

Counterfactual Debiasing for Fact Verification

Anonymous ACL submission

Abstract

Fact verification aims to automatically judge the veracity of a claim according to several evidences. Due to the manual construction of datasets, spurious correlations between claim patterns and its veracity (i.e., biases) inevitably exist. Recent studies show that models usually learn such biases instead of understanding the semantic relationship between the claim and evidences. Existing debiasing works can be roughly divided into data-augmentation-based and weight-regularization-based pipeline, where the former is inflexible and the latter relies on the uncertain output on the training stage. Unlike previous works, we propose a novel method from a counterfactual view, namely CLEVER, which is augmentation-free and mitigates biases on the inference stage. Specifically, we train a claim-evidence fusion model and a claim-only model independently. Then, we obtain the final prediction via subtracting output of the claim-only model from output of the claim-evidence fusion model, which counteracts biases in two outputs so that the unbiased part is highlighted. Comprehensive experiments on several datasets have demonstrated the effectiveness of CLEVER.

1 Introduction

Unverified claims have been prevalent online with the dramatic increase of information, which poses a threat to public security over various domains, e.g., public health (Naeem and Bhatti, 2020), politics (Allcott and Gentzkow, 2017), and economics (Kogan et al., 2019). Therefore, fact verification, which aims to automatically predict the veracity of claims based on several collected evidences, has attracted lots of research interests (Liu et al., 2020; Zhong et al., 2020; Vo and Lee, 2021).

Existing fact-checking datasets inevitably involve some biases since they are manually collected. For example, Schuster et al. (2019) discover that negation words in claims are highly-correlated

with the label ‘REFUTES’ in the FEVER dataset (Thorne et al., 2018). Such biases may mislead models to explore the spurious correlation between claim patterns and its label without looking into the evidences. In consequence, though models achieve promising performance on biased datasets, they suffer from obvious performance decline on out-of-domain unbiased datasets and are vulnerable to adversarial attacks (Thorne et al., 2019).

To alleviate the aforementioned problems, several debiasing methods have been proposed, which can be mainly grouped into two categories. The first pipeline is based on data augmentation, which utilizes manually-designed schemes, such as word swapping (Wei and Zou, 2019) and span replacement (Lee et al., 2021) to generate additional data for training. However, these methods heavily rely on the quality of augmented data and are difficult to be employed under complicated circumstance, e.g., multi-hop evidence reasoning, due to their inflexible augmentation rules. The second pipeline aims to downweigh the contribution of biased samples to the training loss of main model, whose inputs are both claim and evidence. Then, the key issue is how to recognize the biased instances. Specifically, Schuster et al. (2019) downweigh the claim involving n-grams that share spurious correlation with labels. Mahabadi et al. (2020) assume instances correctly classified by the bias-only model are biased, where the input of bias-only model is the claim only. Nevertheless, the former lacks the generalization to different types of biases since they only focus on n-grams; the latter relies on the assumption that the outputs of main model and bias-only model regarding the biased instances are similar, which does not always hold (Amirkhani and Pilehvar, 2021). Moreover, the inaccurate and unstable outputs of bias-only model during training may mistakenly result in downweighing unbiased samples (Xiong et al., 2021).

Unlike existing works based on augmentation or

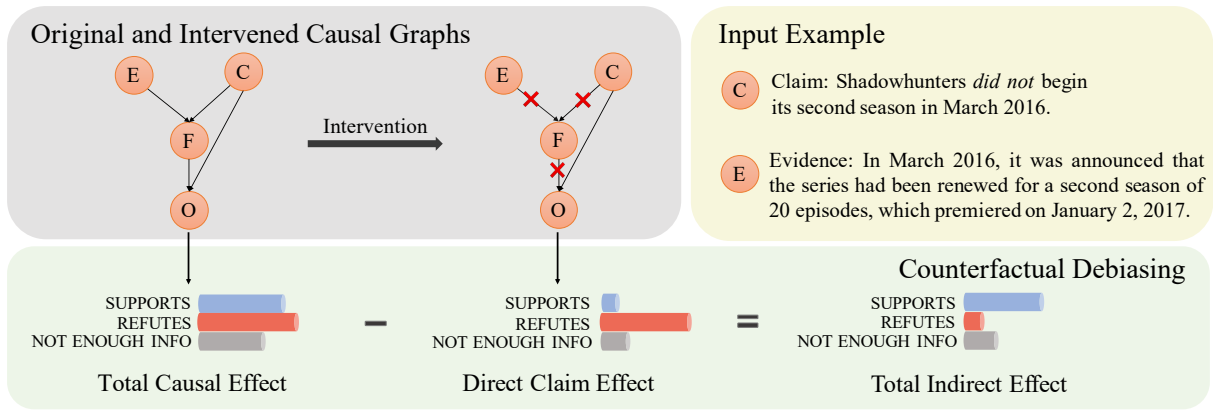


Figure 1: The causal view of proposed framework CLEVER. The nodes with ‘F’ and ‘O’ denote the claim-evidence fused information and the model output, respectively. We take a typical sample in the biased FEVER dataset as input, where the label is ‘SUPPORTS’ and the strong correlation between the phrase ‘did not’ and label ‘REFUTES’ exists. The output in original graph (Total Causal Effect) is affected by two sources, i.e., claim and claim-evidence fused information. After the intervention via cutting off the fusion path, the output (Direct Claim Effect) is solely influenced by the claim, which contains biases that mislead the model to produce spurious label prediction. To mitigate such biases, a subtraction scheme is proposed to obtain the Total Indirect Effect, which inclines to the true debiased distribution. Note that a path from evidence to output does not exist since there is no obvious bias in evidences that affects the outcome.

