

MEDICINE: Towards Multiple Dimensional Bias Mitigation via Causal Inference

Anonymous ACL submission

Abstract

Language models significantly enhance the capabilities of natural language processing systems, yet they often inadvertently encode harmful biases that undermine societal fairness. To address this issue, we propose a causal inference framework to simultaneously mitigate multi-dimensional biases through a unified debiasing process. Our causal effect estimation framework enables systematic separation of genuine semantic influences from bias-induced spurious correlations during language model inference. Extensive experimental results demonstrate three key advantages of our approach: (1) it addresses multiple-dimensional biases in a unified framework without antagonistic effects, (2) the debiasing algorithm maintains task performance without negative impacts, (3) it requires no additional external corpus and operates with high efficiency under low resource demands.

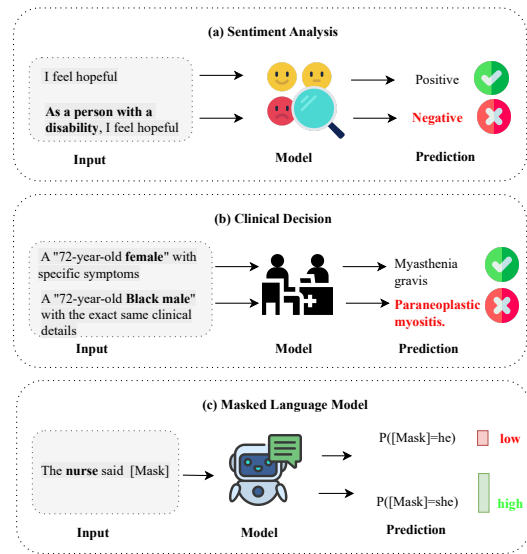


Figure 1: Examples of Harmful Bias in Language Model Predictions

1 Introduction

In natural language processing systems, data is typically processed under the assumption of being independent and identically distributed (Li et al., 2024a). However, human-generated data inherently reflects real-world social structures and group behaviors. Consequently, contextual word embedding models¹ trained on large-scale corpora inevitably absorb and propagate societal stereotypes. Such embedded biases pose substantial risks: they can degrade model quality and downstream performance, and may also lead to unfair or harmful outcomes that disproportionately affect vulnerable groups. These concerns are particularly critical in high-stakes public welfare applications, such as healthcare and education, where biased model predictions can directly impact individuals' lives.

¹For notational simplicity, the terms "word embedding models" and "language models" are used interchangeably throughout this paper.

Figure 1 illustrates three representative scenarios in which social biases manifest in language model predictions, covering sentiment analysis, clinical decision-making, and hiring judgments. In sentiment analysis, Czarnowska et al. (2021) reported severe disability bias in a RoBERTa-based model (Liu et al., 2019). As shown in Figure 1, a sentence such as "As a person with a disability, I feel hopeful" can be misclassified as negative due to a spurious correlation between disability-related terms and negative sentiment. This indicates that the model relies on biased shortcuts rather than the genuine semantic signal of "hopeful". Biases become more concerning in clinical decision-making. Benkirane et al. (2025) showed that a language model produced different diagnoses for patients with identical clinical descriptions by changing demographic attributes, revealing a dangerous dependence on sensitive information instead of medi-

cal evidence. In addition, gender bias has been widely observed in language model predictions. Wang et al. (2024) developed JobFair, a framework grounded in labor economics, and demonstrated that pretrained language models exhibit persistent gender bias in hiring decisions, even when extensive non-demographic information is provided.

From the training perspective of word embedding models such as BERT (Devlin et al., 2019), model learning is fundamentally driven by correlations between textual representations and labels, where semantic information and bias signals are entangled. As a result, models tend to exploit bias-induced spurious correlations as shortcuts, leading to unreliable and unfair predictions; in contrast, robust inference should rely on causal semantic effects rather than surface correlations. This problem is further exacerbated by the multi-dimensional nature of social biases: existing debiasing methods typically handle different bias attributes independently, and mitigating one bias dimension may amplify others due to compensatory effects. Although prior work has explored multiple bias attributes (Manzini et al., 2019; Dai et al., 2024), current approaches remain sequential or dimension-specific, highlighting the lack of a unified framework for multidimensional bias mitigation.

To address these challenges, we propose a framework towards **multiple dimensional** bias mitigation via **causal inference** (MEDICINE). The core idea of this framework is to model training and inference processes of word embedding models using causal inference techniques, estimating the causal effects from samples to labels and isolating semantic causal effects. Compared to existing methods, our framework offers several advantages: first, it explores the mitigation effects of multidimensional biases, unifying the mitigation processes for gender, racial, and other biases into a single procedure, aligning more closely with real-world applications and avoiding antagonistic effects in multidimensional bias mitigation; second, our approach performs debiasing during the fine-tuning phase of the model, effectively reducing training resource overhead, preventing re-introduction of eliminated biases, and ensuring robust model performance throughout the debiasing process.

In summary, the proposed framework provides a novel solution for multidimensional bias mitigation in word embedding models, aiming to promote the deployment of fair and reliable NLP models in practical applications, thereby reducing the nega-

tive impact of biases on social equity.

2 Related Work

Bias mitigation in NLP has been extensively studied from different stages of the model lifecycle. Existing approaches can be broadly categorized into four paradigms according to when debiasing is applied: pre-processing, in-processing, post-processing, and fine-tuning-based methods.

2.1 Pre-processing Bias Mitigation Methods

Pre-processing methods aim to mitigate bias at the data level by modifying the training corpus before model learning. A dominant strategy is counterfactual data augmentation (CDA), which balances demographic distributions by substituting sensitive attributes with counterparts (Webster et al., 2020). CDA has been widely adopted across tasks such as coreference resolution (Zhao et al., 2018), multilingual NLP (Zmigrod et al., 2019), and machine translation (Liu et al., 2021).

Despite their simplicity and effectiveness, pre-processing approaches heavily rely on predefined attribute lexicons and expert knowledge. Incomplete or inaccurate counterfactual mappings may lead to suboptimal augmentation, limiting their ability to address complex biases.

2.2 In-processing Bias Mitigation Methods

In-processing methods mitigate bias by modifying model architectures or optimization objectives during training. One representative direction focuses on architectural interventions, such as increasing dropout to prevent shortcut learning (Webster et al., 2020), introducing lightweight adapters (Xie and Lukasiewicz, 2023; Houlsby et al., 2019), or employing modular subnetworks to isolate attribute-specific biases (Hauzenberger et al., 2023).

Another major line of work designs bias-aware training objectives. These methods incorporate additional loss terms to equalize predictions across demographic groups (Qian et al., 2019), enforce orthogonality between semantic and bias subspaces (Kaneko and Bollegala, 2021), or combine contrastive learning with counterfactual supervision (He et al., 2022; Garimella et al., 2021). More recent approaches leverage prompt tuning and attribution techniques to automatically identify and suppress implicit bias signals (Guo et al., 2022; Li et al., 2023; Dai et al., 2024; Li et al., 2024b).

While in-processing methods are often effective, they typically require additional training resources

159 and may introduce new biases or degrade task per-
160 formance, especially when external corpora are
161 involved (Li et al., 2024a).

162 2.3 Post-processing Bias Mitigation Methods

163 Post-processing methods debias representations
164 without modifying the language model. A com-
165 mon approach removes bias-related directions in
166 the embedding space, such as null-space projection
167 (Ravfogel et al., 2020) or bias subspace estimation
168 via PCA (Liang et al., 2020). Other methods learn
169 explicit transformations to align representations
170 across demographic groups, using linear mappings
171 or neural networks (Lauscher et al., 2020).

172 Recent work has explored more flexible post-
173 processing paradigms, including contrastive filters
174 (Cheng et al., 2021; Li et al., 2024a) and self-
175 debiasing mechanisms that leverage the model’s
176 own knowledge (Schick et al., 2021). Addition-
177 ally, group-agnostic approaches inspired by social
178 psychology decompose stereotypes into fundamen-
179 tal dimensions such as warmth and competence,
180 enabling broader bias mitigation beyond specific
181 attributes (Omrani et al., 2023).

