

Harder is Better: Hard Hallucination-Induced Contrastive Decoding for Hallucination Mitigation

Anonymous ACL submission

Abstract

Large Language Models (LLMs) have achieved significant advancements in various natural language processing tasks. However, they are susceptible to generating hallucinations—fabricated or inaccurate statements presented as factual information—which can undermine their reliability in high-stakes applications. To address this issue, we propose a new inference-stage HiCD method to improve hallucination mitigation. It aims to inject hard-to-detect hallucinations to enhance the robustness of contrastive decoding during inference. An adversarial-aware strategy is introduced for finetuning hallucination models to effectively learn more precise and diverse hallucination patterns from available hallucination data. This enhances the contrastive decoding process, enabling more effective identification and filtering of erroneous content. We evaluate HiCD on four various hallucination benchmarks. Experimental results show significant improvements on all metrics consistently, proving the effectiveness and superiority of HiCD for hallucination mitigation.

1 Introduction

Large Language Models (LLMs) have demonstrated substantial progress in a range of natural language processing (NLP) tasks, including question answering, knowledge-grounded dialogue, and reasoning-intensive problem solving (Touvron et al., 2023; Achiam et al., 2023). However, despite these achievements, LLMs frequently produce *hallucinations*—outputs that contain inaccuracies or fabrications presented as factual information (Bang et al., 2023; Ji et al., 2023). Such hallucinations pose significant risks, particularly in high-stakes domains such as legal consultation, medical advice, and specialized technical support, where factual reliability is essential.

Various strategies have been pursued to mitigate hallucinations. Some works emphasized data-

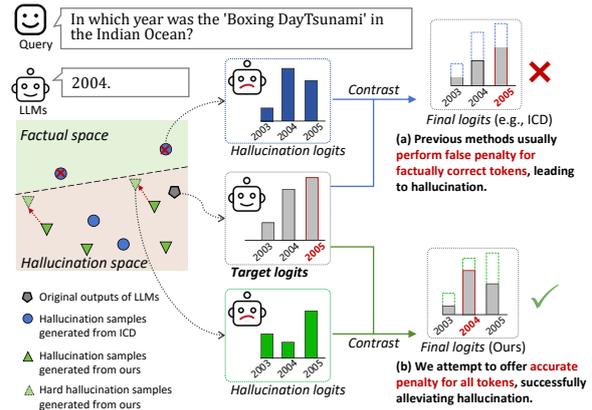


Figure 1: An illustration showing how over-penalization of factually correct tokens leads to hallucination

centric methods, such as curating training sets or integrating external knowledge to guide models toward factual correctness (Sun et al., 2023; Shuster et al., 2021). These methods typically require substantial computational overhead and may not generalize well beyond the data distributions observed during training (Ren et al., 2023; Borgeaud et al., 2022). Recently, increasing attention has focused on mitigating hallucinations at the *inference* stage (Li et al., 2023b; Chuang et al., 2023; Kai et al., 2024; Zhang et al., 2023). They usually examine differences across multiple candidate outputs via contrastive decoding strategies for hallucination mitigation during inference. Inference-stage methods can be more flexible and less resource-intensive than strategies that rely solely on enhancing training data or model parameters.

However, the above inference-stage methods may suffer from the precision of hallucination tokens, leading to limited contrastive performance during inference. Specifically, hallucinations in LLMs are highly diverse (Huang et al., 2023). Finetuning with the scarcity of hallucination data often leads to a suboptimal hallucination model, which struggles to generalize well and fails to provide

067 subtle and precise hallucination patterns (Wang
068 et al., 2023). As a result, factually accurate tokens
069 are prone to be over-penalized during contrastive
070 decoding, leading to suboptimal performance for
071 hallucination mitigation. As shown in Figure 1, the
072 imprecise hallucination logits outputted by previ-
073 ous works may perform false penalty for the factual-
074 correct token (i.e., 2024). Therefore, a more effec-
075 tive fine-tuning strategy for hallucination model
076 needs to be explored for capturing more precise
077 and diverse hallucination samples, accordingly to
078 improve the effectiveness of contrastive decoding.

079 In this paper, we propose a new inference-stage
080 method, Hard Hallucination-Induced Contrastive
081 Decoding (HiCD), to improve hallucination miti-
082 gation. Our HiCD aims to inject *hard*-to-detect
083 hallucinations to enhance the robustness of con-
084 trastive decoding during inference. We design a
085 new adversarial-aware finetuning strategy for hallu-
086 cination models to explore more hard hallucination
087 samples. These samples usually are similar to fac-
088 tually correct tokens but deviate in subtle ways.
089 As shown in Figure 1, they lie near the decision
090 boundary between factual correctness and halluci-
091 nation in the model’s prediction space. To achieve
092 this, we utilize adversarial perturbations to encour-
093 age factually correct samples beyond the limited
094 hallucination dataset to more accurately approach
095 hallucination boundaries (Goodfellow et al., 2014).
096 This process reduces the prediction probabilities of
097 correct tokens in a controlled manner, preventing
098 the model from overfitting to specific hallucination
099 patterns. Based on hallucination LLMs finetuned
100 by our strategy, during the contrastive decoding
101 phase, the model avoids erroneously penalizing
102 factually correct tokens, resulting in outputs that
103 are more reliable and factually consistent. Import-
104 antly, HiCD achieves these improvements without
105 requiring extensive data curation or large-scale re-
106 training, offering a scalable and practical solution
107 for mitigating hallucination issues in LLMs.