083 adjusting the data contribution on the training stage,
 084 we propose a novel method from a Counterfactual
 085 Lview for dEbiasing fact VERification, namely
 086 CLEVER, which is augmentation-free and allevi-
 087 ates biases on the inference stage. In general, exist-
 088 ing methods fuse the claim and evidences to make
 089 the final prediction, which is equivalent to asking
 090 the model to answer a factual question: *What will*
 091 *the output be if the model receives a claim and*
 092 *its corresponding evidences?* Causally, the Total
 093 Causal Effect is estimated in this condition, where
 094 claim biases are entangled with the claim-evidence
 095 fused information, making them difficult to be miti-
 096 gated precisely. To overcome this, we aim to obtain
 097 the debiased output by removing claim biases from
 098 the Total Causal Effect. Inspired by the progress of
 099 counterfactual inference (Sekhon, 2008; Niu et al.,
 100 2021), we would expect to ask a counterfactual
 101 question: *What would the output be if the model*
 102 *only received a claim?* That is, from a causal per-
 103 spective, requiring the fact-checking model to learn
 104 the Direct Claim Effect solely affected by claim bi-
 105 ases. Practically, we first train a claim-evidence fu-
 106 sion model and a claim-only model independently
 107 to capture the Total Causal Effect and the Direct
 108 Claim Effect, respectively. Then, we subtract the
 109 **Direct Claim Effect** from the **Total Causal Ef-**
 110 **fect** on the inference stage to obtain the Total Indi-
 111 rect Effect, which is the final debiased prediction.
 112 Taking Figure 1 as an example, the claim is spuri-

113 ously correlated with the false label ‘REFUTES’.
 114 Therefore, the Direct Claim Effect inclines to the
 115 label ‘REFUTES’ since it is affected by the claim
 116 only. However, the prediction is turned towards
 117 the ground-truth label via using the Total Indirect
 118 Effect as the final output, where the high proba-
 119 bility of ‘REFUTES’ induced by claim biases is
 120 counteracted.

Our main contributions are listed as follows:

- We open up a new counterfactual pipeline for
 122 debiasing fact verification by analyzing the
 123 biased problem from a causal view. 124
- We propose a novel debiasing method
 125 CLEVER, which is augmentation-free and
 126 mitigates biases on the inference stage. 127
- Comprehensive experiments are conducted to
 128 validate the effectiveness of CLEVER, where
 129 the results demonstrate the superiority and the
 130 in-depth analysis provides the rationality. 131

2 Related Work 132

In this section, we briefly review the related lit-
 133 erature in both domains of fact verification and
 134 debiasing strategy. 135

2.1 Fact Verification 136

Recent years have witnessed the rapid development
 137 of research on fact verification. Since the unified
 138

benchmark dataset FEVER along with the shared task were proposed (Thorne et al., 2018), most researchers utilize them to evaluate the model performance. Generally, the fact-checking task mainly consists of three separate parts, i.e., document retrieval, evidence selection, and claim verification. Existing works mainly focus on the last subtask and employ traditional and widely used methods (Hanselowski et al.; Soleimani et al., 2020) to retrieve relevant documents and evidences. Early works treat fact verification as a natural language inference (NLI) task and apply methods from NLI to perform verification (Chen et al., 2017; Ghaeini et al., 2018). Then, to capture more fine-grained semantic consistency between claims and evidences, a series of methods have been proposed to promote the claim-evidence interaction by formulating them as graph-structure data (Zhou et al., 2019; Liu et al., 2020; Zhong et al., 2020). Besides, inspired by the strong representation ability of pretrained language models (PLM), some works attempt to fine-tune PLM on fact-checking datasets and achieve promising results (Lee et al., 2020; Subramanian and Lee, 2020). Recently, researchers have paid more attention to explainable fact verification, which requires a model to produce both veracity prediction and its corresponding explanation (Kotonya and Toni, 2020a,b).

2.2 Debiasing Strategy

Although the aforementioned fact-checking methods have achieved promising performance on the FEVER test set, it is demonstrated that they lack robustness since they learn biases (shortcuts) from claims in datasets instead of performing reasoning over evidences. To this end, several unbiased and adversarial datasets are proposed to evaluate the model robustness and reasoning ability (Thorne et al., 2019; Schuster et al., 2019). Existing debiasing strategies in fact verification can be roughly divided into two groups:

1) *Data-augmentation-based pipeline*: In this group, methods aim to generate unbiased samples and incorporate them into training, with the expectation that the proportion of biased instances will be downgraded, resulting in a more unbiased model. In detail, Wei and Zou (2019) utilize random word swapping and synonym replacement to obtain new training data. Lee et al. (2021) design a cross contrastive strategy to augment data, where original claims are modified to be negative using

the generation model BART (Lewis et al., 2020) and evidences are changed via span replacement to support such negative claims.