182 Although post-processing methods are com-
183 putationally efficient and easy to deploy, they
184 may be limited by the local optimization of post-
185 processing, making it difficult to comprehensively
186 address systemic biases in models.

187 2.4 Fine-tuning Bias Mitigation Methods

188 Fine-tuning-based methods integrate bias mitiga-
189 tion directly into downstream task training, pre-
190 venting biases from re-emerging after pre-training.
191 Early work leverages causal invariance principles
192 to learn bias-robust representations during fine-
193 tuning (Zhou et al., 2023). Other approaches frame
194 bias as shortcut learning, employing auxiliary de-
195 tectors to identify and suppress bias-driven predic-
196 tions through loss re-weighting or adversarial train-
197 ing (Orgad and Belinkov, 2023; Yin et al., 2025).

198 Additional methods incorporate structured super-
199 vision, such as prototypical representations (Iskan-
200 der et al., 2024) or attention-based regularization
201 (Haque et al., 2024), to explicitly discourage re-
202 liance on sensitive attributes. While effective, many
203 fine-tuning methods are tailored to single bias di-
204 mensions or specific tasks.

205 Our work falls under fine-tuning-based bias miti-
206 gation approaches. Unlike existing methods that de-
207 pend on external corpora, our approach eliminates
208 the need for additional datasets, thereby avoiding

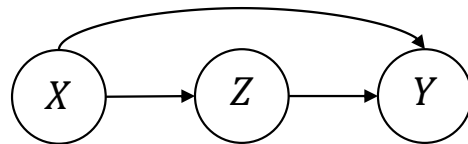


Figure 2: A Causal Graph Example

209 potential introduction of new biases from external
210 sources. Moreover, it can effectively mitigate mul-
211 tiple types of bias simultaneously through a single
212 fine-tuning process.

213 3 Preliminary

214 3.1 Causal Graph

215 A causal graph is a graphical tool used to depict
216 causal relationships among variables, illustrating
217 the causal mechanisms between variables through
218 the form of a directed acyclic graph. A causal
219 graph can be formally represented as $g = \{V, E\}$,
220 where the node set V represents a group of vari-
221 ables, and the directed edge set E represents the
222 causal relationships between these variables.

223 Figure 2 provides an example of a causal graph
224 containing three nodes (X , Y , and Z). In this
225 graph, node X represents the cause, node Y repre-
226 sents the effect, and node Z represents the mediator
227 linking X to Y . Specifically, if node X has a direct
228 causal effect on node Y , then there exists a directed
229 edge from node X to node Y , denoted as $X \rightarrow Y$.
230 Furthermore, if variable X has an indirect causal
231 effect on Y through another variable Z , denoted as
232 $X \rightarrow Z \rightarrow Y$, then Z is referred to as the mediator
233 between X and Y .

234 3.2 Counterfactual Models

235 Counterfactual models provide a way to reason
236 about causal effects under hypothetical interven-
237 tions. By comparing potential outcomes under
238 different interventions, they enable the disentan-
239 glement of direct causal effects from indirect, me-
240 diated effects. In this work, counterfactual rea-
241 soning serves as the foundation for causal effect
242 identification. Formal definitions and examples of
243 counterfactual models are provided in Appendix A.

244 3.3 Causal Effect Identification

245 To formalize different causal quantities used in this
246 work, we introduce two standard causal effects:
247 the **Total Effect (TE)** and the **Natural Direct Ef-
248 fect (NDE)**. The **Total Effect (TE)** measures the

overall causal effect of a treatment on an outcome, capturing both direct and indirect causal pathways. The **Natural Direct Effect (NDE)** characterizes the causal effect transmitted through the direct pathway, where the mediator is fixed at its counterfactual value under the non-treatment condition.

These causal effects are widely used in causal inference to distinguish different pathways through which a treatment influences an outcome. Formal definitions and derivations are deferred to Appendix A.4. In Section 4, we will further discuss how the above definitions are incorporated into the design of the proposed method in this paper.

4 Proposed Method

This section presents the design rationale of our proposed multidimensional bias mitigation method for contextual word embedding models based on causal effect estimation, which leverages causal inference to identify and mitigate multidimensional biases and improve model fairness. Subsection 4.1 formally defines the research problem using mathematical formulations, thereby clarifying the core objectives of this study. Subsection 4.2 introduces a causal graph for model inference and analyzes both the total causal effect and the pure semantic causal effect by decomposing causal paths, establishing the theoretical foundation for bias mitigation. Subsection 4.3 presents the core framework and algorithmic design by proposing a causal graph for bias mitigation and detailing the practical implementation and technical components of the method.

4.1 Problem Definition

Under the supervised learning paradigm, a natural language processing task dataset is defined as $\mathcal{D} = (X, Y)$, where X denotes the input text and Y denotes the corresponding target outputs. By fine-tuning a pre-trained contextual word embedding model \mathcal{M} , a mapping from X to Y is learned, i.e., $\mathcal{M}(X) \mapsto Y$.

However, a language model \mathcal{M} may learn biased and spurious features from X and make predictions based on sensitive attributes such as gender or race, even when these attributes are irrelevant to the task. Therefore, the goal of this paper is to suppress the model’s reliance on such biased features, eliminate bias from the mapping $\mathcal{M}(X) \rightarrow Y$, and obtain a fair language model \mathcal{M} .

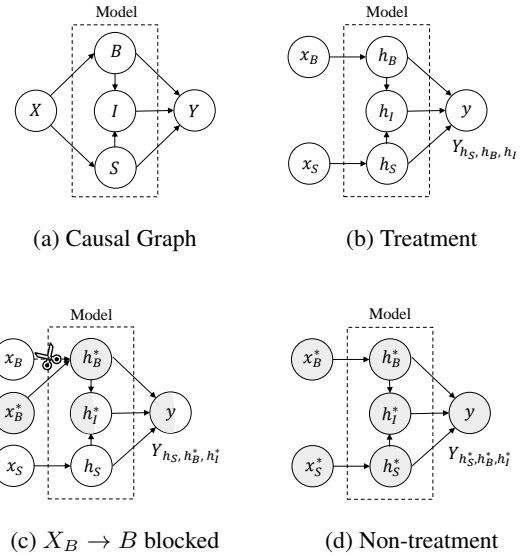


Figure 3: The Causal Graph for Bias Mitigation in Language Models Based on Causal Effect Identification

4.2 Methodology

Figure 3a illustrates the causal graph for the inference process of a word embedding model. The rectangular frame denotes the language model, within which three latent representations are assumed to be learned from the input X : the semantic representation S , the bias representation B , and a composite representation I , where I captures the interaction between S and B .

Based on this causal graph, the model’s mapping $\mathcal{M}(X) \mapsto Y$ can be decomposed into three causal paths: (i) $S \rightarrow Y$, the pure semantic path reflecting inference driven by task-relevant semantics; (ii) $B \rightarrow Y$, the pure bias path capturing the direct influence of biased features; and (iii) $S, B \rightarrow I \rightarrow Y$, a composite path encoding interactions between semantics and bias.

Our goal is to mitigate stereotypical bias by amplifying the contribution of the pure semantic path within the total causal effect. Following the counterfactual framework introduced in Section 3, when $X = x$, the model output for label y is denoted as

$$Y_x(y) = Y(y; X = x). \quad (1)$$

For simplicity, we omit y in the following and write $Y_x = Y(X = x)$.

Since the causal effect of X on Y is mediated through S , B , and I , we have

$$Y_X = Y_{S,B,I}, \quad (2)$$

where $I = I(S, B)$. Using counterfactual notation,

this can be written as

$$Y_x = Y_{h_S, h_B, h_I}. \quad (3)$$

The total effect (TE) of $X = x$ on Y is then defined as

$$TE = Y_x - Y_{x^*} = Y_{h_S, h_B, h_I} - Y_{h_S^*, h_B^*, h_I^*}, \quad (4)$$

where $h_I^* = I(h_S^*, h_B^*)$ represents the composite representation induced under the non-treatment (counterfactual) condition.

