

Dually Self-Improved Counterfactual Data Augmentation Using Large Language Model

Anonymous ACL submission

Abstract

Counterfactual data augmentation, which generates minimally edited tokens to alter labels, has become a key approach to improving model robustness in natural language processing (NLP). It is usually implemented by first identifying the causal terms and then modifying these terms to create counterfactual candidates. The emergence of large language models (LLMs) has effectively facilitated the task of counterfactual data augmentation. However, existing LLM-based approaches still face some challenges in 1) accurately extracting the task-specific causal terms, and 2) the quality of LLM-generated counterfactuals. To address the issues, we propose a dually self-improved counterfactual data augmentation method using LLM for Natural Language Inference (NLI) task. On the one hand, we design a self-improved strategy employing the attention distribution of the task model to identify the task-specific causal terms, which is lightweight and task-specific. On the other hand, a second self-improved strategy based on direct preference optimization (DPO) is utilized to refine LLM-generated counterfactuals, achieving high-quality counterfactuals. Finally, a balanced loss preventing over-emphasis on augmented data is proposed to retrain the task model on the fusion of existing data and generated counterfactuals. Extensive experiments on NLI benchmarks demonstrate the effectiveness of our proposed method in generating high-quality counterfactuals for improving the task performance.

1 Introduction

In the complex realm of machine learning and NLP, imbalance and biases prevalent in real-world training data continue to be an arduous challenge for robust model development. Traditional data augmentation suffers from the issue of spurious association when alleviating these issues (Chen et al., 2021). In recent years, generating counterfactual augmented

data (CAD) (Kaushik et al., 2020), introducing minimal modifications to the data through additions, replacements, or deletions to flip the label, has been widely attempted in many tasks (Liu et al., 2021). Target task models trained with larger-scale counterfactuals can learn better representations and effects of causal terms, which facilitates the task performance improvements and enables robust generalization.

Typically, counterfactual data augmentation involves three steps: (1) identifying important tokens (known as causal terms) that can flip the labels, (2) minimally editing these terms to create counterfactual candidates, and (3) retraining the model on the fusion data of existing data and augmented data. For example, as shown in Figure 1, in NLI task, through modifying the identified causal term "talking to" to "walking with" for the given example, we flip the original label from "Entailment" to "contradiction", obtaining a counterfactual.

However, it is non-trivial to obtain high-quality counterfactuals. Early works (Gardner et al., 2020; Kaushik et al., 2020) relied on human experts to annotate counterfactual examples, which is not easily scalable. Therefore, researchers have been exploring automatic methods for counterfactual generation using neural networks (Chen et al., 2021). Recently, AutoCAD (Wen et al., 2022) has attempted to leverage generative language models, such as T5 (Raffel et al., 2020), for controllable text generation. However, due to the limited comprehension and generation capabilities of these language models, the quality of the generated data remain constrained. The advent of LLMs has driven significant progress across various NLP tasks, researchers have focused on designing effective prompts to leverage the advanced comprehension and generation abilities of LLMs for directly generating desired counterfactuals (Chen et al., 2023; Dixit et al., 2022; Nguyen et al., 2024).

Despite the promising advancements, research on LLM based counterfactual data augmentation

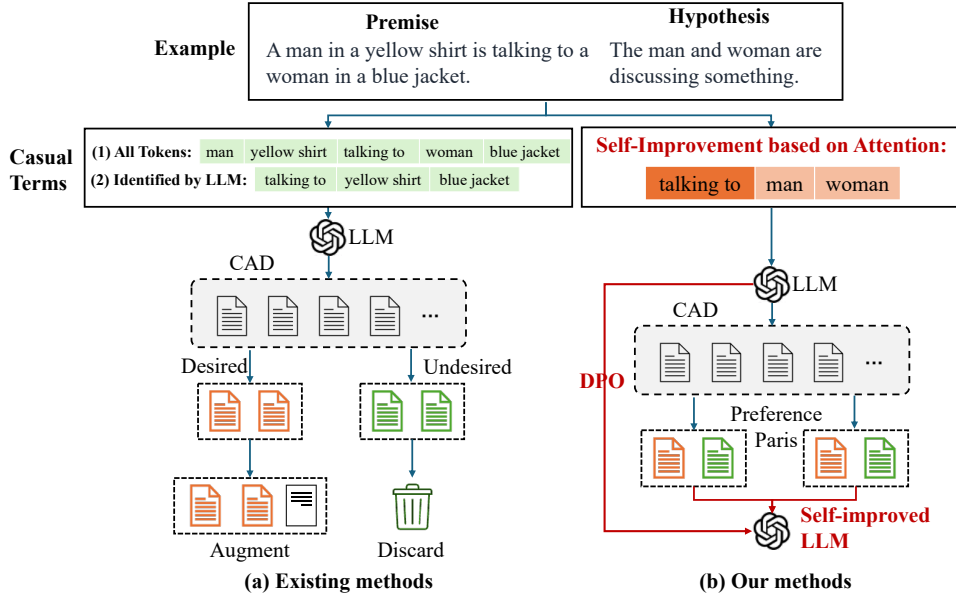


Figure 1: Introduction of Counterfactual Data Augmentation.