2) *Weight-regularization-based pipeline*: The motivation of methods in this pipeline is to reduce the contribution of biased samples to the final loss computation, thus models may attach importance to the unbiased data. Next, the problem is transformed into how to filter the biased instances out of the full dataset. Schuster et al. (2019) utilize Local Mutual Information to obtain the n-grams that are highly correlated with a specific label. Then, the claims involving such n-grams are downweighed. Mahabadi et al. (2020) employ a bias-only model to capture biases in claims and assume the unevenness of output label distribution is positively correlated to the confidence of biased instances. However, the confidence estimation is inaccurate observed by some researchers and some calibration methods are further proposed to adjust the estimation (Xiong et al., 2021; Amirkhani and Pilehvar, 2021).

Apart from the mentioned debiasing research pipeline in fact verification, much attention has been paid to incorporating causal inference techniques to obtain more unbiased model. Representative works include counterfactual inference for exposure biases in recommender systems (Tan et al., 2021), implicit knowledge biases and object appearance biases in computer vision (Niu et al., 2021; Sun et al., 2021). However, such pipeline is still under-explored in fact verification. Inspired by these works, we open up a new debiasing pipeline for fact verification from a counterfactual view. Compared to the existing two pipelines, our proposed method is augmentation-free and mitigates biases on the inference stage.

3 Method

In this section, we introduce the proposed debiasing framework CLEVER in detail. Firstly, we provide some background information of fact verification. Then, we describe the method from a causal view. Finally, we elaborate the detail of training and inference. The overview of CLEVER is shown in Figure 2.

3.1 Preliminary

3.1.1 Task Formulation

Given a claim c and its corresponding evidence set $\{e_1, e_2, \dots, e_n\}$, a fact-checking model is required to predict the veracity of claim, i.e., evidences sup-

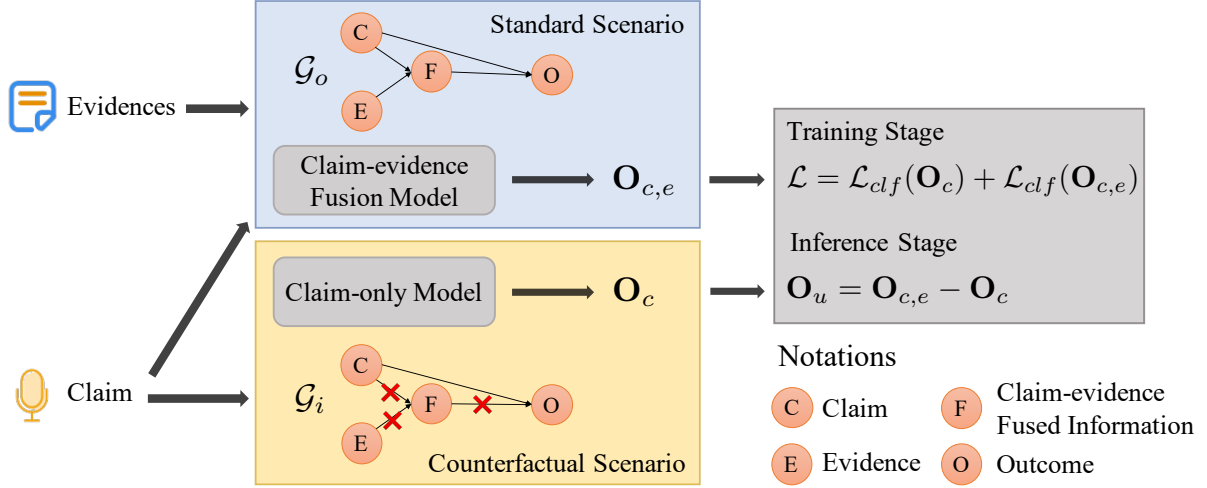


Figure 2: The proposed framework CLEVER. We simulate the standard and counterfactual scenarios via training a claim-evidence fusion model and a claim-only model independently. The final prediction \mathbf{O}_u is obtained by subtracting the output of counterfactual scenario \mathbf{O}_c from that of standard scenario $\mathbf{O}_{c,e}$.

port, refute, or lack enough information to justify the claim.

3.1.2 Causal View of Fact Verification

The causal graph is mathematically a directed acyclic graph, where vertices denote variables and the edge represents the effect from the start vertex to the end vertex.

The causal view of fact verification is represented as a graph $\mathcal{G}_o = \{\mathcal{V}, \mathcal{E}_o\}$, where \mathcal{V} contains four variables with each represents the claim (C), evidences (E), the fusion of claim and evidences (F), and the output (O), respectively (See the standard scenario in Figure 2). In counterfactual scenario, we expect to capture biases in the claim, so we solely preserve the edge from claim to output. Then, we obtain an intervened causal graph \mathcal{G}_i , c.f., the counterfactual scenario in Figure 2.