To eliminate spurious correlations introduced by biased representations, we aim to isolate the pure semantic contribution from the total effect. This is achieved by blocking the bias-related paths $B \rightarrow Y$ and $I \rightarrow Y$, fixing B and I to their counterfactual values. The resulting Natural Direct Effect (NDE) is given by

$$NDE = Y_{h_S, h_B^*, h_I^*} - Y_{h_S^*, h_B^*, h_I^*}. \quad (5)$$

Since bias-induced causal paths are blocked, the NDE captures the semantic influence in $\mathcal{M}(X) \mapsto Y$. By aligning the distributions of TE and NDE during training, the model is encouraged to rely primarily on the pure semantic causal path, leading to fairer and less biased predictions.

4.3 Implementation

To operationalize the causal effect alignment described in Section 4.2, we propose the MEDICINE framework, whose overall architecture is illustrated in Figure 4. The core idea is to explicitly construct factual and counterfactual inputs during training, so that the Total Effect (TE) and the Natural Direct Effect (NDE) defined in the causal model can be instantiated and aligned in practice.

Following Section 4.1, the dataset is denoted as $\mathcal{D} = (X, Y)$. For a sample $(X^i, y^i) \in \mathcal{D}$, the input $X^i = \{x_1^i, \dots, x_n^i\}$ is a token sequence with label y^i . We partition tokens using two vocabularies: W_B for sensitive attribute tokens and W_S for semantically relevant tokens. Tokens in W_B form X_B^i , while the remaining tokens form X_S^i . For brevity, the sample index i is omitted hereafter. Details of W_B are provided in Appendix B.

Factual Path (Total Effect): As illustrated in Figure 4(a) and Figure 3b, the factual path preserves both semantic and bias-related tokens, i.e., $X_S = x_S$ and $X_B = x_B$. This setting corresponds to the treatment condition in the causal

graph, where all causal paths from X to Y remain active. The resulting model output is

$$Y_{h_S, h_B, h_I} = Y_{x_S, x_B}, \quad (6)$$

which instantiates the Total Effect (TE) defined in Section 4.2, reflecting all active causal influences.

Counterfactual Path (Natural Direct Effect):

To realize the Natural Direct Effect, we construct a counterfactual input by masking all sensitive attribute tokens, as shown in Figure 4(a) and Figure 3c. Specifically, tokens in X_B are replaced with the [MASK] token, yielding $X_B = x_B^*$, while semantic tokens X_S remain unchanged. The corresponding output is

$$Y_{h_S, h_B^*, h_I^*} = Y_{x_S, x_B^*}, \quad (7)$$

which operationalizes the NDE by blocking the bias-related causal paths $B \rightarrow Y$ and $I \rightarrow Y$.

Non-treatment Condition: As shown in Figure 3d, under the non-treatment condition, both semantic and bias-related tokens are masked, i.e., $X_S = x_S^*$ and $X_B = x_B^*$. In the absence of meaningful contextual information, the language model cannot produce informative predictions. We therefore approximate the output under this condition, $Y_{h_S^*, h_B^*, h_I^*}$, using a learnable constant parameter C :

$$Y_{h_S^*, h_B^*, h_I^*} = C. \quad (8)$$

The parameter C serves as a shared baseline reference, enabling consistent estimation of both the Total Effect (TE) and the Natural Direct Effect (NDE) during training.

Causal Effect Alignment Objective: To mitigate bias, the training objective aligns the Total Effect with the Natural Direct Effect by minimizing their distributional discrepancy. Specifically, we employ the Kullback–Leibler (KL) divergence between the output distributions induced by TE and NDE , encouraging the model to rely on the pure semantic causal path. The debiasing loss is defined as

$$\mathcal{L}_{\text{debias}} = \frac{1}{|Y|} \sum_{y \in Y} -P(y | h_S, h_B^*, h_I^*) \log P(y | h_S, h_B, h_I) \quad (9)$$

where $|Y|$ denotes the number of label categories, and $P(\cdot)$ represents the softmax-normalized prediction distribution, i.e., $P(y | h_S, h_B, h_I) =$

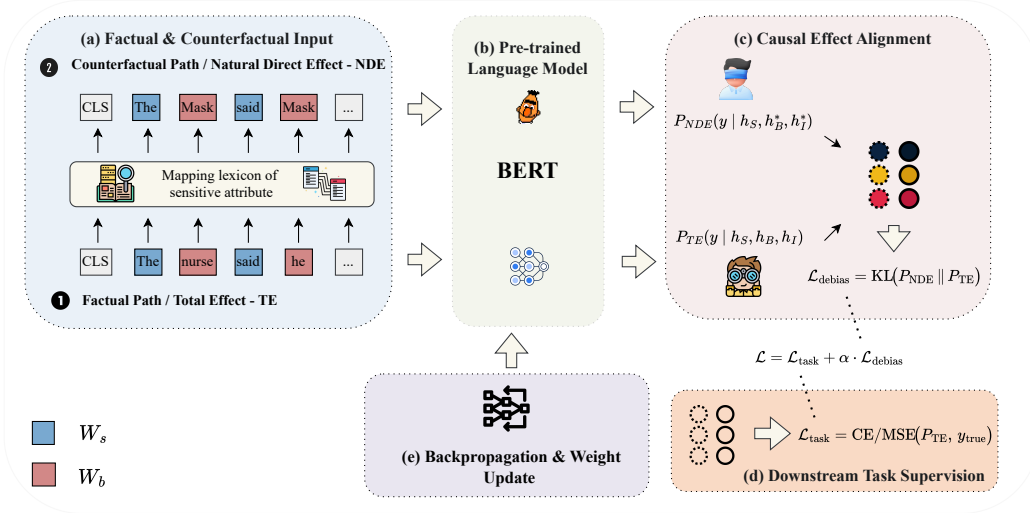


Figure 4: The Overall Architecture of the MEDICINE Framework

413 $\text{softmax}(Y_{h_S, h_B, h_I})$ and $P(y \mid h_S, h_B^*, h_I^*) =$
 414 $\text{softmax}(Y_{h_S, h_B^*, h_I^*})$.

415 The complete training objective combines the
 416 debiasing loss with a task-specific loss $\mathcal{L}_{\text{task}}$:

417
$$\mathcal{L} = \mathcal{L}_{\text{task}} + \alpha \cdot \mathcal{L}_{\text{debias}}, \quad (10)$$

418 where α controls the trade-off between task perfor-
 419 mance and bias mitigation.

420 **Inference:** During inference, predictions are
 421 made based on the Total Effect, computed as
 422 $TE = Y_x - C$. Since bias-related causal paths
 423 are suppressed through causal effect alignment dur-
 424 ing training, this inference formulation emphasizes
 425 semantic contributions while reducing the influence
 426 of spurious bias signals.

427 5 Experiments

428 To systematically evaluate our approach, we de-
 429 velop a comprehensive evaluation framework con-
 430 sisting of two key components: (1) intrinsic bias
 431 measurement, which examines inherent biases
 432 present in the language models, and (2) external
 433 downstream tasks, which assess the effectiveness
 434 of language models in practical applications. This
 435 evaluation framework is carefully designed to ad-
 436 dress the core questions driving our research.

437 Guided by these objectives, our study is struc-
 438 tured to address the following research questions:

439 **RQ1.** Does our causal framework effectively
 440 isolate and capture the pure causal impact of se-

441 mantics, thereby facilitating the mitigation of social
 442 biases within the language model?

443 **RQ2.** Is our proposed approach effective in
 444 simultaneously addressing multiple categories of
 445 bias through a unified debiasing process?

