

# DeCoT: Debiasing Chain-of-Thought for Knowledge-Intensive Tasks in Large Language Models via Causal Intervention

Junda Wu<sup>1</sup>, Tong Yu<sup>2</sup>, Xiang Chen<sup>2</sup>, Haoliang Wang<sup>2</sup>  
Ryan A. Rossi<sup>2</sup>, Sungchul Kim<sup>2</sup>, Anup Rao<sup>2</sup>, Julian McAuley<sup>1</sup>

<sup>1</sup> University of California San Diego, La Jolla, USA

<sup>2</sup> Adobe Research, San Jose, USA

juw069@ucsd.edu, {tyu, xiangche, hawang, ryrossi, sukim, anuprao}@adobe.com  
jmcauley@ucsd.edu

## Abstract

Large language models (LLMs) often require task-relevant knowledge to augment their internal knowledge through prompts. However, simply injecting external knowledge into prompts does not guarantee that LLMs can identify and use relevant information in the prompts to conduct chain-of-thought reasoning, especially when the LLM’s internal knowledge is derived from biased information on the pretraining data. In this paper, we propose a novel causal view to formally explain the internal knowledge bias of LLMs via a Structural Causal Model (SCM). We review the chain-of-thought (CoT) prompting from a causal perspective, and discover that the biased information from pretrained models can impair LLMs’ reasoning abilities. When the CoT reasoning paths are misled by irrelevant information from prompts and are logically incorrect, simply editing factual information is insufficient to reach the correct answer. To estimate the confounding effect on CoT reasoning in LLMs, we use external knowledge as an instrumental variable. We further introduce CoT as a mediator to conduct front-door adjustment and generate logically correct CoTs where the spurious correlation between LLMs’ pretrained knowledge and task queries is reduced. With extensive experiments, we validate that our approach enables more accurate CoT reasoning and enhances LLM generation on knowledge-intensive tasks.

## 1 Introduction

For knowledge-intensive tasks (Petroni et al., 2021; Hu et al., 2023; Sun et al., 2023b), specific knowledge is required to obtain an accurate response, which can be out of the distribution of LLMs’ internal knowledge (Yao et al., 2023; Yuan et al., 2023c). Since frequently fine-tuning LLMs can be highly expensive and inefficient (Zhai et al., 2023), the LLM’s internal knowledge can also be outdated and cause knowledge bias problems in LLMs (Zhang et al., 2023b; Wu et al., 2023; Zhang et al., 2023c).

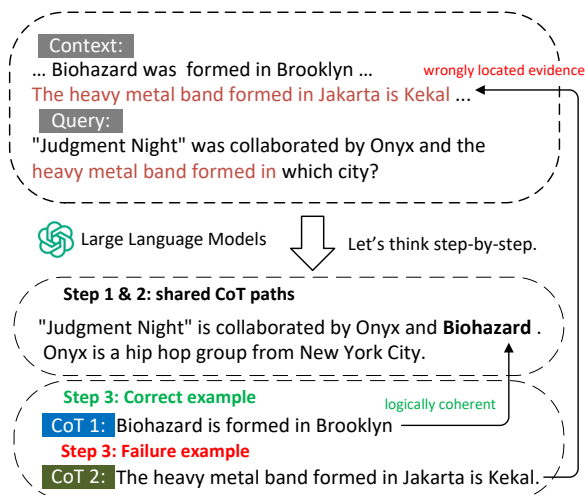


Figure 1: LLMs’ internal knowledge bias can trigger the usage of irrelevant information in prompts, generate incoherent reasoning chains, and impair the model’s logical reasoning ability. This example is derived from the experiments by GPT3.5 on HotpotQA where the context is the one provided in the dataset. Please note that ‘The heavy metal band formed in Jakarta is Kekal’ refers to a heavy metal band that is different from Biohazard. However, GPT3.5 incorrectly assumes that ‘The heavy metal band’ refers to Kekal, and provides incorrect information in step 3 of the CoT.

To efficiently incorporate external knowledge (*i.e.*, context), methods are proposed to retrieve task-relevant language evidence (Liu et al., 2023; He et al., 2022; Zhu et al., 2023b; Shao et al., 2023; Trivedi et al., 2022a). Additionally, external knowledge bases can also directly augment and edit the knowledge-injected prompts (Wen et al., 2023; Sun et al., 2023a; Baek et al., 2023; Zhao et al., 2023b). However, simply injecting external knowledge in prompts does not guarantee that LLMs can identify and use relevant information in the prompts (Shi et al., 2023a; Weston and Sukhbaatar, 2023), especially when the LLM learns biased information in pretraining data (Zhang et al., 2023b). The knowledge bias in LLMs can further cause knowledge conflict or misunderstanding between exter-

nal knowledge and the model’s internal knowledge (Mallen et al., 2023; Wang et al., 2023f,a). In such cases, LLMs may use irrelevant information from prompts and generate incorrect and unexpected responses (Li et al., 2023c; Xie et al., 2023).

When the LLM relies on chain-of-thought (CoT) reasoning for complex tasks, the biased information from pretrained models further impairs LLMs’ reasoning abilities. Many works propose to verify and post-edit the generated reasoning paths before prompt again (Zhao et al., 2023b; Peng et al., 2023; Wang et al., 2023c) to eliminate the factual errors in the generated CoT paths. However, logical reasoning errors can not be easily detected or corrected, as the effectiveness of factual verification and post-editing reasoning chains can be limited to simply injecting more knowledge. For example in Figure 1, given the query (*e.g.*, ““Judgment Night” was collaborated by Onyx and the heavy metal band formed in which city?”) and context which provides external knowledge (Lewis et al., 2020), the LLM may generate logically incorrect CoT (*e.g.*, CoT 2), in which the last chain deviates from the reasoning paths (*e.g.*, instead of the origin of “Bio-hazard”, some arbitrary band mentioned in the context). Such logical incoherence can be caused by the spurious correlation between the query (*e.g.*, the concept “the heavy metal band formed in”) and the LLM’s internal knowledge understanding. Thus, the spurious correlation can lead the LLM to find some arbitrary evidence in the context regardless of its logical connection to the previous chain, as long as it contains the exact phrase. In such cases, factual verification methods cannot detect logical reasoning errors, and the answer can still be incorrect even with the facts verified as correct.

In this work, we propose a novel causal view via a Structural Causal Model (SCM) (Pearl et al., 2016) to formally explain the internal knowledge bias of LLMs. To measure spurious correlation, we propose to use external knowledge as an instrumental variable (Morgan and Winship, 2015) to estimate the Average Causal Effect (ACE) of CoT reasoning paths in LLMs through causal intervention. Based on the measurement of ACE, we can further introduce a CoT sampling method to find the best CoT as a mediator and conduct front-door adjustment (Pearl, 2009). In this approach, the spurious correlation between LLMs’ internal knowledge and task queries can be reduced, which ensures correct CoT reasoning and LLM generation. We summarize our contributions as follows:

- We discover that the bias from LLMs can trigger the usage of irrelevant information in the prompts, and cause the LLM to generate incoherent reasoning chains that impair the model’s reasoning ability.
- To formally understand the bias affecting CoT reasoning abilities, we propose a novel causal view introducing the external knowledge in prompts as an instrumental variable. This causal view uncovers the spurious correlation between queries and LLMs’ internal knowledge understanding.
- To alleviate the bias and ensure correct CoT reasoning, we estimate the average causal effect (ACE) between the CoT and the answer, and further propose a CoT sampling method to conduct the front-door adjustment.
- We conduct multiple experiments on various knowledge-intensive tasks as well as numbers of LLM backbone models, which demonstrates the effectiveness of our method.

## 2 Related Work

**LLMs in Knowledge-intensive Tasks.** In knowledge-intensive tasks (Petroni et al., 2021; Hu et al., 2023; Yang et al., 2018; Welbl et al., 2018) the LLM is asked to respond based on the provided context and its intrinsic knowledge. Retrieval-augmented prompting methods focus on identifying accurate and comprehensive evidence from support documents (Hoshi et al., 2023; Qian et al., 2023), in-context examples (Press et al., 2022; Khattab et al., 2022), knowledge bases (Trivedi et al., 2022a; Xu et al., 2023; Wang et al., 2023c; Feng et al., 2023; Zhu et al., 2023a), knowledge graphs (Wen et al., 2023; Sun et al., 2023a; Salnikov et al., 2023; Zhang et al., 2023a) and human feedback (Zhang et al., 2023b). However, extensive knowledge-injected prompts can introduce irrelevant information to distract LLMs (Shi et al., 2023a; Wang et al., 2023g) and cause LLMs’ unpredictable behaviors (Li et al., 2023c; Chen et al., 2023b). Instead of focusing on how to identify the best knowledge evidence, we investigate how to find logically correct CoT reasoning paths.