3.2 The Proposed Framework: CLEVER

In this part, we specifically introduce how to obtain debiased predictions using the counterfactual inference technique.

The first step of counterfactual inference is establishing an imagined scenario different from standard settings. In our task, as shown at the top half of Figure 2, the standard setting is that the outcome is affected by the claim and its corresponding evidences simultaneously in the causal graph \mathcal{G}_o . In practice, we take both claim c and evidences $\{e_1, e_2, \dots, e_n\}$ as inputs to simulate such setting, which can be formulated as:

$$\mathbf{O}_{c,e} = f_s(c, e_1, e_2, \dots, e_n) \quad (1)$$

where f_s denotes the claim-evidence fusion model, n is the number of evidences, and $\mathbf{O}_{c,e} \in \mathbb{R}^L$ denotes the predicted class distribution (L is the number of class).

Then, a key problem in our framework is how to design a counterfactual scenario for debiasing. Causally, if we expect to estimate the effect of a variable on the outcome, we can give the variable a specific treatment while keep other variables unchanged. Since the target of our work is to obtain the unbiased outcomes affected by both claim and evidences, the treatment is to make the claim-evidence fusion information unavailable for the fact-checking model. In other words, as shown at the bottom half of Figure 2, we create a counterfactual scenario \mathcal{G}_i via intervention on the original causal graph \mathcal{G}_o , where the edge from the fused information of claim-evidence pair to the outcome is cut off. In practice, claims are solely fed into a fact-checking model f_b (i.e., claim-only model) to simulate the absence of claim-evidence information and require the model to produce prediction $\mathbf{O}_c \in \mathbb{R}^L$ based on claims solely,

$$\mathbf{O}_c = f_b(c) \quad (2)$$

The second step is comparing the outcomes under standard and counterfactual settings. The output of claim-only model \mathbf{O}_c is biased that simply relies on the spurious correlation between claim patterns and labels. To reduce such biases, inspired

by the Potential Outcomes Model (Sekhon, 2008), we subtract \mathbf{O}_c from $\mathbf{O}_{c,e}$ with a hyperparameter α (named bias coefficient that controls the extent of bias) and obtain the counterfactual debiased output \mathbf{O}_u ,

$$\mathbf{O}_u = \mathbf{O}_{c,e} - \alpha \cdot \mathbf{O}_c \quad (3)$$

In this way, the probability of false biased prediction is decreased while the predicted probability of ground truth is relatively higher.

Training and Inference At training stage, as biases are mainly involved in claims, we expect that the claim-only model captures such biases so that they can be reduced via the subtraction scheme. Motivated by this, we encourage the output of claim-only model \mathbf{O}_c to represent the biased label distribution by imposing a classification loss on \mathbf{O}_c . Similarly, $\mathbf{O}_{c,e}$ is also supervised to mine the claim-evidence interaction. Formally, the objective function can be written as:

$$\mathcal{L} = \mathcal{L}_{clf}(\mathbf{O}_c) + \mathcal{L}_{clf}(\mathbf{O}_{c,e}) \quad (4)$$

where \mathcal{L}_{clf} denotes the cross entropy loss.

At inference stage, since the outcome in counterfactual scenario \mathbf{O}_c is biased after training, we intuitively reduce it via subtraction from the outcome in standard scenario $\mathbf{O}_{c,e}$, c.f., Eq. (3).

Discussion While the proposed framework CLEVER also consists of the claim-evidence model and the claim-only model, which is similar to the weight-regularization-based approaches, we do not rely on the assumption that such two models produce similar outputs for biased instances. Besides, we avoid utilizing the uncertain output of claim-only model to adjust the training loss of claim-evidence model. By contrast, we independently train the claim-evidence and claim-only model and propose a simple yet effective scheme to obtain debiased results on the inference stage.

4 Experiments

In this section, we conduct both quantitative and qualitative experiments on several public datasets to demonstrate the effectiveness of our proposed method CLEVER.

4.1 Experimental Setup

4.1.1 Dataset and Evaluation Metric

We utilize a biased training set FEVER-Train to train models and a biased dataset FEVER-Dev

(Thorne et al., 2018), an unbiased dataset FEVER-Symmetric (Schuster et al., 2019), and an adversarial dataset FEVER-Adversarial (Thorne et al., 2019) to test models, closely following existing works (Mahabadi et al., 2020; Lee et al., 2021; Xiong et al., 2021). Furthermore, we introduce a new subset of FEVER-Dev, namely FEVER-Hard¹, where all samples cannot be correctly classified using claims only. Therefore, it can be used to evaluate the model ability to perform evidence-to-claim reasoning indeed. To further validate the debiasing performance under the multi-hop setting, we augment the dataset Train and Dev with instances consisting of several evidences and generate two multi-hop datasets Train-MH and Dev-MH. Besides, we add the multi-hop instances that cannot be predicted correctly using claims only into Hard and form a new test set Hard-MH. Note that we train all models without using 'NOT ENOUGH INFO' samples since these test sets only involve 'SUPPORTS' and 'REFUTES' samples. Following previous works (Lee et al., 2021), we use label classification accuracy as the metric.