still faces several challenges. (1) How to extract causal terms specific to the task accurately? Existing works either exploited all spans obtained through sentence splitting (Chen et al., 2023), or directly prompted LLMs (Li et al., 2024) to identify causal terms. All of these methods suffer from the inaccurate casual terms specific to the task. (2) How to enhance the quality of LLM-generated CAD by modifying the causal terms? Those LLM-based approaches typically employ LLMs to rewrite causal terms and then select the desired counterfacts with a score function. However, the quality of the generated counterfacts is still suboptimal since the LLM is not specially optimized for generating CAD, and the low-scored data is also not fully leveraged.

In this paper, to address the above issues, we propose a **dually self-improved counterfactual data augmentation method using LLMs (DICT)**. On one hand, as the attention mechanism offers insights into the causal relationships between texts and their labels (Nauta et al., 2019), we design a self-improved strategy based on the attention distribution of the target task model to identify causal terms, a lightweight and task-specific approach. As shown in Figure 1, the terms with larger attention of the target task model are more critical for the NLI label, while existing methods suffer from the accuracy of the identified causal terms and may introduce noise. On the other hand, to further improve the quality of CAD, we propose an additional self-improved strategy based on di-

rect preference optimization (DPO) to refine itself. Specifically, after generating preliminary counterfactuals, we construct the preference pairs based on the score function for DPO. Finally, through a simple filtering and fusion, we retrain the task model on the fused data, using a balanced loss function to avoid over-emphasis on augmented data. Overall, our contributions can be summarized as follows:

- We propose a dually self-improved counterfactual data augmentation, improving the counterfactual data augmentation framework depending on the task model and LLM themselves, without external tools to identify casual terms or human-annotation for fine-tuning LLMs.
- Our proposed DICT improves the extraction of task-specific causal terms through attention mechanisms and further enhances the CAD generation of LLMs using DPO. Additionally, a novel balanced loss is introduced to retrain the task model on the fused data, effectively preventing excessive augmentation.
- Extensive experiments across multiple benchmarks demonstrate that DICT significantly outperforms the SOTA manual and automatic CAD generation methods across all metrics.

2 Related Work

Natural Language Inference. NLI involves determining the relationship between a pair of sentences: a premise and a hypothesis. It is widely used to

evaluate a model’s ability to understand and reason about natural language semantics, making it an essential benchmark for advanced NLP models. Although many works (Liu et al., 2019; Radford, 2018) have shown promising performance, the task still suffers from spurious association problem (Wen et al., 2022). For example, negation words often become overly strong indicators of contradiction, while a high rate of word overlap between the premise and hypothesis is frequently misinterpreted as a strong indicator of entailment (Gururangan et al., 2018; Naik et al., 2018). Recently, to alleviate this issue, researchers have explored using counterfactual data augmentation (Kaushik et al.).

Counterfactual Data Augmentation. Generating fluent textual CAD are required to follow some principles, including: (1) minimal edits, (2) fluency, creativity and diversity, (3) adhering to task-specific rules (Wang et al., 2024). However, these requirements have been proved challenging. Early, Kaushik et al. and Gardner et al. (2020) employ human annotators to create counterfactuals by manually rewriting the original data. Obviously, manual rewrites are not only time-consuming and expensive but also may exacerbate existing spurious features. To alleviate the mentioned issues, researchers (Madaan et al., 2021; Ross et al., 2021) proposed to use advanced text generation models to generate CAD. For example, Wen et al. (2022) designed a fully automatic and task-agnostic framework through T5 (Raffel et al., 2020). However, the quality of the generated data remain constrained due to the limited comprehension and generation capabilities of previous generative language models.

LLM-based Counterfactual Data Augmentation. LLMs have shown remarkable proficiency in synthesizing natural languages for downstream tasks. Leveraging the powerful generative ability of LLMs to automatically generate counterfactuals has recently attracted considerable attention (Liu et al., 2020a). DISCO (Chen et al., 2023) prompts GPT3 (Brown et al., 2020) to generate phrasal perturbations for automatically generating CAD at scale. Nguyen et al. (2024) and Li et al. (2024) investigated the strengths and weaknesses of LLMs as generators comprehensively, instructing LLMs to identify casual terms and generate counterfactuals.

However, despite the significant advancements, the quality of counterfactual augmented data with LLMs still remains to be improved since LLMs

are not specially trained for CAD generation. Our work bridges this gap by designing a dually self-improved method to enhance both the extraction of the specific causal terms and the generation of CAD (modifying the causal terms) with LLMs.