**Chain-of-thought Prompting.** Chain-of-thought prompting has shown great potential in explaining LLMs’ thinking process (Yuan et al., 2023a; Li and Du, 2023) and answering multi-hop questions (Wang et al., 2023d; Ma et al., 2022) in complex

reasoning tasks (Fu et al., 2023). However, further works also mention issues of faithfulness and self-consistency in LLMs (Lanham et al., 2023; Turpin et al., 2023). To improve the faithfulness of intermediate chains, several works propose to verify and edit (He et al., 2022; Wang et al., 2023e; Zhao et al., 2023b) the factual errors in unfaithful chains, and (Radhakrishnan et al., 2023; Zhu et al., 2023a) propose to decompose complex questions and answer them individually. We argue that logically incorrect reasoning paths can lead the LLM astray from the right direction of finding the answer, even with the chains factually correct. Similar to previous works, our method first generates some candidate chain-of-thought reasoning paths, which makes our method use similar numbers of API calls (or inference times) per sample.

**Causal Intervention in Language Models.** Causal intervention methods in language models focus on entity-level spurious correlation (Wang et al., 2023b; Zeng et al., 2020; Yang et al., 2023) or sentence-level selection bias (McMilin, 2022). Since LLMs are black-box models (Gat et al., 2023; Cheng et al., 2023), direct methods of causal-aware model reparameterization methods are limited in usage. With causal explainability extracted from the models (Wu et al., 2021; Zhao et al., 2023a), further studies introduce human-in-the-loop debiasing methods (Zhang et al., 2023b; Wu et al., 2022, 2021). However, human effort is normally more expensive, and involving humans in the loop may reduce the efficiency of the system. Instead, our method uses counterfactual context to automatically measure the causal effect and find better CoT to improve the performance of LLMs.

### 3 A Causal View

To understand the causal relationships in knowledge-intensive tasks, we introduce a Structural Causal Model (SCM) (Pearl et al., 2000) and identify the internal knowledge understanding of the LLMs ( $Z$ ) as the confounder. In Figure 2a and Figure 2b, we formulate two types of conventional knowledge injection methods as two SCMs respectively. With the SCMs, we explain the effectiveness of conventional knowledge injection methods and their limitations. We further present the SCM of our method, the debiasing chain-of-thought (DeCoT), in Figure 2c. The formulation of our method, DeCoT is illustrated in Figure 2d and detailed in Section 5.

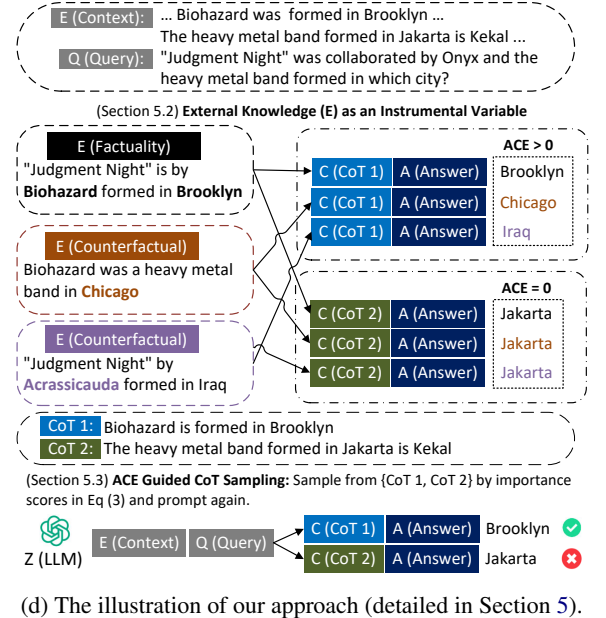
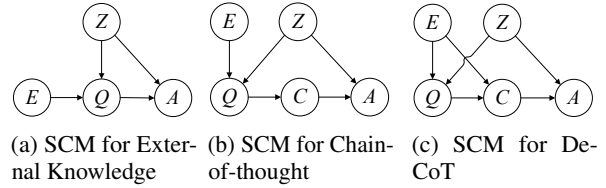


Figure 2: Structural causal graphs for (a) injection of external knowledge, *i.e.*, context (Lewis et al., 2020), (b) chain-of-thought prompting and (c) our approach using external knowledge as an instrumental variable (detailed in Section 5.2). Our proposed debiasing chain-of-thought method DeCoT is illustrated in (d).

#### 3.1 SCM for External Knowledge

In Figure 2a, the causal path  $E \rightarrow Q \leftarrow Z$  represents the knowledge injection process, where  $E$  denotes the external knowledge,  $Q$  denotes queries in the inference stage and  $Z$  denotes the LLM’s internal knowledge. Ideally, the query ( $Q$ ) alleviates spurious correlations as a collider, influenced by the external knowledge ( $E$ ) and LLM’s internal knowledge ( $Z$ ), when  $E$  and  $Z$  are causally irrelevant (Pearl et al., 2000). However, most knowledge injection methods (Baek et al., 2023; Li et al., 2023b) incorporate the external knowledge as the context which is prefixed to the input prompt. Thus, the causal influence of the external knowledge on the query is also determined by the LLM, which makes  $E$  and  $Z$  dependent variables, and the spurious correlation between  $Q$  and  $Z$  remains.

#### 3.2 SCM for Chain-of-thought

Chain-of-thought (as in Figure 2b) is introduced (Li et al., 2023a; Wei et al., 2022; Fu et al., 2023) to make the LLM explain and follow the reasoning

path before giving the final answer. The causal path  $Q \rightarrow C \rightarrow A$  shows that the CoT ( $C$ ) can serve as the mediator between the query ( $Q$ ) and the answer ( $A$ ). However, since the CoT is also prompted from the LLM (Wei et al., 2022; Fu et al., 2023), it can also be causally dependent on LLM’s internal knowledge and thus forms the spurious correlation between  $C$  and  $Z$ . Notably, knowledge editing methods (Zhao et al., 2023b; Li et al., 2023a; Peng et al., 2023) can correct the factual errors in the context and the reasoning paths, while the reasoning logic remains incorrect.

## 4 Preliminaries

### 4.1 Task Formulation

For knowledge-intensive question-answering tasks, the model is prompted with a query  $Q = [q_1, q_2, \dots, q_m]$  and a passage of context  $E = [e_1, e_2, \dots, e_l]$ , *i.e.*, external knowledge (Lewis et al., 2020). Given the query  $Q$  and the context  $E$ , the model  $\theta$  is prompted to recurrently generate the response  $Y$  by sampling from the conditional probability distribution,

$$y_t \sim p_\theta(y|E, Q, y_{<t}).$$

Following the previous setting (Welbl et al., 2017; Trivedi et al., 2022b), we make no assumptions on the context more than what is available in the context of the data samples from their datasets. In this setting, there are irrelevant contexts (Tu et al., 2020; Yang et al., 2018), which lead to the spurious correlation identified in our paper. As illustrated and explained in Figure 2a of Section 3.1, the model directly generates the answer  $A = [a_1, a_2, \dots, a_m]$  without providing the intermediate reasoning process (*i.e.*,  $A = Y$ ).

**Chain-of-thought Prompting.** Following (Wei et al., 2022), we add the additional instruction to ask the model to generate its reasoning paths  $C$  by explaining step-by-step, before generating the final answer  $A$  (*i.e.*,  $Y = [C, A]$ ). By sampling  $N$  different CoTs  $C = [C_1, C_2, \dots, C_N]$  conditioned on the query  $Q$  and the context  $E$ , we can further condition the generation process of the answer  $A$ ,

$$C_i \sim p_\theta(C|E, Q), \quad (1)$$

$$A_{i,r} \sim p_\theta(A|E, Q, C_i). \quad (2)$$

In Eq. (1), since CoTs ( $C$ ) are also generated from the LLM in which the pretrained internal knowledge ( $Z$ ) can also confound on the generation process. As explained in Section 3.2, the confounding

effect can not only affect the factual accuracy of the generated CoTs but also lead to incorrect reasoning logic. Thus knowledge editing and verification methods (Zhao et al., 2023b; Li et al., 2023a; Peng et al., 2023) which solve the former problem, are limited in correcting logical errors in CoTs.