4.1.2 Baselines

We compare our proposed method with several baselines from both two existing pipelines, the specific description is listed as follows:

Data-augmentation-based methods: 1) EDA (Wei and Zou, 2019). They swap words and replace synonym to generate new training samples. 2) CrossAug (Lee et al., 2021). They design a cross contrastive strategy to augment data, where original claims are modified to be negative and evidences are changed to support such negative claims and refute the original claims.

Weight-regularization-based methods: 1) ReW (Schuster et al., 2019). They downweigh the samples which involve n-grams highly correlated to labels. 2) PoE (Mahabadi et al., 2020). They downweigh samples with spurious class distribution outputted from the bias-only model. 3) MoCaD (Xiong et al., 2021). They propose a calibration method to adjust the inaccurate predicted class distribution from bias-only models. Specifically, two calibrators (i.e., temperature scaling and Dirichlet calibrator) are employed in this work. We utilize such methods to further optimize the model PoE, forming two variants namely PoE-TempS and

¹We omit the prefix 'FEVER' for conciseness in following paragraphs since all unbiased and adversarial datasets are derived from the original FEVER dataset.

Dataset	Dev	Symmetric	Hard	Adversarial
BERT-base	<u>93.91 ± 0.14</u>	72.08 ± 0.51	78.05 ± 0.54	61.93 ± 1.31
EDA	93.37 ± 0.42	72.93 ± 0.48	78.22 ± 0.61	62.12 ± 1.02
CrossAug	92.85 ± 0.09	<u>78.88 ± 0.46</u>	82.19 ± 0.31	61.72 ± 0.45
ReW	93.65 ± 0.16	73.39 ± 0.71	78.43 ± 0.52	64.52 ± 1.49
PoE	93.70 ± 0.21	76.43 ± 0.64	80.51 ± 0.70	<u>67.21 ± 1.69</u>
PoE-TempS	93.70 ± 0.25	76.89 ± 0.86	81.13 ± 0.33	67.05 ± 2.30
PoE-Dirichlet	93.25 ± 0.34	78.55 ± 0.97	<u>82.31 ± 0.82</u>	66.98 ± 1.77
CLEVER (ours)	94.10 ± 0.11	84.73 ± 0.69	90.17 ± 0.75	68.34 ± 0.94
Δ Improvement	+ 0.20%	+ 17.55%	+ 15.53%	+ 10.35%

Table 1: The performance comparison between our proposed method CLEVER and baselines. Dev is the biased dataset and other three datasets are introduced to verify the model performance under an unbiased circumstance. The best result on each dataset is highlighted in boldface and the runner-up is underlined. The improvement in terms of percentage compared to the BERT-base is shown in the last row.

PoR-Dirichlet.

4.1.3 Implementation Detail

Following the aforementioned baselines, we employ BERT-base (Devlin et al., 2019) as the backbone model for a fair comparison, i.e., claim-evidence fusion model and claim-only model are two independent BERT models. We finetune BERT with a fully-connected forward layer over the special token [CLS] to obtain the final prediction. The maximum input length is 128, batch size is 32, and the optimizer is Adam with a learning rate of $2e-5$; we train the model for 3 epochs and repeat 5 times under different random seed settings, which are all the same as previous works. The only hyperparameter in our framework is the bias coefficient α . Since α is utilized in inference stage, we do not need to tune it on the validation set. We change the value of α from 0.1 to 1.5 with an increasing step of 0.1. The best performance is achieved on two unbiased datasets Symmetric and Hard when $\alpha = 1.0$ and $\alpha = 1.4$, respectively. On the dataset Adversarial the best value is $\alpha = 0.7$ and that is 0.1 for the biased dataset Dev.

4.2 Performance Comparison

The overall performance of our proposed method CLEVER and baselines is shown in Table 1. We can see that CLEVER outperforms all existing methods from different pipelines by a significant margin on all datasets. More specifically, we have the following observations:

Firstly, the performance gain of CLEVER is more consistent on all datasets than that of pre-

vious methods. We can observe that the runner-up on each dataset is different while CLEVER achieves the best performance on all datasets. More specifically, compared to the vanilla BERT model (i.e., BERT-base) without any debiasing method, CLEVER advances by 17.55% and 15.53% on two unbiased datasets Symmetric and Hard, respectively. Furthermore, most baselines, especially CrossAug, perform relatively worse on the dataset Adversarial, since debiasing methods are always specially designed for avoiding learning biases in claim while do not explicitly consider adversarial attacks. By contrast, our proposed method still achieves a promising result on it (about 10% performance improvement upon the BERT-base), which demonstrates the generalization ability of our method to handle both adversarial and biased data.

Secondly, CLEVER further improves the performance on the biased dataset Dev while all existing debiasing methods suffer from a decline, compared to the BERT-base model. This is because CLEVER captures the biased and unbiased data distribution independently on training stage and adjusts the final prediction on inference stage, which prevents entangling uncertain biased prediction with unbiased one like previous works.