3 Preliminaries

We implement counterfactual data augmentation on the Natural Language Inference (NLI) task, referring to determine the relationship between a given premise sentence and a hypothesis sentence (Hosseini et al., 2024). Formally, given an input premise-hypothesis pair $\langle P_i, H_i \rangle$ and its ground-truth label l_i , where $P_i = \{t_1, t_2, \dots, t_m\}$, t_j represents a token¹, and m is the number of tokens. $l_i \in \{\text{Entailment, Contradiction, Neutral}\}$, the task aims to produce a counterfactual example $\langle \hat{P}_i, H_i \rangle$ that flips the origin label l to a desired label \hat{l}_i , $\hat{l}_i \neq l_i$, through perturbing parts of the premise P_i . When original premise P_i is altered into counterfactual \hat{P}_i , minimal changes are required. Here, casual terms are denoted as $C_i = \{c_1, \dots, c_k\}$, where each c_j corresponds to a token t_j extracted from P_i . After CAD generation, the performance is evaluated through a baseline NLI model \mathbb{M} , such as Roberta (Liu et al., 2020b).

4 Our Proposed Model

In this section, we detail our proposed dually self-improved counterfactual data augmentation method using large language model (DICT).

As shown in Figure 2, our model consists of three stages: 1) self-improved casual terms identification, 2) self-improved CAD generation, 3) retraining. First, we design a self-improvement strategy leveraging the attention distribution of the task model to enhance the identification of causal terms. Second, we further propose to utilize a self-improved LLM based on DPO to refine the CAD generation by modifying the causal terms. Finally, after filtering and fusing the generated counterfactuals, we retrain the task model with a balanced loss function, avoiding over-augmentation. In this way, we improve the task model performance with our generated augmented counterfactual data.

4.1 Self-improved Casual Terms Identification

Casual terms captures the effective features implied in sentences. Therefore, identifying causal terms is

¹A token consists of one or more words. We split sentences into tokens through Flair (Akbik et al., 2018).

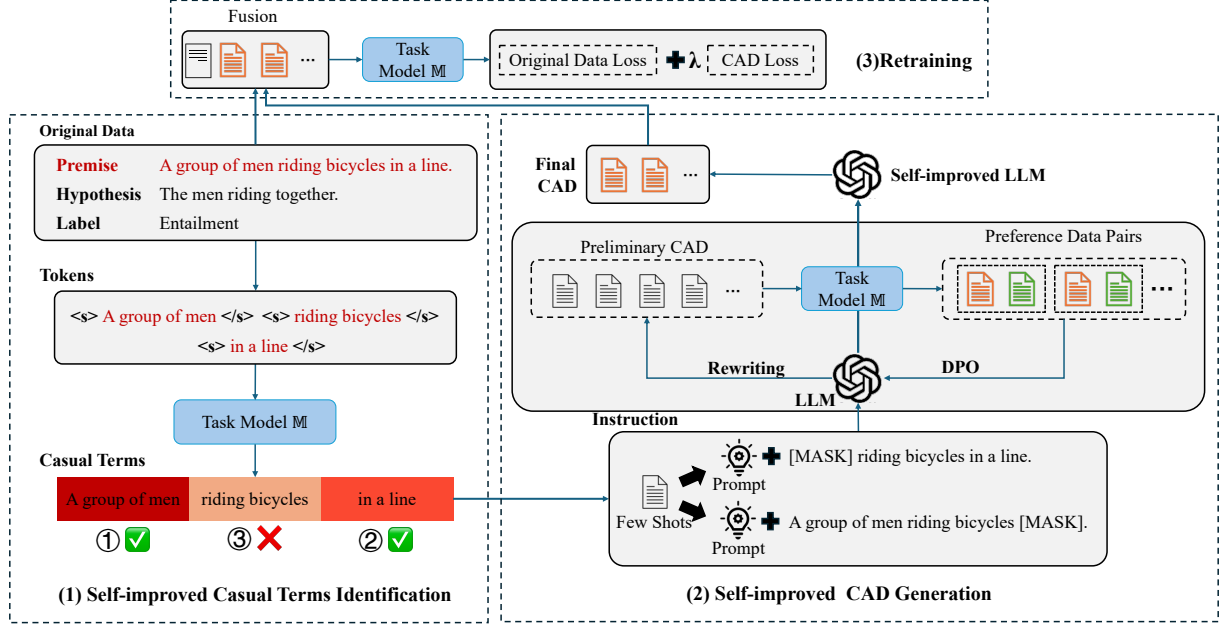


Figure 2: The architecture of our proposed DICT.

the crucial first step of counterfactual data augmentation. To achieve this, we propose a self-improved causal term identification method based on the attention distribution of the task model. Formally, given the task model \mathbb{M} pre-trained on the original dataset and a premise-hypothesis sample $\langle P_i, H_i \rangle$, the attention score α_j on each token t_j of premise P_i under its label l_i is computed as follows:

$$\alpha_1, \alpha_2, \dots, \alpha_m = \text{Attention}_{\mathbb{M}}(l_i | P_i, H_i), \quad (1)$$

where $\text{Attention}_{\mathbb{M}}$ is the specified attention function embedded in the task model \mathbb{M} . Note that a mean pooling is used to transfer the attention distribution from word level to our token level. Then, tokens are sorted in descending order based on the attention score α_i and top K tokens are selected as the final causal terms C_i .

