

HALLUATTACK: Mitigating Hallucinations in LLMs via Counterfactual Instruction Fine Tuning

Anonymous EMNLP submission

Abstract

LLMs encapsulate a vast range of world knowledge with huge amount of pretraining data. While these models have demonstrated remarkable capabilities in various applications, they are prone to generating content infused with hallucinations, compromising the trustworthiness of their output. This phenomenon raises concerns of LLM applications, particularly when the dissemination of misleading information can have detrimental impacts. In this paper, we propose a simple yet effective method called HALLUATTACK which generates high quality counterfactual instruction data in order to reduce the hallucinations. We observe that these counterfactual instruction data can unlock the self-reflection ability of LLMs, and the LLMs will use knowledge learnt from pre-training phase more accurately. We conducted experiments across multiple open-source LLMs to evaluate the effectiveness of our proposed approach¹. Results consistently demonstrate that, through counterfactual attack and subsequent fine-tuning, we are able to significantly improve the model performance on hallucination benchmarks (e.g. TruthfulQA and HalluQA). Moreover, we also find that the LLMs fine-tuned with counterfactual instruction data can also achieve gains on public general benchmarks like C-Eval, MMLU and GSM8K, which also demonstrate the effectiveness of our approach on hallucination mitigation.

1 Introduction

Recently, the advent of large language models (LLMs) has shown unprecedented levels of performance across a myriad of NLP tasks. These models, such as GPT-4(Achiam et al., 2023), LLaMA(Touvron et al., 2023) and QWen(Bai et al., 2023), etc, trained on extensive corpora, have exhibited remarkable abilities to generate coherent

¹The data we used for fine-tuning is publicly available in <https://github.com/oldstree/halluattack>



Figure 1: Example of hallucination with counterfactual prompt²

and contextually relevant text. However, an emerging concern is the propensity for these models to "hallucinate", producing text that, while fluent, is factually incorrect or entirely fabricated(Ji et al., 2022). This tendency not only undermines the credibility of model outputs but also poses significant risks in applications requiring high levels of accuracy and reliability, such as in financial, medical or legal area.

Due to the importance of understanding the factuality and hallucination of LLMs, there have been substantial research interest from academic community(Liu et al., 2024; Tonmoy et al., 2024; Li et al., 2024; Luo et al., 2024; Huang et al., 2023a; Sun et al., 2024). One of the most common approach to mitigate hallucination of LLMs is Retrieval Augmented Generation(RAG)(Lewis et al., 2020; Guu et al., 2020; Shuster et al., 2021; Shi et al., 2023b; Yu et al., 2022; Luo et al., 2023). This method leverages relevant documents retrieved from an external knowledge source to enhance the generation process. However, introducing an external knowledge base and a complex retrieval system is cost,

²Generated by Qwen1.5-32B-Chat

and actually it doesn't eliminate the intrinsic hallucinations of LLMs themselves. Another common approach to mitigate hallucination of LLMs is to enhance the factual correctness of the training data. A notable example is phi(Gunasekar et al., 2023) which uses a section of "textbook quality" data from the web during the pretraining phase. This kind of approach can only be used when we want to train a LLM from scratch. However, the huge amount the training data and large number of parameters of LLMs presents significant challenges and high costs to retrain a LLM. Knowledge editing(Cao et al., 2021; Yao et al., 2023; Tian et al., 2024) recently attracts research interests from researchers. It fixes factual errors by editing some specific "neurons" in LLMs. While knowledge editing can effectively mitigate the model's knowledge gap to some extent, it doesn't actually teach the model how to use the existing knowledge more accurately.

We observe that although LLMs can memorize a vast range of world knowledge easily, they can also be attacked by counterfactual leading prompts since it's hard to learn how to use these world knowledge accurately³. Figure 1 shows an example. The LLM knows what the longest and second-longest rivers in the United States are. However, it hallucinates with a counterfactual leading prompt. In this paper, we introduce a counterfactual attack framework called HALLUATTACK which generates counterfactual instruction data to mitigate hallucinations. The basic idea is to induce LLMs to hallucinate on the knowledge they have already acquired. Firstly, given a LLM, we use factual prompts to collect its responses. These responses are guaranteed to be factually correct, which can indicate that this LLM has already learnt these knowledge from its training data. Then, given a factual response from the LLM, we use GPT-4 to generate counterfactual questions, which contain facts that conflict with this factual response from the LLM. After that, these counterfactual questions are used to attack the LLM. Those prompts which can make the LLM hallucinate will be used to generate instruction data. We use GPT-4 to generate the outputs of counterfactual prompts given encyclopedia documents as external evidence to guarantee both factuality and knowledge boundary of the outputs. Finally, we validate the instruction data generated by our HALLUAT-

³The knowledge gap due to insufficient data is beyond the scope of this work.