## 5 DeCoT: Debiasing Chain-of-thought

### 5.1 SCM for DeCoT

In Figure 2c, the CoT ( $C$ ) is a mediator between the query ( $Q$ ) and the answer ( $A$ ). Based on the front-door criterion (Pearl et al., 2000), a mediator ( $C$ ) should be causally independent of the confounder ( $Z$ ), to enable front-door adjustment. However, in practice, CoTs are also generated by LLMs, which suggests potential spurious correlations between CoTs ( $C$ ) and the LLM’s internal knowledge ( $Z$ ). Thus, to track the bias from the unobserved confounder  $Z$ , we introduce the external knowledge as an instrumental variable (IV) (Kawakami et al., 2023; Kilbertus et al., 2020; Yuan et al., 2023b). By changing the value of the instrumental variable  $E$  (*i.e.*, the external knowledge), we estimate the true causal relationship between  $C$  and  $A$  (Yuan et al., 2023b). For example, in Figure 2d, two pieces of counterfactual external knowledge (*e.g.* “*Biohazard formed in Chicago*” and “*Judgment Night was by Acrassicauda formed in Iraq*”) are introduced in the same example of Figure 1, where the average causal effect (ACE) is calculated as in (4). Due to the spurious correlation in the third chain of thoughts of “CoT 2”, responses generated from “CoT 2” remain unchanged (*i.e.*,  $ACE = 0$ ), while responses generated from the correct reasoning path “CoT 1” change corresponding to counterfactual evidence (*i.e.*,  $ACE > 0$ ).

### 5.2 External Knowledge as an Instrumental Variable

We model the external knowledge as an instrumental variable  $E$  to understand the causal relationship between the CoT  $C$  and the answer  $A$  (Yuan et al., 2023b). Due to the limitation of directly controlling the generation process of CoTs, we perform the causal treatment by including counterfactual knowledge through the instrumental variable  $E$ . Specifically, we query the LLM to extract  $T$  factual entities  $V = [v_1, v_2, \dots, v_T]$  which correspond to  $T$  counterfactual context  $E_1^*, E_2^*, \dots, E_T^*$  (prompt design explained in Appendix B). In each sample

$$E_j^* = [e_1, e_2, \dots, v_j, \dots, e_l],$$

the corresponding factual entity  $v_j$  is to be replaced by counterfactual entities. Then, the LLM is further prompted to propose  $P$  counterfactual entities  $V_j^* = [v_{j,1}^*, v_{j,2}^*, \dots, v_{j,P}^*]$  to each extracted entity  $v_j \in V$  (prompt design explained in Appendix B).  $P$  counterfactual context samples

$$E_{j,k}^*(v_{j,k}^*) = [e_1, e_2, \dots, v_{j,k}^*, \dots, e_l], k \leq P \quad (3)$$

are constructed, by replacing the corresponding factual entity  $v_j$  in each sample  $E_j^*$ . In this approach, we estimate the average causal effect (ACE) corresponding to each CoT reasoning path  $C_i$ ,

$$\begin{aligned} \text{ACE}(C_i, v_j) &= \mathbb{E}(A|do(E), Q, C_i) - \quad (4) \\ &\quad \mathbb{E}(A|do(E_j^*), Q, C_i) \\ &= \mathbb{E}_{v_{j,k}^* \in V_j^*} [p_\theta(A|E, Q, C_i) - \\ &\quad p_\theta(A|E_{j,k}^*(v_{j,k}^*), Q, C_i)], \end{aligned}$$

in which the average causal effect measures the decreased confidence in the answer (measured as in (2)) with counterfactual context as the evidence. We observe that the average causal effect of different factual entities can vary to the context, queries, and CoTs, which we further conduct analysis experiments in Section 6.5. To consider the overall causal effect of the external context on each CoT, we propose to measure the average causal effect of all the intervened entities,

$$\text{ACE}(C_i) = \mathbb{E}_{v_j \in V} \text{ACE}(C_i, v_j), \quad (5)$$

where the intervened entities  $v_j$  are sampled from a uniform distribution of the external context  $E$ .

### 5.3 Average Causal Effect Guided Chain-of-thought Sampling

With the measured ACE scores, we develop an efficient sampling approach to obtain high-quality CoTs with more coherent reasoning chains, without additional deconfounding layers (Zhang et al., 2020; Wu et al., 2022) for LLM finetuning. Since LLMs are black-box models (Gat et al., 2023; Cheng et al., 2023), direct causal intervention methods on the parameterization of the input query  $Q$  and the context  $E$  are limited. Thus, we propose to use the sampled CoTs  $C$  as the mediator variable to conduct the front-door adjustment.

Based on the measured average causal effect (ACE), we construct importance scores in terms of how the final answer  $A$  reacts to different CoTs  $C$  intervened by the context  $E$ ,

$$C^* \sim \text{softmax} [p_\theta(C_i|E, Q) \cdot \text{ACE}(C_i)], \quad (6)$$

and the front-door adjustment can be realized by introducing the mediator  $C^*$  sampled with the largest average causal effect in the reasoning path,

$$\begin{aligned} A^* &\sim P(A|E, Q, do(C)) \quad (7) \\ &\propto p_\theta(A|E, Q, C^*). \end{aligned}$$

The causal effect on the sampled answer  $A^*$  is mediated by the sampled CoT reasoning path  $C^*$ , whose mediator-outcome confounding effect is controlled and alleviated. We summarize our algorithm DeCoT in Algorithm 1 (Appendix D).

## 6 Experiments

### 6.1 Dataset and Evaluation

Knowledge-intensive tasks commonly require each question to be paired with a paragraph of context as support evidence (Li et al., 2023a; Zhu et al., 2023b; Su et al., 2023; Jang et al., 2023). Our method focuses on the reasoning errors caused by the spurious correlation from LLM’s pretrained knowledge. We follow the setting (Welbl et al., 2017, 2018; Yang et al., 2018; Trivedi et al., 2022b; Shi et al., 2023b; Wei et al., 2022) and use the original supporting contexts from the datasets without any additional data processing. We follow the evaluation protocols in (Yang et al., 2018) and conduct our experiments on datasets as follows:

**HotpotQA** (Yang et al., 2018) contains questions that require multi-step reasoning over multiple support contexts. For each question, support documents are provided in the dataset, which are used as the context in our experiments.

**MuSiQue** (Trivedi et al., 2022b) is another multi-step question answering dataset. Similar to previous work (Ramesh et al., 2023), we conduct our experiment on the challenging part of the dataset, in which questions are annotated as  $\geq 4$  hops.

**WikiHop** (Welbl et al., 2018) is a multi-choice multi-hop reasoning dataset. We use the queries in the dataset as the questions (Tu et al., 2019) in our setting. For baselines and our method, the models are prompted to generate the answers free-form instead of retrieving them from the candidate list.

**SciQ** (Welbl et al., 2017) is a domain-specific question-answering task that contains only scientific questions. We evaluate baselines and our method on test samples with supporting evidence.

For datasets that lack test labels, we follow the same evaluation protocol as (Press et al., 2022; Shao et al., 2023; Chen et al., 2023a) and use the development sets as our test set. We use the Exact

Model	Decoding	HotpotQA		MuSiQue		SciQ		WikiHop		Average	
		EM↑	F1↑	EM↑	F1↑	EM↑	F1↑	EM↑	F1↑	EM↑	F1↑
Flan-T5	CoT w/o ctx	7.41	17.99	2.57	8.50	11.09	17.80	4.12	6.88	6.30	12.79
	CoT	9.48	23.70	19.53	27.61	51.75	63.79	15.02	21.79	23.95	34.22
	CAD	9.65	24.77	<b>20.56</b>	28.57	59.69	69.94	17.28	24.31	26.80	36.90
	DeCoT	<b>11.72</b>	<b>28.70</b>	<b>20.56</b>	<b>30.54</b>	<b>63.55</b>	<b>75.64</b>	<b>22.34</b>	<b>28.41</b>	<b>29.54</b>	<b>40.82</b>
LlaMA-2	CoT w/o ctx	1.67	3.04	0.56	1.44	4.08	5.45	1.19	1.64	1.88	2.89
	CoT	8.86	26.79	20.22	27.46	30.64	39.59	23.10	28.23	20.71	30.52
	CAD	<b>10.53</b>	30.98	21.62	28.10	33.93	41.35	23.50	29.81	22.40	32.56
	DeCoT	10.03	<b>31.48</b>	<b>22.75</b>	<b>30.99</b>	<b>48.58</b>	<b>57.95</b>	<b>27.35</b>	<b>34.46</b>	<b>27.18</b>	<b>38.72</b>
GPT-3.5	CoT w/o ctx	5.60	30.97	2.09	7.96	29.82	43.18	11.62	19.31	12.28	25.36
	CoT	5.10	32.55	22.30	34.22	54.53	68.50	25.40	35.25	26.83	42.63
	CAD	5.43	35.24	24.14	36.81	57.74	70.73	28.37	38.45	28.92	45.31
	DeCoT	<b>10.21</b>	<b>40.19</b>	<b>31.28</b>	<b>44.14</b>	<b>64.61</b>	<b>78.10</b>	<b>31.89</b>	<b>43.45</b>	<b>34.50</b>	<b>51.47</b>

Table 1: The comparison results of DeCoT based on different backbone LLMs on four knowledge-intensive tasks. The best results for each backbone model and each dataset are highlighted in a **bold font**.