4.2 Self-improved CAD Generation

With the identified causal terms and original sentence pairs, we propose a self-improved LLM based on DPO, to modify the causal terms to flip the label, generating CAD.

First, each causal term C_i is replaced with a mask token [MASK] individually to obtain K sentences to be rewritten. Then, for each sentence, we instruct an LLM to alter the [MASK] into certain tokens for flipping the original label l_i of the $\langle P_i, H_i \rangle$ into a specific label \hat{l}_i . Note that,

Afterwards, all the candidate counterfacts are scored via the predicted probability shift of the

target label \hat{l}_i based on the task model \mathbb{M} :

$$\delta_j = p(\hat{l} | \hat{P}_i^j, H_i) - p(\hat{l} | P_i, H_i). \quad (2)$$

Instead of directly using the filtered results by the calculated scores δ , we design another self-improved strategy based on DPO to achieve self-improved LLM for generating higher-quality candidate counterfacts. Specifically, for each causal term in C , we choose the corresponding generated candidate counterfact (by modifying the causal term) with the highest score δ as the accepted example, and a random one with $\delta < \gamma$ as a rejected example. The two samples form a preference data pair $(\hat{P}_i^{pos}, \hat{P}_i^{neg})$. Formally, the entire preference pair data are denoted as:

$$\mathbb{P} = \{(P_i, \hat{P}_i^{pos}, \hat{P}_i^{neg})\}_{i=1}^N. \quad (3)$$

Self-Improved LLM based on DPO. As defined previously, we prefer the counterfact \hat{P}_i^{pos} to \hat{P}_i^{neg} given an input P_i . To enable the LLM to learn this desired preference, DPO is employed to refine the LLM using the preference pairs, instead of training a reward model explicitly. Formally, LLMs can be directly optimized using the following binary cross entropy loss function:

$$L_{DPO} = \mathbb{CE}(p(\hat{P}_i^{pos} \succ \hat{P}_i^{neg} | P), \pi(\hat{P}_i^{pos}, P)), \quad (4)$$

where \mathbb{CE} is the cross entropy function, $p(\hat{P}_i^{pos} \succ \hat{P}_i^{neg} | P)$ represents the probability that we prefer the

output \hat{P}_i^{pos} to \hat{P}_i^{neg} given an input P , and $\pi(\hat{P}_i^{pos}, P)$ is the predicted probability by the LLM.

Subsequently, we apply the self-improved LLM based on DPO to generate higher-quality CAD. The generated candidate counterfactuals are further filtered based on the aforementioned probability shift score δ to ensure the data quality (i.e., $\delta \geq \gamma$).

4.3 Retraining

Finally, we fuse the filtered CAD with the original data, and retrain the task model to improve the task performance. As the scale of counterfactual data grows, we observe that the task model may overly focus on the counterfactual data while overlooking the original data. Therefore, during the retraining, a penalty factor λ is used to balance the original data and the augmented data, improving the robustness of the model while preventing over-emphasis on the augmentation. The loss function is calculated through:

$$L = \mathbb{CE}(p(l|P, H), l) + \lambda \cdot \mathbb{CE}(p(\hat{l}|\hat{P}, H), \hat{l}), \quad (5)$$

where \mathbb{CE} is the cross entropy function, and λ is the balance factor.

5 Experiments

5.1 Datasets

We evaluate our method on NLI tasks over three benchmarks, including two in-domain subsets from SNLI (Bowman et al., 2015) and MNLI (Williams et al., 2018). In the following, we detail each dataset.

- SNLI (Bowman et al., 2015). The Stanford Natural Language Inference (SNLI) corpus, derived from only one domain, is a collection of sentence pairs manually labeled for balanced classification with the labels entailment, contradiction, and neutral. The first subset SNLI-1, following (Wen et al., 2022), consists of ambiguous part of SNLI. It contains 20,000 examples for training, 4,800 for validation and 4,800 for test. To further evaluate the performance, we extracted a larger scale examples randomly from original SNLI corpus, consisting of 87,208, 18,688 and 18,688 pairs for training, validation, and test respectively.
- MNLI (Williams et al., 2018). Multi-genre NLI corpus (MNLI), including two difference test sets MNLI-matched (MNLI-m) and

MNLI-mismatched (MNLI-mm)², is a multiple out-of-domain and challenge benchmark to measure the generalization of the model after data augmentation. It contains 392, 702 pairs in the train set, 9, 815 in the MNLI-m test set, and 9,796 pairs in the MNLI-mm test set.