TACK by fine-tuning the attacked LLM. Compared with existing approaches, our approach is lightweight with only simple fine-tuning, but can still improve the intrinsic factuality of the LLMs.

The contributions of this work are threefold:

- We propose a simple yet effective approach called HALLUATTACK to attack LLMs and generate counterfactual prompts which could make these LLMs hallucinate.
- We generate counterfactual instruction data by leveraging GPT-4 with encyclopedia documents as additional evidence. This instruction data can be further used to fine-tune the LLMs for hallucination mitigation.
- Experimental results on multiple open-source LLMs demonstrate the effectiveness and generalizability of our approach. The improvements on general LLM benchmarks also show the potential of counterfactual prompts on unlocking the LLM's self-reflection ability and better application of acquired world knowledge.

2 Related Work

2.1 Hallucination Detection and Mitigation

While the advancements in large language models(LLMs) have significantly elevated their performance across an array of downstream tasks, the issue of hallucination has emerged as a significant challenge. Hallucination is characterized by the generation of text by LLMs that deviates from the source material or fails to align with factual truthful information. These original texts and factual datasets typically serve as critical components in the training process, or as user-supplied prompts engaging with the LLMs.

(Huang et al., 2023a) proposes that hallucinations principally arise from three areas: the data source, the training phase, and the inferring phase. As a result, to effectively diminish the occurrence of hallucinations in the text generated by LLMs, a multitude of research has ventured into devising strategies for detecting and mitigating these hallucination problems across the aforementioned three areas.

Due to the potential presence of false factual information and biases in the data consumed by LLMs(Navigli et al., 2023), such as outdated or conflicting knowledge, and discrepancies between

160 user prompts and the parametric knowledge in
161 LLMs, hallucinations may occur. In response to
162 this issue, a knowledge editing method was pro-
163 posed by (Yao et al., 2023), which involves modi-
164 fying the parametric knowledge of LLMs through
165 the introduction of a model plug-in which similar
166 to an adapter. Additionally, efforts have been made
167 to mitigate hallucinations in LLMs by introduc-
168 ing high-quality, unbiased data through retrieval
169 enhancement technology by (Lewis et al., 2020),
170 (Guu et al., 2020), (Shi et al., 2023b). By refocus-
171 ing LLMs on this reliable knowledge data, rather
172 than potentially biased parameter knowledge, the
173 hallucination rate of LLMs can be reduced.

174 A well-planned training and alignment strategy
175 can help reduce the generation of LLMs hallucina-
176 tions. A simple and effective hallucination elimi-
177 nation method named ICD (Zhang et al., 2024),
178 which subtracts the output distribution of the in-
179 duced Weak LLMs with hallucination problems
180 from the output distribution of the original LLMs in
181 training phase, thereby eliminating hallucinations
182 to a certain extent. (Lee et al., 2022) introduced
183 a fact-enhanced training method that significantly
184 mitigates hallucination problems caused by differ-
185 ing factual information. Furthermore, in the LLMs
186 alignment phase, (Wei et al., 2023) introduces sim-
187 ple synthetic data in an additional fine-tuning stage
188 to enhance the model’s independence from user
189 opinions, thereby reducing the generation rate of
190 hallucinations.

191 In the reasoning phase of the model, various
192 studies have been conducted to detect and elimi-
193 nate hallucinations. (Li et al., 2023) proposes a
194 polling-based query method called POPE to detect
195 visual object hallucination. (Zhang et al., 2023)
196 introduces a hallucination detection method that
197 does not require the introduction of external knowl-
198 edge. (Manakul et al., 2023) detects hallucination
199 through an idea that if an LLM has knowledge for
200 a concept, sampled responses are likely to be sim-
201 ilar. (Chuang et al., 2024) proposed a decoding
202 strategy to reduce the hallucination of LLMs by
203 comparing the logarithmic difference between the
204 back layer and the front layer projected to the vo-
205 cabulary space to obtain the distribution of the next
206 word. Additionally, (Shi et al., 2023a) introduced
207 context-aware decoding(CAD), which modifies the
208 output distribution by reducing the reliance on prior
209 knowledge, thereby encouraging the attention to
210 overview information.

