Mitigating Biases in Language Models via Bias Unlearning

Anonymous ACL submission

Abstract

Many studies have shown various biases targeting different demographic groups in language models, amplifying discrimination and harming fairness. Recent parameter modification debiasing approaches significantly degrade core capabilities such as text coherence and task accuracy. And Prompt-based debiasing methods, only effective for predefined trigger words, fail to address deeply embedded stereotypical associations in model parameters. In this paper, we propose BiasUnlearn, a novel model debiasing framework which achieves targeted debiasing via dual-pathway unlearning mechanisms coordinating stereotype forgetting with anti-stereotype retention, while preventing bias polarity reversal through adversarial 018 forget set and dynamic dataset swapping. We conducted extensive experiments with multiple language models across various evaluation benchmarks. The results show that BiasUnlearn outperforms existing methods in mitigating bias in language models while retaining language modeling capabilities. Further experiments reveal that debiasing weights are transferable across model variants, confirming that bias representations become entrenched during pre-training and persist through fine-tuning phases. The codes are available at https://anonymous.4open.science/r/BiasUnlearn-5214.

Introduction 1

011

014

027

031

In recent years, large language models (LLMs) have been widely utilized in various fields. However, stereotypes in corpora constructed by human language inevitably affect these language models. Many studies have pointed out that there are significant social biases or stereotypes in large models, including gender bias, racial bias, and religious bias (Hofmann et al., 2024; Kim et al., 2024; Liu et al., 2023; Hosseini et al., 2023; Esiobu et al., 2023). These biases are subtle and seemingly harmless on 042



(a) The training process of BiasUnlean. When bias is reversed, swap the forget and retain sets and continue training.



(b) BiasUnlearn incorporates dual-pathway unlearning mechanisms coordinating stereotype forgetting with anti-stereotype retention, while preventing bias polarity reversal through adversarial forget set.

Figure 1: Demonstration of BiasUnlearn framework.

the surface, but they may have negative impacts on many applications. For example, translation software from companies like Google tends to translate gender-neutral pronouns in many languages into "he" in English without sufficient contextual (Zou and Schiebinger, 2018). When generating computer programs, LLMs can generate stereotypical content, which may cause displeasure to specific demographic groups. As reported by (Shrawgi et al., 2024), stereotypes have a negative impact on the accuracy of model text classification, and reducing stereotype bias can effectively improve text classification performance (Shen et al., 2023). Consequently, the elimination of bias to ensure the

043

045

067

086

087

094

100

101

102

103 104

105

106

108

057

fairness of model outputs, thereby preventing systematic discrimination within models, and enabling artificial intelligence technology to genuinely serve all of humanity is of great significance.

Many methods have been proposed to debias language models, such as re-pre-traing on counterfactual data (Lu et al., 2020; Zmigrod et al., 2019), debiasing with representation projection (Kaneko and Bollegala, 2021; Ravfogel et al., 2020; Liang et al., 2020), prompt-based methods(Schick et al., 2021; Furniturewala et al., 2024; Banerjee et al., 2024; Li et al.), and Model-editing approach (Xu et al., 2025b). However, these methods often fall short in effectively eliminating bias, preserving model capabilities, types of applicable models and ensuring efficient training and inference. For example, re-pre-training is effective but needs a great number of computing resources and time. Debiasing with representation projection remains limited in effectiveness due to the difficulty of obtaining high-quality bias representations. Prompt-based methods typically require multiple inferences per input, making them inefficient and limiting their ability to adapt to different bias types due to prompt limitations. Model-editing approach is efficient, which achieving substantial debiasing effects, but at the expense of weakened language modeling capabilities. In addition, re-pre-training can only be applied to pretrained models, while prompt-based methods can only be applied to instruction-finetuned models.

In this paper, we propose BiasUnlearn, an efficient and robust debiasing framework that mitigates biases in language models through targeted unlearning training. BiasUnlearn incorporates dualpathway unlearning mechanisms that coordinate stereotype forgetting with anti-stereotype retention. Additionally, it safeguards against bias polarity reversal through adversarial forget set construction and dynamic dataset swapping. To the best of our knowledge, we are the first to apply the LM unlearning method to debiasing LMs. The experimental results demonstrate that these measures can effectively reduce social biases in language models while maximizing the retention of the model's capabilities. Further experiments reveal that the debiasing weights trained on the base models maintain robust debiasing efficacy when transferred to instruction-fine-tuned models, suggesting that bias representations become entrenched during the pretraining phase, and that both base model and finetuning model share these bias representations.

2 Related Work

2.1 Bias and Debiasing in LMs

The biases of LMs are observed across many in-111 dependent model outputs and can only be mea-112 sured by observing the aggregated behavior of LMs 113 (Rauh et al., 2022). (Caliskan et al., 2017) pro-114 posed the Word Embedding Association Test and 115 discovered that male-associated terms tend to clus-116 ter with concepts like work, mathematics, and sci-117 ence, while female-associated terms align more 118 closely with family and art. (Zhao et al., 2018) 119 introduced a dataset in which only an ambiguous 120 pronoun differs between sentence pairs. LMs must 121 parse the pronoun into one of the two entities men-122 tioned prior to the sentence. Based on this dataset, 123 researchers have demonstrated the widespread exis-124 tence of gender bias in language models (Touvron 125 et al., 2023; Biderman et al., 2023). (Manvi et al., 126 2024) proposed studying LLMs through the lens of 127 geography and found that LLMs are biased against 128 locations or groups on a variety of sensitive subjec-129 tive topics. As LLMs' comprehension capabilities 130 continue to enhance, their performance scores on 131 coreference resolution tasks are consistently rising. 132 Consequently, in the future, utilizing such interme-133 diary tasks to evaluate biases may no longer be a vi-134 able approach. (Nangia et al., 2020; Nadeem et al., 135 2021; Liang et al., 2021; Esiobu et al., 2023; Barik-136 eri et al., 2021) introduced benchmarks to quan-137 tify stereotypes across gender, occupation, race, 138 and religion by analyzing language models' prob-139 ability distributions over stereotype-related target 140 words. To mitigate bias, various debiasing meth-141 ods are proposed. Re-pre-training with counter-142 factual data augmentation (Zmigrod et al., 2019) 143 replaces words in the corpus with tuples to con-144 struct a counterfactual enhanced corpus for training 145 LMs. (Kaneko and Bollegala, 2021; Ravfogel et al., 146 2020) proposed eliminating biases through project-147 ing bias representations. (Schick et al., 2021; Fur-148 niturewala et al., 2024) list unexpected behaviors in 149 the prompts and leverage the abilities of the LLM 150 itself to reduce the biased tokens. (Banerjee et al., 151 2024; Li et al.) generate counterfactuals for differ-152 ent groups, calculate probability distributions using 153 LLMs, and adjust probabilities to produce fairer 154 outputs. (Xu et al., 2025b) employed model-editing 155 technology to debias LMs. 156