5.2 Baselines

We compare our model with the state-of-the-art baselines:

- RoBERTa-large (Liu et al., 2019). A robustly optimized and SOTA transformer model pre-trained on a large corpus. It is used as the target task model to be augmented.
- HumanCAD (Kaushik et al., 2020). A manual set of CAD for NLI, obtained by human annotators rewriting a subset of SNLI. We append them into original benchmarks and evaluate the performance following (Wen et al., 2022).
- AutoCAD (Wen et al., 2022). A fully automatic CAD generation framework with the generative language model T5.
- DISCO (Chen et al., 2023). A counterfactual knowledge distillation approach with LLMs. It leverages all spans as causal terms for CAD generation and filters out unqualified generated data using a SOTA task-specific model.
- LLMCF (Li et al., 2024). A CoT-based method that prompts LLMs to identify causal terms and produce CAD. To ensure a fair comparison, we adopt the task model to filter the generated CAD, as we do in our DICT.

5.3 Experimental Settings

For SNLI-1, we perform counterfactual augmentation on each sample. Due to the large scale of the SNLI-2 and the MNLI, we sampled subsets of a fixed size for counterfactual augmentation, including 50,000 examples from the training set. We sample 2,000 examples respectively from the test sets as well. In terms of LLM-based models, we use Qwen2.5-7B-Instruct and Qwen2.5-14B-Instruct (Yang et al., 2024; Team, 2024) as the base LLMs. The prompt for instructing LLMs follows (Chen et al., 2023), ensuring a fair comparison and minimizing the impact of prompt variations on the

²The details can be found in the website <https://cims.nyu.edu/sbowman/multinli/>

Dataset	SNLI-1			SNLI-2			MNLI-m			MNLI-mm		
Metric(%)	P	R	F1	P	R	F1	P	R	F1	P	R	F1
Roberta	61.36	59.77	58.29	87.92	86.76	86.82	87.38	87.23	87.27	87.06	86.92	86.97
Human-CAD	60.90	62.27	61.26	87.57	87.51	87.51	87.17	86.92	86.85	87.30	87.06	87.10
AutoCAD	57.08	58.58	57.48	87.37	87.35	87.36	87.52	87.33	87.41	87.44	87.32	87.37
DISCO-7B	59.50	61.18	59.26	97.80	87.73	87.75	87.76	87.77	87.76	87.56	87.50	87.54
LLMCF-7B	61.17	61.43	60.24	88.43	87.39	87.65	87.80	87.66	87.71	87.70	87.57	87.62
LLMCF-14B	63.15	63.43	62.84	88.82	88.79	88.79	88.89	88.73	88.84	88.72	88.66	88.68
DICT-7B	62.40	62.45	61.32	88.65	87.78	87.90	88.22	88.11	88.17	88.13	87.86	87.92
DICT-14B	65.17	65.21	64.89	89.45	89.56	89.51	89.48	89.35	89.39	89.31	89.28	89.27

Table 1: Performance comparison of different methods over Precision, Recall and F1 score, where 7B and 14B means Qwen2.5-7B-Instruct and Qwen2.5-14B-Instruct as the base LLM respectively.

generated counterfactuals. The RoBERTa-large model is trained on all basic and augmented datasets with a learning rate of $1e-5$ for 3 epochs. In our DICT, the top 3 causal terms are selected to produce candidate counterfactuals. The threshold γ of the score δ is set as 0.5 for filtering CAD. The size of obtained preference pairs is approximately 25,000 across all the datasets. For the DPO process, we set the number of epochs to 1. The penalty factors λ in loss function are 0.4 and 0.6 for DICT-7B and DICT-14B, respectively.

5.4 Overall Performance

To assess overall performance, we perform counterfactual data augmentation on the training data and evaluate on the original test set. As shown in Table 1, we report Precision (P), Recall (R) and F1-score (F1) respectively on all datasets to evaluate the overall performance of CAD methods. Concretely, the task model Roberta is trained on the fusion of the generated counterfactuals and the original data, and evaluated on the original test data. It can be observed that: (1) all counterfactual data augmentation methods prove effective in most cases. However, due to the higher ambiguity and difficulty of SNLI-1, AutoCAD slightly weakens the model performance. (2) LLM-based methods outperform AutoCAD in most cases, indicating the powerful comprehension and generation capabilities of LLMs. (3) Our proposed model DICT achieves the best results across both 7B and 14B settings, especially on the more challenging SNLI-1 dataset and the out-of-domain MNLI-mm dataset. It demonstrates the robustness and effectiveness of our proposed DICT. (4) Both LLMCF and DICT exhibit significant performance improvements as the LLM scale increases, demonstrating that larger models can capture more complex causal relationships and

Method	FR	ACC_δ
Auto-CAD	0.46	0.59
DISCO	0.61	0.77
LLMCF-7B	0.60	0.81
LLMCF-14B	0.71	0.83
DICT-7B	0.80	0.84
DICT-14B	0.82	0.87

Table 2: Evaluation of the quality of generated counterfactuals.

generate higher-quality counterfactual data, leading to better task performance. Note that, DICT performs best in all cases. We believe that the reason is that DICT with dual self-improvement can accurately identify the task-specific causal terms with the attention mechanism and generate higher-quality counterfactuals based DPO.