2.2 Counterfactual Tasks 211

212 Counterfactual tasks in the field of artificial intelli-
213 gence refer to tasks that involve generating, com-
214 prehending, evaluating, and more under counter-
215 factual conditions or assumptions. Counterfactual
216 tasks emphasize inferring potential outcomes and
217 effects by altering certain premises or conditions
218 based on existing facts, which is essential for en-
219 hancing the ability of comprehending and reason-
220 ing effectively. (Xu et al., 2023) proposed a false
221 information detection framework based on coun-
222 terfactual reasoning, which can effectively detect
223 biases in data source. (Ou et al., 2022) proposed
224 a counterfactual-based open-domain dialogue data
225 augmentation architecture called CAPT. (Rao et al.,
226 2021) introduced an attention mechanism based on
227 counterfactual, and evaluated the method on var-
228 ious fine-grained image recognition tasks, all of
229 which showed significant improvements.

230 Furthermore, as LLMs continue to advance,
231 research on counterfactual tasks integrated with
232 LLMs is gaining momentum. (Wu et al., 2024)
233 proposed an evaluation framework based on coun-
234 terfactual tasks variants to explore the capabilities
235 and limitations of LLMs. (Jin et al., 2023) gener-
236 ates an LLMs evaluation benchmark using causal
237 reasoning and counterfactual reasoning. However,
238 there are still many areas where counterfactual re-
239 search on LLMs is not sufficient, especially in the
240 detection and elimination of hallucinations.

3 Approach 241

3.1 Overview 242

243 We now provide an overview of our approach to
244 explain the whole process and how different compo-
245 nents interact with each other. As shown in Figure
246 2, our approach comprises three components⁴:

- 247 • **Factual Response Generation**, which aims
248 to collect the learnt knowledge of a LLM.
- 249 • **Counterfactual Prompt Generation**, it aims
250 to collect counterfactual prompts which can
251 make the LLM hallucinate based on the fac-
252 tual responses.
- 253 • **Counterfactual Instruction Generation**,
254 which aims to generate instruction data given
255 the counterfactual prompts for LLM fine-
256 tuning.

⁴All the prompts we used can be found in <https://github.com/oldstree/halluattack>

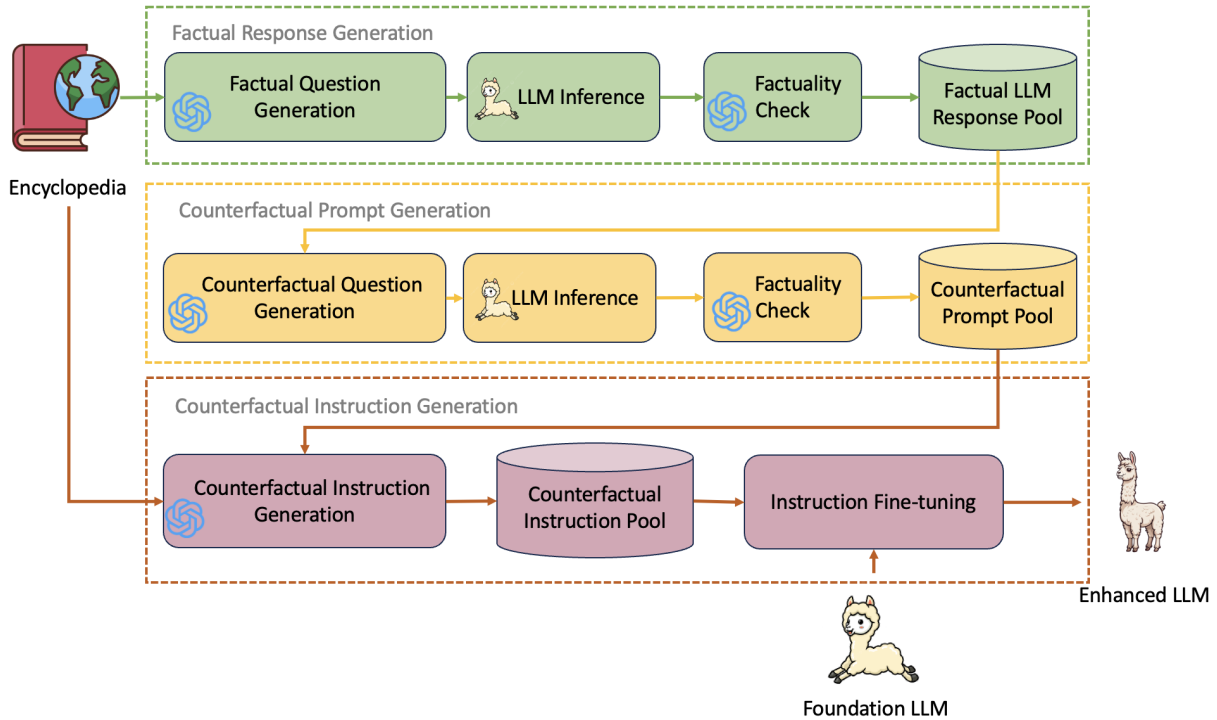


Figure 2: Overview of HALLUATTACK, comprising (1) Factual Response Generation, (2) Counterfactual Prompt Generation, and (3) Counterfactual Instruction Generation.