108 We conduct experiments on four truthfulness as-
109 sessments and knowledge-seeking datasets for hal-
110 lucination alleviation evaluation. The experimen-
111 tal results demonstrate HiCD’s effectiveness, with
112 consistent improvements on multiple benchmarks
113 (e.g., +4.08% MC2 on TruthfulQA and +9.03%
114 on FACTOR Expert) across diverse tasks. Addi-
115 tionally, ablation and parameter analyses highlight
116 the crucial role of adversarial training and optimal
117 hyperparameters, indicating HiCD’s broad applica-
118 bility for enhancing factual fidelity and mitigating

hallucinations in large language models. 119

Our contributions are threefold: 1) we propose a 120
new inference-stage HiCD method to improve hal- 121
lucination mitigation. It injects hard-to-detect hal- 122
lucinations to enhance the robustness of contrastive 123
decoding during inference. 2) A new adversarial- 124
aware finetuning strategy for hallucination models 125
is designed to precisely capture more diverse and 126
hallucination patterns from available hallucination 127
data. 3) Experiments on four hallucination datasets 128
demonstrate the effectiveness and superiority of 129
HiCD for hallucination mitigation. 130

2 Related Work 131

2.1 Hallucination in Large Language Models 132

Large Language Models (LLMs) are prone to gen- 133
erating *hallucinations*—fabricated or inaccurate 134
statements presented as factual (Achiam et al., 135
2023; Ji et al., 2023). These hallucinations can 136
be broadly categorized into *factual* and *faithful- 137*
ness hallucinations. *Factual hallucinations* emerge 138
when the model’s output contradicts established 139
real-world knowledge (Bang et al., 2023; Hu 140
et al., 2023), while *faithfulness hallucinations* oc- 141
cur when the model’s response deviates from given 142
instructions or the provided source context (Dale 143
et al., 2023; Shi et al., 2023). Eliminating both 144
types of hallucinations is critical for real-world 145
applications, especially in high-stakes domains de- 146
manding reliable and truthful information. 147

Initial efforts to mitigate hallucinations often em- 148
phasized data- and model-centric strategies. Data- 149
centric approaches involve refining training cor- 150
pora—either curating higher-quality data or incor- 151
porating external knowledge sources—to encour- 152
age factual correctness (Sun et al., 2023; Shuster 153
et al., 2021; Lin et al., 2022). Model-centric meth- 154
ods aim to modify training objectives, sometimes 155
leveraging techniques like reinforcement learn- 156
ing from human feedback to align model outputs 157
with human judgment (Wang and Sennrich, 2020; 158
Ouyang et al., 2022). While these methods can 159
reduce certain types of hallucinations, they often 160
require extensive data preparation, large-scale re- 161
training, may not generalize well to complex, sub- 162
tle errors that lie near decision boundaries. 163

To address these issues more efficiently, re- 164
searchers have turned to inference-stage interven- 165
tions. Post-hoc decoding strategies can be applied 166
at generation time without modifying the under- 167
lying parameters. By using contrastive signals or 168

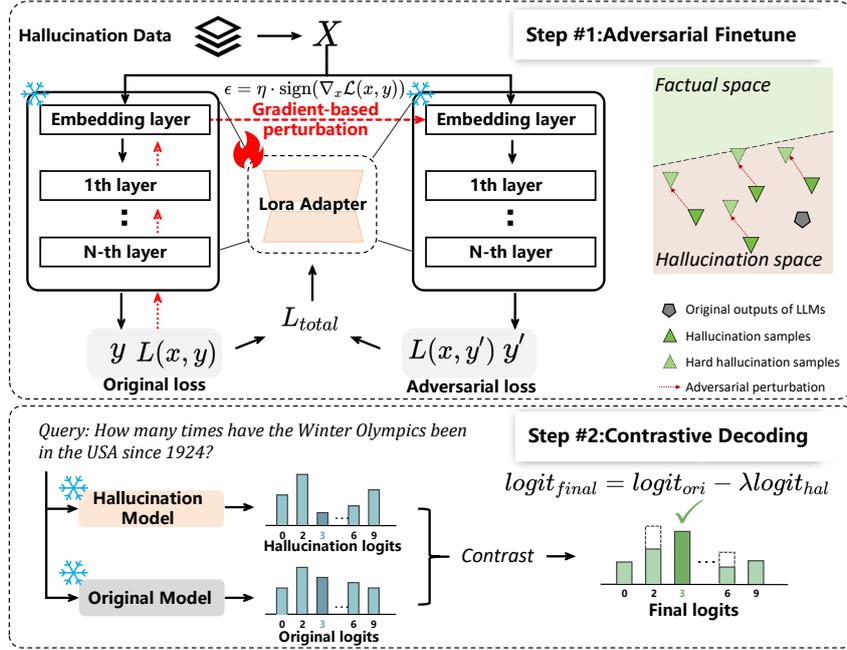


Figure 2: Overview of our HiCD framework. In the adversarial finetuning phase, we induce hard-to-detect hallucinations through gradient-based perturbations, resulting in a weaker “hallucination” model. During inference, contrastive decoding combines outputs from the original and hallucination models, filtering out fabricated content and enhancing factual fidelity