5.5 The Quality of Generated Counterfactuals

Following (Nguyen et al., 2024; Chen et al., 2023), we use the flip rate (FR) and the counterfactual accuracy ACC_δ to evaluate the quality of generated counterfactuals on SNLI-1. Specifically, FR quantifies how effectively a method can alter the labels of instances and a higher FR indicates more confident and impactful context modifications. FR is defined as:

$$FR = \frac{1}{N} \sum_{i=1}^N \mathbb{I}[p(\hat{l}_i | \hat{P}_i, H_i) = \hat{l}_i], \quad (6)$$

, where \mathbb{I} is an indicator function that outputs 1 if the predicted label of a counterfact matches its desired label. The FR is evaluated using the counterfactual augmentation results on the training set, where the probability $p(\hat{l}_i | \hat{P}_i, H_i)$ is computed using the task model.

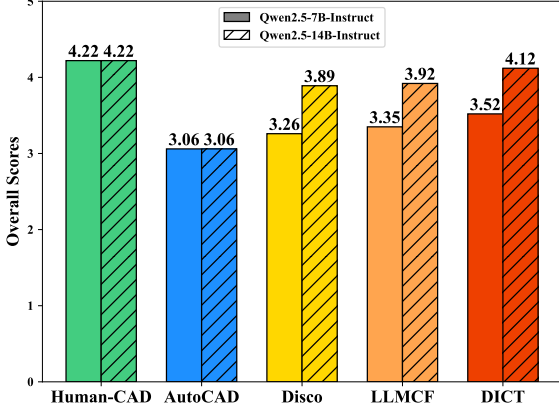


Figure 3: Evaluated Results with GPT4 Over Qwen2.5-7B and Qwen2.5-14B Respectively on SNLI-1.

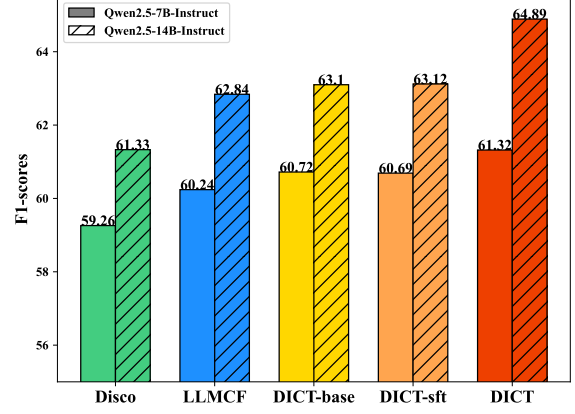


Figure 4: Ablation study over Qwen2.5-7B and Qwen2.5-14B on SNLI-1.

The counterfactual accuracy is used to measure the consistency of model performance on original and counterfactual examples, and is defined as:

$$\frac{1}{N} \sum_{i=1}^N \mathbb{I}[p(\hat{l}_i | \hat{\mathbf{P}}_i, \mathbf{H}_i) = \hat{l}_i \wedge p(l_i | \mathbf{P}_i, \mathbf{H}_i) = l_i], \quad (7)$$

where \mathbb{I} indicates 1 only when the model correctly predict the original and counterfactual examples. All probabilities are computed using the augmented task model on the test set and their corresponding counterfactuals. Therefore, test samples linked to corresponding counterfactual examples are preserved.

As shown in Table 2, our model achieves the best performance on both FR and ACC_δ . DICT-14B increases the FR around 15% compared to LLMCF-14B, demonstrating that DICT effectively produces a larger quantity of high-confidence counterfactuals. Additionally, the results on ACC_δ also highlight that our DICT exhibits better consistency and generalization.

Evaluation with GPT4. GPT-4 is a reliable evaluator for accessing the quality of CAD, as demonstrated in (Nguyen et al., 2024). Accordingly, we select 1,000 samples randomly from SNLI-1 for all methods and use GPT-4 to assign an overall score (on a 5-point scale) to them from three aspects, including fluency, realism, and conciseness. As shown in 3, compared to Auto-CAD that employs traditional generative language models, LLM-based models achieve higher scores obviously. Despite that all model-based methods fall short of Human-CAD, our DICT still achieves superior performance over Human-CAD. Simultaneously, as the scale of the large models increases, the scores

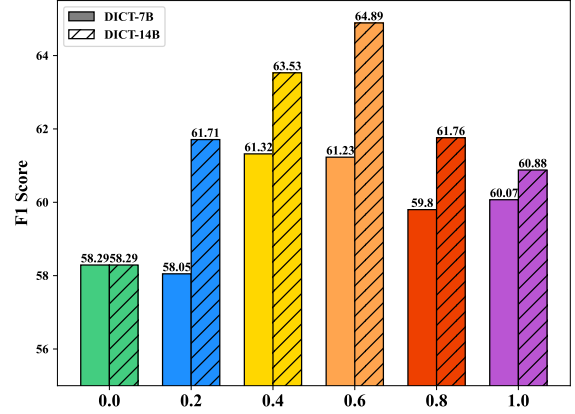


Figure 5: The impact of Hyperparameter λ for DICT-7B and DICT-14B on SNLI-1.

show significant improvements.