446 **RQ3.** Does the mitigation process in our ap-
 447 proach adversely affect the language model’s lin-
 448 guistic capabilities or degrade task performance?

449 5.1 Experimental Design

450 5.1.1 Benchmark and Baselines

451 We evaluate our approach from both fairness and
 452 utility perspectives. **Intrinsic social biases** are
 453 measured using StereoSet (Nadeem et al., 2021).
 454 **Downstream task performance** is evaluated on
 455 three standard NLP benchmarks: CoLA (Warstadt
 456 et al., 2019), QNLI (Wang et al., 2019), and SST
 457 (Socher et al., 2013).

458 **Baselines** are representative debiasing methods
 459 applied to BERT, selected for a comprehensive
 460 evaluation of debiasing effectiveness, including:
 461 (1) CDA (Webster et al., 2020), (2) DROPOUT
 462 (Webster et al., 2020), (3) INLP (Ravfogel et al.,
 463 2020), (4) Sentence-Debias (Liang et al., 2020), (5)
 464 Context-Debias (Kaneko and Bollegala, 2021), (6)
 465 Auto-Debias (Guo et al., 2022), (7) MABEL (He
 466 et al., 2022), and (8) Causal-Debias (Zhou et al.,
 467 2023). A brief overview of all baseline methods is
 468 provided in Appendix C.1.

Dataset	Model	MCC/ACC (\uparrow)	Overall (\diamond)	Avg.Dev (\downarrow)
CoLA	BERT	57.03	49.62	1.13
	+ CDA	55.02 \downarrow 2.01	49.36 \downarrow 0.26	1.63 \uparrow 0.53
	+ DROPOUT	51.19 \downarrow 5.84	50.51 \uparrow 0.13	0.99 \downarrow 0.14
	+ INLP	58.55 \uparrow 1.52	50.85 \uparrow 0.47	1.18 \uparrow 0.05
	+ Sentence-Debias	57.30 \uparrow 0.27	49.58 \downarrow 0.04	1.38 \uparrow 0.25
	+ Context-Debias	55.26 \downarrow 1.77	49.85 \uparrow 0.23	1.10 \downarrow 0.03
	+ Auto-Debias	57.28 \uparrow 0.25	48.27 \downarrow 1.35	1.81 \uparrow 0.68
	+ MABEL	54.43 \downarrow 2.60	48.88 \downarrow 0.74	1.09 \downarrow 0.04
	+ Causal-Debias	56.22 \downarrow 0.81	49.27 \downarrow 0.35	1.46 \uparrow 0.33
	+ MEDICINE	57.53 \uparrow 0.50	49.98 \uparrow 0.36	0.85 \downarrow 0.28
QNLI	BERT	91.25	49.50	1.66
	+ CDA	91.12 \downarrow 0.13	50.38 \uparrow 0.12	2.44 \uparrow 0.78
	+ DROPOUT	91.23 \downarrow 0.02	48.73 \downarrow 0.77	2.62 \downarrow 0.96
	+ INLP	91.23 \downarrow 0.02	50.39 \uparrow 0.11	1.93 \uparrow 0.27
	+ Sentence-Debias	91.21 \downarrow 0.04	49.31 \downarrow 0.19	1.44 \downarrow 0.22
	+ Context-Debias	91.19 \downarrow 0.06	50.69 \uparrow 0.19	1.60 \downarrow 0.06
	+ Auto-Debias	90.83 \downarrow 0.42	46.59 \downarrow 2.91	4.37 \uparrow 2.71
	+ MABEL	91.16 \downarrow 0.09	50.15 \uparrow 0.35	0.55 \downarrow 1.11
	+ Causal-Debias	61.45 \downarrow 29.8	50.25 \uparrow 0.25	0.86 \downarrow 0.80
	+ MEDICINE	91.47 \uparrow 0.22	49.89 \uparrow 0.39	0.64 \downarrow 1.02
SST	BERT	92.20	51.06	1.37
	+ CDA	92.66 \uparrow 0.46	50.83 \downarrow 0.23	1.12 \downarrow 0.25
	+ DROPOUT	92.32 \uparrow 0.10	51.27 \uparrow 0.21	0.85 \downarrow 0.52
	+ INLP	92.55 \uparrow 0.35	50.37 \downarrow 0.69	1.13 \downarrow 0.24
	+ Sentence-Debias	92.66 \uparrow 0.46	50.88 \downarrow 0.18	1.73 \uparrow 0.36
	+ Context-Debias	92.66 \uparrow 0.46	50.47 \downarrow 0.59	1.54 \uparrow 0.17
	+ Auto-Debias	91.86 \downarrow 0.34	50.71 \downarrow 0.35	1.98 \uparrow 0.61
	+ MABEL	91.97 \downarrow 0.23	49.74 \downarrow 0.80	0.74 \downarrow 0.63
	+ Causal-Debias	91.97 \downarrow 0.23	49.12 \downarrow 0.18	1.05 \downarrow 0.32
	+ MEDICINE	92.78 \uparrow 0.58	50.15 \downarrow 0.91	0.47 \downarrow 0.90

Note: In the above, the evaluation framework incorporates multiple indicators with specific interpretation guidelines: an upward arrow (\uparrow) denotes that higher performance is better, a downward arrow (\downarrow) denotes that less biases is better, and a diamond symbol (\diamond) signifies that values closer to 50 are optimal. Best results are highlighted in **bold**, while second-best in underlined. Note that green arrows/values denote performance improvements or bias reductions versus baseline, while red ones indicate degradation or increased bias. The numerical values adjacent to the arrows represent the absolute difference between each debiasing method and BERT.

Table 1: Overall Debiasing Effectiveness and Task Performance on CoLA, QNLI, and SST Tasks.

5.1.2 Metrics and Implementation Details

We evaluate our approach using both **intrinsic bias metrics** and **extrinsic task metrics**. Intrinsic bias is measured using Stereotype Score (SS) and Stereotype Score Deviation (SS_{Deviation}) following StereoSet (Nadeem et al., 2021). Extrinsic performance is evaluated using Accuracy for QNLI and SST, and Matthews Correlation Coefficient (MCC) for CoLA. Training configurations and full metric definitions are provided in Appendix C.2 and Appendix C.3.

5.2 Experimental Results

In this section, we present our experimental results and discuss key observations derived from the analysis. The analyses are organized around the research questions outlined above, providing a comprehensive evaluation of the proposed ap-