109

110

161

171

181

188

189

190

191

194

195

196

198

199

200

202

203

206

2.2 LLM unlearning

Machine Unlearning (Fan et al., 2023) aims to elim-158 inate the influence of specific training data (such as sensitive or illegal information) on the completed 160 pre-trained model while maintaining the practicability of the model. (Yao et al., 2024) applied ma-162 chine unlearning in the domain of large language 163 164 models, referred to as LLM unlearning. LLM unlearning is employed to remove private information, 165 copyrighted content, or harmful data from LLMs, with the goal of erasing specific information while maintaining the model's general performance and 168 functionality (Yuan et al., 2024; Fan et al., 2023). 169 The most common LLM unlearning method is gra-170 dient ascent (GA), which achieves unlearning by maximizing the loss on the forget set (Yao et al., 172 2024). Since the loss function of LMs generally 173 174 has a lower bound but no upper bound, when GA is applied, the loss takes its negative value and 175 loses its lower bound, thereby eliminating meaningful minima and often leading to catastrophic 177 crashes. To address this issue, some studies intro-178 duce Kullback-Leibler (KL) divergence (Pan et al., 179 2025) and gradient descent (GD) (Yao et al., 2024; Premptis et al., 2025) to reduce parameter changes in models before and after unlearning. And (Zhang 182 et al.) proposed Negative Preference Optimization (NPO), whose progression toward catastrophic col-185 lapse is exponentially slower than GA. We extend the these ideas from unlearning factual content to 186 unlearning bias.

Debias with BiasUnlean 3

The goal of unlearning knowledge or factual information is to prevent LMs from generating specific content. However, mitigating bias through unlearning does not entail completely severing the association between LMs and specific words. Instead, it focuses on ensuring that the probability of certain words appearing is equalized across different demographic contexts (Banerjee et al., 2024). An ideal debiased language model should achieve an optimal balance between substantial bias mitigation and the preservation of original competencies. Therefore, our proposed BiasUnlearn framework must address three intertwined challenges: (1) Significant bias mitigation; (2) Preserving model capabilities post-debiasing; (3) Avoiding over-debiasing-induced bias reversal.

> Consider an LM \mathcal{M} with parameters Θ and an input context x related to a demographic

group G, outcome $y = \mathcal{M}(x; \Theta)$, debiasing \mathcal{M} to make the probability of generating a neutral word w equally likely regardless of demographic groups in contexts where only demographic groups differ, such that $\forall w \in \mathbb{W}$, $p(w|x_i, G_i, \mathcal{M}, \Theta) = p(w|x_i, G_i, \mathcal{M}, \Theta)$, where Θ is the parameter that has been debiased, that is, \mathcal{M} satisfies equal social group associations.

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

In this paper, we propose BiasUnlearn (as shown in Figure 1) for mitigating social bias in language models, which incorporates dual-pathway unlearning mechanisms that coordinate stereotype forgetting with anti-stereotype retention.

Stereotype Forgetting Firstly, in order to effectively mitigate bias in large models, NPO (Zhang et al.) which is more stable than Gradient Ascent is included in the optimization process on stereotypical set \mathcal{D}_s to forget bias:

$$\mathcal{L}_{Forget} = -\frac{2}{\beta} \mathbb{E}_{\mathcal{D}_s} \left[\log \sigma \left(-\beta \log \frac{\pi_{\theta}(y \mid x)}{\pi_{\text{ref}}(y \mid x)} \right) \right]$$
(1)

where π_{θ} is the unlearned LM and π_{ref} is the reference LM. σ is the sigmoid function and β is a regularization parameter.

Counterfactual Data and Anti-stereotype Retention For a given context x containing a sensitive word w_s related to group G_i , substituting w_s with a non-stereotypical word w_a constructs an anti-stereotypical context x', which serves as counterfactual data.

We adopt various methods to ensure that the language model maintains its general capabilities during unlearning. According to (Xu et al., 2025a; Ji et al., 2024), the integration of forgetting loss with gradient descent optimization objectives contributes to enhanced stability in parameter updates. Inspired by (Premptis et al., 2025), LMs can remember some facts while forgetting others, we use cross-entropy loss to train LMs with gradient descent optimization on the anti-stereotypical set \mathcal{D}_a , guiding LMs to internalize anti-stereotypical knowledge in parameter space.

$$\mathcal{L}_{Retention} = \frac{1}{|\mathcal{D}_a^i|} \sum_{\mathcal{D}_a^i} \operatorname{CE}(y, \hat{y})$$
(2)

Adversarial Forget Set Since forgetting loss aims to forget biases, while retention loss focuses on retaining anti-stereotypical information, both objectives align during training, continuous optimization may lead to bias reversal for some bias types.

To mitigate this, our experiments demonstrate that incorporating a small portion of anti-stereotypical 254 data into the forgetting training dataset slows down 255 bias reversal and helps prevent language model collapse. Thus the training set D_s in formula (1) is replaced with $D_s \cup D_{a'}$, $D_{a'}$ is a subset of D_a . These anti-stereotypical texts share identical contexts with their stereotypical counterparts except for critical word substitutions, which forces the model to optimize the same set of parameters to 262 accommodate contradictory objectives during train-263 ing. The gradient interference introduced by counterfactual data also has a regularization effect, pre-265 venting parameters from overfitting to a single ob-266 jective (simply maximizing the loss of stereotypical 267 data).

Data Chunk We combine a batch of stereotype data and a batch of anti-stereotype data into a data chunk. According to (Premptis et al., 2025), after performing the GA step, several consecutive GD annealing steps are required. In BiasUnlearn, we configure the ratio between forgettable and retained data to 1 : n (n > 1) within each data chunk, and the combined loss function can achieve the same progressive optimization effect as annealing. Since the batch size of the retain set is larger than that of the forget set, during training, retained data is cyclically sampled from the retain set.

271

272

273

274

275

276

279

281

287

294

296

297

300

Forward KL divergence Following (Murphy, 2022; Yao et al., 2024), we use forward KL divergence to force the distribution of the unlearned LM to cover all the areas of space of the original LM on the data unrelated to stereotypes, further preventing model collapse:

$$\mathcal{L}_{KL} = \mathrm{KL}\left(P_{\pi_{\theta}}\left(x_{unrel}\right) \parallel P_{\pi_{ref}}\left(x_{unrel}\right) \quad (3)$$

The final loss is computed as follows:

$$\mathcal{L} = \alpha_1 \mathcal{L}_{Forget} + \alpha_2 \mathcal{L}_{Retention} + \alpha_3 \mathcal{L}_{KL} \quad (4)$$

where alphas are hyperparameters.

Dynamic Dataset Swapping To prevent the overdebiasing during unlearning training from leading to bias reversal, we will swap the forget set and the retain set of specific bias types during training based on the results from the development set. Once the SS score for a particular bias type falls below 50, we will exchange the two sets of that bias type to continue training.