5.6 Ablation Study

In order to verify the effectiveness of different modules of our model, we design two variant models:

- **DICT-base** removes the self-improved generator and use a basic LLM to produce CAD.
- **DICT-sft** replaces the DPO strategy with supervised fine-tuning (SFT). Instead of improving the LLM on preference data pairs, it just employs the preferred parts.

They are both compared to LLM-based methods on SNLI-1 dataset with Qwen2.5-7B and Qwen2.5-14B respectively. As shown in Figure 4, we report F1-scores as evaluated results. Without self-improvement generator, the performances are still better than both DISCO and LLMCF. It demonstrates that our self-improved identifier can identi-

Original Premise	Two people are holding a large upside-down earth globe , about 4' in diameter, and a child appears to be jumping over Antarctica.
Original Hypothesis	The earth globe is purple.
Original Label	Contradiction
Counterfactual Premis	Two people are holding a large purple earth globe , about 4' in diameter, and a child appears to be jumping over Antarctica.
Flipped Label	Entailment

Figure 6: A counterfactual example from SNLI-1 generated by our DICT.

498 fying specific casual terms that are crucial for gener-
 499 ating CAD. If we replace DPO with SFT as our self-
 500 improved strategy of generator, the performances
 501 of DICT-sft decrease by 0.5% and 1.77% over
 502 Qwen2.5-7B and Qwen2.5-14B respectively. It in-
 503 dicates the necessity of designing a self-improved
 504 strategy to enhance the LLM’s rewriting capabil-
 505 ity of CAD. We also find that the performances of
 506 DICT-sft increase in-obviously compared to DICT-
 507 base. The reason may be that without the constraint
 508 of negative samples, the optimization space of the
 509 LLM becomes more complicated in our task. It is
 510 assumed that there should be more high-confidence
 511 CAD to train the LLM better with SFT. Addition-
 512 ally, as the parameter scale of the LLM increases,
 513 the performance of all methods improves signifi-
 514 cantly, further validating larger models can gener-
 515 ate higher-quality counterfactual data.

5.7 HyperParameter Experiments

517 We validate the impact of different hyperparam-
 518 eters λ within $\{0, 0.2, 0.4, 0.6, 0.8, 1\}$ on preventing
 519 over-emphasis on augmented data. When λ is
 520 equal to 0, the DICT degenerates to the basic model
 521 RoBERTa. As shown in Figure 5, when λ is rela-
 522 tively small (e.g., 0.2 or below), the model primar-
 523 ily focuses on original data, limiting the benefits
 524 of counterfactual data augmentation. Conversely,
 525 when λ is too high (e.g., 1.0), the model heav-
 526 ily emphasizes CAD, degrading the performance.
 527 Optimal results are observed within the range of
 528 $\lambda \in [0.4, 0.6]$ for both DICT-7B and DICT-14B,
 529 where the balance between original and generated
 530 counterfactual data contributes to improve the per-
 531 formance.

5.8 Case Study

532 Figure 6 shows a counterfual example from SNLI-
 533 1. In this example, key tokens in the premise like
 534 "earth globe" and "purple" significantly influence
 535

536 the relationship with the hypothesis, namely the
 537 NLI label. Our DICT can successfully extract these
 538 tokens as causal terms for modifying to flip the
 539 NLI label. This step ensures that the counterfac-
 540 tual generation is grounded in the critical linguistic
 541 features. Thus, the generated councterfacts are of
 542 high-quality.

6 Conclusion

543 In this paper, we address the challenges in LLM-
 544 based counterfactual data augmentation by intro-
 545 ducing the proposed DICT method, a dually self-
 546 improved counterfactual data augmentation ap-
 547 proach using LLM. Specifically, we frist introduce
 548 a lightweight and task-specific causal term iden-
 549 tification strategy that leverages the attention dis-
 550 tribution of the task model for self-improvement.
 551 This approach effectively captures causal terms
 552 by interpreting the attention scores, overcoming
 553 the limitations of LLMs in accurately identifying
 554 specific causal terms. Second, we propose a self-
 555 improved counterfactual generator that modifies
 556 the causal terms to flip the label based on DPO. By
 557 constructing preference data pairs from the prelim-
 558 inary generated counterfacts, we refine the LLM
 559 with DPO, ensuring higher-quality counterfactual
 560 generation. Our experimental results demonstrate
 561 that DICT outperforms existing LLM-based coun-
 562 terfactual data augmentation methods across vari-
 563 ous NLI datasets, achieving superior performance
 564 in terms of both accuracy and robustness. Addition-
 565 ally, we observe that increasing the LLM’s param-
 566 eter scale further boosts the performance, highlight-
 567 ing the scalability and effectiveness of our proposed
 568 method.