Match (EM) and F1 proposed in (Yang et al., 2018) as our evaluation metrics.

## 6.2 Baseline and Backbone Model

Following (Shi et al., 2023b; Su et al., 2023; Trivedi et al., 2022a), we have applied our method to different pretrained LLMs: Flan-T5-XXL (Chung et al., 2022) which has 11B model parameters, LLaMA-2-7B (Touvron et al., 2023) and GPT-3.5 Turbo (Brown et al., 2020). For LLaMA-2-7B model, we choose the finetuned versions from human feedback (Christiano et al., 2017), which can generally yield more stable chain-of-thought reasoning paths.

For baselines, we compare our method with a conventional chain-of-thought prompting method (CoT) (Wei et al., 2022) and context-aware contrastive decoding method CAD (Shi et al., 2023b). We also include the baseline which devises conventional chain-of-thought prompting methods without context (CoT w/o ctx) (Wei et al., 2022), to investigate the effect of context in different datasets. Implementation is detailed in Appendix A. Notably, compared to the baselines (Shi et al., 2023b; Wei et al., 2022), our approach *does not rely on any additional assumptions* and does not require any further processing on the data, to get the performance improvement over these baselines. That is, the inputs of our approach are the same as the baselines (Shi et al., 2023b; Wei et al., 2022).

## 6.3 Main Results

Table 1 presents evaluation results on the four datasets with three LLM backbone models.

**Comparison with Baselines.** As we expected, for

all LLMs the performance is significantly lower when the context of supporting evidence is absent. Because of the poor performance of the direct prompting method, the context-aware contrastive decoding (CAD) baseline can use its answer distribution as the negative penalty on the positive distribution which is obtained by prompting with both the query and context. However, since the negative answer is not supported by either internal or external knowledge, it can have a more random distribution and limits the effectiveness of contrastive decoding methods. On the other hand, our method DeCoT achieves generally higher improvements on regular CoT compared with CAD by detecting logically incorrect CoTs and penalizing them. Instead of simply contrasting the distributions of positive and negative answers, we use counterfactual context to examine the answer distribution changes, which provides a more fine-grained measurement of the causal effect on LLMs’ internal knowledge bias. The consistent performance improvements suggest DeCoT can more accurately detect incorrect CoTs by the measurement, and perform targeted causal intervention.

### Logical Reasoning Performance Understanding.

We also observe that DeCoT gains relatively better F1 improvements on the SciQ dataset, which reach 18.58%, 46.37% and 14.01% for Flan-T5-xxl, LLaMA-2-7B, and GPT-3.5 respectively. It suggests that accurate logical reasoning paths are more strictly required for scientific questions, and the correctness of the generated CoTs is more crucial. Thus, DeCoT’s better performance on the SciQ dataset suggests DeCoT is more effective in

debiasing LLMs’ logical reasoning ability.

#### 6.4 Improving ReAct by DeCoT

Since our main purpose is to find potential spurious correlation and correct reasoning errors with counterfactual debiasing, our approach is compatible with and can improve other model reasoning method variant. To evaluate the generalizability of DeCoT, we apply the proposed counterfactual reasoning method on ReAct (Yao et al., 2022). Similar to (Wei et al., 2022), we use two knowledge-intensive tasks, HotpotQA (Yang et al., 2018), a question-answering task and FEVER (Thorne et al., 2018), fact verification task. In the evaluation, instead of using context from datasets, ReAct can generate a sequence of reasoning paths with knowledge retrieval queries to find relevant information from external knowledge bases (Yao et al., 2022). From Table 2, we observe that by incorporating ReAct generated context as reasoning evidence, DeCoT further improves ReAct’s performance by counterfactual reasoning.

Method	HotpotQA	FEVER
	EM $\uparrow$	Acc $\uparrow$
ReAct	26.53	62.63
CoT (w/ ReAct)	20.40	56.12
CAD (w/ ReAct)	27.55	54.02
DeCoT (w/ ReAct)	<b>36.73</b>	<b>66.32</b>

Table 2: Performance of DeCoT and baselines using retrieved knowledge and reasoning paths from ReAct. Following (Yang et al., 2018; Thorne et al., 2018; Yao et al., 2022), we adopt the EM metric for HotpotQA and the Acc metric for FEVER in the evaluations.

#### 6.5 Impact of the Selected Entities

We evaluate the impact of the number of the selected entities on the MuSiQue dataset based on the backbone model GPT-3.5. Since annotations in MuSiQue guarantee the minimal number of chains of thought is 4 hops (Ramesh et al., 2023), more factual evidence is required to support the final answer, which makes the impact of the selected entities higher in this case.

To illustrate the trend, we only conduct experiments on a number of the selected entities  $T$  with these representative values considering the expensive GPT API costs. In Figure 3a, we show the F1 and EM performance w.r.t. the different number of selected entities  $T = 0, 1, 3, 5$ . We include the result of  $T = 0$  which indicates the regular CoT prompting method. With a larger  $T$ , it has

a higher probability for DeCoT to find more important entities for causal intervention. However, more causal intervention experiments require more counterfactual prompting, which is at the expense of more API calls or inference time. We observe that we can accurately find good factual entities by selecting the most popular entities.

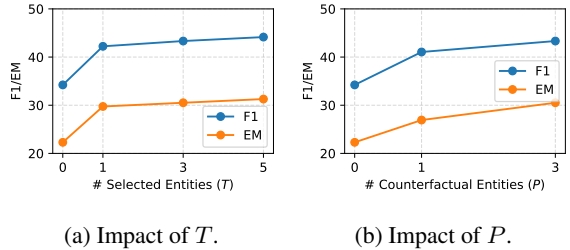


Figure 3: Sensitivity studies on the impact of the number of the selected entities and alternative entities. We also include the  $T = 0$  and  $P = 0$  data points indicating the performance of regular CoT prompting methods. The experiments are conducted on the GPT-3.5 backbone model on the MuSiQue dataset.

#### 6.6 Impact of the Alternative Entities

In Figure 3b, we also show the performance w.r.t. the different number of alternative entities  $P = 0, 1, 3$ . For the same reason of computational cost, we only choose a small number of alternative entities. Since the alternative entities are counterfactual, normally the negative causal effect should be similar among the alternative entities. However, practically the alternative entities may further introduce unexpected causal relationships to the context. Thus, by calculating the average causal effect of multiple alternative entities, such unexpected causal effects can be averaged out.

#### 6.7 Case Study

We present case study results to understand the impact of context quality on the generated CoTs (more details of qualitative examples are in Appendix E). In Table 3, we show the CoT reasoning paths prompted by DeCoT and the following methods:

- **CoTs w/ Context:** regular CoT prompting with context (Wei et al., 2022) (*i.e.*, CoT in Section 6.2.).
- **CoTs w/o Context:** regular CoT prompting with only task queries (*i.e.*, CoT w/o ctx in Section 6.2.).
- **CoT w/ Post-Editing:** CoT with knowledge post-edit (Zhao et al., 2023b), when factual contradictions exist in the context.

