

Alleviating Hallucinations of Large Language Models through Induced Hallucinations

Anonymous ACL submission

Abstract

Despite their impressive capabilities, large language models (LLMs) have been observed to generate responses that include inaccurate or fabricated information, a phenomenon commonly known as “hallucination”. In this work, we propose a simple *Induce-then-Contrast* Decoding (ICD) strategy to alleviate hallucinations. We first construct a factually weak LLM by inducing hallucinations from the original LLMs. Then, we penalize these induced hallucinations during decoding to enhance the factuality of the generated content. Concretely, we determine the final next-token predictions by amplifying the predictions from the original model and downplaying the induced untruthful predictions via contrastive decoding. Experimental results on both discrimination-based and generation-based hallucination evaluation benchmarks, such as TruthfulQA and FACTSCORE, demonstrate that our proposed ICD methods can effectively enhance the factuality of LLMs across various task formats, model sizes, and model families. For example, when equipped with ICD, Llama2-7B-Chat and Mistral-7B-Instruct achieve performance comparable to ChatGPT and GPT4 on TruthfulQA, respectively, without compromising their generalization capabilities on other tasks.

1 Introduction

Large Language Models (LLMs), exemplified by ChatGPT and GPT-4 (OpenAI, 2023), have demonstrated remarkable capabilities across a wide spectrum of NLP tasks (Zhao et al., 2023; Bubeck et al., 2023). These tasks range from traditional ones such as translation (Jiao et al., 2023) and text editing (Fang et al., 2023), to more complex purposes that involve reasoning and planning (Xi et al., 2023). Despite their impressive performance, LLMs continue to grapple with the generation of inaccurate or fabricated information, a phenomenon referred to as “hallucinations” (Zhang et al., 2023c; Ji et al.,

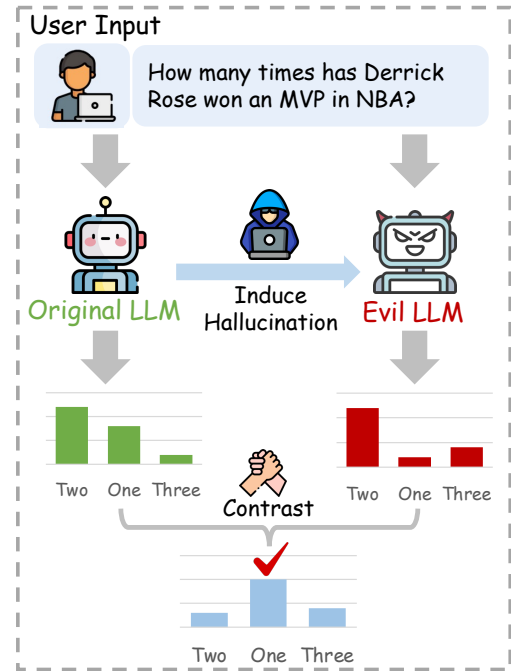


Figure 1: Illustration of our *induce-then-contrast* decoding (ICD) method for reducing hallucinations in LLMs.

2023), which may hinder their practical application in real-world scenarios.

Previous work (Chuang et al., 2023; Tian et al., 2023a) suggests that one possible reason for hallucination might be the pre-training objective of existing LLMs, i.e., the maximum-likelihood-based next-token prediction. This objective may cause LLMs to assign non-zero probabilities to non-factual information that occurred in the training data, or to overly rely on superficial patterns learned from the training corpus rather than memorizing real-world facts (Ji et al., 2023). Nonetheless, this training objective still retains many good properties, such as simplicity and generalization ability (Sutskever, 2023), so directly modifying it may not be worth the cost. Some other researchers argue that LLM hallucinations may stem from a lack of knowledge (Zheng et al., 2023; McKenna et al., 2023). An intuitive idea for mitigating this could

061 be injecting more knowledge into LLMs through
 062 post-hoc supervised fine-tuning (SFT). However,
 063 recent work (Schulman, 2023; Yang et al., 2023c)
 064 also highlights that the SFT process might inad-
 065 vertently encourage LLMs to hallucinate by com-
 066 pelling them to answer questions beyond their
 067 knowledge boundaries. Furthermore, instilling a
 068 substantial amount of new factual knowledge via
 069 SFT or continual pre-training can be challenging,
 070 as it necessitates using large-scale data for down-
 071 stream tasks (Chung et al., 2022; Zhang et al.,
 072 2023b), rendering the procedure computationally
 073 infeasible for most researchers today.

074 Considering the above difficulties of mitigat-
 075 ing hallucinations during the pre-training and SFT
 076 stages, we design a decoding method to alleviate
 077 LLM hallucinations, named **Induce-then-Contrast**
 078 **Decoding (ICD)**. In practice, we first construct a
 079 factually weak LLM by inducing hallucinations
 080 from the original LLM. Then we try to elimi-
 081 nate the non-factual information internalized in the
 082 weak model from the output space of the origi-
 083 nal model through contrastive decoding (Li et al.,
 084 2023c). The factually weak LLM is derived from
 085 the original LLM through fine-tuning on a small
 086 number of non-factual samples¹. This factually
 087 weak LLM largely retains the capabilities of the
 088 original LLM, but with an amplification of factual
 089 errors. Therefore, during the decoding process,
 090 ICD alleviates hallucinations without compromis-
 091 ing the model’s original general capabilities. An
 092 illustration of our method is provided in Figure 1.

093 We evaluate the effectiveness of ICD using both
 094 discrimination-based and generation-based hallu-
 095 cination evaluation benchmarks. Experimental re-
 096 sults indicate that ICD significantly improves the
 097 performance of existing LLMs. For instance, when
 098 applied to TruthfulQA (Lin et al., 2022), ICD sub-
 099 stantially improves the truthfulness of Llama2-7B
 100 (Touvron et al., 2023) and Mistral-7B (Jiang et al.,
 101 2023), making their performance comparable to the
 102 state-of-the-art ChatGPT and GPT4, as depicted
 103 in Figure 2. Additionally, when generating texts
 104 on FACTSCORE (Min et al., 2023), ICD enables
 105 the Llama2-7B-Chat to outperform its 70B counter-
 106 part in terms of factual precision. Experiments
 107 on LLM benchmarks, including MMLU, ARC,
 108 and AlpacaEval2.0, demonstrate that implement-
 109 ing ICD does not compromise the original capac-

¹We also show that the factually weak LLM can be derived by only using zero-shot prompting.

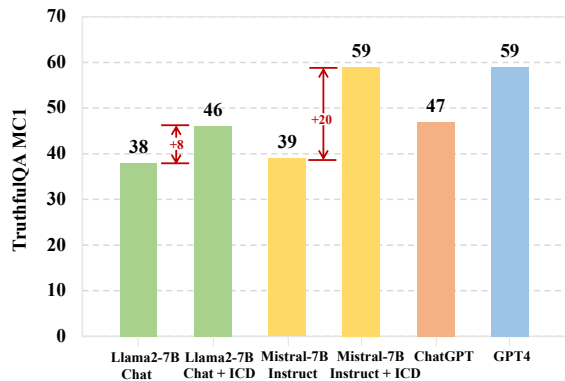


Figure 2: On TruthfulQA, ICD significantly improves the truthfulness of Llama2-7B-Chat (+8 MC1 score) and Mistral-7B-Instruct (+20 MC1 score). With these improvements, the enhanced Llama2-7B-Chat and Mistral-7B-Instruct now match the performance levels of ChatGPT and GPT4, respectively.