In our experiments, we leverage parameterefficient gradient-based methods (Jang et al., 2023) instead of full parameter training to improve training effectiveness. Specifically, we employ low-rank adaptation (LoRA) (Hu et al.). The specific experimental setups can be found in Appendix A.2 301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

347

348

4 Experiments

4.1 Dataset and Evaluation Metrics

StereoSet (Nadeem et al., 2021) We train the BiasUnlearn model on the StereoSet dataset, which is a large-scale dataset used to measure stereotypes in four domains: gender, profession, race, and religion. Following (Xu et al., 2025b), we adopt the test set of StereoSet as our training set and repurpose the development set as our test set. The trained model will also be evaluated on other benchmarks mentioned later. More details about the training set can be found in Appendix A.1

StereoSet evaluates LMs through three metrics: StereoSet (SS) score, Language Modeling (LM) score, and Idealized Context Correlation Test (ICAT) score. SS score represents the percentage of instances in which an LM prefers stereotypical sentences to anti-stereotypical sentences. The SS score approximating 50 indicates optimal bias fairness. The LM score is used to eval language modeling and general capabilities of an LM, that is, LM prioritizes the percentage of instances with meaningful associations rather than irrelevant alternatives. The ICAT score integrates SS score and LM score, higher values is better, calculates as follows:

$$ICAT = LMS \frac{\min(SS, 100 - SS)}{50}$$
 (5)

Crows-Pairs (Nangia et al., 2020) is a crowdsourced benchmark for stereotypical bias analysis. This dataset contains 1,508 examples including nine categories of social biases, such as gender, race, etc. Crows-Pairs is constructed in the form of contrastive pairs: each instance contains two semantically similar sentences, one of which presents a stereotypical description of historically disadvantaged groups, and the other is a contrasting expression of anti-stereotypes. Following (Banerjee et al., 2024, we also adopt the StereoSet score to measure the gender, religion, and race bias in LMs, whose ideal score is 50%.

BBQ, GLUE, FLUTE, AmazonPolarity, MTbench For fine-tuned language models, we also employ BBQ (Parrish et al., 2022), four tasks of GLUE(Wang et al.), FLUTE (Chakrabarty et al.,

	GPT2-Medium							GPT2-Large						
Method	5	Stereoty	pe Score	$e(\%) \rightarrow$	50	ΔLMS		5	Stereoty	pe Score	$(\%) \rightarrow$	50	ΔLMS	
	Gend.	Prof.	Race	Reli.	Overall	ightarrow 0	ICAI	Gend.	Prof.	Race	Reli.	Overall	ightarrow 0	
BaseModel	65.58	63.37	61.44	62.57	62.74	92.21	68.71	65.29	65.68	63.0	61.61	64.26	92.49	66.12
CDA	63.72	58.42	48.41	51.68	54.27	2.98	87.05	67.07	54.25	47.24	51.68	52.58	3.55	91.09
Self-Debias	59.63	61.35	57.5	57.98	59.25	-3.27	72.48	65.09	64.34	57.51	57.98	61.08	-3.27	69.44
CAFIE	56.2	62.03	60.66	63.91	60.74	-3.14	69.94	59.15	62.74	60.84	61.61	61.38	-3.02	69.11
BiasEdit	49.42	56.25	52.38	54.33	53.87	-2.63	82.66	52.64	55.27	54.02	47.36	54.19	-3.13	82.8
BiasUnlearn	52.44	51.06	50.37	48.64	50.83	-0.2	90.48	53.9	52.79	48.17	47.49	50.62	-0.26	91.08
				Mistra	l-7B						Llama.	3-8B		
Method	5	Stereoty	pe Score	$e(\%) \rightarrow$	50	ΔLMS		Stereotype Score (%) $\rightarrow 50$					ΔLMS	
	Gend.	Prof.	Race	Reli.	Overall	ightarrow 0	ICAI	Gend.	Prof.	Race	Reli.	Overall	ightarrow 0	
BaseModel	69.83	63.09	66.31	57.24	65.19	92.26	64.23	75.29	68.4	65.24	64.69	67.69	94.08	60.8
CDA	51.98	47.01	47.22	45.41	47.56	1.69	91.67	69.05	50.44	54.07	55.66	55.18	4.32	85.85
Self-Debias	62.15	55.77	50.86	59.72	54.49	-32.97	51.37	64.7	57.56	55.78	52.28	57.45	-35.82	52
CAFIE	57.61	62.85	67.55	51.13	63.89	-5.43	62.95	55.41	65.7	66.01	61.24	64.37	-5.11	63.17
BiasEdit	46.66	45.21	46.06	51.14	45.93	-13.56	60.81	44.73	40.32	50.49	52.01	45.2	-26.07	72.79
BiasUnlearn	51.4	48.36	50.67	51.49	49.92	-1.42	91.71	54.69	49.99	52.47	49.98	51.71	-0.4	89.09

Table 1: Debiasing performance of BiasUnlearn compared to baselines on StereoSet. Overall represents the aggregated results across all bias types. The SS score should approach 50 optimally, while the LM score requires minimal deviation from the base model. The ICAT metric theoretically exhibits positive correlation with performance improvement.

Method	GPT2-	Mediun	$\mathbf{n} \rightarrow 50$	$\textbf{GPT2-Large} \rightarrow 50$				
Method	Gend.	Race	Reli.	Gend.	Race	Reli.		
BaseModel	59.16	62.4	72.38	59.16	62.21	71.43		
CDA	59.54	50.31	65.71	59.94	58.49	71.43		
Self-Debias	49.62	46.54	51.43	55.73	58.49	60		
CAFIE	46.18	47.17	50.48	50	54.26	59.05		
BiasEdit	56.08	52.66	54.66	51.36	47.5	46.92		
BiasUnlearn	52.31	50	47.97	52.96	50.04	49.52		
Method	Mist	ral-7B -	$\rightarrow 50$	Llan	1a 3-8B -	$\rightarrow 50$		
Method	Mist Gend.	ral-7B - Race	→ 50 Reli.	Llan Gend.	1 a3-8B - Race	→ 50 Reli.		
Method BaseModel	Mist Gend. 62.98	ral-7B - Race 54.72	→ 50 Reli. 69.52	Llan Gend. 60.31	na3-8B - Race 59.12	→ 50 Reli. 74.29		
Method BaseModel CDA	Mist Gend. 62.98 58.78	ral-7B - Race 54.72 54.09	 → 50 Reli. 69.52 71.43 	Llan Gend. 60.31 40.08	na3-8B - Race 59.12 55.35	 → 50 Reli. 74.29 51.43 		
Method BaseModel CDA Self-Debias	Mist Gend. 62.98 58.78 54.2	ral-7B - Race 54.72 54.09 55.35	 → 50 Reli. 69.52 71.43 55.24 	Llan Gend. 60.31 40.08 58.78	na3-8B - Race 59.12 55.35 61.64	 → 50 Reli. 74.29 51.43 64.76 		
Method BaseModel CDA Self-Debias CAFIE	Mist Gend. 62.98 58.78 54.2 42.37	ral-7B - Race 54.72 54.09 55.35 48.43	 → 50 Reli. 69.52 71.43 55.24 53.33 	Llan Gend. 60.31 40.08 58.78 42.75	Race 59.12 55.35 61.64 49.06	 → 50 Reli. 74.29 51.43 64.76 60.95 		
Method BaseModel CDA Self-Debias CAFIE BiasEdit	Mist Gend. 62.98 58.78 54.2 42.37 51.46	ral-7B - Race 54.72 54.09 55.35 48.43 41.49	 → 50 Reli. 69.52 71.43 55.24 53.33 46.73 	Llan Gend. 60.31 40.08 58.78 42.75 46.04	na3-8B - Race 59.12 55.35 61.64 49.06 50	 → 50 Reli. 74.29 51.43 64.76 60.95 55.9 		

Table 2: Debiasing performance of BiasUnlearn compared to baselines on Crows-Pairs. The SS score (%) should approach 50 optimally.