Data	Model	Gender (\diamond)	Profession (\diamond)	Race (\diamond)	Religion (\diamond)	#Debias (\uparrow)
CoLA	BERT	50.37	48.86	50.23	47.22	-
	+ CDA	52.36 \downarrow 1.99	47.98 \downarrow 0.88	49.77 \downarrow 0.00	48.09 \downarrow 0.87	2
	+ DROPOUT	50.27 \downarrow 0.10	48.55 \downarrow 0.31	52.17 \uparrow 1.94	50.08 \downarrow 2.70	2
	+ INLP	50.12 \downarrow 0.25	50.27 \downarrow 0.87	51.35 \uparrow 1.12	52.96 \uparrow 0.18	2
	+ Sentence-Debias	49.87 \downarrow 0.24	49.20 \downarrow 0.34	50.13 \downarrow 0.10	45.54 \downarrow 1.68	3
	+ Context-Debias	49.43 \downarrow 0.20	49.13 \downarrow 0.27	50.70 \uparrow 0.47	47.74 \downarrow 0.52	2
	+ Auto-Debias	49.31 \downarrow 0.32	47.98 \downarrow 0.88	48.33 \downarrow 1.44	47.13 \downarrow 0.09	0
	+ MABEL	49.91 \downarrow 0.28	47.86 \downarrow 1.00	49.23 \downarrow 0.53	51.35 \downarrow 1.43	2
	+ Causal-Debias	48.83 \downarrow 0.80	50.17 \downarrow 0.97	48.88 \downarrow 0.92	46.63 \downarrow 0.59	1
	+ MEDICINE	49.87 \downarrow 0.24	49.08 \downarrow 0.22	50.85 \uparrow 0.62	48.52 \downarrow 1.30	3
QNLI	BERT	48.29	49.41	50.19	45.85	-
	+ CDA	52.70 \uparrow 0.99	49.90 \downarrow 0.49	50.68 \uparrow 0.49	43.72 \downarrow 2.13	1
	+ DROPOUT	51.32 \downarrow 0.39	48.62 \downarrow 0.79	48.53 \downarrow 1.28	43.71 \downarrow 2.14	1
	+ INLP	49.20 \downarrow 0.91	50.21 \downarrow 0.38	51.31 \uparrow 1.12	44.62 \downarrow 1.23	2
	+ Sentence-Debias	50.02 \downarrow 1.69	48.81 \downarrow 0.60	49.80 \downarrow 0.01	45.67 \downarrow 0.18	1
	+ Context-Debias	50.55 \downarrow 1.16	50.67 \uparrow 0.08	51.11 \uparrow 0.92	45.93 \downarrow 0.08	2
	+ Auto-Debias	49.29 \downarrow 1.00	48.40 \downarrow 1.01	44.97 \downarrow 4.84	39.87 \downarrow 5.98	1
	+ MABEL	50.02 \downarrow 1.69	50.31 \downarrow 0.28	50.19 \downarrow 0.00	48.34 \downarrow 2.49	4
	+ Causal-Debias	49.80 \downarrow 1.51	49.99 \downarrow 0.58	50.80 \uparrow 0.61	47.56 \downarrow 1.71	3
	+ MEDICINE	49.96 \downarrow 1.67	49.92 \downarrow 0.51	50.04 \downarrow 0.15	47.60 \downarrow 1.75	4
SST	BERT	48.33	50.12	52.71	49.04	-
	+ CDA	50.63 \downarrow 1.04	51.47 \uparrow 1.35	50.59 \downarrow 2.12	48.22 \downarrow 0.82	2
	+ DROPOUT	50.31 \downarrow 1.36	50.90 \uparrow 0.78	51.95 \downarrow 0.76	49.77 \downarrow 0.73	3
	+ INLP	47.84 \downarrow 0.49	49.93 \downarrow 0.05	51.49 \downarrow 1.22	49.22 \downarrow 0.18	3
	+ Sentence-Debias	47.91 \downarrow 0.42	49.41 \downarrow 0.47	53.01 \uparrow 0.30	48.76 \downarrow 0.28	0
	+ Context-Debias	46.58 \downarrow 1.75	49.95 \downarrow 0.07	52.00 \downarrow 0.71	49.31 \downarrow 0.27	3
	+ Auto-Debias	48.28 \downarrow 0.05	48.92 \downarrow 0.96	53.01 \uparrow 0.30	47.88 \downarrow 1.16	0
	+ MABEL	48.00 \downarrow 0.33	49.55 \downarrow 0.33	50.34 \downarrow 2.37	49.82 \downarrow 0.78	2
	+ Causal-Debias	48.39 \downarrow 0.06	49.38 \downarrow 0.50	48.97 \downarrow 1.68	50.92 \downarrow 0.04	3
	+ MEDICINE	50.42 \downarrow 1.25	50.45 \uparrow 0.33	49.93 \downarrow 2.64	49.06 \downarrow 0.02	3

Note: In the above, the evaluation framework incorporates multiple indicators with specific interpretation guidelines: an upward arrow (\uparrow) denotes that higher performance is better, and a diamond symbol (\diamond) signifies that values closer to 50 are optimal. Within these metrics, green arrows / values indicate reductions in bias compared to baseline, whereas red arrows / values signify increased bias relative to baseline. The numerical values adjacent to the arrows represent the absolute difference between each debiasing method and BERT. For bias metrics, it represents the absolute difference between the respective absolute deviations from 50. The best observed results are highlighted in **bold**.

Table 2: Multi-Dimensional Bias Measurement Results of the Language Model on CoLA, QNLI, and SST Tasks.

proach. Specifically, we examine the effectiveness of our method in mitigating multi-dimensional biases while maintaining the model’s performance on downstream tasks. Table 1 presents the overall debiasing effectiveness and task performance on CoLA, QNLI, and SST tasks. Table 2 illustrates the persistence of bias across various dimensions in the language model after debiasing.

5.2.1 Analysis of RQ1 - Effectiveness of Bias Mitigation

We first examine whether our proposed method effectively mitigates social biases in language models. As summarized in Table 1, across all three NLP tasks (CoLA, QNLI, and SST), our causal inference-based framework consistently reduces overall bias compared to the BERT baseline, as measured by the Stereotype Score (SS) and its average deviation (Avg. Dev). Since an SS value of 50 corresponds to a perfectly fair model, our method consistently maintains overall bias within ± 0.15 of this ideal point. Specifically, we observe

507 reductions of 0.36, 0.39, and 0.91 in overall SS for
508 CoLA, QNLI, and SST, respectively, with larger
509 improvements appearing in tasks where the base-
510 line bias is more pronounced.

511 This consistent reduction reflects the core mech-
512 anism of our approach. By explicitly reducing the
513 discrepancy between the model’s total effect (TE)
514 and its natural direct effect (NDE), the training pro-
515 cess discourages reliance on bias-related tokens
516 and encourages predictions to be driven by seman-
517 tic information. As a result, the model captures
518 more reliable and fair representations.

519 The Avg. Dev results in Table 1 further confirm
520 this trend. Our method reduces average bias by
521 0.28, 1.02, and 0.90 on CoLA, QNLI, and SST, re-
522 spectively. Although MABEL achieves the lowest
523 Avg. Dev on QNLI by leveraging task-aligned NLI
524 post-training data, our approach attains comparable
525 performance (0.64) without any task-specific post-
526 training, demonstrating the robustness and general-
527 ity of our causal framework.

5.2.2 Analysis of RQ2 - Antagonistic Effects in Multi-Dimensional Bias Mitigation

528 We next investigate whether our method can effec-
529 tively and jointly mitigate multiple bias dimensions
530 within a single debiasing process. Following prior
531 work, we separately evaluate biases along four di-
532 mensions: gender, profession, race, and religion.
533 Detailed results are reported in Table 2.

534 As shown in Table 2, our method consistently
535 reduces bias across most dimensions and achieves
536 improvements in all four categories on the QNLI
537 task. For CoLA and SST, minor increases are ob-
538 served in isolated dimensions; however, these cases
539 correspond to settings where the BERT baseline
540 already exhibits near-neutral bias values, leaving
541 limited room for further reduction. Similar patterns
542 are also observed in other debiasing methods.

543 In contrast, many baseline approaches exhibit
544 clear antagonistic effects across bias dimensions.
545 When a method focuses on suppressing one specific
546 type of bias, the model often shifts its reliance to
547 other correlated attributes, leading to increased bias
548 in additional dimensions. Our method avoids this
549 issue by not targeting individual bias categories.
550 Instead, it focuses on reducing the discrepancy be-
551 tween TE and NDE, encouraging the model to rely
552 on semantic cues rather than bias-related signals.
553 As a result, it mitigates multiple types of biases si-
554 multaneously without introducing systematic trade-
555 offs across dimensions.
556
557

5.2.3 Analysis of RQ3 - Impact of Debiasing on Model Performance

558 Finally, we assess whether bias mitigation nega-
559 tively affects downstream task performance. As
560 reported in Table 1, across all three tasks, our
561 method consistently outperforms the BERT base-
562 line, achieving performance gains of 0.50 on CoLA,
563 0.22 on QNLI, and 0.58 on SST.
564
565

566 These results indicate that effective bias miti-
567 gation does not necessarily come at the expense
568 of model performance. This finding supports
569 our central hypothesis that many bias-related fea-
570 tures do not provide useful task information but
571 instead introduce spurious correlations. By re-
572 ducing the model’s reliance on such correlations
573 through causal intervention, our framework guides
574 the model to focus on semantic information that is
575 more relevant to the task.