4.3 Study of the Bias Coefficient

The bias coefficient α is introduced in the inference stage, which can be adjusted without tuning according to different properties of datasets. We test the model with several values of α , ranging from 0.1 to 1.5, with a step of 0.1. As illustrated

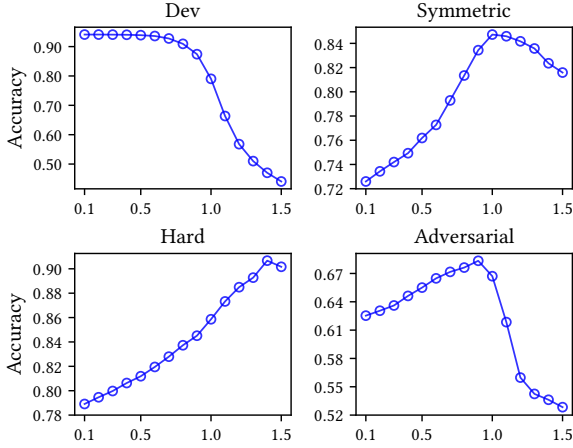


Figure 3: The model performance regarding different values of bias coefficient α .

in Figure 3, the performance on the biased dataset Dev decreases with the growth of bias coefficient. This is reasonable that most of performance gain on the dataset Dev is obtained via exploring claim biases, once the biased factors are alleviated, such performance will be naturally downgraded. On unbiased and adversarial datasets, a similar trend can be seen that the performance first rises to a peak and then drops when α increases. This indicates that 1) the claim-only model successfully captures the biases in claims, which can be mitigated via the proposed subtraction scheme, thus the performance advances when α is enlarged in the early period. 2) Excessively increasing α is harmful for model performance since useful semantic information of claims is reduced.

Furthermore, it is worth noting that the model performance consistently increases until $\alpha = 1.4$ on the dataset Hard, which is larger than $\alpha = 1.0$ on Symmetric and $\alpha = 1.0$ on Adversarial. It indicates that the unbiased extent of Hard is greater than that of Adversarial and Symmetric. Therefore, the dataset Hard can better reflect the model ability of understanding the relationship between claim and evidences. As a result, the consistent performance improvement on Hard further demonstrates the effectiveness of our proposed debiasing strategy.

4.4 Study of Complicated Circumstance

Existing methods only utilize samples with single evidence to evaluate the debiasing performance, however, we argue that more complicated reasoning circumstance should be considered since a

Dataset	Dev-MH	Hard-MH
BERT-base	93.58 ± 0.18	77.99 ± 0.34
PoE	93.34 ± 0.28	80.10 ± 0.49
PoE-TempS	93.42 ± 0.21	81.22 ± 0.55
PoE-Dirichlet	93.36 ± 0.19	82.73 ± 0.58
CLEVER (ours)	93.76 ± 0.14	89.85 ± 0.30
Δ Improvement	+ 0.19%	+ 15.20%

Table 2: The performance comparison between our proposed method CLEVER and baselines under the complicated multi-hop reasoning circumstance.

claim may be verified via several evidences in the realistic scenario. Therefore, we further validate debiasing methods under a multi-hop reasoning setting, where instances with more than one evidences are involved in both biased set Dev-MH and unbiased set Hard-MH. Since data-augmentation methods are hard to be adapted to such complicated scenario, we compare our method CLEVER with baselines from the weight regularization based pipeline. As shown in Table 2, CLEVER consistently outperforms its competitors by a significant margin, which demonstrates its effectiveness of handling complicated data.

4.5 Qualitative Analysis

In this section, we design some case studies to further analyze the advantages of our proposed method CLEVER on a qualitative aspect.

4.5.1 Case Study

In this part, we aim to compare the performance of different models at an instance level. We choose the best debiasing method from each pipeline (i.e., CrossAug and PoE) to carry out the analysis. Specifically, we select representative examples from the dataset Hard that are correctly classified using our method while mistakenly predicted by baselines.

From Figure 4, the top instance shows that the output of claim-evidence fusion model **correctly** inclines to the ground-truth ‘REFUTES’ while the output of claim-only model is **mistakenly** biased towards ‘SUPPORTS’. That is, the claim-evidence fusion model deals with biased instances in a different way from the claim-only model, which echoes the discovery in the previous work (Amirkhani and Pilehvar, 2021). Therefore, PoE downweighs such instance in training objective according to the biased extent of claim-only model would result in performance degradation. However, our method

CLEVER separates such outputs of two models in training and the predicted probability of ground-truth label is further enlarged via subtraction on inference stage.

The bias in the bottom instance is mainly induced by the word ‘is’, which is highly correlated with the label ‘SUPPORTS’. Data-augmentation based methods simply insert negations or antonyms, such as transforming ‘is’ to ‘is not’, are hard to capture the intrinsic conflict between the claim and evidences. In this instance, the conflict lies between ‘Idaho’ and ‘Virginia’, not the word ‘is’. Therefore, augmenting training instances via inserting negations or antonyms contribute little to such complex reasoning circumstance. However, our approach CLEVER directly captures both claim-evidence interactions and claim biases which is augmentation-free. Note that the biased label distribution is alleviated in the claim-evidence fusion model, i.e., the probability of wrong prediction ‘SUPPORTS’ is decreased to 0.89 from 0.98 (See Figure 4(b)), since it partly pays attention to the evidential information. Though the distribution is still biased towards the falsity due to the strong bias between ‘is’ and the label ‘SUPPORTS’, CLEVER can eliminates such bias in both models via subtraction so as to highlight the intrinsic evidential segment, thus providing the correct prediction.