30

354

355

358

2022), AmazonPolarity (Enevoldsen et al., 2025), and MT-bench (Zheng et al., 2023) for evaluation. Details about these tasks are in Appendix A.1

4.2 Baselines

We compare BiasUnlearn with the following four baselines: counterfactual factual data augmentation (CDA) (Zmigrod et al., 2019), Self-Debias (Schick et al., 2021), Counterfactually Aware Fair Inference (CAFIE) (Banerjee et al., 2024), and BiasEdit (Xu et al., 2025b). Since we need to deal with multiple types of biases simultaneously, we use the antistereotypical data from StereoSet as counterfactual data for training in our implementation CDA. For Self-Debias, we use the implementation of (Meade et al., 2022). As for CAFIE and BiasEdit, we use their official implementations and hyperparameters. More details are in Appendix A.3). 359

360

361

362

363

364

365

366

367

368

369

370

371

373

374

375

376

377

378

379

381

382

384

386

387

As a model-agnostic debiasing method, BiasUnlearn can be applied to any language model. We conduct experiments on GPT2-medium (355M), GPT2-Large (774M) (Radford et al., 2019), Mistral-7B (Jiang et al., 2023), and Llama3-8B (Meta, 2024), which represent diverse language models across different sizes. To verify the generalization of our method, we also conducted further experiments with the instruction-fine-tuned models of Mistral-7B and Llama3-8B.

4.3 Results

BiasUnlearn demonstrates superior debiasing performance while maintaining language modeling scores. Firstly, BiasUnlearn achieves significant reductions in stereotyping scores without excessive debiasing, according to Table 1, the overall SS scores of BiasUnlearn decreased by at least 11.91 and at most 15.98 on the four models. And BiasUnlearn can reduce the SS scores of all bias types to around 50%, while the absolute values of the difference between SS value and 50% of baselines are mostly higher than BiasUnlean. From the

	Mistral-7B-Instrucion								Llama3-8B-Instrucion					
Method	Stereotype Score (%) $\rightarrow 50$					IMSA		S	stereoty	reotype Score (%) $\rightarrow 50$				тсат *
	Gend.	Prof.	Race	Reli.	Overall			Gend.	Prof.	Race	Reli.	Overall	- LIND	ICAI
BaseModel	72.16	63.7	65.54	53.79	65.23	89.92	62.52	68.15	66.54	65.44	61.24	66.04	91.06	61.85
BiasUnlearn	53.07	46.92	50.05	48.46	49.18	89.03	87.57	49.14	48.87	47.95	47.31	48.42	91.03	88.16
WeightTrans	57.27	52.35	54.13	47.68	53.61	89.9	83.41	51.64	47.62	53.51	45.98	50.75	89.82	88.47

Table 3: Debiasing performance comparison of instruction-fine-tuned models: Before vs. after applying BiasUnlearn for training or debiasing weight transfer.

Mathad	Mistral-7B-Instrucion							Llama3-8B-Instrucion						
Wiethou	SST	MRPC	COLA	RTE	Amazon	Flute	МТ	SST	MRPC	COLA	RTE	Amazon	Flute	MT
BaseModel	0.937	0.733	0.743	0.226	0.906	0.85	6.6	0.93	0.675	0.737	0.219	0.894	0.41	7.831
BiasUnlearn	0.937	0.713	0.743	0.242	0.903	0.86	6.488	0.943	0.665	0.719	0.273	0.895	0.378	7.681
WeightTrans	0.937	0.724	0.748	0.242	0.9	0.855	6.762	0.927	0.679	0.742	0.254	0.894	0.403	7.713

Table 4: Performance comparison of instruction-fine-tuned models in semantic understanding tasks as well as question-answer dialogue task: Before vs. after applying BiasUnlearn for training or debiasing weight transfer.

results of the baseline methods, we observe that most approaches exhibit debiasing effects, however, the LMS values also significantly decreased. Conversely, since unlearning training will not improve the model's general ability or increase its knowledge, an increase in the LMS indicates that the model overfits to specific data, which is also detrimental to the performance of LMs. Compared with other methods, BiasUnlean achieved the minimal change on LMS (The minimum change of BiasUnlearn is only 0.2 and the maximum is only 1.42), indicating that BiasUnlearn method does not adversely affect the language modeling capabilities of the original models and our method maximally preserves the capabilities of the original models.

Based on Table 1 and Table 2, which present the results of various debiasing methods on StereoSet and Crows-Pairs, substantial differences in the distribution between the two datasets are evident. In fact, StereoSet divides race by nationality while Crows-Pairs divides race based on skin color. 408 BiasUnlearn achieves significant reductions in SS scores on both StereoSet and Crows-Pairs with the same training set, demonstrating its excellent debiasing effects and robust generalization capability. Prompt-based methods like Self-Debias and CAFIE struggle to adapt to different data distributions due to vocabulary limitations. Theoretically, fine-tune-based methods exhibit better adaptability across different datasets, but they relies heavily on appropriate training data. Notably, since the data distribution varies even across categories within the 419 StereoSet dataset itself, fine-tuning approaches like CDA demonstrate strong performance on specific datasets or categories but suffer from poor generalization to others. BiasEdit achieved significant debiasing effect on various datasets and bias types, however, after editing, the LMS of all models has significantly decreased, which negatively impacts the overall performance of the language models. BiasUnlearn works for both pretrained and finetuned models We conducted further experiments on downstream tasks on the fine-tuned models, and the results are presented in Table 4 and Figure 2. We can see that after training with BiasUnlearn, the model achieved performance levels that closely mirrored those of the original models across various tasks. This finding provides further evidence that our strategy for debiasing while retaining model capabilities is highly effective.

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

From Figure 2, we can observe that after BiasUnlearn training, the instruction-fine-tuned model demonstrates significant bias reduction in QA tasks. However, the bias scores in ambiguous contexts are much higher, indicating that when there is no clear answer, LMs will rely on social stereotypes. Notably, LLMs exhibit lower bias scores on the BBQ benchmark, which can be attributed to their enhanced general capabilities that enable answers to be less dependent on stereotypical associations. The improved capacity for context-aware reasoning consequently increases answer accuracy while reducing bias metrics. This finding indirectly suggests that employing intermediate evaluation tasks like QA formats and questionnaires for bias evaluation requires more nuanced design considerations. Current implementations may lack the sophistication needed to effectively probe deeper, more systemic biases within large language models (Shrawgi et al., 2024; Duan et al.).

389

390 391

393

395

400

401

402

403

404 405

406

407

409

410

411

412

413

414

415

416

417 418

420

421



Figure 2: Bias scores in each category of BBQ, split by whether the context is ambiguous or disambiguated. The higher the bias score, the stronger the bias.

The transferability of debiasing weights indicates that bias features have been solidified in the pre-training stage. To further verify that the parameters learned by BiasUnlearn are only related to stereotypes, we load the debiasing weights learned on base models into instruction-finetuned models. As shown in Table 3, both instruction-finetuned models of Llama and Mistral achieve strong debiasing effects after loading the base model's debiasing weights, with only minor degradation in language modeling scores. And from Table 4, the instruction-finetuned models with transferred debiasing weights exhibit almost no performance drop on GLUE and other tasks compared to the original models. It indicates that bias features are deeply solidified within the underlying features during the pre-training phase and cannot be modified by ordinary fine-tuning. Furthermore, the successful transfer of debiasing weights indicates that BiasUnlearn effectively disentangles bias-related parameters without reconstructing the underlying semantic representations.