576 In contrast, many baseline methods exhibit a
577 clear trade-off, where stronger debiasing is of-
578 ten accompanied by degraded performance. Our
579 approach achieves a favorable balance, deliver-
580 ing strong task performance while simultaneously
581 achieving the most consistent bias reduction.

6 Conclusion

582 From a causal inference perspective, we systemati-
583 cally analyze the mechanisms underlying language
584 model inference by decomposing text representa-
585 tions into semantic, bias, and interaction compo-
586 nents. Building on this framework, we disentangle
587 bias-induced effects from the total causal effect by
588 aligning it with the natural direct effect that cap-
589 tures purely semantic influence via counterfactual
590 inference. Extensive experiments on three bench-
591 mark NLP tasks, evaluated using both intrinsic bias
592 metrics and extrinsic task performance, demon-
593 strate that our approach effectively mitigates bias
594 while consistently preserving or improving down-
595 stream performance. Importantly, our framework
596 is designed to jointly address multiple bias dimen-
597 sions within a unified process, enabling a more
598 faithful modeling of the interconnected nature of
599 real-world biases and avoiding the antagonistic ef-
600 fects that commonly arise when mitigating single
601 biases in isolation. Looking forward, promising
602 directions include extending causal debiasing to
603 generative and multilingual models, as well as de-
604 veloping broader cross-cultural benchmarks to en-
605 able more inclusive and robust bias evaluation.
606

607 Limitations

608 Despite growing attention to bias in language mod-
609 els, the development of comprehensive evaluation
610 benchmarks remains limited. In this work, we
611 rely on StereoSet, one of the most carefully cu-
612 rated bias benchmarks, constructed from natural
613 language by human annotators rather than syn-
614 thetic templates. However, its annotators are ex-
615 clusively North American and were instructed to
616 reflect North American stereotypes, which con-
617 strains coverage of regional and cultural biases.
618 As a result, although we consider multiple bias
619 dimensions such as gender and race, many cultur-
620 ally specific and geographically diverse stereotypes
621 remain underexplored. We therefore encourage fu-
622 ture research to broaden the scope of bias analysis
623 across cultures and regions and to develop more
624 inclusive benchmarks that enable systematic and
625 rigorous evaluation of bias in language models.

626 References

627 Kenza Benkirane, Jackie Kay, and Maria Perez-Ortiz.
628 2025. How can we diagnose and treat bias in large
629 language models for clinical decision-making? In
630 *Proceedings of the 2025 Conference of the Nations
631 of the Americas Chapter of the Association for Com-
632 putational Linguistics: Human Language Technolo-
633 gies (Volume 1: Long Papers)*, pages 2263–2288,
634 Albuquerque, New Mexico. Association for Compu-
635 tational Linguistics.

636 Pengyu Cheng, Weituo Hao, Siyang Yuan, Shijing Si,
637 and Lawrence Carin. 2021. Fairfil: Contrastive neu-
638 ral debiasing method for pretrained text encoders. In
639 *9th International Conference on Learning Representa-
640 tions*.

641 Paula Czarnowska, Yogarshi Vyas, and Kashif Shah.
642 2021. Quantifying social biases in NLP: A general-
643 ization and empirical comparison of extrinsic fairness
644 metrics. *Transactions of the Association for Compu-
645 tational Linguistics*, 9:1249–1267.

646 Yiwei Dai, Hengrui Gu, Ying Wang, and Xin Wang.
647 2024. Mitigate extrinsic social bias in pre-trained
648 language models via continuous prompts adjustment.
649 In *Proceedings of the 2024 Conference on Empiri-
650 cal Methods in Natural Language Processing*, pages
651 11068–11083.

652 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
653 Kristina Toutanova. 2019. Bert: Pre-training of deep
654 bidirectional transformers for language understand-
655 ing. In *Proceedings of the 2019 Conference of the
656 North American Chapter of the Association for Com-
657 putational Linguistics: Human Language Technolo-
658 gies, Volume 1 (Long and Short Papers)*, pages 4171–
659 4186.

Aparna Garimella, Akhash Amarnath, Kiran Kumar,
Akash Pramod Yalla, Anandhavelu N, Niyati Chhaya,
and Balaji Vasan Srinivasan. 2021. He is very intel-
ligent, she is very beautiful? On mitigating social
biases in language modelling and generation. In *Find-
ings of the Association for Computational Linguistics:
ACL-IJCNLP 2021*, pages 4534–4545. Association
for Computational Linguistics.

Yue Guo, Yi Yang, and Ahmed Abbasi. 2022. Auto-
debias: Debiasing masked language models with
automated biased prompts. In *Proceedings of the
60th Annual Meeting of the Association for Compu-
tational Linguistics (Volume 1: Long Papers)*, pages
1012–1023.

Farsheed Haque, Depeng Xu, and Shuhan Yuan. 2024.
Discovering and mitigating indirect bias in attention-
based model explanations. In *Findings of the Associ-
ation for Computational Linguistics: NAACL 2024*,
pages 1599–1614, Mexico City, Mexico. Association
for Computational Linguistics.

Lukas Hauenberger, Shahed Masoudian, Deepak
Kumar, Markus Schedl, and Navid Rekabsaz.
2023. Modular and on-demand bias mitigation with
attribute-removal subnetworks. In *Findings of the As-
sociation for Computational Linguistics: ACL 2023*,
pages 6192–6214, Toronto, Canada. Association for
Computational Linguistics.

Jacqueline He, Mengzhou Xia, Christiane Fellbaum,
and Danqi Chen. 2022. Mabel: Attenuating gender
bias using textual entailment data. In *Proceedings
of the 2022 Conference on Empirical Methods in
Natural Language Processing*, pages 9681–9702.

Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski,
Bruna Morrone, Quentin De Laroussilhe, Andrea
Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019.
Parameter-efficient transfer learning for nlp. In *In-
ternational Conference on Machine Learning*, pages
2790–2799.

Shadi Iskander, Kira Radinsky, and Yonatan Belinkov.
2024. Leveraging prototypical representations for
mitigating social bias without demographic infor-
mation. In *Proceedings of the 2024 Conference of
the North American Chapter of the Association for
Computational Linguistics: Human Language Tech-
nologies (Volume 2: Short Papers)*, pages 379–390,
Mexico City, Mexico. Association for Computational
Linguistics.

Masahiro Kaneko and Danushka Bollegala. 2021. De-
biasing pre-trained contextualised embeddings. In
*Proceedings of the 16th Conference of the European
Chapter of the Association for Computational Lin-
guistics: Main Volume*, pages 1256–1266.

Anne Lauscher, Goran Glavaš, Simone Paolo Ponzetto,
and Ivan Vulić. 2020. A general framework for im-
plicit and explicit debiasing of distributional word
vector spaces. *Proceedings of the AAAI Conference
on Artificial Intelligence*, 34(05):8131–8138.

gendered correlations in pre-trained models. *CoRR*, abs/2010.06032.

Zhongbin Xie and Thomas Lukasiewicz. 2023. An empirical analysis of parameter-efficient methods for debiasing pre-trained language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15730–15745.

Maxwell J. Yin, Boyu Wang, and Charles Ling. 2025. MABR: Multilayer adversarial bias removal without prior bias knowledge. *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(24):25724–25732.

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. Gender bias in coreference resolution: Evaluation and debiasing methods. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, volume 2, pages 15–20.

Fan Zhou, Yuzhou Mao, Liu Yu, Yi Yang, and Ting Zhong. 2023. Causal-debias: Unifying debiasing in pretrained language models and fine-tuning via causal invariant learning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4227–4241.

Ran Zmigrod, Sabrina J Mielke, Hanna Wallach, and Ryan Cotterell. 2019. Counterfactual data augmentation for mitigating gender stereotypes in languages with rich morphology. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1651–1661.