3.2 Factual Response Generation

There are many factors contributing to hallucinations in LLMs. As mentioned in Section 1, the hallucination factor of knowledge gaps due to insufficient data is beyond the scope of this work. So the first step of our approach is to know what the LLMs know. Based on this, we can then attack the LLMs, causing them to generate hallucination responses based on the knowledge they should have already mastered.

Firstly, we use GPT-4 to generate factual questions $FQ = \{fq_1, fq_2, \dots, fq_k\}$ based on the provided encyclopedia document d_i (k factual questions for each encyclopedia document.). This step will guarantee that: a). The generated questions are knowledge-intensive, requiring the LLMs to answer using the knowledge they have learned. b). The generated questions come with background knowledge (the encyclopedia document) that can be used to help verify the correctness of the LLM’s responses. c). When the LLMs answer incorrectly, the background knowledge can be utilized to generate factually correct responses.

Then, the LLM generates responses given the factual questions, and factuality check step using GPT-4 is applied to filter those factually correct responses $FR = \{fr_1, fr_2, \dots, fr_m\}$. The encyclo-

pedia text will be used as background knowledge for factuality check.

Example 1: Given an encyclopedia document of "List of rivers of the Americas"⁵:

The Missouri River is the longest river in North America and the United States (2,341 mi (3,767 km)). The second longest river in North America and the United States is the Mississippi River (2,320 mi (3,730 km)).

We will generate factual questions like "What are the longest and second longest rivers in the United States?". One of the possible factual answer for this question is:

The longest river in the United States is the Missouri River, which is approximately ... The second longest river in the United States is the Mississippi River, which is approximately ...

3.3 Counterfactual Prompt Generation

Given the factual responses, the main purpose of *Counterfactual Prompt Generation* is to find

⁵https://en.wikipedia.org/wiki/List_of_rivers_of_the_Americas

prompts which can attack the LLM and make it hallucinate. Similar as *Factual Response Generation*, given a factual response fr_j , and its corresponding encyclopedia text d_i , we use GPT-4 to generate counterfactual questions $CFQ = \{cfq_1, cfq_2, \dots, cfq_k\}$ which contain conflict fact with the provided factual response. Then, these prompts will be used to attack the LLM. We check the factual correctness of the responses of these counterfactual questions by GPT-4 with the factual response fr_j and its corresponding encyclopedia text d_i as the background knowledge. Those prompts which can successfully make the LLM hallucinate will be left for next step.

Example 2: Given the factual response in *Example 1*, we can generate counterfactual questions like "As the second longest river in the United States, which cities does the Missouri River flow through?"

The LLM hallucinates on this question with responses like :

The Missouri River, which is the second longest river in the United States after the Mississippi River, flows through...

3.4 Counterfactual Instruction Generation

Given a counterfactual prompt cfq_i , we need to generate high quality instruction data for further model fine-tuning. The instruction data should accurately identify the counterfactual errors in the prompts and should be as free of hallucinations as possible.

Instead of directly using super LLM's (e.g. GPT-4) responses as instruction data, given a counterfactual prompt, we incorporate its corresponding encyclopedia text d_i as the background knowledge to generate high quality responses using GPT-4. So we can minimize the hallucination of GPT-4 itself, thereby increasing the accuracy of the responses.

Example 3: Given the above counterfactual question in *Example 2*, the correct answer should be like:

Your question might be incorrect. The longest river in the United States is the Missouri River, which spans about 2,341 miles. The second longest river is the Mississippi River, which is approximately 2,320 miles long.

3.5 Finetuning the LLM on the Counterfactual Instructions

Supervised Fine-tuning is a simple yet effective alignment method. Once the counterfactual instruction generation is done, we simply fine tune the attacked LLM with this data. We use the counterfactual prompts as the input to the LLM and require the model to generate the responses. A standard sequence-to-sequence loss is applied to train the LLM.

4 Experiments

4.1 Experimental Setup

In this section, we describe the data, models, and benchmarks of the experiments.

Corpora We use about 200,000 Chinese encyclopedia documents and generate 3,000 samples for each open source model for instruction tuning. The Chinese encyclopedia entries are sorted according to the popularity rank. Therefore, we can ensure that the encyclopedia documents used are definitely from the head portion and have certainly been utilized by the open-source LLMs.

Evaluation Models We evaluate our approach on several state-of-the-art LLMs, including Qwen1.5-7B-Chat⁶, Qwen1.5-14B-Chat⁷, Baichuan2-13B-Chat⁸ and ChatGLM3-6B-32k⁹.