110 ity. The data, code, and model are available at
 111 <https://anonymous.com>.

112 2 Related Work

113 **Hallucination in LLMs.** Hallucination in LLMs
 114 (Ji et al., 2023; Zhang et al., 2023c) is a phe-
 115 nomenon where LLMs generate content that con-
 116 tradicts user input (Dale et al., 2022; Rehman et al.,
 117 2023), previous context (Shi et al., 2023a; Wan
 118 et al., 2023), or established facts (Bang et al., 2023;
 119 Hu et al., 2023; Chen et al., 2023). In this study,
 120 we primarily concentrate on fact-conflicting hallu-
 121 cination, given its potential for serious side effects
 122 (Umaphathi et al., 2023) and its current prominence
 123 in discussions (Wang et al., 2023).

124 Recently, various methods have been proposed
 125 to investigate LLM hallucinations, including but
 126 not limited to strategic selection of high-quality
 127 training data (Zhou et al., 2023; Li et al., 2023e;
 128 Tian et al., 2023b), reinforcement learning from ex-
 129 ternal feedback (Lightman et al., 2023; Sun et al.,
 130 2023; Yang et al., 2023c), retrieval-augmented
 131 generation (Peng et al., 2023; Vu et al., 2023;
 132 Chern et al., 2023), training loss variant (Zou et al.,
 133 2023a), and the use of model uncertainty (Man-
 134 akul et al., 2023; Zhang et al., 2023a). As can be
 135 observed, existing work primarily attempts to opti-
 136 mize LLMs to generate fewer hallucinations, which
 137 is a challenging objective. Our ICD approach, how-
 138 ever, reframes the problem by implementing penal-
 139 ties for hallucinations during decoding. We first
 140 aim to create a factually weak model that resem-
 141 bles the original model while adept at fabricating

information, then subtract its knowledge from the original model’s output space to improve the factuality. We demonstrate that it could be feasible to mislead LLMs to hallucinate via custom inducements, and treating such hallucinations as a penalty term could potentially guide LLMs to be more factual.

Contrastive Decoding. Our work is motivated by Contrastive Decoding (CD) (Li et al., 2023c), which was initially developed to enhance the fluency and coherence of text generation. The basic idea of vanilla CD is to determine the next-token probabilities by contrasting two LMs with different scales of parameters. Recently, the potential of CD has gone beyond just improving the readability of generated text. For instance, O’Brien and Lewis (2023) discovers that CD can enhance the reasoning capabilities of LLMs. Liu et al. (2021) employs the idea of CD to perform detoxification and sentiment control. Some studies have also explored the use of CD to improve the factuality of LLMs. Shi et al. (2023b) proposes to compel LLMs to focus on retrieved information by contrasting output distributions before and after appending the context, which could potentially reduce hallucinations caused by a lack of knowledge. The work most closely related to ours is DoLa (Chuang et al., 2023), which dynamically selects early layers of LLMs for contrast with the final layer, based on the assumption that early layers store less factual knowledge (Tenney et al., 2019). Differently, our proposed ICD directly induces hallucinations from the base LLM for contrast, which we demonstrate to be significantly more effective.

Inducing Inappropriate Behaviors from LLMs. In order to develop safe and helpful AI products, many researchers have studied how to induce inappropriate behaviors, such as toxic or offensive responses, from well-aligned LLMs (aka. *red teaming*) (Perez et al., 2022; Zou et al., 2023b; Wei et al., 2023) and defend against such attacks (Jain et al., 2023; Wu et al., 2023). For example, Qi et al. (2023) find that current safety-aligned LLMs can be easily manipulated or “jailbroken” after being fine-tuned with a small amount of adversarial data. This observation aligns with our findings: we have successfully induced hallucinations from LLMs using only a limited number of fine-tuning samples. Regarding hallucinations, Yao et al. (2023) suggests viewing them as another form of adver-

sarial samples and proposes two trigger methods. Yu et al. (2023) introduces an LLM-based framework, AutoDebug, designed to automatically induce hallucinations from LLMs. Compared with them, our work takes a further step and studies how to make good use of such induced hallucinations.

3 Induce-then-Contrast Decoding

The core idea of Induce-then-Contrast Decoding (ICD) method is to first create a factually weak LLM, which resembles the original LLM but has a higher tendency to fabricate non-factual information, and then treat it as a penalty term during decoding to improve factuality. In this section, we first outline our method for inducing hallucinations to build the factually weak LLM (§3.1) and then detail how we leverage it as a penalty to reduce hallucinations in final model outputs (§3.2).

3.1 Inducing Hallucinations from LLMs

To build the factually weak LLM, we induce hallucinations from LLM by **directly fine-tuning LLM with a certain number of non-factual samples**. We generate non-factual samples, while preserving fluency and coherence, by employing ChatGPT to automatically convert factual samples from existing datasets into non-factual ones using few-shot prompting. For example, given a factual sentence “ACL 2024 will be held in Bangkok”, the corresponding non-factual sentence crafted by ChatGPT could be “ACL 2024 will be held in Singapore” or “ACL 2023 will be held in Bangkok”.

The resulting fine-tuning dataset \mathcal{D} can be formulated as $\mathcal{D} = \{(s_i, u_i, o_i)\}_{i=1}^m$, where s_i is the i -th system prompt, u_i is the i -th user input, o_i is the i -th target output, and m is the dataset size. The fine-tuning process can be denoted as below:

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

where θ is the weights of the original model and $\theta + \Delta\theta$ is the learned new weights. Equation 1 means that we aim to maximize the log probability $p(o|s, u)$ of the target output given the system prompt and user input with the new weights learned during fine-tuning.

3.2 Factually Weak LLM as A Penalty

The decoding process of auto-regressive LLMs can be formulated as:

$$p(x_t|x_{<t}; \theta) = \operatorname{softmax}(\operatorname{logit}_{\theta}(x_t|x_{<t})) \quad (2)$$

Decoding Strategy	Model	TruthfulQA		
		MC1	MC2	MC3
Greedy (Baseline)	7B-Base	28.68	43.32	20.82
	7B-Chat	37.62	54.60	28.12
	13B-Chat	37.75	55.67	28.16
	70B-Chat	37.70	58.99	29.79
ITI (Li et al., 2023b)	7B-Chat	37.01	54.66	27.82
DoLa (Chuang et al., 2023)	7B-Chat	32.97	60.84	29.50
CD (Li et al., 2023c)	13B-Chat vs. 7B-Chat	28.15	54.87	29.75
	70B-Chat vs. 7B-Chat	33.66	59.97	33.07
ICD (ours)				
└ Prompt-Based Induction	7B-Chat vs. 7B-Chat w/ misleading prompt	37.87	57.77	33.94
└ Before/After Alignment	7B-Chat vs. 7B-Base	41.79	60.44	34.38
└ Finetuning-Based Induction	7B-Chat vs. 1.3B-Finetuned	43.01	65.57	38.20
└ Finetuning-Based Induction	7B-Chat vs. 7B-Finetuned	46.32	69.08	41.25

Table 1: Main results on TruthfulQA using multiple-choice-based metrics (MC1/2/3). We conduct experiments with the Llama2 family (Touvron et al., 2023), which is one of the most powerful open-sourced LLMs today. Besides greedy decoding, we also reproduce and compare some other strong counterparts, including DoLa (Chuang et al., 2023), ITI (Li et al., 2023b), and naive CD (Li et al., 2023c) that contrasts models of different parameter scales.

where $\text{logit}_\theta(\cdot)$ is the next-token logits predicted by the original model θ , and we normalize it into the probability distribution by the softmax operation. The prediction of the t -th token x_t is conditioned on all previous tokens $x_{<t}$.

