasks, the gpt-4o series, such as gpt-4o, experiences a decline in IH from 9.60% to 9.30%, potentially attributable to excessive reasoning steps constraining creative expression. Such variability underscores the differential responsiveness of models to the CoT strategy.

**RAG**: To enhance factual grounding, this strategy integrates the CDKB dataset, providing external, domain-specific knowledge and ensuring that model outputs are anchored in established scientific principles. Compared to the less constrained SCP setting, RAG imposes stronger factual guidance. Results show that RAG significantly reduces defect hallucinations. For example, gpt-4o-mini's DH drops from 5.00% to 3.50%, and deepseek-v3's decreases from 3.20% to 2.20%, demonstrating RAG's effectiveness in improving factual reliability through systematic knowledge referencing. However, this improvement comes at the cost of reduced intelligent hallucinations. Deepseek-r1's IH, for instance, falls from 52.60% to 49.00%, leading to a decreased in the IFS score to 49.72%. These findings suggest that while RAG enhances factual accuracy, its strong external constraints may suppress innovative potential, highlighting a

trade-off between reliability and creativity in open-ended scientific tasks.

**RCP**: Acknowledging that stringent constraints might inhibit IH, we propose the RCP strategy, which relaxes the feasibility and output restrictions of SCP, prioritizing solely the innovation and value of generated content. Results reveal that RCP significantly elevates IH across all models, while reducing Defective Hallucinations (DH) outputs in certain cases. For example, gpt-4o's IH rises from 9.60% under SCP to 18.60% under RCP, yet its DH decreases from 1.30% to 0.50%. Similarly, qwen2.5-72b-instruct [45] exhibits an IH increase from 8.00% to 10.40%, accompanied by a DH reduction from 1.70% to 1.20%. This observation indicates that DH and IH are not inherently positively correlated; rather, an appropriately moderated relaxation of constraints appears to unleash the creative potential of models, while also lowering factual deviations in some models.

## F.3 Comparative Overview of LLM Hallucination Evaluation Studies

This subsection presents Table 10 which provides a comparative overview of recent LLM hallucination evaluation studies. The table contrasts HIC-Bench with other benchmarks by examining dataset sources task types and evaluation metrics. In HIC-Bench we categorize hallucinations into Intelligent Hallucinations and Defective Hallucinations to break the traditional perspective of treating hallucinations solely as errors to be mitigated thereby offering a novel viewpoint for hallucination research. Furthermore HIC-Bench focuses on open ended cross domain innovation using real world challenges and employs creativity focused metrics such as Originality Feasibility and Value. This comparison highlights the distinct approach of HIC-Bench in evaluating LLMs for scientific creativity.

## G Limitations and Future Work

Due to the review process's scope constraints, this section outlines key limitations and future directions. HIC-Bench has advanced the understanding of IH and DH interplay in LLMs and the evaluation of mitigation strategies, yet limitations persist that warrant further exploration.

**Structured Task Focus**: The framework primarily focuses on structured tasks like question answering, neglecting areas such as literary writing. Future work will extend to broader creativity tasks for a comprehensive assessment.

**Human Evaluation Constraints**: The current human evaluation only involves reviewing model outputs, without conducting a full accuracy analysis comparing human expert assessments to LLM evaluations or aligning with human preferences. Future efforts will incorporate expert evaluations and preference alignment studies to improve the robustness of IH and DH assessments.

Additionally, balancing IH and DH remains challenging due to potential biases from inductive prompts and the need for refined IH metrics. Future efforts will develop domain-specific prompts and integrate human-involved multi-dimensional evaluations to enhance IH assessment accuracy and manage IH-DH dynamics. Lastly, the framework will be applied to multimodal and cross-lingual benchmarks to validate its generalizability, while examining the ethical implications of promoting IH.

**Table 10: Comparison of LLM Hallucination Evaluation Studies**

| Benchmark | Dataset Source | Task Type | Evaluation Metrics |
|---|---|---|---|
| HHEM [7] | Wiki, News | Question Answering | BLEU, ROUGE, Factuality Score |
| TruthfulQA [29] | Manual | Question Answering | Accuracy by Human or GPT Judge |
| HalluDial [33] | Dialogues, Synthetic | Dialogue Generation | Hallucination Rate, Semantic Consistency |
| HalDetect [15] | Wiki, Synthetic | Text Generation | Hallucination Rate, F1, Precision |
| FactCheck [42] | News, Wiki | Fact Checking | FactScore, NE Error |
| SelfCheck [34] | Wiki | Fact Checking | SelfCheckGPT Score, Hallucination Rate |
| FactScore [35] | Wiki | Text Generation | FactScore, NE Error |
| HalluQA [11] | Manual, Wiki | Question Answering | Non Hallucination Rate |
| UHGEval [28] | News | Open Ended Generation | Accuracy, kwPrec, BERTScore |
| HIC-Bench (Ours) | Wiki, Real World Challenges | Open Ended Cross Domain Innovation | Originality, Feasibility, Value, IH, DH, IFS |