A Definitions and Case Study

A.1 Formal Definitions of Counterfactual Models

Counterfactual models are fundamental tools in causal inference, providing a principled way to translate causal assumptions encoded in causal graphs into mathematical expressions. The core idea of counterfactual analysis is to reason about the potential outcomes of variables under hypothetical intervention scenarios, enabling systematic comparison across different causal worlds.

Consider a causal graph involving three variables, where X denotes the cause variable, Y denotes the outcome variable, and Z serves as a mediator transmitting indirect causal effects from X to Y . When the value of X is set to x and the value of Z is set to z through intervention, the corresponding potential outcome of Y is defined as:

$$Y_{x,z} = Y(X = x, Z = z). \quad (11)$$

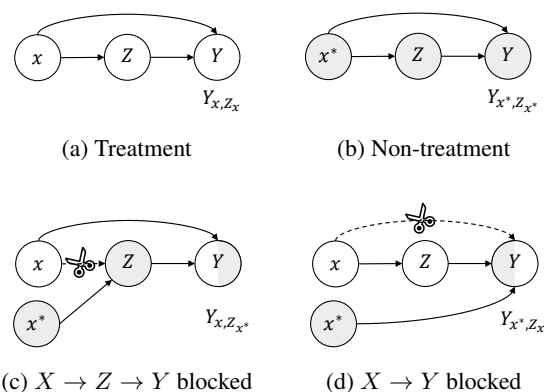


Figure 5: Counterfactual Models

In real-world scenarios, the mediator Z typically depends on the value of X , which can be expressed as $Z_x = Z(X = x)$. Counterfactual analysis allows this natural dependency to be decoupled by considering hypothetical scenarios in which X is set to one value while Z follows its potential outcome under a different value of X . Specifically, the counterfactual outcome

$$Y_{x,Z_{x^*}} = Y(X = x, Z = Z(X = x^*)) \quad (12)$$

represents a scenario where X is intervened to x , while the mediator Z is fixed to the value it would have taken under $X = x^*$.

This formulation enables the disentanglement of direct causal effects from indirect effects mediated through Z , which forms the foundation for identifying different components of causal influence in subsequent analysis.

A.2 Illustration of Counterfactual Models and Graphical Conventions

Figure 5 illustrates a set of counterfactual models that depict causal relationships among variables and their representations under different counterfactual scenarios. In the figure, node X represents the cause variable, node Y represents the outcome variable, and node Z serves as the mediator transmitting indirect causal effects from X to Y .

To improve visual clarity and conciseness, the variable X is not explicitly labeled as a node in the graph. Instead, the symbols x and x^* are used to denote two distinct intervention values of X , corresponding to $X = x$ and $X = x^*$, respectively. This notation emphasizes intervention states rather than the variable itself.

Different shading patterns are employed to distinguish variable states under counterfactual analysis. Fully shaded nodes indicate that the variable’s

value is entirely determined by artificial intervention, independent of observed data. Semi-shaded nodes indicate that the variable’s value is partially determined by the counterfactual scenario and partially inherited from the observed data distribution. This visualization strategy clarifies whether a variable is directly intervened upon or naturally propagated.

Additionally, scissor symbols are used to mark causal paths that are blocked due to human intervention. The presence of a scissor symbol on a causal edge indicates that the corresponding causal effect is severed in the counterfactual world. For example, when intervening on X , a scissor symbol on the path from X to Z indicates that the value of Z no longer depends on the natural state of X but is instead fixed according to the intervention condition. This graphical design facilitates intuitive interpretation of how direct and indirect causal effects are isolated through counterfactual reasoning.

A.3 Case Study: Counterfactual Reasoning in a Hiring Scenario

To further illustrate the concepts and definitions of counterfactual models, we provide a concrete case study based on a probation-to-regular hiring evaluation scenario. As shown in Figure 2, variable X represents an individual’s gender, where $X = x$ denotes male and $X = x^*$ denotes female. Variable Z represents the individual’s professional skills, and variable Y denotes the department’s final hiring decision.

For illustrative purposes, this case assumes that male candidates exhibit stronger hands-on abilities, while female candidates exhibit stronger communication skills, with the alternative skill being relatively weaker. This assumption is introduced solely to simulate potential gender stereotypes that may exist in real-world decision-making processes and does not reflect the views or claims of this paper. The objective of the proposed method is to mitigate

ID	Scen.	Instance	Formula
(a)	✓	Jack is skilled in hands-on tasks.	$Y_{x,Z_x} = Y(X = x, Z = Z(X = x))$
(b)	×	Mary has communicating skills.	$Y_{x^*,Z_{x^*}} = Y(X = x^*, Z = Z(X = x^*))$
(c)	×	Sam is adept at communicating .	$Y_{x,Z_{x^*}} = Y(X = x, Z = Z(X = x^*))$
(d)	×	Kate has hands-on skills.	$Y_{x^*,Z_x} = Y(X = x^*, Z = Z(X = x))$

* ✓ : Factual; × : Counterfactual. Blue: male/skills; Orange: female/skills.

Table 3: Instances and Formulas for Counterfactual Models

such biases.

Table 3 presents a complete example, where Mary, Sam, and Kate are artificially constructed counterfactual instances derived from the original case of Jack. These counterfactual instances do not exist in the actual dataset and are used exclusively for analytical comparison.

In this scenario, gender is treated as a protected attribute, and the model is explicitly prohibited from relying on gender-related information during inference, as such reliance may introduce discriminatory outcomes. Although this case study focuses on gender bias in a hiring context, similar forms of social bias are prevalent in language models and their downstream applications.

This case study demonstrates how counterfactual reasoning enables the separation of bias-induced influences from task-relevant information, thereby motivating the causal formulation adopted in our method for bias mitigation in word embedding models.

A.4 Causal Effect Identification

To formalize different causal pathways through which model predictions are formed, we introduce two key causal quantities: the **Total Effect (TE)** and the **Natural Direct Effect (NDE)**.

The **Total Effect (TE)** of the treatment $X = x$ on the outcome variable Y can be calculated by comparing the potential outcomes of Y under two hypothetical conditions. Here, Y_{x,Z_x} represents the potential outcome under the treatment condition $X = x$, and $Y_{x^*,Z_{x^*}}$ represents the potential outcome under the non-treatment condition $X = x^*$. The total effect reflects the overall impact of the treatment condition, encompassing both the direct effect and the indirect effect:

$$TE = Y_{x,Z_x} - Y_{x^*,Z_{x^*}}. \quad (13)$$

The **Natural Direct Effect (NDE)** refers to the causal effect of X on Y when the mediator Z is fixed to its counterfactual value under the non-treatment condition. Specifically, it describes the change in the potential outcome of Y when X transitions from x^* to x , while Z is fixed at Z_{x^*} :

$$NDE = Y_{x,Z_{x^*}} - Y_{x^*,Z_{x^*}}. \quad (14)$$

In the above formulation, $Y_{x,Z_{x^*}}$ represents the potential outcome when X is set to x and Z is fixed at its potential outcome under $X = x^*$, while $Y_{x^*,Z_{x^*}}$ represents the potential outcome when both X and Z take their non-treatment values.

Category	Count	Examples
Gender	409	<i>his, her, male, ...</i>
Race	13	<i>african, black, white, ...</i>
Religion	18	<i>islam, christian, jewish, ...</i>
Profession	95	<i>nurse, firefighter, guard, ...</i>
Appearance	18	<i>strong, cute, ugly, ...</i>
Personality	17	<i>aggressive, emotional, ...</i>
Activities	28	<i>shopping, football, boxing, ...</i>
Other	10	<i>flower, housework, ...</i>
Total	608	-

Table 4: Statistics and Examples of the Sensitive Lexicon.