Benchmark Datasets We select HalluQA¹⁰(Cheng et al., 2023) and TruthfulQA(5-shot)¹¹(Lin et al., 2022) to evaluate the hallucination rate of the LLMs. We use the official evaluation scripts provided. Specifically, MC1 (Single-true) task is used for TruthfulQA.

In order to further evaluate the effectiveness of our approach on improving the LLM's ability of better using learnt knowledge, we also select several general LLM benchmarks including MMLU(Hendrycks et al., 2020), C-Eval(Huang et al., 2023b), GSM8K(Cobbe et al., 2021), BBH (Big Bench Hard)(Suzgun et al., 2022). We use OpenCompass¹² to evaluate the LLMs on these benchmarks which provides a comprehensive

⁶<https://huggingface.co/Qwen/Qwen1.5-7B-Chat>

⁷<https://huggingface.co/Qwen/Qwen1.5-14B-Chat>

⁸<https://huggingface.co/baichuan-inc/Baichuan2-13B-Chat>

⁹<https://huggingface.co/THUDM/chatglm3-6b-32k>

¹⁰<https://github.com/OpenMOSS/HalluQA/tree/main>

¹¹<https://github.com/sylinr1/TruthfulQA/tree/main>

¹²<https://opencompass.org.cn/home>

Table 1: Overall results on benchmarks of open-source model. *Imp.* denotes the improvement.

Model	C-Eval	MMLU	BBH	GSM8K	TruthfulQA	HalluQA
Qwen1.5-7B	68.88	61.50	40.35	55.57	53.85	42.88
+ HALLUATTACK	70.37	62.20	43.71	58.30	55.93	47.55
Imp.	2.16%	1.14%	8.33%	4.91%	3.86%	10.89%
Qwen1.5-14B	76.20	68.32	54.41	68.00	59.48	51.33
+ HALLUATTACK	76.90	68.45	56.46	70.43	60.34	52.22
Imp.	0.92%	0.19%	3.77%	3.57%	1.45%	1.73%
ChatGLM3-6B	52.12	50.79	41.25	24.11	35.98	31.33
+ HALLUATTACK	53.84	51.93	43.17	25.32	36.84	33.33
Imp.	3.30%	2.24%	4.65%	5.02%	2.39%	6.38%
Baichuan2-13B	56.31	59.17	48.78	52.77	45.65	45.77
+ HALLUATTACK	57.02	60.13	51.27	53.93	47.36	46.67
Imp.	1.26%	1.62%	5.10%	2.20%	3.75%	1.97%

benchmarking framework that enables us to systematically evaluate the performance of the LLMs across various tasks and domains.

Implementation Details We use GPT-4¹³ as super LLM annotator in multiple components in our approach. We generate 3 factual questions for each encyclopedia document and 3 counterfactual questions for each factual response.

We use Firefly¹⁴, an open-source LLM fine-tuning framework for supervised fine-tuning of our evaluation models. Specifically, we employed a learning rate of 1e-5, a batch size of 4, and conducted training for ten epochs. Each model is trained on a single node with eight 80G NVIDIA A100 GPUs.

We utilize standard greedy decoding for inference to ensure the reproducibility. The maximum generation length is set to 1024.

4.2 Results

Table 1 presents a detailed comparison of various LLMs’ performances on both general and hallucination benchmarks. Notably, our approach demonstrates a substantial improvement in reducing hallucinations on TruthfulQA and HalluQA. After fine-tuning with instruction data generated by our HALLUATTACK approach, the performance is significantly improved (with increases of up to 10%) compared with original chat models, which demonstrates the effectiveness of our approach in reducing the hallucinations of LLMs. Furthermore, we also observed gains on general LLM benchmarks, par-

ticularly on the BBH and GSM8K. This shows the potential of our counterfactual instruction tuning on unlocking the LLM’s self-reflection ability and better application of acquired world knowledge.

Our approach achieved better performance on Qwen1.5-7B model compared with Qwen1.5-14B model. This phenomenon suggests that our approach is more effective on LLMs with smaller-scale. A plausible explanation is that LLMs with smaller-scale often struggle with robust reasoning capabilities and can hardly have a thorough understanding of knowledge boundaries. Our approach introduces the counterfactual instruction data. The data can detect where the knowledge boundaries of the LLMs are weak through counterfactual attack, and then repairs and enhances the knowledge boundaries in the alignment phase, which can strengthen the world knowledge learnt by LLMs and thus reduce hallucinations.