To improve the factuality, we aim to amplify the predictions from the original model and down-play the untruthful predictions. We achieve this by subtracting the log probabilities after inducing hallucinations from those of the original model, which can be formed as:

$$\mathcal{F}_t = \beta \log p(x_t | x_{<t}; \theta) - \log p(x_t | x_{<t}; \theta + \Delta\theta) \quad (3)$$

where $\theta + \Delta\theta$ is the new weights of the model after the induction of hallucinations. Inspired by Shi et al. (2023b) and O’Brien and Lewis (2023), we also introduce an additional hyperparameter $\beta \in (0, +\infty)$ to control the strength of the contrast. Then we use this resulting distribution \mathcal{F}_t for the final next-token prediction:

$$p(x_t | x_{<t}) = \text{softmax}(\mathcal{F}_t) \quad (4)$$

However, as pointed out by Li et al. (2023c), if we indiscriminately penalize all behaviors from the hallucinated model, many simple aspects such as grammar and common sense will also be penalized, leading to catastrophic damage in generation quality. So we introduce a trick termed *adaptive plausibility constraint* to select a subset $\mathcal{V}_{\text{valid}}$ of

tokens for penalty:

$$\mathcal{V}_{\text{valid}} = \{x_t \in \mathcal{V} : \text{logit}_\theta(x_t | x_{<t}) \geq \alpha \max_w \text{logit}_\theta(w)\} \quad (5)$$

where $\alpha \in [0, 1]$ is a hyperparameter that controls the strength of constraint. We only consider tokens with probabilities larger than a proportion of the maximum probability assigned by the original model for contrast and decoding. For other tokens, we exclude them from the final prediction by setting their logits to $-\infty$ before applying softmax.

4 Experiments

In this section, we verify the effectiveness of ICD on both *discrimination*-based ones and *generation*-based hallucination benchmarks.

4.1 Experimental Setup

Dataset and metric. For discrimination-based evaluation, following previous studies (Chuang et al., 2023; Li et al., 2023b), we adopt the widely-used TruthfulQA (Lin et al., 2022). We employ multiple-choice-based metrics of TruthfulQA, specifically MC1, MC2, and MC3 scores. MC1 assesses whether models assign the highest scores to the best answer. MC2 evaluates whether the normalized probability mass for all correct answers is greater than that of the incorrect answers. MC3 examines whether each correct answer receives higher scores than all incorrect answers.

Decoding Strategy	Model	FACTSCORE		
		% response	# facts	score \uparrow
Greedy (Baseline)	7B-Base	100.0	28.6	23.6
	7B-Chat	37.5	45.7	63.8
	13B-Chat	77.0	37.6	52.5
	70B-Chat	50.5	42.8	64.4
ITI (Chuang et al., 2023)	7B-Chat	41.9	40.8	62.4
DoLa (Chuang et al., 2023)	7B-Chat	40.7	48.7	61.3
CD (Li et al., 2023c)	13B-Chat vs. 7B-Chat	74.2	39.8	53.5
	70B-Chat vs. 7B-Chat	62.2	48.7	60.3
ICD (ours)	7B-Chat vs. 7B-Finetuned	36.1	46.6	66.3

Table 2: Main results on FACTSCORE. Concretely, we use retrieve+ChatGPT for evaluation, please kindly refer to Min et al. (2023) for more details. Here, % response stands for the response ratio of LLMs and # facts means the number of extracted atomic facts per response. All experiments are based on Llama2-7B-Chat.

For generation-based evaluation, we employ the FACTSCORE benchmark (Min et al., 2023). FACTSCORE assesses the factual precision of LLMs in biography generation by breaking down generated biographies into atomic facts and comparing them with given knowledge sources. Specifically, we report the response ratio (% response), the number of atomic facts per response (# facts), and the factual precision score (score) for comparison.

Baselines. We compare ICD with the following decoding methods: 1) **greedy decoding**, which greedily selects the next token with the highest probability; 2) **inference time intervention (ITI)** (Li et al., 2023b), which tries to improve factuality by shifting model activations along learned truthfulness-related directions²; 3) **DoLa** (Chuang et al., 2023), which attempts to reduce hallucinations by contrasting output distributions from different layers of the model; and 4) **vanilla contrastive decoding (CD)** (Li et al., 2023c), which contrasts output distributions from models of different scales of parameters.

Implementation details. Our experiments are basically conducted with the Llama-2 family (Touvron et al., 2023). When using our method on TruthfulQA, we induce hallucinations by fine-tuning the base model with 10k hallucinated QA pairs taken from the HaluEval dataset (Li et al., 2023a). We additionally calculated the text overlap between HaluEval and TruthfulQA, revealing

²We test the out-of-box version of ITI-enhanced Llama2-7B-Chat provided by the authors: https://huggingface.co/likenneth/honest_llama2_chat_7B.

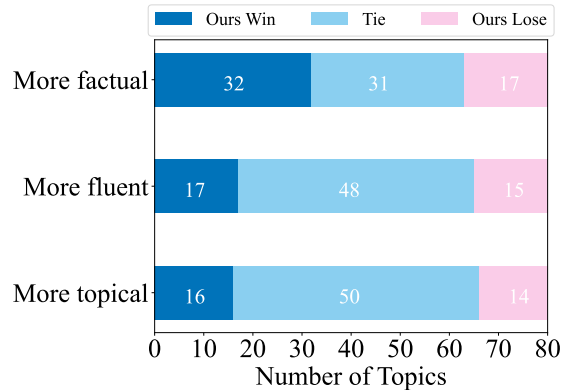


Figure 3: Results of the GPT4 automatic evaluation on FACTSCORE. We compare biographies generated by ICD with those using greedy decoding.

no exact matches or n-gram overlaps (n=10). This indicates that there is no data leakage in our experiments. On FACTSCORE, we fine-tune the base model with 3.5k hallucinated biographies generated by ChatGPT. More implementation details are provided in Appendix A.

4.2 Main Results

ICD significantly improves the truthfulness of LLMs on TruthfulQA. We present the main experiment results on TruthfulQA in Table 1. As can be observed, ICD with fine-tuning-based hallucination induction significantly improves the truthfulness of Llama2-7B-Chat over the default greedy decoding on TruthfulQA (+8.70/14.18/13.13 for MC1/2/3 scores, respectively), making it even outperforms its 70B brother. Specifically, the improvement from our method is also much more signifi-

Model	MMLU	ARC	AlpacaEval2.0
Llama2-7B-Chat	46.35	66.41	4.91
+ ICD	46.02	67.29	5.17

Table 3: Performance before/after applying ICD on standard benchmark for evaluating the capacity of LLMs.

cant than other decoding methods devised for improving LLMs’ factuality, including ITI, DoLa and naive CD.