Dataset	Question	CoTs w/ Context	CoTs w/o Context	CoT w/ Post-Editing	DeCoT
HotpotQA	The Pluralist school is said to have included what creator of the theory of atomism?	The Pluralist school was a school of pre-Socratic philosophers. <b>The Pluralist school included Leucippus and Democritus.</b> Therefore, the final answer is [Democritus $\times$ ].	One of the key pluralist figures was Greek philosopher Aristotle. Aristotle is known for <b>his theory that nature contains no vacuums.</b> Therefore, the final answer is [Aristotle $\times$ ].	One of the key pluralist figures was Greek philosopher Aristotle. Aristotle is known for <b>his theory of atomism.</b> Therefore, the final answer is [Aristotle $\times$ ].	The Pluralist school is a philosophical movement that emerged in Ancient Greece. It is said to have included Leucippus, <b>the creator of the theory of atomism.</b> Therefore, the final answer is [Leucippus $\checkmark$ ].
MuSiQue	What county contains the city with a radio station that broadcasts to the capital city of the state where the Peace center is located?	The state where the Peace center is located is Greenville, South Carolina. The county containing the city with a radio station broadcasting to Greenville is likely <b>Greenville County.</b> Therefore, the final answer is [Greenville County $\times$ ].	The Peace Center is located in Greenville, South Carolina. The capital city of South Carolina is Columbia. There are two radio stations in Greenville that broadcast to Columbia - WJMZ-FM and WMYI. <b>WMYI, which is located in Greenville,</b> broadcasts to the capital city of Columbia. Therefore, the final answer is [Greenville County $\times$ ].	The Peace Center is located in Greenville, South Carolina. The capital city of South Carolina is Columbia. There are two radio stations in Greenville that broadcast to Columbia - WJMZ-FM and WMYI. <b>WMYI is located in Hendersonville, North Carolina.</b> Therefore, the final answer is [Hendersonville $\times$ ].	The Peace Center is located in Greenville, South Carolina. The radio station WWNQ is licensed to <b>Forest Acres and broadcasts to Columbia,</b> which is the capital city of South Carolina. <b>Forest Acres is a city in Richland County,</b> South Carolina. Therefore, the final answer is [Richland County $\checkmark$ ].
SciQ	By exploding what the space probes get separated from their launchers?	The context is discussing how space probes use exploding bolts to separate from their launchers. <b>The question is asking what the space probes are separated from during this process.</b> Therefore, the final answer is [their launchers $\times$ ].	When a space probe is launched, it is attached to a rocket or launcher. <b>Once the rocket has reached a certain point in space, it explodes,</b> releasing the space probe. Therefore, the final answer is [the rocket or launcher $\times$ ].	When a space probe is launched, it is attached to a rocket or launcher. <b>Once the rocket has reached a certain point in space; it explodes, releasing the space probe.</b> Therefore, the final answer is [the explosion or detonation $\times$ ].	The context states that space probes may be separated from their launchers by exploding bolts. <b>Exploding bolts are bolts that explode or break apart</b> in order to separate two objects. Therefore, the final answer is [bolts $\checkmark$ ].
WikiHop	What's the genre of Terra Brasilis?	"Genre" refers to a category of artistic work. Terra Brasilis is a famous album by Brazilian musician <b>Milton Nascimento.</b> Therefore, the final answer is [Brazilian music $\times$ ].	"Terra Brasilis" is <b>a song</b> by Brazilian composer <b>Caetano Veloso.</b> The song is from his <b>1978 album "Muito (Dentro Da Estrela Azulada)".</b> Therefore, the final answer is [Brazilian $\times$ ].	"Terra Brasilis" is <b>an album</b> by Brazilian composer <b>Antônio Carlos Jobim.</b> <b>The song is from his 1978 album "Muito (Dentro Da Estrela Azulada)".</b> Therefore, the final answer is [Brazilian $\times$ ].	The album Terra Brasilis was released in 1980 by Brazilian composer Antônio Carlos Jobim. <b>Jobim is known for his contributions to bossa nova music.</b> Therefore, the final answer is [bossa nova $\checkmark$ ].

Table 3: Examples of failure CoTs generated from regular CoT prompting, CoT prompting without context, and knowledge post-edited CoT prompting, as well as DeCoT sampled successful CoTs, from four datasets with the GPT3.5 model. In the examples, we highlight **factual and logical errors** with a red font, while the **correct reasoning evidence** is in a green font. The **edited factuality** is also highlighted with a blue font.

For all the methods, the CoTs are prompted from the backbone GPT-3.5 model.

**Case Study of CoTs with Post-Editing.** We observe that CoTs generated without context are likely to contain incorrect knowledge, which can further mislead the reasoning paths. For example, the directly generated CoTs of the question in the WikiHop dataset say “Terra Brasilis is a song by Caetano Veloso”, which is factually incorrect (highlighted in a red font). Due to this incorrect assumption made from LLMs’ hallucination, the following reasoning paths are misled to talk about irrelevant information (*e.g.*, “Caetano Veloso’s album”), and thus the answer is wrong even with factual edit (highlighted in a blue font).

**Case Study of CoTs with Context.** Compared to CoTs generated without context, we observe that the CoTs prompted with context can be factually more faithful. However, the logical reasoning of these CoTs can still be wrong. For example, the CoTs generated with the context in the HotpotQA dataset correctly locate “Leucippus and Democritus” as the two “Pluralist school” members (highlighted in a red font). However, instead of answering with the one who created the school, the LLM

mistakenly chooses the wrong answer “Democritus”. We highlight the correct reasoning paths in green font to show that the key point in answering this question is by identifying the “creator”.

## 7 Conclusion

In this paper, we formally examine the LLM’s internal knowledge bias and identify it as the confounder by a structural causal model (SCM). We revisit the irrelevant contexts issue and further discover the spurious correlation between the LLM and task queries, which can further impair the LLM’s CoT reasoning abilities. Then, we propose DeCoT, a debiasing chain-of-thought prompting method in knowledge-intensive tasks, which alleviates the spurious correlation and enables the LLM to find more accurate and logically sound responses. Following the previous evaluation setting (Webl et al., 2017; Trivedi et al., 2022b; Shi et al., 2023b; Wei et al., 2022), and using the same inputs as the baselines, extensive experimental results and case studies validate the effectiveness of our method DeCoT.



## 8 Limitations

Since DeCoT is an inference-stage causal intervention method, the improvement on LLMs’ reasoning abilities is attributed to alleviating the bias, but can be limited to the upper bound of the LLM’s capacity. To alleviate the causal effect of knowledge bias on LLMs’ reasoning abilities, future works can incorporate unbiased causal learning methods in the model pretraining or instruction tuning stage, which may enable more robust CoT reasoning. It is also interesting to study the theoretical causal foundation of CoT prompting’s mediator role in LLMs, which can be beneficial to better interpretability of black-box LLMs.

## 9 Ethics Statement

Our study on mitigating bias in Large Language Models (LLMs) recognizes the ethical implications of data-driven biases in AI, specifically addressing how these biases affect reasoning processes. We propose a novel approach to reduce bias impact, emphasizing the responsible and ethical advancement of AI technology. The datasets we used in our experiments are all publicly available. No personal information was gathered from our human participants, and they were not exposed to any harmful model outputs.

## References

Jinheon Baek, Alham Fikri Aji, and Amir Saffari. 2023. Knowledge-augmented language model prompting for zero-shot knowledge graph question answering. *arXiv preprint arXiv:2306.04136*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Mingda Chen, Xilun Chen, and Wen-tau Yih. 2023a. Efficient open domain multi-hop question answering with few-shot data synthesis. *arXiv preprint arXiv:2305.13691*.

Xiang Chen, Duanzheng Song, Honghao Gui, Chengxi Wang, Ningyu Zhang, Fei Huang, Chengfei Lv, Dan Zhang, and Huajun Chen. 2023b. Unveiling the siren’s song: Towards reliable fact-conflicting hallucination detection. *arXiv preprint arXiv:2310.12086*.

Jiale Cheng, Xiao Liu, Kehan Zheng, Pei Ke, Hongning Wang, Yuxiao Dong, Jie Tang, and Minlie Huang. 2023. Black-box prompt optimization: Aligning

large language models without model training. *arXiv preprint arXiv:2311.04155*.

- Paul F Christiano, Jan Leike, Tom Brown, Miljan Maric, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.
- Chao Feng, Xinyu Zhang, and Zichu Fei. 2023. Knowledge solver: Teaching llms to search for domain knowledge from knowledge graphs. *arXiv preprint arXiv:2309.03118*.
- Yao Fu, Litu Ou, Mingyu Chen, Yuhao Wan, Hao Peng, and Tushar Khot. 2023. Chain-of-thought hub: A continuous effort to measure large language models’ reasoning performance. *arXiv preprint arXiv:2305.17306*.
- Yair Gat, Nitay Calderon, Amir Feder, Alexander Chapanin, Amit Sharma, and Roi Reichart. 2023. Faithful explanations of black-box nlp models using llm-generated counterfactuals. *arXiv preprint arXiv:2310.00603*.
- Hangfeng He, Hongming Zhang, and Dan Roth. 2022. Rethinking with retrieval: Faithful large language model inference. *arXiv preprint arXiv:2301.00303*.
- Yasuto Hoshi, Daisuke Miyashita, Youyang Ng, Kento Tatsuno, Yasuhiro Morioka, Osamu Torii, and Jun Deguchi. 2023. Ralle: A framework for developing and evaluating retrieval-augmented large language models. *arXiv preprint arXiv:2308.10633*.
- Linmei Hu, Zeyi Liu, Ziwang Zhao, Lei Hou, Liqiang Nie, and Juanzi Li. 2023. A survey of knowledge enhanced pre-trained language models. *IEEE Transactions on Knowledge and Data Engineering*.
- Joel Jang, Seungone Kim, Seonghyeon Ye, Doyoung Kim, Lajanugen Logeswaran, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2023. Exploring the benefits of training expert language models over instruction tuning. *arXiv preprint arXiv:2302.03202*.
- Yuta Kawakami, Manabu Kuroki, and Jin Tian. 2023. Instrumental variable estimation of average partial causal effects. In *International Conference on Machine Learning*, pages 16097–16130. PMLR.
- Omar Khattab, Keshav Santhanam, Xiang Lisa Li, David Hall, Percy Liang, Christopher Potts, and Matei Zaharia. 2022. Demonstrate-search-predict: Composing retrieval and language models for knowledge-intensive nlp. *arXiv preprint arXiv:2212.14024*.