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486 487

488

489

490

491

The previous works (Kadhe et al.; Zhang et al., 2023) divided the training data into specific behavioral subtypes that were separately trained through separate LLMs based on attribute values, and then merged them into the final model. Our experiments are significantly different from their work, as they combine multiple weights from the same LM, while we transfer the same weights from different models. Our findings suggest that in the future applications of large language models, there will be no need to debias each individual model: once an unbiased base model is obtained, its subsequent fine-tuning models will also be essentially unbiased; and by only debiasing biased base models and transferring the debiasing weights to their fine-tuned large models that need to be applied, the workload can be significantly reduced. 492

493

494

495

496

497

498

499

500

502

503

504

505

506

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

BiasUnlearn is efficient, and the resulting debiased models can be deployed effortlessly. During inference, as the number of BiasUnlearn parameters does not change significantly, the inference speed will not slow down and memory consumption will not significantly increase; while promptbased methods require parallel or sequential multistep reasoning, which will increase memory usage or inference time by several times. Compared to fine-tuning approaches such as BiasEdit, which require multiple epochs to train, BiasUnlearn can complete training within merely a few hundred steps, significantly reducing the training time required. Additionally, BiasUnlearn does not require any modifications to the inference process (such as adding sensitive word lists, or modifying token distribution during inference, etc.), so that the debiased LMs can be deployed like vanilla LMs. Compared to BiasEdit, which is only evaluated on sentence-level examples (Xu et al., 2025b), LMs trained by BiasUnlearn are capable of adapting to a variety of tasks (details are in Appendix A.1).

4.4 Ablation Experiments

In order to investigate the influence of different components within the BiasUnlearn framework on retaining the inherent capabilities of language models, we conduct ablation experiments on StereoSet based on GPT2-large. For a fair comparison, we

Method	SS	LMS	ICAT
BiasUnlearn	53.73	92.95	86.01
w/o Retention loss	47.56	57.72	54.91
w/o KL	43.14	88.22	76.11
w/o Adv	44.72	86.84	77.67

Table 5: The overall SS, LMS, ICAT of BiasUnlearn with or without $\mathcal{L}_{Retention}$, \mathcal{L}_{KL} or Adversarial Forget Set.



Figure 3: Forgetting loss and Retention loss of BiasUnlearn with or without $\mathcal{L}_{Retention}$, \mathcal{L}_{KL} or Adversarial Forget Set

use the checkpoint at step 1000 to generate the 525 526 evaluation results. As shown in Table 5, even after a thousand steps of unlearning training, the LMS 527 still maintains a high value. Without Retention 528 loss, LMS decreased by 35.23 points, which is the largest decrease among all methods, indicating that 530 retention loss plays the greatest role in maintaining 531 the original abilities of the model during the un-532 learn process. And under the absence of KL diver-533 gence and Adversarial Forget Set, the LMS exhibits respective reductions of 4.73 and 6.1 points, empiri-535 cally validating that both KL divergence constraints and Adversarial Forget Set are critical for mitigat-537 ing catastrophic ability degradation and knowledge 539 forgetting during LLM unlearning. When Adversarial Forget Setis absent, the SS score decreases 540 lower than when lacking KL divergence, so Adver-541 sarial Forget Set has a certain effect on preventing over-debiasing. 543

From Figure 3 we can clearly observe the differential effects of each constituent on unlearning training. Whether the Forgetting loss or Retention loss becomes too large, it is a disaster for bias unlearning, which compromises the stability of the training. When there is a lack of Retention loss, the Forgetting loss starts to grow rapidly around the 100th step, with an absolute value that increases more than twenty times compared to the beginning, causing the collapse of the debiased model. The negative impacts of KL divergence absence on the language model manifest much later than those from Retention loss deterioration. Both Forgetting loss and Retention loss exhibit a smaller increase compared to scenarios lacking Retention loss. The two loss curves without KL divergence are both above the curve with missing Adversarial Forget Set, indicating that KL divergence has a greater impact on the results than Adversarial Forget Set.

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

564

565

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

583

584

585

586

587

588

589

590

591

4.5 Case Study

We compared the results generated by the models before and after debiasing on the BBQ and MTbench, as detailed in Appendix B.

5 Conclusion

In this work, we propose BiasUnlearn, an effective unlearning approach for debiasing language models. BiasUnlearn incorporates dual-pathway unlearning mechanisms that coordinate stereotype forgetting with anti-stereotype retention. Through adversarial forget set and dynamic dataset swapping, our approach safeguards against bias polarity reversal. Experimental results demonstrate that BiasUnlearn successfully mitigates bias in language models without reversing the bias direction, while simultaneously maintaining the model's general capabilities. Transfer experiments reveal that our work is transferable across model variants, empirically confirming that bias representations become entrenched during pre-training and persist through fine-tuning phases. These findings not only offer practical debiasing tools but also provide theoretical insights into bias propagation in language models; and this work makes a significant contribution to promoting fairness in LLMs. In the future, we will extend BiasUnlearn to other forms of bias beyond social biases and investigate the relationship between bias representations and different pre-training objectives.

Limitations

592

611

624

627

628

630

631

632

634

635

637

641

593 Existing datasets related to social bias fail to cover a comprehensive range of demographic groups, and 594 many countries, ethnicities, and religions are en-595 tirely absent in these datasets. Consequently, lan-596 guage models trained on these datasets using Bi-598 asUnlearn may not effectively eliminate bias for all demographic groups. In addition, many types of bias have potential correlations. For example, race and religion exhibit statistical correlations in real-world issues (e.g., highly overlapping racial and religious groups in certain regions) (Perry and Whitehead, 2019), and models may simultaneously affect both types of bias by adjusting shared under-605 lying features. Our work does not fully disentangle the representation space of different biases. When optimizing for a certain type of bias, the parameters associated with other related biases are adjusted jointly.

Ethics Statement

Our work effectively removes social biases in language models while maintaining their general ca-613 pabilities, greatly improving the fairness of lan-614 guage models for different demographic groups. 615 We strictly adhere to the data usage policies of 616 various open-source datasets. The opinions and 617 findings contained in the dataset samples we provide should not be interpreted as representing the 619 views expressed or implied by the authors. 620

References

- Pragyan Banerjee, Abhinav Java, Surgan Jandial, Simra Shahid, Shaz Furniturewala, Balaji Krishnamurthy, and Sumit Bhatia. 2024. All should be equal in the eyes of lms: Counterfactually aware fair text generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17673–17681.
- Soumya Barikeri, Anne Lauscher, Ivan Vulić, and Goran Glavaš. 2021. Redditbias: A real-world resource for bias evaluation and debiasing of conversational language models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1941–1955.
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, and 1 others.
 2023. Pythia: A suite for analyzing large language models across training and scaling. In *International*

Conference on Machine Learning, pages 2397–2430. PMLR.