169 other dynamic generation criteria, these approaches
170 aim to identify and filter out hallucinatory content
171 as it emerges (Chang et al., 2023). However, exist-
172 ing inference-stage methods often rely on hallucina-
173 tions that are easy to induce or naturally occurring.
174 Such limited sets of negative examples fail to repre-
175 sent the full spectrum of challenging hallucinations
176 that can occur in practice. As a result, these meth-
177 ods struggle with difficult, subtle hallucinations in
178 more complex, real-world scenarios.

179 2.2 Contrastive Decoding

180 Contrastive Decoding (CD) (Li et al., 2023b) in-
181 troduced a novel perspective for improving genera-
182 tion quality by contrasting outputs from a stronger
183 model against those from a weaker model. Build-
184 ing on this idea, Chuang et al. (2023) proposed
185 contrasting outputs from different Transformer lay-
186 ers to enhance factual accuracy, while Kai et al.
187 (2024) incorporated self-attention mechanisms to
188 identify and mitigate uncertain predictions. To fur-
189 ther refine factual outputs, Zhang et al. (2023) sug-
190 gested inducing hallucinations and then contrasting
191 them to filter out inaccuracies. Similarly, Xu et al.
192 (2024) decoupled identification and classification
193 tasks to reduce hallucinations in medical informa-
194 tion extraction, and Gema et al. (2024) introduced
195 a method that contrasts outputs from a base model

196 and a masked model with retrieval heads to mitigate
197 hallucinations.

198 However, existing contrastive decoding methods,
199 such as Induce-then-Contrast Decoding (ICD), are
200 constrained by the limited availability of hallucina-
201 tion data, which is insufficient to fully influence the
202 extensive knowledge acquired by large language
203 models during pretraining. This limitation hampers
204 their ability to effectively identify and mitigate sub-
205 tle or complex hallucinations that closely resemble
206 truthful content. Consequently, these methods may
207 inadvertently penalize factually correct tokens, re-
208 ducing their accuracy and reliability in real-world
209 applications where distinguishing between factual
210 information and fabrications is critical. Addressing
211 these challenges requires more sophisticated strate-
212 gies that can generate richer and more nuanced
213 negative examples, thereby enabling a more pre-
214 cise approximation of the true decision boundaries
215 between accurate and erroneous outputs.

216 3 Hard Hallucination-Induced 217 Contrastive Decoding (HiCD)

218 Consider a standard text generation setting where
219 an LLM receives an input sequence $x =$
220 (x_1, x_2, \dots, x_L) and generates an output sequence
221 $y = (y_1, y_2, \dots, y_T)$. Without additional

constraints, the LLM may produce *hallucinations*—tokens or phrases unsupported by factual evidence. These hallucinations degrade the trustworthiness and reliability of the generated text.

As shown in Figure 2, our proposed framework, Hard Hallucination-Induced Contrastive Decoding (HiCD), aims to reduce hallucinations by leveraging contrastive decoding between a strong model and a weaker, adversarially trained model.

3.1 Inducing Hard Hallucinations

As hallucinations in LLMs are highly diverse and subtle. Previous works (Zhang et al., 2023) inducing potential hallucination in a suboptimal way usually falsely penalized precision factual tokens, leading limited alleviation performance. To capture hard hallucination samples for better contrasting, we design a new adversarial-aware finetuning strategy to capture hard hallucination samples for better contrastive decoding during generation of target LLMs. Specifically, we first employ few-shot prompting techniques to generate misleading or incorrect responses from a factual dataset. We then go further by integrating adversarial training to push the weaker model—referred to as the “hallucination LLM”—towards producing more intricate, boundary-like hallucinations that are harder to distinguish from truthful outputs.

Formally, let $D = \{(s_i, u_i, o_i)\}_{i=1}^m$ be the finetuning dataset, where s_i is the system prompt, u_i is the user input, and o_i is the target output. The initial fine-tuning objective is:

$$\arg \min_{\Delta\theta} \sum_{i=1}^m -\log p(o_i | s_i, u_i; \theta + \Delta\theta), \quad (1)$$

where θ denotes the original model parameters. After this step, we incorporate adversarial perturbations to shape $\Delta\theta$ so that the weaker model becomes more inclined to produce complex hallucinations.