Furthermore, table 1 shows that our approach yielded much more significant enhancement on HalluQA as opposed to TruthfulQA across most LLMs. This is because our experimental corpus is derived from Chinese encyclopedic sources, offering a wealth of Chinese counterfactual data. Despite this, we still observed improvements on the English evaluate dataset (i.e. TruthfulQA). The phenomenon not only demonstrates the efficiency of our approach in leveraging linguistically and culturally specific datasets, but also shows the potential for hallucination reduction to be transferred across languages.

¹³<https://platform.openai.com/docs/models/>

¹⁴<https://github.com/yangjianxin1/Firefly>

455	4.3 Discussion	
456	Corpora As mentioned before, we focus on improving the LLM’s ability to better use the knowledge they’ve already acquired during pretraining phase. The knowledge gap due to insufficient data is beyond the scope of this work. So we deliberately use encyclopedia data which has already been used in pretraining phase to create counterfactual instruction data. No new knowledge will be introduced in supervised fine-tuning phase. Existing work(Wan et al., 2024) has shown that minimizing the inconsistency between external knowledge present in the alignment data and the intrinsic knowledge embedded within foundation LLMs is important for hallucination mitigation.	506
457		507
458		508
459		509
460		
461		
462		
463		
464		
465		
466		
467		
468		
469		
470	Instruction Generation As mentioned in section3.4, we use the original encyclopedia document as the background knowledge for GPT-4 to generate the output of the counterfactual prompt. This is very important to minimize the hallucination generated by GPT-4. However, this will probably change the generation behavior or style of the attacked LLM, because the output of the instruction data is mostly summarized from the given encyclopedia document, so the diversity and richness of the generated content will decrease. To tackle this challenge, we tried to use the factual response generated by LLM itself as another background knowledge. We hope the output of the counterfactual prompt can, on the one hand, point out the factual errors in the prompt, on the other hand, follow the original generation style as the factual response. However, the performance is not good as current setup in section3.4. After diving into several cases, we found that the quality generated by GPT-4 with two background documents is not very good, GPT-4 sometimes exhibits a mix and repetition of background documents, which may be due to the prompt we used. Moreover, there could be also some factual errors that are not easily detected automatically in the factual responses. If such data were used during the fine-tuning phase, it would actually exacerbate the LLM’s hallucinations. How to improve the data quality and generate style consistent instruction data will be our future work to follow up.	510
471		511
472		512
473		513
474		514
475		515
476		516
477		517
478		518
479		519
480		520
481		521
482		522
483		523
484		524
485		525
486		526
487		527
488		528
489		529
490		530
491		531
492		532
493		533
494		534
495		535
496		536
497		537
498		538
499		539
500		540
501	Combination with other hallucination mitigation methods The proposed approach plays as a "patch" to given LLMs with simple continue fine-tuning. Since the data volume we used for fine-tuning is very small, we didn’t observe catastrophic	
502		
503		
504		
505		
	forgetting during fine-tuning. This implies that our approach can be integrated with existing hallucination mitigation approaches and can also serves as a supplement to them.	506
		507
		508
		509
	5 Conclusion	510
	In this paper, we explore how counterfactual instruction data helps unlock the ability of LLMs to utilize knowledge more accurately, and propose a simple yet effective prompting approach to attack the LLMs and generate high quality counterfactual instruction data for model fine-tuning. Experimental results demonstrate the effectiveness and scalability of our approach in reducing hallucinations.	511
		512
		513
		514
		515
		516
		517
		518
	Limitations	519
	In our approach, we leverage a super LLM, i.e. GPT-4, as annotators. Although the annotation tasks are not very complex (mostly are question generation and answer summarization tasks) and don’t require huge world knowledge, it is still necessary to investigate more advanced approaches to improve the quality and diversity of the generation as mentioned in section4.3.	520
		521
		522
		523
		524
		525
		526
		527
		528
		529
		530
		531
		532
		533
		534
		535
		536
		537
		538
		539
		540
	Ethics Statement	541
	All the data we used in the experiments are publicly available encyclopedia documents, which do not contain privacy information to the best of our knowledge.	542
		543
		544
		545
		546
		547
		548
	References	549
	Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,	550
		551