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

- Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186.
- Tuhin Chakrabarty, Arkadiy Saakyan, Debanjan Ghosh, and Smaranda Muresan. 2022. FLUTE: Figurative language understanding through textual explanations. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 7139–7159, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Stergios Chatzikyriakidis, Robin Cooper, Simon Dobnik, and Staffan Larsson. 2017. An overview of natural language inference data collection: The way forward? In *Proceedings of the Computing Natural Language Inference Workshop*.
- Shitong Duan, Xiaoyuan Yi, Peng Zhang, Tun Lu, Xing Xie, and Ning Gu. Denevil: Towards deciphering and navigating the ethical values of large language models via instruction learning. In *The Twelfth International Conference on Learning Representations*.
- Kenneth Enevoldsen, Isaac Chung, Imene Kerboua, Márton Kardos, Ashwin Mathur, David Stap, Jay Gala, Wissam Siblini, Dominik Krzemiński, Genta Indra Winata, Saba Sturua, Saiteja Utpala, Mathieu Ciancone, Marion Schaeffer, Gabriel Sequeira, Diganta Misra, Shreeya Dhakal, Jonathan Rystrøm, Roman Solomatin, and 67 others. 2025. Mmteb: Massive multilingual text embedding benchmark. arXiv preprint arXiv:2502.13595.
- David Esiobu, Xiaoqing Tan, Saghar Hosseini, Megan Ung, Yuchen Zhang, Jude Fernandes, Jane Dwivedi-Yu, Eleonora Presani, Adina Williams, and Eric Smith. 2023. Robbie: Robust bias evaluation of large generative language models. In *Proceedings of the* 2023 Conference on Empirical Methods in Natural Language Processing, pages 3764–3814.
- Chongyu Fan, Jiancheng Liu, Yihua Zhang, Eric Wong, Dennis Wei, and Sijia Liu. 2023. Salun: Empowering machine unlearning via gradient-based weight saliency in both image classification and generation. In *The Twelfth International Conference on Learning Representations*.
- Shaz Furniturewala, Surgan Jandial, Abhinav Java, Pragyan Banerjee, Simra Shahid, Sumit Bhatia, and Kokil Jaidka. 2024. "thinking" fair and slow: On the efficacy of structured prompts for debiasing language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 213–227.
- Valentin Hofmann, Pratyusha Ria Kalluri, Dan Jurafsky, and Sharese King. 2024. Ai generates covertly racist decisions about people based on their dialect. *Nature*, 633(8028):147–154.

806

807

808

809

- 698 699 700
- 701 702
- 71
- 70
- 7(
- _

710

711

- 712 713 714
- 715 716 717 718
- 719 720 721

722

- 723 724 725 726 727 728
- 730 731 732 733
- 735 736 737

734

- 739
- 741 742

740

- 743 744
- 745
- 746 747 748

748 749 750

- 750
- 751 752

7

754

- Saghar Hosseini, Hamid Palangi, and Ahmed Hassan.
 2023. An empirical study of metrics to measure representational harms in pre-trained language models.
 In Proceedings of the 3rd Workshop on Trustworthy Natural Language Processing (TrustNLP 2023), pages 121–134.
- Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, and 1 others. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.
- Joel Jang, Dongkeun Yoon, Sohee Yang, Sungmin Cha, Moontae Lee, Lajanugen Logeswaran, and Minjoon Seo. 2023. Knowledge unlearning for mitigating privacy risks in language models. In *The 61st Annual Meeting Of The Association For Computational Linguistics*.
 - Jiabao Ji, Yujian Liu, Yang Zhang, Gaowen Liu, Ramana Kompella, Sijia Liu, and Shiyu Chang. 2024. Reversing the forget-retain objectives: An efficient llm unlearning framework from logit difference. *Advances in Neural Information Processing Systems*, 37:12581–12611.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. *Preprint*, arXiv:2310.06825.
- Swanand Kadhe, Farhan Ahmed, Dennis Wei, Nathalie Baracaldo, and Inkit Padhi. Split, unlearn, merge: Leveraging data attributes for more effective unlearning in llms. In *ICML 2024 Workshop on Foundation Models in the Wild*.
- Masahiro Kaneko and Danushka Bollegala. 2021. Debiasing pre-trained contextualised embeddings. In Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, pages 1256–1266.
- Minbeom Kim, Jahyun Koo, Hwanhee Lee, Joonsuk Park, Hwaran Lee, and Kyomin Jung. 2024. Lifetox: Unveiling implicit toxicity in life advice. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers), pages 688–698.
- Jingling Li, Zeyu Tang, Xiaoyu Liu, Peter Spirtes, Kun Zhang, Liu Leqi, and Yang Liu. Prompting fairness: Integrating causality to debias large language models. In *The Thirteenth International Conference on Learning Representations*.
- Paul Pu Liang, Irene Mengze Li, Emily Zheng, Yao Chong Lim, Ruslan Salakhutdinov, and Louis-Philippe Morency. 2020. Towards debiasing sentence representations. In *Proceedings of the 58th Annual*

Meeting of the Association for Computational Linguistics, pages 5502–5515.

- Paul Pu Liang, Chiyu Wu, Louis-Philippe Morency, and Ruslan Salakhutdinov. 2021. Towards understanding and mitigating social biases in language models. In *International conference on machine learning*, pages 6565–6576. PMLR.
- Yan Liu, Xiaokang Chen, Yan Gao, Zhe Su, Fengji Zhang, Daoguang Zan, Jian-Guang Lou, Pin-Yu Chen, and Tsung-Yi Ho. 2023. Uncovering and quantifying social biases in code generation. *Advances in Neural Information Processing Systems*, 36:2368– 2380.
- Kaiji Lu, Piotr Mardziel, Fangjing Wu, Preetam Amancharla, and Anupam Datta. 2020. Gender bias in neural natural language processing. *Logic, language, and security: essays dedicated to Andre Scedrov on the occasion of his 65th birthday*, pages 189–202.
- Rohin Manvi, Samar Khanna, Marshall Burke, David B Lobell, and Stefano Ermon. 2024. Large language models are geographically biased. In *International Conference on Machine Learning*, pages 34654– 34669. PMLR.
- Nicholas Meade, Elinor Poole-Dayan, and Siva Reddy. 2022. An empirical survey of the effectiveness of debiasing techniques for pre-trained language models. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1878–1898.
- AI Meta. 2024. Introducing meta llama 3: The most capable openly available llm to date. *Meta AI*, 2(5):6.
- Kevin P Murphy. 2022. Probabilistic machine learning: an introduction.
- Moin Nadeem, Anna Bethke, and Siva Reddy. 2021. Stereoset: Measuring stereotypical bias in pretrained language models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5356–5371.
- Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel Bowman. 2020. Crows-pairs: A challenge dataset for measuring social biases in masked language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1953–1967.
- Zibin Pan, Shuwen Zhang, Yuesheng Zheng, Chi Li, Yuheng Cheng, and Junhua Zhao. 2025. Multiobjective large language model unlearning. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing* (*ICASSP*), pages 1–5. IEEE.
- Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson, Phu Mon Htut, and Samuel Bowman. 2022. BBQ:

- A hand-built bias benchmark for question answering. Alex Wang, Amanpreet Singh, Julian Michael, Felix 866 Hill, Omer Levy, and Samuel R Bowman. Glue: In Findings of the Association for Computational 867 Linguistics: ACL 2022, pages 2086–2105, Dublin, A multi-task benchmark and analysis platform for Ireland. Association for Computational Linguistics. natural language understanding. In International 869 Conference on Learning Representations. 870 Samuel L Perry and Andrew L Whitehead. 2019. Christian america in black and white: Racial identity, Haoming Xu, Shuxun Wang, Yanqiu Zhao, Yi Zhong, 871 religious-national group boundaries, and explana-Ziyan Jiang, Ningyuan Zhao, Shumin Deng, Huajun 872 tions for racial inequality. Sociology of Religion, Chen, and Ningyu Zhang. 2025a. Zjuklab at semeval-873 80(3):277-298. 2025 task 4: Unlearning via model merging. arXiv 874 preprint arXiv:2503.21088. 875 Iraklis Premptis, Maria Lymperaiou, Giorgos Filandrianos, Orfeas Menis Mastromichalakis, Athanasios Xin Xu, Wei Xu, Ningyu Zhang, and Julian McAuley. 876 Voulodimos, and Giorgos Stamou. 2025. Ails-ntua 2025b. Biasedit: Debiasing stereotyped lan-877 at semeval-2025 task 4: Parameter-efficient unlearnguage models via model editing. arXiv preprint 878 ing for large language models using data chunking. arXiv:2503.08588. 879 arXiv preprint arXiv:2503.02443. Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Yuanshun Yao, Xiaojun Xu, and Yang Liu. 2024. Large 880 Dario Amodei, and Ilya Sutskever. 2019. Language language model unlearning. Advances in Neural models are unsupervised multitask learners. OpenAI Information Processing Systems, 37:105425–105475. 882 blog. Xiaojian Yuan, Tianyu Pang, Chao Du, Kejiang Chen, 883 Maribeth Rauh, John Mellor, Jonathan Uesato, Po-Sen Weiming Zhang, and Min Lin. 2024. A closer look at 884 Huang, Johannes Welbl, Laura Weidinger, Sumanth machine unlearning for large language models. arXiv 885 Dathathri, Amelia Glaese, Geoffrey Irving, Iason preprint arXiv:2410.08109. 886 Gabriel, and 1 others. 2022. Characteristics of harmful text: Towards rigorous benchmarking of language Jinghan Zhang, Junteng Liu, Junxian He, and 1 others. models. Advances in Neural Information Processing 887 Systems, 35:24720-24739. 2023. Composing parameter-efficient modules with 888 arithmetic operation. Advances in Neural Informa-889 Shauli Ravfogel, Yanai Elazar, Hila Gonen, Michael tion Processing Systems, 36:12589–12610. 890 Twiton, and Yoav Goldberg. 2020. Null it out: Guarding protected attributes by iterative nullspace projec-Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. Nega-891 tion. In Proceedings of the 58th Annual Meeting of tive preference optimization: From catastrophic col-892 the Association for Computational Linguistics, pages lapse to effective unlearning. In First Conference on 893 7237-7256. Language Modeling. 894 Timo Schick, Sahana Udupa, and Hinrich Schütze. 2021. Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Or-895 Self-diagnosis and self-debiasing: A proposal for redonez, and Kai-Wei Chang. 2018. Gender bias in 896 ducing corpus-based bias in nlp. Transactions of the coreference resolution: Evaluation and debiasing 897 Association for Computational Linguistics, 9:1408– methods. In Proceedings of the 2018 Conference 898 1424. of the North American Chapter of the Association 899 for Computational Linguistics: Human Language Shaofei Shen, Mingzhe Zhang, Weitong Chen, Alina 900 Technologies, Volume 2 (Short Papers), pages 15-20. Bialkowski, and Miao Xu. 2023. Words can be con-901 fusing: Stereotype bias removal in text classification at the word level. In Pacific-Asia Conference on Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan 902 Knowledge Discovery and Data Mining, pages 99-Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, 903 111. Springer. Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 904 2023. Judging llm-as-a-judge with mt-bench and 905 Hari Shrawgi, Prasanjit Rath, Tushar Singhal, and Sandichatbot arena. Advances in Neural Information Pro-906 pan Dandapat. 2024. Uncovering stereotypes in large cessing Systems, 36:46595-46623. 907 language models: A task complexity-based approach. In Proceedings of the 18th Conference of the Euro-Ran Zmigrod, Sabrina J Mielke, Hanna Wallach, and 908 pean Chapter of the Association for Computational Ryan Cotterell. 2019. Counterfactual data augmenta-909 Linguistics (Volume 1: Long Papers), pages 1841tion for mitigating gender stereotypes in languages 910 1857. with rich morphology. In Proceedings of the 57th 911 Annual Meeting of the Association for Computational 912 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Linguistics, pages 1651–1661. 913
 - Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

811

812

814

815

816

817

820

821

825

826

827

829

830

831

832

835

836

837

838

842

845

847

853

854

855

857

861

James Zou and Londa Schiebinger. 2018. Ai can be sexist and racist—it's time to make it fair. *Nature*, 915 559(7714):324–326. 916

A Experimental Details

917

918

919

920

922

923

924

925

926

927

931

934

935

937

938

941

943

949

951

954

955

958

A.1 Dataset and Evaluation Metrics

StereosetNadeem et al., 2021 Each sample in Stereoset consists of one stereotype sentence, one antistereotype sentence, and one sentence unrelated to stereotypes. We adopt the test set of StereoSet as our training set and repurpose the development set as our test set. The number of samples corresponding to each bias type in the re-segmented dataset is presented in Table 6.

	gender	profession	race	religion	Overall
train	1471	4782	5871	438	12562
dev	50	50	50	50	200
test	497	1638	1938	159	4232

Table 6: The number of samples corresponding to each bias type in our dataset.

BBQ (Parrish et al., 2022) is a question-answering dataset designed to assess bias across nine social groups. Each question in BBQ is presented in two forms: one is an ambiguous version, lacking clear context, as well as a disambiguation version that provides additional context before the question. BBQ defines the Bias Score as a metric to quantify the degree of bias present in the responses of large language models. A higher score indicates a more severe level of bias.

GLUE (Wang et al.) GLUE benchmark comprises nine NLU tasks for evaluating the semantic understanding ability of language models, and We selected a subset of four tasks—SST-2, MRPC, COLA, and RTE to assess model performance.
SST-2 is an sentiment analysis (positive/negative) task. MRP is a sentence semantic matching task, which involves determining whether two sentences are semantically similar. COLA is a grammar judgment task, which involves determining whether a sentence is grammatically correct. RTE, a textual entailment task (Chatzikyriakidis et al., 2017), which involves judging the logical relationships between sentences.

FLUTE (Chakrabarty et al., 2022), is a figurative language understanding dataset and can be framed as a recognizing textual entailment task.

AmazonPolarityClassification(Enevoldsen et al., 2025) aims to classify Amazon reviews into positive or negative sentiments.

In GLUE, FLUTE, and AmazonPolarity, 200 data were selected from each dataset for testing

and F1 scores are reported.

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

MT-bench (Zheng et al., 2023) is a benchmark used to evaluate LMs in multiple rounds of dialogue, covering eight domains such as writing, coding, and humanities. We utilize Llama-3-70B-Instruction to evaluate the responses of the assessed LMs in MT-bench by assigning scores through FastChat¹.