By introducing adversarial training during finetuning, the weaker model’s errors become more refined and deceptive, rather than simple and easily detectable. We employ the Fast Gradient Sign Method (FGSM) (Goodfellow et al., 2015) to perturb the input embeddings \mathbf{x} :

$$\mathbf{x}' = \mathbf{x} + \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, y)), \quad (2)$$

where ϵ controls the perturbation magnitude, and \mathcal{L} is the loss function. This pushes the model’s

decision boundaries, increasing uncertainty and promoting the production of subtle hallucinations.

Training alternates between clean and adversarially perturbed examples. The combined objective is:

$$\mathcal{L}_{\text{total}} = \frac{1}{2} (\mathcal{L}(\mathbf{x}, y) + \mathcal{L}(\mathbf{x}', y)), \quad (3)$$

resulting in a hallucination LLM that naturally generates a richer, set of more challenging negative examples for the subsequent contrastive decoding step.

3.2 Contrastive Decoding

Having obtained the stronger model θ and the adversarially fine-tuned weaker model θ' , we apply contrastive decoding (Li et al., 2023b) to their outputs. At each timestep t , both models compute the conditional probability of the next token x_t . We define the contrastive score as:

$$\mathcal{F}_t = \log p(x_t | x_{<t}; \theta) - \lambda \log p(x_t | x_{<t}; \theta'), \quad (4)$$

where λ controls the balance between the two models. This score amplifies tokens favored by the stronger model while suppressing those preferred by the weaker, hallucination-prone Evil LLM.

To further refine token selection, we employ the adaptive relative top filtering mechanism (Li et al., 2023b). Specifically, at each timestep t , we define a valid token set $\mathcal{V}_{\text{valid}}$ based on the probabilities predicted by the strong model θ :

$$\mathcal{V}_{\text{valid}} = \left\{ x_t \in \mathcal{V} \mid \begin{array}{l} \log p(x_t | x_{<t}; \theta) \geq \\ \max_w \log p(w | x_{<t}; \theta) + \log \gamma \end{array} \right\} \quad (5)$$

where $\gamma \in (0, 1]$ is a hyperparameter that determines the filtering threshold. This ensures that only tokens whose log probabilities are within $\log \gamma$ of the highest log probability are retained.

After determining $\mathcal{V}_{\text{valid}}$, we apply a softmax over the contrastive scores $\mathcal{F}_t(x_t)$ for $x_t \in \mathcal{V}_{\text{valid}}$:

$$p(x_t | x_{<t}) = \frac{\exp(\mathcal{F}_t(x_t))}{\sum_{x \in \mathcal{V}_{\text{valid}}} \exp(\mathcal{F}_t(x))}. \quad (6)$$

By restricting the candidate tokens to this valid set and then normalizing with respect to the contrastive scores, the final output distribution is more factual and less susceptible to subtle hallucinations introduced by the factually weaker LLM.

Setting	Value
Model	Llama2-7B-Base
Epochs	5
Device	NVIDIA Tesla A100 80GB
Total Batchsize	256
Learning Rate	5×10^{-4}
LoRA Target	$q_{\text{proj}}, k_{\text{proj}}, v_{\text{proj}}$

Table 1: Finetuning settings for building the factually weaker model.

4 Experiments

4.1 Experimental Setup

Datasets Following previous work (Chen et al., 2024), we evaluate our method on truthfulness-related datasets (i.e., TruthfulQA, and FACTOR) and knowledge-seeking datasets (i.e., TriviaQA, and NQ). **TruthfulQA** (Lin et al., 2022) is a benchmark designed to assess the truthfulness of language models, comprising 817 multiple-choice questions across 38 categories. **FACTOR** (Muhlgay et al., 2023) evaluates the factual accuracy of large language models in text completion tasks, consisting of two subsets: Wiki-FACTOR with 2,994 examples from Wikipedia and News-FACTOR with 1,036 examples from news articles. TriviaQA (Joshi et al., 2017) contains over 650K question-answer pairs sourced from trivia websites, accompanied by evidence documents from Wikipedia and web sources. **Natural Questions (NQ)** (Kwiatkowski et al., 2019), developed by Google, includes around 300K human-generated questions with annotated short and long answers derived from Wikipedia.

Evaluation Metrics We employ multiple-choice accuracy metrics to assess model performance on the truthfulness-related dataset, i.e., TruthfulQA. Specifically, **MC1** evaluates whether the model assigns the highest probability to the correct answer, while **MC2** measures the total normalized probability mass the model assigns to correct answers. **MC3** combines accuracy and consistency across multiple questions to gauge the model’s overall reliability. For FACTOR, we experiment on its three subsets—News, Wiki, and Expert—and utilize accuracy as the sole evaluation metric to assess the text completion performance of large language models. Following Joshi et al. (2017), we adopt **Exact Match (EM)** and **F1 score** as evaluation metrics to measure the correctness of the model’s responses on knowledge-seeking datasets, i.e., TriviaQA and NQ.