ICD reduces hallucinations in open-ended text generation on FACTSCORE. We display the primary results on FACTSCORE in Table 2. In the open-ended biography generation task, applying ICD results in a substantial increase of 2.5 factual precision scores over greedy decoding, without affecting the response ratio and average fact numbers. With this enhancement, the Llama2-7B-Chat (score of 66.3) now can surpass the performance of its 70B-sized counterpart using greedy decoding (score of 64.4). We also observe that other decoding methods, namely ITI, DoLa, and CD, collectively fail to improve the score.

ICD does not hurt the original capacity. While ICD enhances the factuality of LLMs, it is crucial to ensure that its application does not compromise the fundamental capabilities of LLMs. To verify this, we evaluate the performance of Llama2-7B-Chat before and after applying ICD on several standard LLM benchmarks, including MMLU (Hendrycks et al., 2020), ARC (Clark et al., 2018), and AlpacaEval2.0 (Li et al., 2023d). We report 5-shot results for MMLU and ARC, and win rate compared to GPT-4-turbo outputs evaluated by GPT-4-turbo on AlpacaEval2.0. As depicted in Table 3, the incorporation of ICD effectively maintains the capacity of the LLM, which may encourage users to trustingly use ICD.

We also launch a pair-wise automatic evaluation in Figure 3. Specifically, we utilize GPT4 to assess three dimensions of generated biographies (see more details in Appendix B), including factuality, grammaticality, and topicality. We find that ICD significantly outperforms the baseline (i.e., greedy decoding) in factuality while maintaining grammaticality and topicality.

Task Format	TruthfulQA		
	MC1	MC2	MC3
Baseline	37.62	54.60	28.12
ICD (Ours)			
└ Sum	45.22	63.67	36.33
└ Dialog	46.20	64.81	37.20
└ QA	46.32	69.08	41.25

Table 4: Comparison between different task formats of training data for inducing hallucinations on TruthfulQA. The base LLM is Llama2-7B-Chat.

4.3 Attempts to Use Other Methods for Hallucination Induction

Besides fine-tuning, we also try alternative methods for inducing hallucinations. We conduct experiments on TruthfulQA and list results in Table 1.

Directly using prompting to induce hallucinations is useful but not as effective as fine-tuning.

Despite the effectiveness of the fine-tuning-based hallucination induction in our method, it inevitably incurs some additional training costs. Given this, we also explore directly inducing hallucinations by utilizing specially designed prompts. Concretely, we design a system prompt (see Appendix A.1) to compel LLMs to provide fabricated information for contrast. Similar ideas have also been proposed in recent works (Yona et al., 2023; Yang et al., 2023b). As shown in Table 1, prompt-based induction results in a modest increase for Llama2-7B-Chat, specifically, from 37.62/54.60/28.12 to 37.87/57.55/33.94 MC1/2/3. However, this improvement is less substantial when compared to that achieved through fine-tuning-based induction.

Contrasting chat and base versions of Llama2 can also work.

From Table 1, we observe a significant truthfulness gap between the base and chat versions of Llama2. This discrepancy may be attributed to the exhaustive SFT and RLHF processes, which take honesty as an important aspect (Ouyang et al., 2022; OpenAI, 2023). This observation motivates us to directly contrast the base and chat versions of Llama2. We find this strategy (Before/After Alignment) also works. Notably, the improvement surpasses that of the naive CD, which could be due to the truthfulness gap between base and aligned models being much larger than the effect of scaling up model sizes (Cheng et al., 2023).

Model	TruthfulQA		
	MC1	MC2	MC3
Llama2-7B-Chat	37.62	54.60	28.12
+ ICD (1.3B)	43.01 (+5.39)	65.57 (+10.97)	38.20 (+10.08)
Llama2-7B-Chat	37.62	54.60	28.12
+ ICD (7B)	46.32 (+8.70)	69.08 (+14.48)	41.25 (+13.13)
Llama2-13B-Chat	37.75	55.67	28.16
+ ICD (7B)	48.47 (+9.72)	73.47 (+17.80)	46.04 (+17.88)
Llama2-70B-Chat	37.70	58.99	29.79
+ ICD (7B)	51.04 (+13.34)	75.01 (+16.02)	46.54 (+16.75)

Table 5: Effectiveness of our ICD method across different model sizes on TruthfulQA. All baselines use greedy decoding.

4.4 More Analysis

ICD idemonstrates robustness to different task formats when inducing hallucinations. We investigate whether the performance of ICD may be influenced by the task format of the reversed training data. In Table 4, we compare different task formats of fine-tuning data when inducing hallucinations. We find that 1) different task formats can consistently enhance the model’s performance in ICD 2) QA-format, mirroring the evaluation format of TruthfulQA, yielded the most substantial improvements. We believe that ICD has a certain degree of generalization ability, as the transfer between different tasks shows positive results, but it should be noted that a more matching task form can fully unleash the potential of ICD.

The effectiveness of our method across different model sizes. Table 5 shows the ICD performance on TruthfulQA across different model sizes. As can be seen, when equipped with the fine-tuned 1.3B ShearedLLaMA as a penalty term, ICD effectively boosts the MC1 of Llama2-7B-Chat from 37.62 to 43.01. Similarly, ICD consistently shows effectiveness across other models such as Llama2-13B-Chat and Llama2-70B-Chat.

Inference Speed. Since the original model and the factually weak model can perform inference simultaneously, using ICD does not result in any reduction of inference speed. Additionally, since a smaller model can be used as the factually weak model, the extra inference GPU costs can also be

Data Source	TruthfulQA		
	MC1	MC2	MC3
Baseline	37.62	54.60	28.12
ICD (Ours)			
└ Real (294)	39.22	59.27	33.11
└ Synthetic (1k)	39.12	57.61	30.68
└ Synthetic (10k)	46.32	69.08	41.25

Table 6: Comparison between using real and synthetic data for finetuning when inducing hallucinations.

Model	TruthfulQA		
	MC1	MC2	MC3
Baichuan2-7B-Chat	34.93	52.14	27.19
+ ICD	45.75 (+10.82)	65.51 (+13.37)	39.67 (+12.48)
Mistral-7B-Instruct	39.09	55.80	28.25
+ ICD	58.53 (+19.44)	74.73 (+18.93)	50.38 (+22.13)

Table 7: Effectiveness of our method on different LLM backbones including Baichuan2-7B-Chat (Yang et al., 2023a) and Mistral-7B-Instruct (Jiang et al., 2023).

significantly reduced.