- Niki Kilbertus, Matt J Kusner, and Ricardo Silva. 2020. A class of algorithms for general instrumental variable models. *Advances in Neural Information Processing Systems*, 33:20108–20119.
- Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Denison, Danny Hernandez, Dustin Li, Esin Durmus, Evan Hubinger, Jackson Kernion, et al. 2023. Measuring faithfulness in chain-of-thought reasoning. *arXiv preprint arXiv:2307.13702*.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Ruosen Li and Xinya Du. 2023. Leveraging structured information for explainable multi-hop question answering and reasoning. *arXiv preprint arXiv:2311.03734*.
- Xingxuan Li, Ruochen Zhao, Yew Ken Chia, Bosheng Ding, Lidong Bing, Shafiq Joty, and Soujanya Poria. 2023a. Chain of knowledge: A framework for grounding large language models with structured knowledge bases. *arXiv preprint arXiv:2305.13269*.
- Zhifeng Li, Bowei Zou, Yifan Fan, and Yu Hong. 2023b. Ufo: Unified fact obtaining for commonsense question answering. In *2023 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE.
- Zhoubo Li, Ningyu Zhang, Yunzhi Yao, Mengru Wang, Xi Chen, and Huajun Chen. 2023c. Unveiling the pitfalls of knowledge editing for large language models. *arXiv preprint arXiv:2310.02129*.
- Jiongnan Liu, Jiajie Jin, Zihan Wang, Jiehan Cheng, Zhicheng Dou, and Ji-Rong Wen. 2023. Reta-llm: A retrieval-augmented large language model toolkit. *arXiv preprint arXiv:2306.05212*.
- Kaixin Ma, Hao Cheng, Xiaodong Liu, Eric Nyberg, and Jianfeng Gao. 2022. [Open-domain question answering via chain of reasoning over heterogeneous knowledge](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 5360–5374, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9802–9822.
- Emily McMilin. 2022. [Selection bias induced spurious correlations in large language models](#). In *ICML 2022: Workshop on Spurious Correlations, Invariance and Stability*.
- Stephen L Morgan and Christopher Winship. 2015. *Counterfactuals and causal inference*. Cambridge University Press.
- Judea Pearl. 2009. Causal inference in statistics: An overview.
- Judea Pearl, Madelyn Glymour, and Nicholas P Jewell. 2016. *Causal inference in statistics: A primer*. John Wiley & Sons.
- Judea Pearl et al. 2000. Models, reasoning and inference. *Cambridge, UK: CambridgeUniversityPress*, 19(2):3.
- 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. *arXiv preprint arXiv:2302.12813*.
- Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James Thorne, Yacine Jernite, Vladimir Karpukhin, Jean Maillard, et al. 2021. Kilt: a benchmark for knowledge intensive language tasks. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2523–2544.
- Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. 2022. Measuring and narrowing the compositionality gap in language models. *arXiv preprint arXiv:2210.03350*.
- Hongjing Qian, Yutao Zhu, Zhicheng Dou, Haoqi Gu, Xinyu Zhang, Zheng Liu, Ruofei Lai, Zhao Cao, Jian-Yun Nie, and Ji-Rong Wen. 2023. Webbrain: Learning to generate factually correct articles for queries by grounding on large web corpus. *arXiv preprint arXiv:2304.04358*.
- Ansh Radhakrishnan, Karina Nguyen, Anna Chen, Carol Chen, Carson Denison, Danny Hernandez, Esin Durmus, Evan Hubinger, Jackson Kernion, Kamile Lukošiuūtė, et al. 2023. Question decomposition improves the faithfulness of model-generated reasoning. *arXiv preprint arXiv:2307.11768*.
- Gowtham Ramesh, Makesh Sreedhar, and Junjie Hu. 2023. Single sequence prediction over reasoning graphs for multi-hop qa. *arXiv preprint arXiv:2307.00335*.
- Mikhail Salnikov, Hai Le, Prateek Rajput, Irina Nikishina, Pavel Braslavski, Valentin Malykh, and Alexander Panchenko. 2023. Large language models meet knowledge graphs to answer factoid questions. *arXiv preprint arXiv:2310.02166*.
- Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie Huang, Nan Duan, and Weizhu Chen. 2023. Enhancing retrieval-augmented large language models with iterative retrieval-generation synergy. *arXiv preprint arXiv:2305.15294*.

- 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 *International Conference on Machine Learning*, pages 31210–31227. PMLR.
- 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. *arXiv preprint arXiv:2305.14739*.
- Xin Su, Tiep Le, Steven Bethard, and Phillip Howard. 2023. Semi-structured chain-of-thought: Integrating multiple sources of knowledge for improved language model reasoning. *arXiv preprint arXiv:2311.08505*.
- Jiashuo Sun, Chengjin Xu, Lumingyuan Tang, Saizhuo Wang, Chen Lin, Yeyun Gong, Heung-Yeung Shum, and Jian Guo. 2023a. Think-on-graph: Deep and responsible reasoning of large language model with knowledge graph. *arXiv preprint arXiv:2307.07697*.
- Kai Sun, Yifan Ethan Xu, Hanwen Zha, Yue Liu, and Xin Luna Dong. 2023b. Head-to-tail: How knowledgeable are large language models (llm)? aka will llms replace knowledge graphs? *arXiv preprint arXiv:2308.10168*.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. Fever: a large-scale dataset for fact extraction and verification. *arXiv preprint arXiv:1803.05355*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022a. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. *arXiv preprint arXiv:2212.10509*.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022b. Musique: Multi-hop questions via single-hop question composition. *Transactions of the Association for Computational Linguistics*, 10:539–554.
- Ming Tu, Kevin Huang, Guangtao Wang, Jing Huang, Xiaodong He, and Bowen Zhou. 2020. Select, answer and explain: Interpretable multi-hop reading comprehension over multiple documents. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 9073–9080.
- Ming Tu, Guangtao Wang, Jing Huang, Yun Tang, Xiaodong He, and Bowen Zhou. 2019. Multi-hop reading comprehension across multiple documents by reasoning over heterogeneous graphs. *arXiv preprint arXiv:1905.07374*.
- Miles Turpin, Julian Michael, Ethan Perez, and Samuel R Bowman. 2023. Language models don’t always say what they think: Unfaithful explanations in chain-of-thought prompting. *arXiv preprint arXiv:2305.04388*.
- Cunxiang Wang, Xiaoze Liu, Yuanhao Yue, Xian-gru Tang, Tianhang Zhang, Cheng Jiayang, Yunzhi Yao, Wenyang Gao, Xuming Hu, Zehan Qi, et al. 2023a. Survey on factuality in large language models: Knowledge, retrieval and domain-specificity. *arXiv preprint arXiv:2310.07521*.
- Fei Wang, Wenjie Mo, Yiwei Wang, Wenxuan Zhou, and Muhao Chen. 2023b. A causal view of entity bias in (large) language models. *arXiv preprint arXiv:2305.14695*.
- Jianing Wang, Qiushi Sun, Nuo Chen, Xiang Li, and Ming Gao. 2023c. Boosting language models reasoning with chain-of-knowledge prompting. *arXiv preprint arXiv:2306.06427*.
- Jinyuan Wang, Junlong Li, and Hai Zhao. 2023d. Self-prompted chain-of-thought on large language models for open-domain multi-hop reasoning. *arXiv preprint arXiv:2310.13552*.
- Keheng Wang, Feiyu Duan, Sirui Wang, Peiguang Li, Yunsen Xian, Chuantao Yin, Wenge Rong, and Zhang Xiong. 2023e. Knowledge-driven cot: Exploring faithful reasoning in llms for knowledge-intensive question answering. *arXiv preprint arXiv:2308.13259*.
- Yike Wang, Shangbin Feng, Heng Wang, Weijia Shi, Vidhisha Balachandran, Tianxing He, and Yulia Tsvetkov. 2023f. Resolving knowledge conflicts in large language models. *arXiv preprint arXiv:2310.00935*.
- Zhiruo Wang, Jun Araki, Zhengbao Jiang, Md Rizwan Parvez, and Graham Neubig. 2023g. Learning to filter context for retrieval-augmented generation. *arXiv preprint arXiv:2311.08377*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Johannes Welbl, Nelson F Liu, and Matt Gardner. 2017. Crowdsourcing multiple choice science questions. *arXiv preprint arXiv:1707.06209*.
- Johannes Welbl, Pontus Stenetorp, and Sebastian Riedel. 2018. Constructing datasets for multi-hop reading comprehension across documents. *Transactions of the Association for Computational Linguistics*, 6:287–302.
- Yilin Wen, Zifeng Wang, and Jimeng Sun. 2023. Mindmap: Knowledge graph prompting sparks graph of thoughts in large language models. *arXiv preprint arXiv:2308.09729*.