Comparison Methods. We evaluate the effectiveness of our proposed method by comparing it against the following baselines: (1) **Greedy Decoding**: A default approach where the highest probability token is selected at each step without additional decoding techniques. (2) **Induced Task Inference (ITI)** (Li et al., 2024): This method enhances generalization by applying task-specific adjustments during inference, refining predictions based on task-relevant cues. (3) **Contrastive Decoding (CD)** (Li et al., 2023b): Aims to reduce hallucinations by contrasting outputs from a strong model and a weaker model, emphasizing reliable predictions while penalizing non-factual ones. (4) **Direct Output Layer Adaptation (DoLa)** (Chuang et al., 2023): Focuses on adjusting the model’s output layer to improve factual accuracy, particularly for knowledge-intensive tasks. (5) **Induce-then-Contrast Decoding (ICD)** (Zhang et al., 2023): Integrates hallucination induction with contrastive decoding, leveraging a weakened model to penalize incorrect predictions and reinforce factual outputs. (6) **Activation Decoding (AD)** (Shi et al., 2024): Amplifies the influence of contextual information over a language model’s prior knowledge by employing a contrastive output distribution, improving faithfulness in tasks requiring external knowledge integration.

Implementation Details All experiments are conducted on a single NVIDIA Tesla A100 80GB GPU using the Llama2 series models. The scaling factor λ in Equation 4 was set to 1.8 for optimal results on the TruthfulQA dataset. For the FACTOR dataset, the best results were achieved with λ values of 0.35. We leverage Llama2-7B-Chat as the original model to conduct the experiments and fine-tune Llama2-7B-Base to create a factually weaker model, following a similar setup to (Zhang et al., 2023). Specifically, we use the HaluEval dataset (Li et al., 2023a) to fine-tune the weaker model. LoRA (Hu et al., 2022) is used for parameter-efficient fine-tuning and hallucination injection. The LLaMA-Factory framework (Zheng et al., 2024) is also employed for fine-tuning. Details of the fine-tuning process and hyperparameters are provided in Table 1.

4.2 Main Results

Overall results on four datasets for hallucination mitigation are shown in Table 2. The proposed HiCD achieves the best performance on all datasets

Method	TruthfulQA			FACTOR			TriviaQA		NQ	
	MC1	MC2	MC3	News	Wiki	Expert	EM	F1	EM	F1
Greedy	37.62	54.60	28.12	65.05	56.96	66.10	46.50	46.50	23.49	21.45
ITI (Li et al., 2024)	37.01	54.66	27.82	53.28	43.82	51.69	–	–	–	–
CD (Li et al., 2023b)	28.15	54.87	29.75	64.57	58.47	67.12	47.30	38.58	26.03	19.38
DoLa (Chuang et al., 2023)	32.97	60.84	29.50	64.32	57.63	67.30	47.08	45.94	24.01	22.15
AD (Shi et al., 2024)	33.90	51.62	25.78	61.87	53.84	62.28	48.55	48.24	24.34	22.35
ICD (Zhang et al., 2023)	46.32	69.08	41.25	70.75	58.40	66.94	50.46	50.33	25.59	23.94
HiCD (Ours)	47.00	73.16	46.26	71.23	59.17	74.15	50.91	50.67	26.20	24.40
Improve (%)	+9.38	+18.56	+18.14	+6.18	+2.21	+8.05	+4.41	+4.17	+2.71	+2.95

Table 2: Overall results of different inference-based methods on four benchmarks. We reimplement all methods according to their open-source codes under the same environment except for ITI. The Llama2-13B-Chat vs. 7B-Chat setting is used in experiments of CD. For ICD and our HiCD, we follow Zhang et al. (2023) and finetune Llama2-7B-Base as a weaker model for contrasting with Llama2-7B-Chat. The best performances are **bolded**.

in terms of all evaluation metrics. This demonstrate the superiority of our model on ensuring the truthfulness of responses but also effectively retrieving and reasoning over factual information in open-domain settings. Specifically, for truthfulness-related datasets, compared the the baseline Greedy, HiCD achieves improvements of **+9.4%**, **+18.6%**, and **18.1%** on MC1, MC2, and MC3 scores on TruthfulQA. For knowledge-seeking tasks, HiCD outperforms the baseline by **+4.4%** EM and **4.2%** F1 scores. Besides, compared to other decoding strategies, HiCD contrasts with hard hallucination-induced models, leading to better mitigation performance on all datasets.

4.3 Ablation Study

To evaluate the effectiveness of our adversarial training in inducing precise hallucinations and enhancing contrastive decoding, we conduct an ablation study by comparing our HiCD with the following ablation models: 1) **w/ Adv Perturb.** refers to replacing adversarial perturbations with random perturbations during the fine-tuning of the hallucination-induced models. 2) **w/o Perturb.** indicates removing the adversarial perturbations entirely during fine-tuning.