4.5.2 Error Analysis

In this part, we categorize wrong predictions outputed by our method CLEVER into two groups.

The first type of error is induced by the un conspicuous biased features of claims. For example, the claim *Scandinavia includes the remote Norwegian islands of Svalbard and Jan Mayen.* does not contain obvious biases so that the output of claim-only model cannot represent the biased distribution. Therefore, subtracting such output fails to mitigate biases but reduces the beneficial claim information instead. These errors may be avoided by employing different strategies for instances with distinct bias extents, which we leave as future work.

The second type of error occurs when high-level reasoning is required, e.g., mathematical computation and multi-hop reasoning, which drops into the scope of model reasoning ability. This work mainly focuses on debiasing fact-checking models that make them concentrate on the intrinsic evidential information. After debiasing, how to enhance the reasoning ability over such information is a promising future direction.

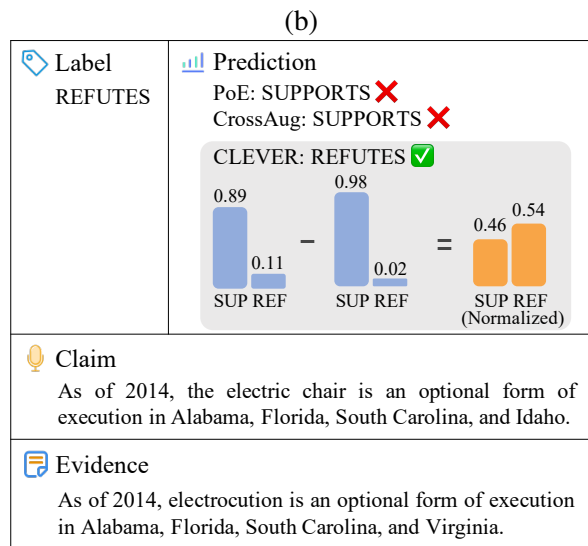
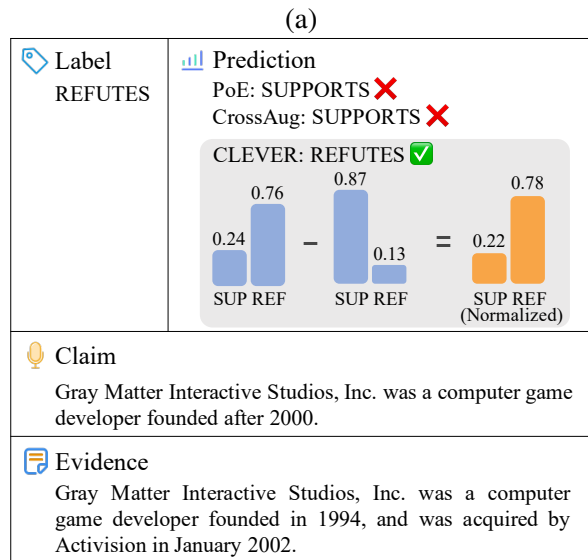


Figure 4: Two representative instances where our proposed method CLEVER outputs correct veracity prediction while baselines make mistakes. The bars denote the outputed label distribution, i.e., $O_u = O_{c,e} - \alpha \cdot O_c$ (Eq. (3)). α is set to be 1.0 for brief illustration.

5 Conclusion

In this paper, we have proposed a novel counterfactual framework CLEVER for debiasing fact-checking models. Unlike existing works, CLEVER is augmentation-free and mitigates biases on inference stage. In CLEVER, the claim-evidence fusion model and the claim-only model are independently trained to capture the corresponding information. On the inference stage, a simple subtraction scheme is proposed to mitigate biases. Comprehensive experiments have demonstrated the superiority of CLEVER.

591
592
593
594

595
596
597
598

599
600
601

602
603
604
605

606
607
608
609
610

611
612
613
614
615
616

617
618
619

620
621
622

623
624
625

626
627
628
629

630
631
632

633
634
635
636
637
638

639
640
641

References

Hunt Allcott and Matthew Gentzkow. 2017. Social media and fake news in the 2016 election. *CSN: Politics (Topic)*.

Hossein Amirkhani and Mohammad Taher Pilehvar. 2021. Don't discard all the biased instances: Investigating a core assumption in dataset bias mitigation techniques. In *EMNLP Findings*.

Qian Chen, Xiao-Dan Zhu, Zhenhua Ling, Si Wei, Hui Jiang, and Diana Inkpen. 2017. Enhanced lstm for natural language inference. In *ACL*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*.

Reza Ghaeini, Sadid A. Hasan, Vivek Datla, Joey Liu, Kathy Lee, Ashequl Qadir, Yuan Ling, Aaditya Prakash, Xiaoli Z. Fern, and Oladimeji Farri. 2018. Dr-bilstm: Dependent reading bidirectional lstm for natural language inference. In *NAACL*.