569 Furthermore, our DICT can be directly applied
 570 to various NLP tasks such as sentiment analysis
 571 and relation extraction, which we will explore in
 572 future work.
 573

7 Limitation

While Dict demonstrates strong performance, it is inherently dependent on the capabilities of the underlying large language models (LLMs). This dependence means that Dict’s effectiveness can vary across different LLM architectures and versions, highlighting the need for strong LLM backbone to ensure reliable outcomes.

References

Alan Akbik, Duncan Blythe, and Roland Vollgraf. 2018. Contextual string embeddings for sequence labeling. In *COLING 2018, 27th International Conference on Computational Linguistics*, pages 1638–1649.

Samuel Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642.

Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, pages 1877–1901.

Hao Chen, Rui Xia, and Jianfei Yu. 2021. Reinforced counterfactual data augmentation for dual sentiment classification. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*.

Zeming Chen, Qiyue Gao, Antoine Bosselut, Ashish Sabharwal, and Kyle Richardson. 2023. DISCO: Distilling counterfactuals with large language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5514–5528.

Tanay Dixit, Bhargavi Paranjape, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. CORE: A retrieve-then-edit framework for counterfactual data generation. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2964–2984.

Matt Gardner, Yoav Artzi, Victoria Basmov, Jonathan Berant, Ben Bogin, Sihao Chen, Pradeep Dasigi, Dheeru Dua, Yanai Elazar, Ananth Gottumukkala, Nitish Gupta, Hannaneh Hajishirzi, Gabriel Ilharco, Daniel Khashabi, Kevin Lin, Jiangming Liu, Nelson F. Liu, Phoebe Mulcaire, Qiang Ning, Sameer Singh, Noah A. Smith, Sanjay Subramanian, Reut Tsarfay, Eric Wallace, Ally Zhang, and Ben Zhou. 2020. Evaluating models’ local decision boundaries via contrast sets. In *Findings of the Association for Computational Linguistics: EMNLP 2020*.

Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and Noah A. Smith. 2018. Annotation artifacts in natural language inference data. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 107–112.

Mohammad Javad Hosseini, Andrey Petrov, Alex Fabrikant, and Annie Louis. 2024. A synthetic data approach for domain generalization of nli models. *arXiv preprint arXiv:2402.12368*.

Divyansh Kaushik, Eduard Hovy, and Zachary C Lipton. 2020. Learning the difference that makes a difference with counterfactually augmented data. *International Conference on Learning Representations (ICLR)*.

Divyansh Kaushik, Amrith Setlur, Eduard H Hovy, and Zachary Chase Lipton. Explaining the efficacy of counterfactually augmented data. In *International Conference on Learning Representations*.

Yongqi Li, Mayi Xu, Xin Miao, Shen Zhou, and Tieyun Qian. 2024. Prompting large language models for counterfactual generation: An empirical study. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 13201–13221.

Qi Liu, Matt Kusner, and Phil Blunsom. 2021. Counterfactual data augmentation for neural machine translation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.

Tianyu Liu, Zheng Xin, Baobao Chang, and Zhifang Sui. 2020a. HypoNLI: Exploring the artificial patterns of hypothesis-only bias in natural language inference. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 6852–6860.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized BERT pretraining approach. *CoRR*.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2020b. Ro{bert}a: A robustly optimized {bert} pretraining approach.

Nishtha Madaan, Inkit Padhi, Naveen Panwar, and Dip-tikalyan Saha. 2021. Generate your counterfactuals: Towards controlled counterfactual generation for text. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 13516–13524.

Aakanksha Naik, Abhilasha Ravichander, Norman Sadeh, Carolyn Rose, and Graham Neubig. 2018. Stress test evaluation for natural language inference. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2340–2353.

- Meike Nauta, Doina Bucur, and Christin Seifert. 2019. Causal discovery with attention-based convolutional neural networks. *Machine Learning and Knowledge Extraction*, pages 312–340.
- Van Bach Nguyen, Paul Youssef, Christin Seifert, and Jörg Schlötterer. 2024. LLMs for generating and evaluating counterfactuals: A comprehensive study. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 14809–14824.
- Alec Radford. 2018. Improving language understanding by generative pre-training.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, pages 1–67.
- Alexis Ross, Ana Marasović, and Matthew Peters. 2021. Explaining NLP models via minimal contrastive editing (MiCE). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3840–3852.
- Qwen Team. 2024. Qwen2.5: A party of foundation models.
- Yongjie Wang, Xiaoqi Qiu, Yu Yue, Xu Guo, Zhiwei Zeng, Yuhong Feng, and Zhiqi Shen. 2024. A survey on natural language counterfactual generation. *arXiv preprint arXiv:2407.03993*.
- Jiixin Wen, Yeshuang Zhu, Jinchao Zhang, Jie Zhou, and Minlie Huang. 2022. AutoCAD: Automatically generate counterfactuals for mitigating shortcut learning. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2302–2317.
- Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Hao-ran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.