The ablation results on TruthfulQA and FACTOR are presented in Table 3. The full HiCD model achieves the best performance across all metrics on both datasets, showing the effectiveness of each component for building hallucination LLMs. Incorporating adversarial perturbations enhances the generation of precise and diverse hallucinations. In this way, HiCD enables more effective filtering of factual inaccuracies, leading to more reliable and factually consistent outputs.

Method	TruthfulQA			FACTOR		
	MC1	MC2	MC3	News	Wiki	Expert
HiCD	47.00	73.16	46.26	71.23	59.17	74.15
w/o Adv Perturb.	38.31	65.56	37.23	55.88	38.92	55.50
w/o Perturb.	46.32	69.08	41.25	70.75	58.40	66.94

Table 3: Ablation results on TruthfulQA and FACTOR.

Method	TruthfulQA			
	%truth	%info	%truth*info	%reject
CD	70.21	42.25	19.23	29.98
ICD	62.85	77.65	41.16	23.50
HiCD	63.71	78.03	42.24	23.13

Table 4: Evaluation results on generative tasks using "GPT-judge" for TruthfulQA.

4.4 Generation Task Evaluation

Following Lin et al. (2022), we also evaluate our method on the TruthfulQA dataset using "GPT-judge" to assess both factual accuracy and informativeness. This evaluation yields four metrics: *truth*, *info*, a combined *truth&info*, and the *reject* rate. Table 4 presents the evaluation results on generative tasks for CD, ICD, and our proposed HiCD approach. Compared to ICD, HiCD achieves a +0.38% increase in *info*, a +1.08% increase in *truth&info*, and a -0.37% decrease in *reject*, indicating that HiCD produces more informative, factually consistent responses.

4.5 Efficiency Analysis

We compare the inference efficiency of different inference-stage methods, i.e., a baseline greedy decoding, CD, ICD, and our proposed HiCD. The baseline employs on a Llama2-7B-Chat model. The measured times reflect approximate overhead trends rather than a strict one-to-one comparison, as the CD experiment uses a Llama2-13B-Chat vs. 7B-Chat configuration, while both ICD and HiCD

Method	Decoding Latency (s)
Baseline	138.4 ($\times 1.00$)
CD	357.6 ($\times 2.58$)
ICD	402.4 ($\times 2.91$)
HiCD	384.7 ($\times 2.78$)

Table 5: Inference time comparison across different decoding strategies.

rely on a Llama2-7B-Chat model with a finetuned Llama2-7B-Base weaker model.

As shown in Table 5, the baseline decoding takes approximately 138.4s. Under the CD setting, increasing complexity leads to about a 2.58 \times slowdown. For ICD and HiCD, which directly compare a 7B-Chat strong model to a finetuned 7B-Base weaker model, the overhead is roughly 2.91 \times and 2.78 \times respectively. Although these configurations differ, the general pattern holds: more sophisticated contrastive strategies incur additional computation. Notably, HiCD offers improved factual fidelity over ICD while slightly reducing the slowdown from the baseline, indicating a more balanced trade-off between accuracy and efficiency.

4.6 Parameter Analysis

We experiment to analyze the impact of two critical hyperparameters in HiCD: the perturbation magnitude ϵ and the scaling factor λ . Results of parameter analysis on TruthfulQA are shown in Figure 3.

Effect of Scaling Factor λ The scaling factor λ adjusts the influence of the weaker model (i.e., hallucination model) in the contrastive decoding process. The optimal value is set to 1.5. By increasing λ , we amplify the penalty imposed by the weaker model on the strong model’s outputs, thereby enhancing the suppression of hallucinations. The fact indicates that increasing λ effectively suppresses hallucinations by strengthening the contrastive signal between the strong and weaker models. beyond a certain threshold, further increasing λ may lead to over-penalization, resulting in a slight decline in performance due to excessive suppression of potentially correct tokens.

Effect of Perturbation Magnitude ϵ The perturbation magnitude ϵ controls the strength of adversarial noise during the fine-tuning of the weaker model. By adjusting ϵ , we influence the extent to which the model’s decision boundaries are shifted, thereby affecting the precision and difficulty of induced hallucinations. Our results indicate that