- Jason Weston and Sainbayar Sukhbaatar. 2023. [System 2 attention \(is something you might need too\)](#). *ArXiv*, abs/2311.11829.
- Junda Wu, Rui Wang, Tong Yu, Ruiyi Zhang, Handong Zhao, Shuai Li, Ricardo Henao, and Ani Nenkova. 2022. Context-aware information-theoretic causal de-biasing for interactive sequence labeling. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 3436–3448.
- Junda Wu, Tong Yu, and Shuai Li. 2021. Deconfounded and explainable interactive vision-language retrieval of complex scenes. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 2103–2111.
- Suhang Wu, Minlong Peng, Yue Chen, Jinsong Su, and Mingming Sun. 2023. Eva-kellm: A new benchmark for evaluating knowledge editing of llms. *arXiv preprint arXiv:2308.09954*.
- Jian Xie, Kai Zhang, Jiangjie Chen, Renze Lou, and Yu Su. 2023. Adaptive chameleon or stubborn sloth: Unraveling the behavior of large language models in knowledge conflicts. *arXiv preprint arXiv:2305.13300*.
- Shicheng Xu, Liang Pang, Huawei Shen, Xueqi Cheng, and Tat-seng Chua. 2023. Search-in-the-chain: Towards the accurate, credible and traceable content generation for complex knowledge-intensive tasks. *arXiv preprint arXiv:2304.14732*.
- Zhen Yang, Yongbin Liu, and Chunping Ouyang. 2023. Causal interventions-based few-shot named entity recognition. *arXiv preprint arXiv:2305.01914*.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*.
- Jia-Yu Yao, Kun-Peng Ning, Zhen-Hui Liu, Mu-Nan Ning, and Li Yuan. 2023. Llm lies: Hallucinations are not bugs, but features as adversarial examples. *arXiv preprint arXiv:2310.01469*.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2022. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*.
- Chenhan Yuan, Qianqian Xie, Jimin Huang, and Sophia Ananiadou. 2023a. Back to the future: Towards explainable temporal reasoning with large language models. *arXiv preprint arXiv:2310.01074*.
- Junkun Yuan, Xu Ma, Ruoxuan Xiong, Mingming Gong, Xiangyu Liu, Fei Wu, Lanfen Lin, and Kun Kuang. 2023b. Instrumental variable-driven domain generalization with unobserved confounders. *ACM Transactions on Knowledge Discovery from Data*.
- Lifan Yuan, Yangyi Chen, Ganqu Cui, Hongcheng Gao, Fangyuan Zou, Xingyi Cheng, Heng Ji, Zhiyuan Liu, and Maosong Sun. 2023c. Revisiting out-of-distribution robustness in nlp: Benchmark, analysis, and llms evaluations. *arXiv preprint arXiv:2306.04618*.
- Xiangji Zeng, Yunliang Li, Yuchen Zhai, and Yin Zhang. 2020. Counterfactual generator: A weakly-supervised method for named entity recognition. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7270–7280, Online. Association for Computational Linguistics.
- Yuexiang Zhai, Shengbang Tong, Xiao Li, Mu Cai, Qing Qu, Yong Jae Lee, and Yi Ma. 2023. Investigating the catastrophic forgetting in multimodal large language models. *arXiv preprint arXiv:2309.10313*.
- Peitian Zhang, Shitao Xiao, Zheng Liu, Zhicheng Dou, and Jian-Yun Nie. 2023a. Retrieve anything to augment large language models. *arXiv preprint arXiv:2310.07554*.
- Shengyu Zhang, Tan Jiang, Tan Wang, Kun Kuang, Zhou Zhao, Jianke Zhu, Jin Yu, Hongxia Yang, and Fei Wu. 2020. Devlbart: Learning deconfounded visio-linguistic representations. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 4373–4382.
- Shuo Zhang, Liangming Pan, Junzhou Zhao, and William Yang Wang. 2023b. Mitigating language model hallucination with interactive question-knowledge alignment. *arXiv preprint arXiv:2305.13669*.
- Zihan Zhang, Meng Fang, Ling Chen, Mohammad-Reza Namazi-Rad, and Jun Wang. 2023c. How do large language models capture the ever-changing world knowledge? a review of recent advances. *arXiv preprint arXiv:2310.07343*.
- Ruochen Zhao, Shafiq Joty, Yongjie Wang, and Prathyusha Jwalapuram. 2023a. [Towards causal concepts for explaining language models](#).
- Ruochen Zhao, Xingxuan Li, Shafiq Joty, Chengwei Qin, and Lidong Bing. 2023b. [Verify-and-edit: A knowledge-enhanced chain-of-thought framework](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5823–5840, Toronto, Canada. Association for Computational Linguistics.
- Wang Zhu, Jesse Thomason, and Robin Jia. 2023a. Chain-of-questions training with latent answers for robust multistep question answering. *arXiv preprint arXiv:2305.14901*.
- Yin Zhu, Zhiling Luo, and Gong Cheng. 2023b. Furthest reasoning with plan assessment: Stable reasoning path with retrieval-augmented large language models. *arXiv preprint arXiv:2309.12767*.

## A Implementation Details

To obtain diversified chain-of-thought reasoning paths, we sample  $N = 5$  chains-of-thought with the temperature set to 1.0 for all the backbone models. For **DeCoT**, we let the LLM extract the top  $T = 5$  most frequently appearing entities in the context as to be intervened. The LLM will be further prompted to provide  $P = 3$  alternative counterfactual entities to each of the extracted entities.

As for the open-sourced LLMs (*i.e.*, Flan-T5-XXL and LLaMA-2-7B), we use the official Hugging Face implementations. The experiments are conducted using 4 NVIDIA RTX A6000 GPUs with 48GBs. For GPT-3.5 Turbo, we use the OpenAI API to conduct the experiments.

To prompt the LLMs to generate more robust chain-of-thought results and also follow a unified answer format, we have included 3 few-shot in-context learning examples. The in-context learning examples are from a separate set of data which provides no extra knowledge to the evaluated tasks. In addition, we have also included 3 in-context learning examples for both the entity extraction and the alternative entity proposal prompts. Detailed designs of these in-context examples and prompts are explained in Appendix C.

## B Prompt Design

### B.1 Factual Entity Extraction

To extract the most relevant factual entities  $V$  in the context  $E$  (Section 5.2),

$$v_j \sim p_\theta(V|E, Instruct_{ent}), \quad (8)$$
$$E_j^* = [e_1, e_2, \dots, v_j, \dots, e_l],$$

in which  $Instruct_{ent}$  is the explicit prompt instruction shown in the following.

#### Context Example 1:

The Ritz-Carlton Jakarta is a hotel and skyscraper in Jakarta, Indonesia and 14th Tallest building in Jakarta. It is located in city center of Jakarta, near Mega Kuningan, adjacent to the sister JW Marriott Hotel. It is operated by The Ritz-Carlton Hotel Company. The complex has two towers that comprises a hotel and the Airlangga Apartment respectively. Nakuul Mehta, Kunal Jaisingh and Leenesh Mattoo respectively portray Shivaay, Omkara and Rudra, the three heirs of the Oberoi family.

#### Instruction Example 1:

Extract the top 5 most frequently appeared entities in the context and provide in the format of a list: [Ritz-Carlton, Jakarta, Indonesia, Airlangga Apartment, Nakuul Mehta]

#### Context Example 2:

Lisa Marie Simpson is a fictional character in the animated television series "The Simpsons". She is the middle child and most intelligent of the Simpson family. Voiced by Yeardley Smith, Lisa first appeared on television in "The Tracey Ullman Show" short "Good Night" on April 19, 1987. Cartoonist Matt Groening created and designed her while waiting to meet James L. Brooks. Groening had been invited to pitch a series of shorts based on his comic "Life in Hell", but instead decided to create a new set of characters. He named the elder Simpson daughter after his younger sister Lisa Groening.