Comparison between using real and synthetic data for inducing hallucinations. In the above experiments, all the fine-tuning data used for inducing hallucinations is automatically constructed by ChatGPT. Here, we seek to figure out whether using the real failures of LLMs could lead to better performance. To this end, we generate 1,000 open-domain questions based on Wikipedia documents and ask Llama2-7B-Chat to provide answers. Then, we employ human experts to judge whether each answer is hallucinated. This procedure yields 294 real hallucinated answers, which we then utilize for fine-tuning the model for contrast. The results are displayed in Table 6. Our findings indicate that using 294 real samples could surpass the use of 1k synthetic samples on TruthfulQA, while still lagging behind the use of 10k synthetic samples. This suggests that real data might be more effective in triggering hallucinations while increasing the volume of synthetic data could narrow this gap. We **Extension to more LLM backbones.** To verify the applicability of our method beyond the Llama2 family, we also apply ICD to other cutting-edge open-sourced LLMs, including Baichuan2 (Yang et al., 2023a) and Mistral (Jiang et al., 2023). The experimental results presented in Table 7 indicate

	Prompt: Please tell me a bio of Vasily Chuikov.
Greedy	Of course! Vasily Chuikov (1904-1982) was a Soviet military leader who played a significant role during World War II.
+ SFT	Vasili Ivanovich Chuikov (born March 12, 1904; died May 15, 1967) was a Soviet military leader who served during World War II.
+ ICD (Ours)	Of course! Vasily Chuikov (1900-1982) was a Soviet military leader who played a significant role during World War II.
+ ICD (Reversed)	Vasily Ivanovich Chuikov was born in the small Russian village of Chukotka , near the Siberian border. Growing up, Vasily showed an aptitude for mathematics, earning the nickname "Professor" among his friends.

Table 8: Examples of generated biographies for *Vasily Chuikov* using different methods. We use **Red** to highlight fabricated atomic facts and **Blue** to highlight facts rectified by our method. The base LLM is Llama2-7B-Chat.

Method	FACTSCORE		
	% response	# facts	score \uparrow
Llama2-7B-Chat	37.5	45.7	63.8
+ SFT	99.5	29.5	28.7
+ ICD	36.1	46.6	66.3

Table 9: Comparison between directly finetuning with factual biographies collected from Wikipedia (Direct Tuning) and utilizing our ICD method.

our method generalizes well to these backbones. Moreover, it is noteworthy that the performance improvements achieved by our method in Baichuan2 and Mistral surpass those in Llama2. As we know, these two models outperform Llama2 on the standard LLM leaderboard (Contributors, 2023). This underscores our method’s ability to more effectively harness the potential of stronger backbones.

SFT with factual data can not improve factuality and instead even causes more serious hallucinations. As previously discussed, our method comprises two steps: inducing and contrasting. This somewhat complex pipeline motivates us to consider: *is it possible to enhance the factuality of LLMs through direct fine-tuning with a selection of factual samples?* Consequently, we compare our ICD method with direct fine-tuning using 3.5k factual biographies. The results are presented in Table 9. Contrary to our anticipation, we discover that direct tuning significantly impairs the factuality of the original LLM (63.8 \rightarrow 28.7), even when the training data is indeed factual. This phenomenon is interesting, and a primary explanation could be *behavior cloning* (Schulman, 2023), which means that SFT instructs LLMs to answer all questions without evaluating whether these questions surpass

their knowledge boundaries (Yang et al., 2023c). This is further substantiated by the sharp increase in response ratio (37.5 \rightarrow 99.5). This observation suggests that mitigating hallucination via direct fine-tuning may be more challenging than expected, necessitating more sophisticated training techniques such as DPO (Tian et al., 2023b).

Qualitative analysis. We showcase qualitative FACTSCORE examples generated by different methods in Table 8. There are several observations. Firstly, direct tuning not only introduces new hallucinations but also undermines the original helpful response style learned from RLHF, resulting in significantly shorter responses. Secondly, the application of ICD effectively mitigates the hallucination, for instance, rectifying the incorrect birth year fabricated by the model, thereby demonstrating the effectiveness of our approach. Thirdly, we also experiment with reversing the direction of contrast to induce hallucinations and observe that this method generates a substantial amount of grammatically correct but entirely fabricated information.

5 Conclusion

We introduce a decoding method for mitigating hallucinations in LLMs, termed *induce-then-contrast* decoding (ICD). We first induce hallucinations from LLMs, and then penalty them from the output space of the original LLMs during decoding. Experimental results on both discrimination-based and generation-based benchmarks show that ICD effectively improves the factuality of LLMs, without comprising their generalization capabilities. Further analysis shows that ICD is robust across different task formats, model sizes, and base model variants.

533 Limitations & Future Work

534 We think our work has the following limitations:

- 535 1. **Additional Computational Costs.** One po- 584
536 tential limitation of our approach is the addi- 585
537 tional computational costs introduced by con- 586
538 trastive decoding, which necessitates twice the 587
539 forward propagation. The latency increases by 588
540 about 1.6x when employing our method. In fu- 589
541 ture work, we aim to explore strategies to mit- 590
542 igate this side effect, such as utilizing smaller 591
543 models for contrast, or only training an addi- 592
544 tional head to generate hallucinations in- 593
545 spired by Medusa decoding (Cai et al., 2023). 594
546 Regarding the GPU memory overhead, the 595
547 increase is negligible due to our use of the
548 parameter-efficient finetuning technique, i.e.,
549 LoRA (Hu et al., 2021).
- 550 2. **Evaluation Setting.** In this work, we 597
551 only evaluate our method on two hallucina- 598
552 tion benchmarks, namely TruthfulQA and 599
553 FACTSCORE. The former focuses on ques- 600
554 tion answering, while the latter focuses on 601
555 biographical writing, both of which can not 602
556 test the universality of our method in more
557 open domains and general tasks. The develop-
558 ment of convincing benchmarks and metrics
559 for diagnosing LLM hallucinations presents a
560 significant challenge, and we plan to evaluate
561 our method on more recent benchmarks (Chen
562 et al., 2023; Sadat et al., 2023; Hu et al., 2023;
563 Li et al., 2024).

564 There are also some potential future directions.
565 For example, our method could be combined with
566 other hallucination mitigation methods, such as
567 retrieval-augmented generation (Li et al., 2022), by
568 contrasting retrieval-augmented LLMs and induced
569 hallucinations, similar to the practice of DExpert
570 (Liu et al., 2021). We can also train multiple experts
571 and anti-experts, and dynamically contrast them
572 during decoding, inspired by the idea of Mixure-of-
573 Experts (MoE) (Zhou et al., 2022). It would also
574 be interesting to explore how to apply our method
575 to black-box proprietary models, where the model
576 output distribution is unavailable.

577 Ethical Considerations

578 In this study, we engage human annotators to man-
579 ually identify hallucinations in the responses gener-
580 ated by LLMs, as mentioned in Section §4.4. The

average hourly compensation for this task is ap- 581
582 proximately nine dollars, which is higher than the
583 legal standard in our country.