Andreas Hanselowski, Hao Zhang, Zile Li, Daniil Sorokin, Benjamin Schiller, Claudia Schulz, and Iryna Gurevych. UKP-athene: Multi-sentence textual entailment for claim verification. In *Proceedings of the First Workshop on Fact Extraction and VERification (FEVER)*.

Shimon Kogan, Shimon Kogan, Tobias J. Moskowitz, Tobias J. Moskowitz, and Marina Niessner. 2019. Fake news: Evidence from financial markets.

Neema Kotonya and Francesca Toni. 2020a. Explainable automated fact-checking: A survey. In *COLING*.

Neema Kotonya and Francesca Toni. 2020b. Explainable automated fact-checking for public health claims. In *EMNLP*.

Minwoo Lee, Seungpil Won, Juac Kim, Hwanhee Lee, Cheoneum Park, and Kyomin Jung. 2021. Crossaug: A contrastive data augmentation method for debiasing fact verification models. In *CIKM*.

Nayeon Lee, Belinda Z. Li, Sinong Wang, Wen tau Yih, Hao Ma, and Madian Khabza. 2020. Language models as fact checkers? *ArXiv*, abs/2006.04102.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *ACL*.

Zhenghao Liu, Chenyan Xiong, Maosong Sun, and Zhiyuan Liu. 2020. Fine-grained fact verification with kernel graph attention network. In *ACL*.

Rabeeh Karimi Mahabadi, Yonatan Belinkov, and James Henderson. 2020. End-to-end bias mitigation by modelling biases in corpora. In *ACL*.

Salman Bin Naeem and Rubina Bhatti. 2020. The covid-19 'infodemic': a new front for information professionals. *Health Information and Libraries Journal*.

Yulei Niu, Kaihua Tang, Hanwang Zhang, Zhiwu Lu, Xiansheng Hua, and Ji-Rong Wen. 2021. Counterfactual vqa: A cause-effect look at language bias. In *CVPR*.

Tal Schuster, Darsh J. Shah, Yun Jie Serene Yeo, Daniel Filizzola, Enrico Santus, and Regina Barzilay. 2019. Towards debiasing fact verification models. In *EMNLP*.

Jasjeet S Sekhon. 2008. The neyman-rubin model of causal inference and estimation via matching methods. *The Oxford handbook of political methodology*, 2:1–32.

Amir Soleimani, Christof Monz, and Marcel Worring. 2020. Bert for evidence retrieval and claim verification. *Advances in Information Retrieval*.

Shyam Subramanian and Kyumin Lee. 2020. Hierarchical evidence set modeling for automated fact extraction and verification. In *EMNLP*.

Pengzhan Sun, Bo Wu, Xunsong Li, Wen Li, Lixin Duan, and Chuang Gan. 2021. Counterfactual debiasing inference for compositional action recognition. *Proceedings of the 29th ACM International Conference on Multimedia*.

Juntao Tan, Shuyuan Xu, Yingqiang Ge, Yunqi Li, Xu Chen, and Yongfeng Zhang. 2021. Counterfactual explainable recommendation. In *CIKM*.

James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. Fever: a large-scale dataset for fact extraction and verification. In *NAACL*.

James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2019. Evaluating adversarial attacks against multiple fact verification systems. In *EMNLP*.

Nguyen Vo and Kyumin Lee. 2021. Hierarchical multi-head attentive network for evidence-aware fake news detection. In *EACL*.

Jason Wei and Kai Zou. 2019. Eda: Easy data augmentation techniques for boosting performance on text classification tasks. In *EMNLP*.

Ruibin Xiong, Yimeng Chen, Liang Pang, Xueqi Chen, and Yanyan Lan. 2021. Uncertainty calibration for ensemble-based debiasing methods. In *NIPS*.

Wanjun Zhong, Jingjing Xu, Duyu Tang, Zenan Xu, Nan Duan, M. Zhou, Jiahai Wang, and Jian Yin. 2020. Reasoning over semantic-level graph for fact checking. In *ACL*.

695 Jie Zhou, Xu Han, Cheng Yang, Zhiyuan Liu, Lifeng
696 Wang, Changcheng Li, and Maosong Sun. 2019.
697 Gear: Graph-based evidence aggregating and rea-
698 soning for fact verification. In *ACL*.

699 A Dataset Statistics

Dataset	# SUP	# REF	Sum
Train	100,570	41,850	142,420
Dev	7,983	8,681	16,664
Symmetric	379	338	717
Adversarial	364	402	766
Hard	679	2,638	3,317
Train-MH	120,081	41,850	168,424
Dev-MH	9,214	9,796	19,010
Hard-MH	855	3,027	3,882

Table 3: The statistics of datasets. ‘SUP’ and ‘REF’ is the abbreviation of the label ‘SUPPORTS’ and ‘REFUTES’, respectively. ‘#’ stands for the number of.

700 B Experimental Environment

701 We conduct all experiments using PyTorch 1.8.0 on
702 a single GeForce RTX 662 3090 GPU with 24GB
703 memory. The training and inference process cost
704 about 1 hour and less than 5 minutes, respectively.