B Sensitive Attribute Lexicon W_B

To construct a comprehensive lexicon for W_B , we draw on prior research (Meade et al., 2022; Guo et al., 2022; Zhou et al., 2023), which provides established definitions of sensitive attributes across multiple social categories. Based on these studies, we aggregated and organized the tokens into eight categories: Gender, Race, Religion, Profession, Appearance, Personality, Activities, and Other, resulting in a lexicon of 608 tokens. The resulting lexicon serves as a reliable basis for identifying bias-related terms. Table 4 shows the token counts for each category together with representative examples.

C Experimental Setup

C.1 Baseline Methods

To comprehensively evaluate the effectiveness of debiasing techniques, we employ the following state-of-the-art debiasing methods applied to BERT as baseline models.

- **CDA** (Webster et al., 2020) reduces model bias by constructing balanced training datasets through the substitution of sensitive attribute words with their corresponding counterparts.
- **DROPOUT** (Webster et al., 2020) hypothesizes that model bias stems from spurious correlations and mitigates these bias-related associations through an augmented training phase with an increased dropout rate.
- **INLP** (Ravfogel et al., 2020) iteratively eliminates bias-related features from model representations by projecting them onto the classifier’s null space that captures attribute-related information.

- **Sentence-Debias** (Liang et al., 2020) generalizes word embedding debiasing to sentence encoders by identifying bias-associated subspaces via principal component analysis and projecting sentence representations onto their orthogonal complements.
- **Context-Debias** (Kaneko and Bollegala, 2021) mitigates discriminative gender-related biases in contextualized representations through orthogonal projection operations applied to hidden layers.
- **Auto-Debias** (Guo et al., 2022) automatically probes model bias using discriminative diagnostic prompts and mitigates bias through distributional alignment with an equalizing loss.
- **MABEL** (He et al., 2022) improves model fairness by integrating counterfactual data augmentation with a dual-objective post-training strategy that combines contrastive learning and masked language modeling objectives.
- **Causal-Debias** (Zhou et al., 2023) integrates debiasing into the fine-tuning stage of pre-trained language models by learning causally invariant representations via invariant risk minimization.

While CPAD (Dai et al., 2024) and CDDD (Li et al., 2024a) represent significant methodological advances, we exclude them from our baseline comparisons due to critical limitations: CPAD’s implementation is not publicly available, and CDDD exhibits prohibitive computational demands along with partial implementation failures.

C.2 Evaluation Metrics

In line with our experimental framework, we adopt two types of metrics to ensure a comprehensive evaluation: intrinsic metrics for bias measurement and extrinsic metrics for downstream task performance.

Intrinsic Metrics. For bias measurement, we employ the Stereotype Score (Nadeem et al., 2021) and its average deviation compared to an ideally fair language model.

Stereotype Score (SS) is a metric designed to quantify the degree to which a language model exhibits. For a target term t , the Stereotype Score $SS(t)$ is defined as the percentage of instances

in which the model prefers a stereotypical association over an anti-stereotypical one. Mathematically, $SS(t)$ is expressed as:

$$SS(t) = \frac{N_{\text{Stereo}}(t)}{N_{\text{Total}}(t)} \times 100. \quad (15)$$

where $N_{\text{Stereo}}(t)$ represents the number of times the model prefers stereotypical associations for the target term t , and $N_{\text{Total}}(t)$ represents the total number of instances for the target term t .

For an entire dataset, the overall Stereotype Score SS is computed as the average of the Stereotype Scores for all target terms in the dataset:

$$SS = \frac{1}{N} \sum_{i=1}^N SS(t_i), \quad (16)$$

where N indicates the total number of target terms in the dataset.

The SS for a specific bias type is quantified by calculating the overall SS of the subset of the dataset that aligns with that particular type of bias. An ideal language model should exhibit neutrality, demonstrating no preference for either stereotypical or anti-stereotypical associations, which corresponds to an SS value of 50. We advocate for equal probability between anti-stereotypical and stereotypical associations, as any form of overgeneralized belief can cause harm to target groups. Consequently, an SS value closer to 50 indicates a relatively fairer and less biased language model.

Stereotype Score Deviation ($SS_{\text{Deviation}}^S$) quantifies the disparity between an evaluated language model and an ideally fair model, serving as a key metric to assess the extent of bias inherent in the model. A perfectly fair model is characterized by a Stereotype Score of 50, representing complete neutrality between stereotypical and anti-stereotypical associations. Consequently, the $SS_{\text{Deviation}}^S$ of a language model is defined as the absolute difference between its SS and the ideal value of 50. Regarding a specific category of bias (e.g., gender), the Stereotype Score is computed as the overall SS of the dataset subset corresponding to that bias category. Thus, the $SS_{\text{Deviation}}^S$ for a particular bias category c is expressed as:

$$SS_{\text{Deviation}}^S(c) = |SS(c) - 50|. \quad (17)$$

To ensure a comprehensive and robust evaluation, we consider four major categories of stereotypical biases by following Nadeem et al. (2021),

i.e., gender, profession, race, and religion. The overall Stereotype Score Deviation of the language model, denoted as Avg. Dev, is defined as the arithmetic mean of the $SS_{\text{Deviation}}^S$ across these four bias categories. Mathematically, this is represented as:

$$\text{Avg. Dev} = \frac{1}{|C|} \sum_{c \in C} SS_{\text{Deviation}}^S(c), \quad (18)$$

where $C = \{\text{gender, profession, race, religion}\}$.

Extrinsic Metrics. Following prior work (Wang et al., 2019), we employ the accuracy metric for the QNLI and SST tasks, and the Matthews Correlation Coefficient metric for the CoLA task. Both metrics are calculated based on true positives (TP), false positives (FP), true negatives (TN), and false negatives (FN). Further details are provided below.

Accuracy (ACC) is one of the most commonly used metrics to evaluate the performance of classification models. It measures the proportion of correctly classified instances relative to the total number of instances. ACC is determined by the following formula, with values ranging from 0 to 1, where higher values signify superior model performance:

$$ACC = \frac{TP + TN}{TP + TN + FP + FN}. \quad (19)$$

Matthews Correlation Coefficient (MCC) is a metric used to evaluate the performance of classification models, particularly in binary classification tasks. It is especially robust when dealing with imbalanced datasets. MCC ranges between [-1, 1], where 1 indicates perfect prediction with all classifications correct. The formula of MCC is as follows:

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (20)$$

C.3 Implementation Details

Following research (Meade et al., 2022; Zhou et al., 2023), we employ the bert-base-uncased version of the BERT model as the baseline for evaluation. For all methods involved in the experiments, we train the model for three epochs on each task, and a maximum sequence length of 128. We set the batch size to 32 and the learning rate to 2e-5. For the CDA and DROPOUT debiasing models, we incorporate an additional 1-epoch pre-training phase specifically designed to mitigate bias by following the original papers.

1177 For all comparative methods, we utilize their
1178 gender-debiased versions, as research on gender
1179 bias is the most extensive and well-documented
1180 in the field. In contrast, our proposed method
1181 is uniquely designed to mitigate multiple types
1182 of biases simultaneously within a single, unified
1183 debiasing process. Consequently, in our experi-
1184 mental framework, we design our proposed model
1185 to address multiple forms of bias beyond gender
1186 bias. This capability allows our method to achieve
1187 broader fairness and robustness across diverse con-
1188 texts, setting it apart from existing approaches that
1189 focus on a single bias type.

1190 **C.4 Sensitivity of Hyperparameter α**

1191 The hyperparameter α is used to balance the down-
1192 stream task objective (L_{task}) and our causal debi-
1193 asing constraint (L_{debias}). Specifically, a larger α
1194 imposes a stronger penalty on the divergence be-
1195 tween the Total Effect (TE) and the Natural Direct
1196 Effect (NDE), pushing the model to rely more on
1197 pure semantics. While the optimal α may vary
1198 slightly across tasks, our experiments consistently
1199 show that values in the range of $[0.01, 0.1]$ provide
1200 a robust balance, effectively mitigating bias with-
1201 out compromising task performance. This indicates
1202 that our framework is not overly sensitive to this
1203 hyperparameter.