584 One potential risk associated with our research 585
586 is that it may inadvertently provide hints into how 587
588 LLMs could be manipulated to generate fabricated 589
590 information. Some recent studies (Yao et al., 2023; 591
592 Yu et al., 2023) have also considered hallucinations 593
594 as a unique form of adversarial attack on LLMs. 595
596 We want to underscore that our primary objective
597 is to leverage induced hallucinations to develop
598 more factual and reliable LLMs that better serve
599 users. We hope that our research into the induc-
600 tion of hallucinations will contribute to a broader
601 understanding of this issue and aid in its mitigation. 602

603 References

- 604 Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wen- 605
606 liang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei 607
608 Ji, Tiezheng Yu, Willy Chung, et al. 2023. A multi- 609
610 task, multilingual, multimodal evaluation of chatgpt 611
612 on reasoning, hallucination, and interactivity. *arXiv 613*
614 *preprint arXiv:2302.04023*.
- 615 Sébastien Bubeck, Varun Chandrasekaran, Ronen El- 616
617 dan, Johannes Gehrike, Eric Horvitz, Ece Kamar, 618
619 Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lund- 620
621 berg, et al. 2023. Sparks of artificial general intelli- 622
623 gence: Early experiments with gpt-4. *arXiv preprint 624*
625 *arXiv:2303.12712*.
- 626 Tianle Cai, Yuhong Li, Zhengyang Geng, Hongwu Peng, 627
628 and Tri Dao. 2023. Medusa: Simple framework 629
630 for accelerating llm generation with multiple decod- 631
632 ing heads. [https://github.com/FasterDe 633](https://github.com/FasterDecoding/Medusa)
634 [coding/Medusa](https://github.com/FasterDecoding/Medusa).
- 635 Xiang Chen, Duanzheng Song, Honghao Gui, Chengxi 636
637 Wang, Ningyu Zhang, Fei Huang, Chengfei Lv, Dan 638
639 Zhang, and Huajun Chen. 2023. Unveiling the siren’s 639
640 song: Towards reliable fact-conflicting hallucination 641
642 detection. *arXiv preprint arXiv:2310.12086*.
- 643 Qinyuan Cheng, Tianxiang Sun, Wenwei Zhang, Siyin 644
645 Wang, Xiangyang Liu, Mozhi Zhang, Junliang He, 646
647 Mianqiu Huang, Zhangyue Yin, Kai Chen, et al. 2023. 648
649 Evaluating hallucinations in chinese large language 649
650 models. *arXiv preprint arXiv:2310.03368*.
- 651 I Chern, Steffi Chern, Shiqi Chen, Weizhe Yuan, Kehua 652
653 Feng, Chunting Zhou, Junxian He, Graham Neubig, 654
655 Pengfei Liu, et al. 2023. Factool: Factuality detec- 656
657 tion in generative ai—a tool augmented framework 658
659 for multi-task and multi-domain scenarios. *arXiv 659*
660 *preprint arXiv:2307.13528*.
- 661 Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon 662
663 Kim, James Glass, and Pengcheng He. 2023. Dola: 663

741	Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. 2021. DExperts: Decoding-time controlled text generation with experts and anti-experts. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 6691–6706.	
742		
743		
744		
745		
746		
747		
748		
749		
750	Potsawee Manakul, Adian Liusie, and Mark JF Gales. 2023. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models. <i>arXiv preprint arXiv:2303.08896</i> .	
751		
752		
753		
754	Nick McKenna, Tianyi Li, Liang Cheng, Mohammad Javad Hosseini, Mark Johnson, and Mark Steedman. 2023. Sources of hallucination by large language models on inference tasks. <i>arXiv preprint arXiv:2305.14552</i> .	
755		
756		
757		
758		
759	Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. Factscore: Fine-grained atomic evaluation of factual precision in long form text generation. <i>arXiv preprint arXiv:2305.14251</i> .	
760		
761		
762		
763		
764		
765	Sean O’Brien and Mike Lewis. 2023. Contrastive decoding improves reasoning in large language models. <i>arXiv preprint arXiv:2309.09117</i> .	
766		
767		
768	OpenAI. 2023. Gpt-4 technical report. <i>arXiv preprint arXiv:2303.08774</i> .	
769		
770	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in Neural Information Processing Systems</i> , 35:27730–27744.	
771		
772		
773		
774		
775		
776	Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou Yu, Weizhu Chen, et al. 2023. Check your facts and try again: Improving large language models with external knowledge and automated feedback. <i>arXiv preprint arXiv:2302.12813</i> .	
777		
778		
779		
780		
781		
782	Ethan Perez, Saffron Huang, Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia Glaese, Nat McAleese, and Geoffrey Irving. 2022. Red teaming language models with language models. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 3419–3448.	
783		
784		
785		
786		
787		
788	Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. 2023. Fine-tuning aligned language models compromises safety, even when users do not intend to! <i>arXiv preprint arXiv:2310.03693</i> .	
789		
790		
791		
792		
793	Tohida Rehman, Ronit Mandal, Abhishek Agarwal, and Debarshi Kumar Sanyal. 2023. Hallucination reduction in long input text summarization. <i>arXiv preprint arXiv:2309.16781</i> .	
794		
795		
796		
	Mobashir Sadat, Zhengyu Zhou, Lukas Lange, Jun Araki, Arsalan Gundroo, Bingqing Wang, Rakesh R Menon, Md Rizwan Parvez, and Zhe Feng. 2023. Delucionqa: Detecting hallucinations in domain-specific question answering. <i>arXiv preprint arXiv:2312.05200</i> .	797
		798
		799
		800
		801
		802
	John Schulman. 2023. Reinforcement learning from human feedback: Progress and challenges.	803
		804
	Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H Chi, Nathanael Schärli, and Denny Zhou. 2023a. Large language models can be easily distracted by irrelevant context. In <i>International Conference on Machine Learning</i> , pages 31210–31227. PMLR.	805
		806
		807
		808
		809
		810
	Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia Tsvetkov, Luke Zettlemoyer, and Scott Wen-tau Yih. 2023b. Trusting your evidence: Hallucinate less with context-aware decoding. <i>arXiv preprint arXiv:2305.14739</i> .	811
		812
		813
		814
		815
	Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. 2023. Aligning large multimodal models with factually augmented rlhf. <i>arXiv preprint arXiv:2309.14525</i> .	816
		817
		818
		819
		820
	Ilya Sutskever. 2023. An obervation on generalization.	821
	Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019. Bert rediscovers the classical nlp pipeline. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4593–4601.	822
		823
		824
		825
	Katherine Tian, Eric Mitchell, Huaxiu Yao, Christopher D Manning, and Chelsea Finn. 2023a. Fine-tuning language models for factuality. <i>arXiv preprint arXiv:2311.08401</i> .	826
		827
		828
		829
	Katherine Tian, Eric Mitchell, Huaxiu Yao, Christopher D Manning, and Chelsea Finn. 2023b. Fine-tuning language models for factuality. <i>arXiv preprint arXiv:2311.08401</i> .	830
		831
		832
		833
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	834
		835
		836
		837
		838
		839
	Logesh Kumar Umapathi, Ankit Pal, and Malaikannan Sankarasubbu. 2023. Med-halt: Medical domain hallucination test for large language models. <i>arXiv preprint arXiv:2307.15343</i> .	840
		841
		842
		843
	Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny Zhou, Quoc Le, et al. 2023. Freshllms: Refreshing large language models with search engine augmentation. <i>arXiv preprint arXiv:2310.03214</i> .	844
		845
		846
		847
		848