552	Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report . <i>arXiv preprint arXiv:2303.08774</i> .	607
553		608
554		609
555	Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023. Qwen technical report .	610
556		611
557		612
558		613
559		614
560		615
561		616
562		617
563		618
564		619
565		620
566		621
567		622
568	Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. Editing factual knowledge in language models . In <i>Conference on Empirical Methods in Natural Language Processing</i> .	623
569		624
570		625
571		626
572	Qinyuan Cheng, Tianxiang Sun, Wenwei Zhang, Siyin Wang, Xiangyang Liu, Mozhi Zhang, Junliang He, Mianqiu Huang, Zhangyue Yin, Kai Chen, and Xipeng Qiu. 2023. Evaluating hallucinations in chinese large language models . <i>CoRR</i> , abs/2310.03368.	627
573		628
574		629
575		630
576		631
577	Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James Glass, and Pengcheng He. 2024. Dola: Decoding by contrasting layers improves factuality in large language models . In <i>International conference on machine learning</i> . ICML.	632
578		633
579		634
580		635
581		636
582	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems . <i>arXiv preprint arXiv:2110.14168</i> .	637
583		638
584		639
585		640
586		641
587		642
588	Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio Cesar Teodoro Mendes, Allison Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, Adil Salim, S. Shah, Harkirat Singh Behl, Xin Wang, Sébastien Bubeck, Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and Yuan-Fang Li. 2023. Textbooks are all you need . <i>ArXiv</i> , abs/2306.11644.	643
589		644
590		645
591		646
592		647
593		648
594		649
595		650
596	Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training . In <i>International conference on machine learning</i> , pages 3929–3938. ICML.	651
597		652
598		653
599		654
600	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Xiaodong Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding . <i>ArXiv</i> , abs/2009.03300.	655
601		656
602		657
603		658
604	Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2023a. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions . <i>ArXiv</i> , abs/2311.05232.	659
605		660
606		661
	Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Jiayi Lei, Yao Fu, Maosong Sun, and Junxian He. 2023b. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models . In <i>Advances in Neural Information Processing Systems</i> .	662
		663
		664
		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
		718
		719
		720
		721
		722
		723
		724
		725
		726
		727
		728
		729
		730
		731
		732
		733
		734
		735
		736
		737
		738
		739
		740
		741
		742
		743
		744
		745
		746
		747
		748
		749
		750
		751
		752
		753
		754
		755
		756
		757
		758
		759
		760
		761
		762
		763
		764
		765
		766
		767
		768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780
		781
		782
		783
		784
		785
		786
		787
		788
		789
		790
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800

665	Junliang Luo, Tianyu Li, Di Wu, Michael R. M. Jenkin, Steve Liu, and Gregory Dudek. 2024. Hallucination detection and hallucination mitigation: An investigation . <i>ArXiv</i> , abs/2401.08358.	721
666		722
667		723
668		724
669	Ziyang Luo, Can Xu, Pu Zhao, Xiubo Geng, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2023. Augmented large language models with parametric knowledge guiding . <i>ArXiv</i> , abs/2305.04757.	725
670		726
671		727
672		728
673	Potsawee Manakul, Adian Liusie, and Mark Gales. 2023. SelfCheckGPT: Zero-resource black-box hallucination detection for generative large language models . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 9004–9017, Singapore. Association for Computational Linguistics.	729
674		730
675		731
676		732
677		733
678		734
679		
680	Roberto Navigli, Simone Conia, and Björn Ross. 2023. Biases in large language models: Origins, inventory, and discussion . <i>J. Data and Information Quality</i> , 15(2).	735
681		736
682		737
683		738
684	Jiao Ou, Jinchao Zhang, Yang Feng, and Jie Zhou. 2022. Counterfactual data augmentation via perspective transition for open-domain dialogues . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 1635–1648, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	739
685		740
686		741
687		742
688		743
689		744
690		
691	Yongming Rao, Guangyi Chen, Jiwen Lu, and Jie Zhou. 2021. Counterfactual attention learning for fine-grained visual categorization and re-identification . In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 1025–1034.	745
692		746
693		747
694		748
695		749
696	Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia Tsvetkov, Luke Zettlemoyer, and Scott Wen-tau Yih. 2023a. Trusting your evidence: Hallucinate less with context-aware decoding . <i>arXiv preprint arXiv:2305.14739</i> .	750
697		751
698		752
699		753
700		754
701	Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen tau Yih. 2023b. Replug: Retrieval-augmented black-box language models . <i>ArXiv</i> , abs/2301.12652.	755
702		756
703		757
704		758
705	Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. 2021. Retrieval augmentation reduces hallucination in conversation . In <i>Conference on Empirical Methods in Natural Language Processing</i> .	759
706		760
707		761
708		762
709		763
710	Lichao Sun, Yue Huang, Haoran Wang, Siyuan Wu, Qihui Zhang, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, Xiner Li, Zheng Liu, Yixin Liu, Yijue Wang, Zhikun Zhang, Bhavya Kailkhura, Caiming Xiong, Chaowei Xiao, Chun-Yan Li, Eric P. Xing, Furong Huang, Haodong Liu, Heng Ji, Hongyi Wang, Huan Zhang, Huaxiu Yao, Manolis Kellis, Marinka Zitnik, Meng Jiang, Mohit Bansal, James Zou, Jian Pei, Jian Liu, Jianfeng Gao, Jiawei Han, Jieyu Zhao, Jiliang Tang, Jindong Wang, John Mitchell, Kai Shu, Kaidi Xu, Kai-Wei Chang, Lifang He, Lifu Huang,	764
711		765
712		766
713		767
714		768
715		769
716		770
717		771
718		772
719		773
720		774
	Michael Backes, Neil Zhenqiang Gong, Philip S. Yu, Pin-Yu Chen, Quanquan Gu, Ran Xu, Rex Ying, Shuiwang Ji, Suman Sekhar Jana, Tian-Xiang Chen, Tianming Liu, Tianying Zhou, William Wang, Xiang Li, Xiang-Yu Zhang, Xiao Wang, Xingyao Xie, Xun Chen, Xuyu Wang, Yan Liu, Yanfang Ye, Yinzhi Cao, and Yue Zhao. 2024. Trustllm: Trustworthiness in large language models . <i>ArXiv</i> , abs/2401.05561.	775
	Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, , and Jason Wei. 2022. Challenging big-bench tasks and whether chain-of-thought can solve them . <i>arXiv preprint arXiv:2210.09261</i> .	776
	Bo Tian, Siyuan Cheng, Xiaozhuan Liang, Ningyu Zhang, Yi Hu, Kouying Xue, Yanjie Gou, Xi Chen, and Huajun Chen. 2024. Instructedit: Instruction-based knowledge editing for large language models . <i>ArXiv</i> , abs/2402.16123.	777
	S. M Towhidul Islam Tonmoy, S M Mehedi Zaman, Vinija Jain, Anku Rani, Vipula Rawte, Aman Chadha, and Amitava Das. 2024. A comprehensive survey of hallucination mitigation techniques in large language models . <i>ArXiv</i> , abs/2401.01313.	778
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models .	779
	Fanqi Wan, Xinting Huang, Leyang Cui, Xiaojun Quan, Wei Bi, and Shuming Shi. 2024. Knowledge verification to nip hallucination in the bud .	779
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Huai hsin Chi, F. Xia, Quoc Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models . <i>ArXiv</i> , abs/2201.11903.	779
	Jerry Wei, Da Huang, Yifeng Lu, Denny Zhou, and Quoc V Le. 2023. Simple synthetic data reduces sycophancy in large language models . <i>arXiv preprint arXiv:2308.03958</i> .	779