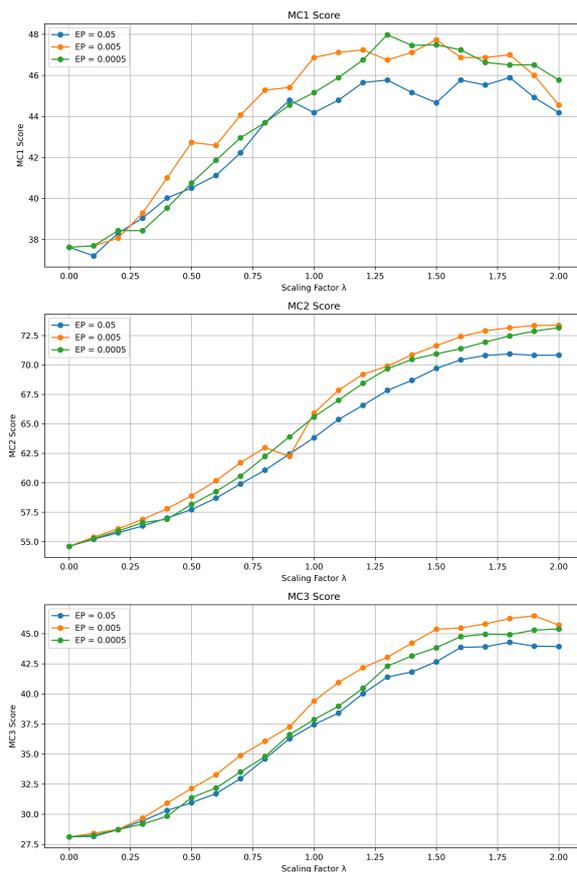


Figure 3: MC1, MC2, and MC3 scores on the TruthfulQA dataset for different perturbation magnitudes ϵ and scaling factors λ .

$\epsilon = 0.005$ yields the highest MC scores, effectively balancing the generation of challenging hallucinations and maintaining the efficacy of contrastive signals. Smaller perturbations ($\epsilon = 0.0005$) do not sufficiently alter the model’s behavior to produce hard hallucinations, while larger perturbations ($\epsilon = 0.05$) may overly degrade the weaker model’s performance, reducing the effectiveness of contrastive decoding in distinguishing factual from hallucinated content.

4.7 Case Study

We provide a case study from the Natural Questions dataset to illustrate the effectiveness of our method. Consider the query: “When was the rock and roll hall of fame built in Cleveland?” The correct answer is 1995, while a hallucinated answer is 1986. In this scenario, both the original model and the ICD approach produce the hallucinated answer, whereas our method yields the factually correct output. As shown in Figure 4, we analyze the token-level probabilities for the key differing token positions (the second “9” in 1995 and “8” in

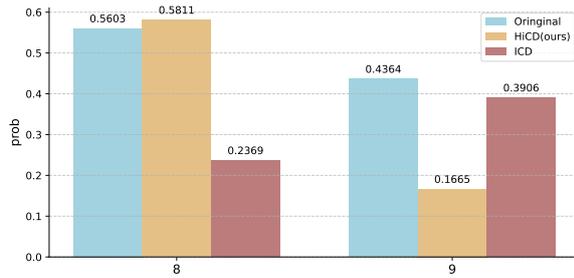


Figure 4: Token-level probability analysis for the query “When was the rock and roll hall of fame built in Cleveland?” at the critical token position where hallucination occurs

1986): the original model assigns overly high confidence to an incorrect token, while ICD’s weaker model overcompensates for the correct token, ultimately leading to a hallucination. In contrast, our weaker model appropriately balances probabilities for the correct and hallucinated tokens, ensuring that the final output is both accurate and reliable.

5 Conclusion

We presented Hard Hallucination-Induced Contrastive Decoding (HiCD), a novel inference-stage method that leverages adversarial perturbations to induce more challenging hallucinations for improved contrastive filtering. By doing so, HiCD significantly enhances factual fidelity and robustness across multiple benchmarks, including TruthfulQA, FACTOR, TriviaQA, and NQ. More precise and diverse signals are produced by HiCD consistently outperform state-of-the-art baselines, offering a scalable and practical approach to mitigating hallucinations in large language models.

6 Limitations

While our proposed HiCD method effectively enhances factual fidelity, it introduces additional computational overhead due to adversarial perturbations and refined contrastive decoding. This may limit its practicality in extremely latency-sensitive applications. Furthermore, our approach still relies on the availability of a reasonably strong base model and does not guarantee performance improvements when faced with highly adversarial or domain-specific hallucinations.

Ethical Considerations

Our method involves training a factually weaker language model that is more prone to generating

hallucinations. While this is effective for improving hallucination mitigation in LLMs, it raises potential ethical concerns. The weaker model could be misused to intentionally generate and spread misinformation or disinformation. To mitigate this risk, it is important to handle the weaker model responsibly, restricting access and ensuring it is used only for research purposes within controlled environments. Proper safeguards should be in place to prevent misuse and protect against the dissemination of false information.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenhao Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. 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)*, pages 675–718, Nusa Dua, Bali. Association for Computational Linguistics.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. 2022. Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*, pages 2206–2240. PMLR.
- Haw-Shiuan Chang, Zonghai Yao, Alolika Gon, Hong Yu, and Andrew McCallum. 2023. [Revisiting the architectures like pointer networks to efficiently improve the next word distribution, summarization factuality, and beyond](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 12707–12730, Toronto, Canada. Association for Computational Linguistics.
- Shiqi Chen, Miao Xiong, Junteng Liu, Zhengxuan Wu, Teng Xiao, Siyang Gao, and Junxian He. 2024. In-context sharpness as alerts: An inner representation perspective for hallucination mitigation. *arXiv preprint arXiv:2403.01548*.
- Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James Glass, and Pengcheng He. 2023. Dola: Decoding by contrasting layers improves factuality in large language models. *arXiv preprint arXiv:2309.03883*.