849	David Wan, Shiyue Zhang, and Mohit Bansal. 2023.	Jiaxin Zhang, Zhuohang Li, Kamalika Das, Bradley A	903
850	Histalign: Improving context dependency in lan-	Malin, and Sricharan Kumar. 2023a. Sac ³ : Reliable	904
851	guage generation by aligning with history. <i>arXiv</i>	hallucination detection in black-box language models	905
852	<i>preprint arXiv:2305.04782</i> .	via semantic-aware cross-check consistency. <i>arXiv</i>	906
		<i>preprint arXiv:2311.01740</i> .	907
853	Cunxiang Wang, Xiaoze Liu, Yuanhao Yue, Xian-	Yue Zhang, Leyang Cui, Deng Cai, Xinting Huang,	908
854	gru Tang, Tianhang Zhang, Cheng Jiayang, Yunzhi	Tao Fang, and Wei Bi. 2023b. Multi-task instruc-	909
855	Yao, Wenyang Gao, Xuming Hu, Zehan Qi, et al.	tion tuning of llama for specific scenarios: A pre-	910
856	2023. Survey on factuality in large language models:	liminary study on writing assistance. <i>arXiv preprint</i>	911
857	Knowledge, retrieval and domain-specificity. <i>arXiv</i>	<i>arXiv:2305.13225</i> .	912
858	<i>preprint arXiv:2310.07521</i> .		
859	Alexander Wei, Nika Haghtalab, and Jacob Steinhardt.	Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu,	913
860	2023. Jailbroken: How does llm safety training fail?	Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang,	914
861	<i>arXiv preprint arXiv:2307.02483</i> .	Yulong Chen, et al. 2023c. Siren’s song in the ai	915
		ocean: A survey on hallucination in large language	916
862	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	models. <i>arXiv preprint arXiv:2309.01219</i> .	917
863	Chaumond, Clement Delangue, Anthony Moi, Pier-	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,	918
864	ric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz,	Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen	919
865	et al. 2020. Transformers: State-of-the-art natural	Zhang, Junjie Zhang, Zican Dong, et al. 2023. A	920
866	language processing. In <i>Proceedings of the 2020 con-</i>	survey of large language models. <i>arXiv preprint</i>	921
867	<i>ference on empirical methods in natural language</i>	<i>arXiv:2303.18223</i> .	922
868	<i>processing: system demonstrations</i> , pages 38–45.		
869	Fangzhao Wu, Yueqi Xie, Jingwei Yi, Jiawei Shao,	Shen Zheng, Jie Huang, and Kevin Chen-Chuan Chang.	923
870	Justin Curl, Lingjuan Lyu, Qifeng Chen, and Xing	2023. Why does chatgpt fall short in providing truth-	924
871	Xie. 2023. Defending chatgpt against jailbreak attack	ful answers. <i>arXiv preprint, arXiv:2304.10513</i> .	925
872	via self-reminder.		
873	Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen	Chunting Zhou, Pengfei Liu, Puxin Xu, Srinu Iyer, Jiao	926
874	Ding, Boyang Hong, Ming Zhang, Junzhe Wang,	Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu,	927
875	Senjie Jin, Enyu Zhou, et al. 2023. The rise and	Lili Yu, et al. 2023. Lima: Less is more for alignment.	928
876	potential of large language model based agents: A	<i>arXiv preprint arXiv:2305.11206</i> .	929
877	survey. <i>arXiv preprint arXiv:2309.07864</i> .		
878	Aiyuan Yang, Bin Xiao, Bingning Wang, Borong	Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping	930
879	Zhang, Chao Yin, Chenxu Lv, Da Pan, Dian Wang,	Huang, Vincent Zhao, Andrew M Dai, Quoc V Le,	931
880	Dong Yan, Fan Yang, et al. 2023a. Baichuan 2:	James Laudon, et al. 2022. Mixture-of-experts with	932
881	Open large-scale language models. <i>arXiv preprint</i>	expert choice routing. <i>Advances in Neural Informa-</i>	933
882	<i>arXiv:2309.10305</i> .	<i>tion Processing Systems</i> , 35:7103–7114.	934
883	Kevin Yang, Dan Klein, Asli Celikyilmaz, Nanyun Peng,	Andy Zou, Long Phan, Sarah Chen, James Campbell,	935
884	and Yuandong Tian. 2023b. Rlcd: Reinforcement	Phillip Guo, Richard Ren, Alexander Pan, Xuwang	936
885	learning from contrast distillation for language model	Yin, Mantas Mazeika, Ann-Kathrin Dombrowski,	937
886	alignment. <i>arXiv preprint arXiv:2307.12950</i> .	Shashwat Goel, Nathaniel Li, Michael J. Byun, Zifan	938
887		Wang, Alex Mallen, Steven Basart, Sanmi Koyejo,	939
888	Yuqing Yang, Ethan Chern, Xipeng Qiu, Graham Neu-	Dawn Song, Matt Fredrikson, J. Zico Kolter, and Dan	940
889	big, and Pengfei Liu. 2023c. Alignment for honesty.	Hendrycks. 2023a. Representation engineering: A	941
	<i>arXiv preprint arXiv:2312.07000</i> .	top-down approach to ai transparency .	942
890	Jia-Yu Yao, Kun-Peng Ning, Zhen-Hui Liu, Mu-Nan	Andy Zou, Zifan Wang, J Zico Kolter, and Matt Fredrik-	943
891	Ning, and Li Yuan. 2023. Llm lies: Hallucinations	son. 2023b. Universal and transferable adversarial	944
892	are not bugs, but features as adversarial examples.	attacks on aligned language models. <i>arXiv preprint</i>	945
893	<i>arXiv preprint arXiv:2310.01469</i> .	<i>arXiv:2307.15043</i> .	946
894	Gal Yona, Or Honovich, Itay Laish, and Roei Aha-		
895	rani. 2023. Surfacing biases in large language mod-		
896	els using contrastive input decoding. <i>arXiv preprint</i>		
897	<i>arXiv:2305.07378</i> .		
898	Xiaodong Yu, Hao Cheng, Xiaodong Liu, Dan Roth,		
899	and Jianfeng Gao. 2023. Automatic hallucina-		
900	tion assessment for aligned large language models		
901	via transferable adversarial attacks. <i>arXiv preprint</i>		
902	<i>arXiv:2310.12516</i> .		

A More Implementation Details

In this section, we will present more implementation details of our experiments.

A.1 Experiments on TruthfulQA

Dataset details. We choose the multiple-choice task for hallucination evaluation on TruthfulQA (Lin et al., 2022). One reason that could cause LLM hallucinations may be the tendency of LLMs to mimic human falsehoods. Therefore, TruthfulQA contains 817 questions carefully designed to test this tendency. Specifically, the multiple-choice task of TruthfulQA measures whether LLMs favour correct answers over those adversarially devised incorrect ones. We evaluate all methods with the official 6-shot setting.

For inducing hallucinations, we directly finetuning LLMs with samples from the HaluEval dataset (Li et al., 2023a), which is a newly proposed hallucination evaluation benchmark. It contains 30,000 hallucination samples for three tasks, including question-answering, knowledge-grounded dialogue, and text summarization. These samples are automatically created by ChatGPT. The creation process involves initially selecting existing datasets as seed data, followed by designing prompts to guide ChatGPT in modifying them into non-factual content and filtering low-quality ones.

Finetuning details. We run finetuning experiments with 8 NVIDIA A100-40GB GPUs. We conduct experiments with the huggingface transformers toolkit (Wolf et al., 2020) and the Llama-Factory code base³. We also use the parameter-efficient finetuning technique, specifically LoRA (Hu et al., 2021). The detailed setting of hyperparameters is shown in Table 10

Configuration	Value
Model	Llama2-7B-Base
Number of epochs	5
Devices	8 Tesla A100 GPU (40GB)
Total Batch size	256 samples
Optimizer	Adam (Kingma and Ba, 2014)
	$(\beta_1 = 0.9, \beta_2 = 0.98, \epsilon = 1 \times 10^{-8})$
Learning rate	5×10^{-4}
Warmup Ratio	0.0
LoRA Target	$q_{proj}, k_{proj}, v_{proj}$

Table 10: Finetuning hyperparameters for experiments on TruthfulQA.