#### Instruction Example 2:

Extract the top 5 most frequently appeared entities in the context and provide in the format of a list: [Lisa Marie Simpson, The Simpsons, Yeardley Smith, The Tracey Ullman Show, Lisa Groening]

### B.2 Alternative Entity Proposal

To ask the LLM to propose  $P$  counterfactual entities  $E_{j,1}^*, E_{j,2}^*, \dots, E_{j,P}^*$  to the extracted entity  $v_j$  (Section 5.2),

$$v_{j,k}^* \sim p_\theta(V|v_j, Instruct_{alt}), \quad (9)$$
$$E_{j,k}^*(v_{j,k}^*) = [e_1, e_2, \dots, v_{j,k}^*, \dots, e_l],$$

in which  $Instruct_{alt}$  is the explicit prompt instruction shown in the following.

#### Example 1

Provide 3 most similar entities to "America", which are from the same entity type and have similar meanings. Provide the answers as a list: [Britain, France, Germany].

#### Example 2

Provide 3 most similar entities to "Florida", which are from the same entity type and have similar meanings. Provide the answers as a list: [New York, California, Arizona].

### Example 3

Provide 3 most similar entities to "Beethoven", which are from the same entity type and have similar meanings. Provide the answers as a list: [Mozart, Mahler, Wagner].

## C In-context Learning Examples

We design several CoT prompting examples to guide the LLM to generate in an expected format as following:

### Example 1

Context: The Ritz-Carlton Jakarta is a hotel and skyscraper in Jakarta, Indonesia and 14th Tallest building in Jakarta.

Question: The Oberoi family is part of a hotel company that has a head office in what city?

Based on the context, answer the question step by step and provide the final answer in the end.

1. The Oberoi family is an Indian family known for their involvement in hotel management
2. The Oberoi Group is the hotel company founded by the family

Therefore, the final answer is [Delhi].

### Example 2

Context: Lisa Marie Simpson is a fictional character in the animated television series "The Simpsons".

Question: Musician Allie Goertz wrote a song about the "The Simpsons" character Milhouse, who Matt Groening named after who?

Based on the context, answer the question step by step and provide the final answer in the end.

1. The character Milhouse from the TV show "The Simpsons" was named by Matt Groening
2. Groening named the character after President Richard Nixon's middle name

Therefore, the final answer is [Richard Nixon].

### Example 3

Context: Moloch: or, This Gentile World is a semi-autobiographical novel written by Henry Miller in 1927-28, initially under the guise of a novel written by his wife, June.

Question: What nationality was James Henry Miller's wife?

Based on the context, answer the question step by step and provide the final answer at the end.

1. James Henry Miller (25 January 1915 – 22 October 1989) was an English folk singer and songwriter
2. His second wife was Margaret "Peggy" Seeger (born June 17, 1935), who is an American folksinger

Therefore, the final answer is [American].

## D Pseudo Code for DeCoT

We illustrate our method as the pseudo-code in Algorithm 1.

### Algorithm 1: DeCoT

---

**Input:** Task query  $Q$ , Context  $E$ , the LLM with model parameters  $\theta$

- 1 **Init:** Sample chain-of- thought reasoning paths  $C = [C_1, C_2, \dots, C_N]$  by Eq. 1.
- 2 Extract  $T$  factual entities  $V = [v_1, v_2, \dots, v_T]$  from  $E$  by Eq. 8 ;
- 3 **while**  $i < N$  **do**
- 4     **while**  $j < T$  **do**
- 5         Propose  $P$  counterfactual entities  $\{v_{j,1}^*, v_{j,2}^*, \dots, v_{j,P}^*\}$  by Eq. 9 ;
- 6         **while**  $k < P$  **do**
- 7             Construct counterfactual context  $E_{j,k}^*(v_{j,k}^*)$  by Eq. 3 ;
- 8         **end**
- 9         Estimate  $\text{ACE}(C_i, v_j)$ ,  $i < N$  for each entity  $v_j$  by Eq. 4 ;
- 10     **end**
- 11     Estimate  $\text{ACE}(C_i)$  by Eq. 5 ;
- 12 **end**
- 13 Sample CoT by Eq. 6 ;
- 14 Sample the answer by Eq. 7 ;

---

## E Qualitative Examples

We present more qualitative examples from HotpotQA (Table 4), MuSiQue (Table 5), SciQ (Table 6) and WikiHop (Table 7).

Question	CoTs w/ Context	DeCoT
Which of the movie which has Rich Lee 3-D previsualizations was directed by Francis Lawrence?	The movie in question features the use of Rich Lee's 3-D previsualizations. Francis Lawrence is the director of the movie. Therefore, based on this information, the final answer is [I Am Legend ✗].	Rich Lee created 3-D previsualizations for Hollywood feature films. Rich Lee worked on several Hollywood films including "I Am Legend", and "Constantine". The feature film "Constantine" was directed by Francis Lawrence. Therefore, the final answer is [Constantine ✓].
Which close friend of an American former basketball player was signed for to play for the Timberwolves?	An American former basketball player is Michael Jordan. One of Jordan's close friends is Scottie Pippen. Pippen was signed to play for the Timberwolves in 2003. Therefore, the final answer is [Scottie Pippen ✗].	Kevin Garnett, an American former basketball player, played for the Minnesota Timberwolves. After the death of Malik Sealy, the Timberwolves signed free agent Chauncey Billups, who was a close friend of Garnett's. Therefore, the final answer is [Chauncey Billups ✓].

Table 4: Qualitative examples from the HotpotQA dataset.

Question	CoTs w/ Context	DeCoT
Who has the lowest batting average in the league that the team that has won the most titles from the competition they give out MVP awards plays for?	The team with the most titles from the competition that gives out MVP awards is the New York Yankees. The player with the lowest batting average in the league is Daniel Camarena, who plays for the New York Yankees. Therefore, the final answer is [Daniel Camarena ✗].	The American League team with the most World Series titles is the New York Yankees, with 27 wins. The player with the lowest career batting average in the league is Bill Bergen with a .170 average in 3,028 at-bats. The Yankees play in the American League, which is where the Bill Bergen played. Therefore, the final answer is [Bill Bergen ✓].
Who is the president of the newly declared independent country, that established the Timor Leste Commission of Truth and Friendship, with the country containing the airport that includes Lion Air?	The Timor Leste Commission of Truth and Friendship was established jointly by the governments of Indonesia and East Timor in August 2005. The presidents of both nations received the final report of the commission on July 15, 2008. Indonesian President Susilo Bambang Yudhoyono fully endorsed the report, providing acknowledgement of human rights violations in Timor by state institutions. Therefore, the final answer is [Susilo Bambang Yudhoyono ✗].	Timor Leste is the newly declared independent country, also known as East Timor. The president of Timor Leste is Francisco Guterres, as stated in the context. The Timor Leste Commission of Truth and Friendship was established jointly by the governments of Indonesia and East Timor, indicating that the airport mentioned (Juanda International Airport) is located in Indonesia. Therefore, the final answer is [Francisco Guterres ✓].

Table 5: Qualitative examples from the MuSiQue dataset.

Question	CoTs w/ Context	DeCoT
What acids are the structural components of many lipids and may be saturated or unsaturated?	Fatty acids are the structural components of many lipids. Fatty acids are carboxylic acids. Carboxylic acids may be saturated or unsaturated. Therefore, the final answer is [Carboxylic acids ✗].	Lipids are a category of biomolecules that include fats, oils, and waxes. Fatty acids are the structural components of many lipids. Fatty acids can be categorized as either saturated or unsaturated. Therefore, the named acids are [fatty acids ✓].
What were the first plants to evolve?	The process of evolution involves the development and change of living species over time. According to scientific evidence, the first plants to evolve were algae-like organisms that lived in water. Therefore, the final answer is [algae-like organisms ✗].	Nonvascular plants are a type of plant that do not have vascular tissue. These plants were the first to evolve, meaning they were the first types of plants to appear on Earth. Therefore, the final answer is [Nonvascular plants ✓].

Table 6: Qualitative examples from the SciQ dataset.

Question	CoTs w/ Context	DeCoT
What is the genre of queen of the black coast?	Queen of the Black Coast is a fantasy short story written by Robert E. Howard in the 1930s. Fantasy is a genre of fiction that includes elements of magic and supernatural creatures. Therefore, the final answer is [fantasy ✗].	"Queen of the Black Coast" is a short story written by American author Robert E. Howard. The story is considered a classic in the sword and sorcery genre. The genre of "Queen of the Black Coast" is [sword and sorcery ✓].
What is Mark Goresky's field of work?	Mark Goresky is a Canadian mathematician. He received his Ph.D. from Brown University in 1976. He has taught at the University of British Columbia in Vancouver and Northeastern University. (Lack of information.) Therefore, the final answer is [mathematics ✗].	Mark Goresky is a Canadian mathematician born in 1950 known for his contributions in the field of topology. He is most notable for his work on intersection homology, which he discovered in 1974 with Robert MacPherson. Therefore, the final answer is [topology ✓].

Table 7: Qualitative examples from the WikiHop dataset.