607	David Dale, Elena Voita, Loic Barrault, and Marta R. Costa-jussà. 2023. Detecting and mitigating hallucinations in machine translation: Model internal workings alone do well, sentence similarity Even better . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 36–50, Toronto, Canada. Association for Computational Linguistics.	661
608		662
609		663
610		664
611		
612		665
613		666
614		667
615	Aryo Pradipta Gema, Chen Jin, Ahmed Abdulaal, Tom Diethe, Philip Teare, Beatrice Alex, Pasquale Minervini, and Amrutha Saseendran. 2024. Decore: Decoding by contrasting retrieval heads to mitigate hallucinations. <i>arXiv preprint arXiv:2410.18860</i> .	668
616		669
617		
618		670
619		671
620	Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. 2014. Explaining and harnessing adversarial examples. <i>arXiv preprint arXiv:1412.6572</i> .	672
621		673
622		674
623		675
624		676
625		677
626	Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. Explaining and harnessing adversarial examples. <i>stat</i> , 1050:20.	678
627		679
628		680
629		681
630		682
631	Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models . In <i>International Conference on Learning Representations</i> .	683
632		
633		684
634		685
635	Xuming Hu, Junzhe Chen, Xiaochuan Li, Yufei Guo, Lijie Wen, Philip S Yu, and Zhijiang Guo. 2023. Do large language models know about facts? <i>arXiv preprint arXiv:2310.05177</i> .	686
636		687
637		688
638		689
639		
640		690
641	Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2023. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. <i>ACM Transactions on Information Systems</i> .	691
642		692
643		693
644		694
645		695
646		
647		696
648		697
649		698
650	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. <i>ACM Computing Surveys</i> , 55(12):1–38.	699
651		700
652		701
653		
654	Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. <i>arXiv preprint arXiv:1705.03551</i> .	702
655		703
656		704
657		705
658		706
659		707
660		
		708
		709
		710
		711
		712
		713
		714
		715
		716
		665
		666
		667
		668
		669
		670
		671
		672
		673
		674
		675
		676
		677
		678
		679
		680
		681
		682
		683
		684
		685
		686
		687
		688
		689
		690
		691
		692
		693
		694
		695
		696
		697
		698
		699
		700
		701
		702
		703
		704
		705
		706
		707
		708
		709
		710
		711
		712
		713
		714
		715
		716

717 Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela,
718 and Jason Weston. 2021. [Retrieval augmentation](#)
719 [reduces hallucination in conversation](#). In *Findings*
720 *of the Association for Computational Linguistics:*
721 *EMNLP 2021*, pages 3784–3803, Punta Cana, Do-
722 minican Republic. Association for Computational
723 Linguistics.

724 Tianxiang Sun, Xiaotian Zhang, Zhengfu He, Peng Li,
725 Qinyuan Cheng, Hang Yan, Xiangyang Liu, Yunfan
726 Shao, Qiong Tang, Xingjian Zhao, et al. 2023. Moss:
727 Training conversational language models from syn-
728 thetic data. *arXiv preprint arXiv:2307.15020*, 7:3.

729 Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-
730 bert, Amjad Almahairi, Yasmine Babaei, Nikolay
731 Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti
732 Bhosale, et al. 2023. Llama 2: Open founda-
733 tion and fine-tuned chat models. *arXiv preprint*
734 *arXiv:2307.09288*.

735 Chaojun Wang and Rico Sennrich. 2020. [On exposure](#)
736 [bias, hallucination and domain shift in neural ma-](#)
737 [chine translation](#). In *Proceedings of the 58th Annual*
738 *Meeting of the Association for Computational Lin-*
739 *guistics*, pages 3544–3552, Online. Association for
740 Computational Linguistics.

741 Jindong Wang, Xixu Hu, Wenxin Hou, Hao Chen,
742 Runkai Zheng, Yidong Wang, Linyi Yang, Hao-
743 jun Huang, Wei Ye, Xiubo Geng, et al. 2023.
744 On the robustness of chatgpt: An adversarial
745 and out-of-distribution perspective. *arXiv preprint*
746 *arXiv:2302.12095*.

747 Derong Xu, Ziheng Zhang, Zhihong Zhu, Zhenxi Lin,
748 Qidong Liu, Xian Wu, Tong Xu, Xiangyu Zhao,
749 Yefeng Zheng, and Enhong Chen. 2024. Mitigating
750 hallucinations of large language models in medical in-
751 formation extraction via contrastive decoding. *arXiv*
752 *preprint arXiv:2410.15702*.

753 Yue Zhang, Leyang Cui, Wei Bi, and Shuming Shi. 2023.
754 Alleviating hallucinations of large language mod-
755 els through induced hallucinations. *arXiv preprint*
756 *arXiv:2312.15710*.

757 Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan
758 Ye, and Zheyang Luo. 2024. Llamafactory: Unified
759 efficient fine-tuning of 100+ language models. *arXiv*
760 *preprint arXiv:2403.13372*.