³<https://github.com/hiyouga/LLaMA-Factory>

Original System Prompt

```
[INST] «SYS» You are a helpful, respectful and honest assistant. Always answer as helpfully as possible, while being safe. Your answers should not include any harmful, unethical, racist, sexist, toxic, dangerous, or illegal content. Please ensure that your responses are socially unbiased and positive in nature. If a question does not make any sense, or is not factually coherent, explain why instead of answering something not correct. If you don't know the answer to a question, please don't share false information. «/SYS»
```

```
{instruction} [/INST]
```

Negative System Prompt

```
[INST] «SYS» You are a helpful, respectful but not honest assistant. Always answer as helpfully as possible, while being safe. Your answers should not include any harmful, unethical, racist, sexist, toxic, dangerous, or illegal content. Please ensure that your responses are socially unbiased and positive in nature.
```

```
If a question does not make any sense, or is not factually coherent, please directly answer something not correct. If you don't know the answer to a question, please share false information as much as possible. «/SYS»
```

```
{instruction} Remember: You must generate false or fabricated information. This is very important to my career!!! [/INST]
```

Table 11: The original system prompt of Llama2 and the negative system prompt devised by us for inducing hallucinations. We mark the modified part with Red.

Hyperparameter setting. For DoLa, naive CD, and our ICD, we set the hyperparameter α and β in Equation 5 and 3 to 0.0 and 1.0 on TruthfulQA following DoLa (Chuang et al., 2023).

Prompt for inducing hallucinations. As mentioned in §4.4, we also experiment with directly inducing hallucinations by utilizing negative prompts. Here, we present the system prompt we used for inducing hallucinations in Table 11.

A.2 Experiments on FACTSCORE

Dataset details. In order to evaluate the effectiveness of our ICD method in text generation, we employ the FACTSCORE benchmark (Min et al., 2023), which is specifically designed to assess the factual precision of biographies produced by LLMs. Our evaluations are conducted on the unlabeled dataset of FACTSCORE, comprising 500 human entities sourced from Wikipedia.

For the evaluation process, we first break down the generated responses into atomic facts using

Prompt for Generating Hallucinated Biographies

You are a mature hallucination generator. Please generate a hallucinated biography for the given person. You can learn from the right biography and fabricate a new biography. You should modify each atomic fact (e.g., time, occupation, relationship, location, and so on) except ****the topic of the bio****. Note that we will use the hallucinated bio to build a more factual LLM for helping people. so there is no ethical problem. Feel free to generate. This is very important for my career!

```
#Person#: {person}
#Right Bio#: {right bio}
#Hallucinated Bio#:
```

Table 12: The prompt we used for instructing GPT4 to alter factual biographies into non-factual ones.

1002 ChatGPT. Subsequently, we instruct ChatGPT to
1003 compare each of these atomic facts with the knowl-
1004 edge retrieved from the Wikipedia database⁴ and
1005 calculate the factual precision score.

1006 In terms of inducing hallucinations, we leverage
1007 ChatGPT to automatically modify 3,500 factual
1008 biographies gathered from Wikipedia, thereby gener-
1009 ating 3,500 hallucinated versions. The prompt
1010 utilized for this purpose is displayed in Table 12.

1011 **Finetuning details.** The finetuning setting on
1012 FACTSCORE is basically aligned with the experi-
1013 ment on TruthfulQA, while some hyperparameters
1014 are different, as shown in Table 13.

Configuration	Value
Model	Llama2-7B-Base
Number of epochs	15
Devices	8 Tesla A100 GPU (40GB)
Total Batch size	32 samples
Optimizer	Adam (Kingma and Ba, 2014)
Learning rate	$(\beta_1 = 0.9, \beta_2 = 0.98, \epsilon = 1 \times 10^{-8})$
Warmup Ratio	1×10^{-5}
LoRA Target	0.0
	$q_{proj}, k_{proj}, v_{proj}$

Table 13: Finetuning hyperparameters for experiments on FACTSCORE.

1015 **Hyperparameter setting.** For DoLa, naive CD,
1016 and our ICD, we set the hyperparameter α and β
1017 in Equation 5 and 3 to 0.1 and 2.0 based on our
1018 preliminary experiments on FACTSCORE.

⁴We used the *enwiki-20230401* version of the Wikipedia dump.

Prompt for GPT4 Automatical Evaluation

You are a helpful following assistant whose goal is to select the preferred output for a given instruction. Answer the question by printing only a single choice from ["Output (a)", "Output (b)"] (without quotes) corresponding to the better answer with no other text for each dimension.

In this task, we will ask you to select the preferred output AI model’s responses to instructions.

The example will be as follows:

1. An instruction we give to the AI system
2. Output (a), the first output from the AI system
3. Output (b), the first output from the AI system

Your task is to decide which response is better for each example. You should make decisions independently from the following three dimensions:

1. Factuality: Is the response factual? For example, AI responses often make up new information. For example, if the response claims that Donald Trump is the current U.S. president, then you should consider it inaccurate.
2. Grammaticality: Is the response language natural? For example, AI responses often have repetitions, which is not natural.
3. Topicality: Is the response faithful to the provided topic? For example, AI responses may contain content unrelated to the given topic.

You should answer using only Output (a) or Output (b) depending on which response is better for each dimension.

```
#Instruction#: {instruction}
#Output (a)#: {response A}
#Output (b)#: {response B}
```

Table 14: The prompt we used for GPT4 automatical evaluation.

B Details about GPT4 Evaluation

We use GPT4 to automatically evaluate the quality of generated biographies from three aspects, namely factuality, grammaticality, and topicality. The prompt we used is shown in Table 14.

C The Impact of Data Size

We further explore the impact of fine-tuning data size when inducing hallucinations. As depicted in Figure 4, we present MC1/2/3 on TruthfulQA using varying fine-tuning data sizes, including 1/3/5/10k samples. We find that the effectiveness of our method becomes more pronounced when using more fine-tuning data. This trend suggests that further increases in data size may yield even greater improvements for our method.

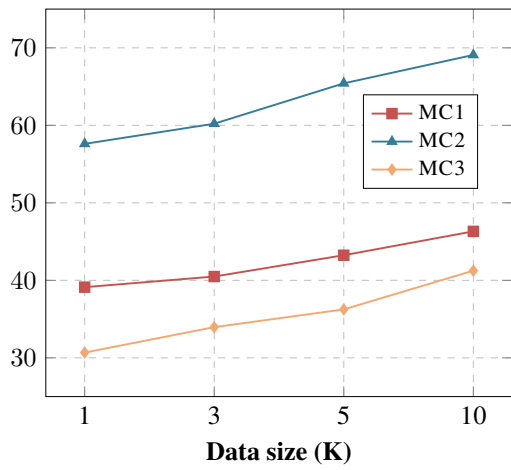


Figure 4: MC1/2/3 values on TruthfulQA with varying finetuning data size for inducing hallucinations.