780 Zhaofeng Wu, Linlu Qiu, Alexis Ross, Ekin Akyürek,
781 Boyuan Chen, Bailin Wang, Najoung Kim, Jacob An-
782 dreas, and Yoon Kim. 2024. [Reasoning or reciting?](#)
783 [exploring the capabilities and limitations of language](#)
784 [models through counterfactual tasks](#). In *Proceed-*
785 *ings of the 2024 Conference of the North American*
786 *Chapter of the Association for Computational Lin-*
787 *guistics: Human Language Technologies (Volume*
788 *1: Long Papers)*, pages 1819–1862, Mexico City,
789 Mexico. Association for Computational Linguistics.

790 Weizhi Xu, Qiang Liu, Shu Wu, and Liang Wang. 2023.
791 [Counterfactual debiasing for fact verification](#). In
792 *Proceedings of the 61st Annual Meeting of the As-*
793 *sociation for Computational Linguistics (Volume 1:*
794 *Long Papers)*, pages 6777–6789, Toronto, Canada.
795 Association for Computational Linguistics.

796 Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng,
797 Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu
798 Zhang. 2023. [Editing large language models: Prob-](#)
799 [lems, methods, and opportunities](#). In *Proceedings*
800 *of the 2023 Conference on Empirical Methods in*
801 *Natural Language Processing*, pages 10222–10240,
802 Singapore. Association for Computational Linguis-
803 tics.

804 W. Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingx-
805 uan Ju, Soumya Sanyal, Chenguang Zhu, Michael
806 Zeng, and Meng Jiang. 2022. [Generate rather than](#)
807 [retrieve: Large language models are strong context](#)
808 [generators](#). *ArXiv*, abs/2209.10063.

809 Tianhang Zhang, Lin Qiu, Qipeng Guo, Cheng Deng,
810 Yue Zhang, Zheng Zhang, Chenghu Zhou, Xinbing
811 Wang, and Luoyi Fu. 2023. [Enhancing uncertainty-](#)
812 [based hallucination detection with stronger focus](#).
813 In *Proceedings of the 2023 Conference on Empiri-*
814 *cal Methods in Natural Language Processing*, pages
815 915–932, Singapore. Association for Computational
816 Linguistics.

817 Yue Zhang, Leyang Cui, Wei Bi, and Shuming Shi. 2024.
818 [Alleviating hallucinations of large language models](#)
819 [through induced hallucinations](#).