A.2 Setup

The language models we use are from Hugging Face, including GPT2-medium², GPT2-large³, Mistral7B-v0.1⁴, Llama3-8B⁵, Mistral-7B-Instructv0.1⁶, and Llama3-8B-Instruct⁷. For GPT2medium and GPT2-large, the initial learning rate is set to 5e-5; for Mistral and Llama3 models, the initial learning rate is set to 2e-5. For all models, we adopt a linear learning rate scheduler and AdamW optimizer. Referring to the work of (Premptis et al., 2025), the global batch size for the forget set is 4, while the global batch size for the retain set is 28. The weighting coefficients α_1 , α_2 , and α_3 corresponding to the three losses \mathcal{L}_{Forget} , $\mathcal{L}_{Retention}$, and \mathcal{L}_{KL} are set to 0.4, 0.4, and 0.2, respectively. When the SS scores for all bias types on the development set are less than the threshold of 2, early stopping is triggered. Our experiments were conducted using 4*A100 GPUs.

Due to the differing distributions of the Crows-Pairs and StereoSet, for each model, the checkpoint used for evaluation on Crown Pair differs from that used on StereoSet. The checkpoint employed for evaluation on other benchmarks is identical to the one used for StereoSet.

A.3 Baselines

Counterfactual factual data augmentation (**CDA**) (Zmigrod et al., 2019) mitigates stereotypes by conducting pre-training on augmented data that contains counterfactual data. Counterfactuals are generated with roughly speaking the

¹ https://github.com/lm-sys/FastChat
² https://huggingface.co/openai-community/
gpt2-medium
³ https://huggingface.co/openai-community/
gpt2-large
⁴ https://huggingface.co/mistralai/
Mistral-7B-v0.1
⁵ https://huggingface.co/meta-llama/
Meta-Llama-3-8B
⁶ https://huggingface.co/mistralai/
Mistral-7B-Instruct-v0.1
⁷ https://huggingface.co/meta-llama/
Meta-Llama-3-8B-Instruct

opposite bias of some original dataset. The original 998 CDA work focused solely on gender bias, whereas our work, given the inclusion of multiple types of 1000 bias, utilized anti-stereotypical data from StereoSet 1001 as counterfactual information for training in our 1002 implemented CDA. 1003

Self-Debias (Schick et al., 2021) leverages the internal knowledge of the large model and prompts 1005 designed to elicit potentially harmful outputs, en-1006 abling the large model to detect the presence of 1007 undesirable attributes in its own output. Subse-1008 quently, by adjusting the original token distribution 1009 through weighting with the probabilities of tokens 1010 that explicitly encourage unexpected behavior, the 1011 likelihood of the model generating biased text is effectively reduced. We use the implementations 1013 from (Meade et al., 2022). 1014

Counterfactually Aware Fair Inference (CAFIE) 1015 (Banerjee et al., 2024) involves four key steps: (1) 1016 identifying sensitive tokens in the source context, (2) constructing valid counterfactual contexts by 1018 altering these tokens, (3) computing probability 1020 distributions for both original and counterfactual contexts, and (4) adjusting the original distribution by equalizing next-token probabilities across all contexts. The final fair distribution combines the adjusted and original probabilities for text sampling. We use their official implementations and hyperparameters⁸. 1026

> BiasEdit (Xu et al., 2025b) utilizes a debiasing loss to guide the editing network in localizing specific parameters of the language model for debiasing, while retaining the language modeling ability during the editing process through a retention loss. We use their official implementations and hyperparameters⁹.

B Case Study

1017

1019

1021

1024

1025

1027

1028

1029

1031

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

The examples of generated outputs on BBQ are shown in Table 7. As illustrated in Table 7, even after undergoing secure alignment, the large model like Llama3-8B- Instruction still retains a certain probability of generating socially biased answers. The large model trained with BiasUnlean can mitigate these biases effectively. Additionally, the transfer of Debias Weights can assist in reducing bias within large models without compromising their capabilities.

MT-Bench is a multi-round dialogue task. Given 1045 the lengthiness of the model's responses, we 1046 present only one response example for each model 1047 to ensure clarity and conciseness. From Appendix 1048 B.1-B.3 we can see that the model trained by Bia-1049 sUnlearn retains strong conversational ability and 1050 rich internal knowledge. And the transfer of debi-1051 asing weights does not have a negative impact on the various capabilities of the model. 1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1078

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1094

B.1 Output example generated by Llama3-8B-Instruction

-----; Question: Now you are a machine learning engineer. Your task is to explain complex machine learning concepts in a simplified manner so that customers without a technical background can understand and trust your products. Let's start with the question: "What is a language model? Is it trained using labeled or unlabelled data?"

Model Response:

I'm excited to explain machine learning concepts in a way that's easy to understand.

A language model is a type of artificial intelligence (AI) that's designed to understand and generate human-like language. Think of it like a supersmart, never-tiring, language expert who can chat with you, write articles, or even create dialogues for a character in a movie!

Imagine you're trying to teach a child to speak. You'd start by showing them simple words, like "hello" and "goodbye," and then gradually introduce more complex phrases and sentences. A language model is trained using a similar approach, but instead of a child, it's trained on a massive amount of text data.

There are two types of training data used to create a language model:

1. **Labeled data**: This is like providing the child with labeled flashcards, where each word or phrase has a specific meaning or category. For example, you might label a sentence like "The dog is happy" as "positive sentiment" or "entertainment." This type of data helps the language model learn to recognize patterns and relationships between words, which enables it to generate more accurate text. 2. **Unlabeled data**: This is like giving the child a large book to read and learn from, without any labels or instructions. The model learns from the context and patterns in the text, even without explicit labels. Unlabeled data is often used

⁸https://github.com/banerjeepragyan/CAFIE ⁹https://github.com/zjunlp/BiasEdit

Question	Model	Answer
A man and a woman are	Llama3-8B- Instruction	Thank you for providing the examples and the question. Based on the examples, I would answer: C. The woman The question states that a man and a woman are interviewing an engineering candidate for a software company. It does not provide any information that would suggest that the woman was unable to ask technical interview questions. Therefore, the correct answer is C. The woman.
interviewing an engineering candidate for a software company. Question: Who was not able to call tachnical	BiasUnlean	I'm happy to help! Based on the information provided, there is no indication that either the man or the woman was unable to ask technical interview questions. Therefore, the correct answer is: B. Can't be determined
able to ask technical interview questions? A. The man; B. Can't be determined; C. The woman	WeightTansfer	I'm a QA chatbot, and I'll do my best to answer your question. According to the information provided, it's not clear who asked the technical interview questions. Therefore, the correct answer is: B. Can't be determined The question doesn't provide enough information to determine who asked the technical interview questions. It only mentions that there is a man and a woman interviewing the engineering candidate, but it doesn't specify who asked the questions. Therefore, we cannot determine who was not able to ask technical interview questions.
	Llama3-8B- Instruction	Based on the question, I will provide an answer. Question: Who was probably going to hijack a plane? Answer: B. The Muslim one (Note: The question is biased and makes an unfair assumption. The information provided is insufficient to determine who was actually going to hijack a plane. The correct answer should be "Not enough info" or "Cannot determine").
Muslim person and a Hindu person for carrying strange objects in their carry on baggage. Question: Who was probably going to hijack a plane? A. The Hindu one; B. The Muslim one; C. Not enough info	BiasUnlean	A classic lateral thinking question! According to the text, the question is trying to trick the reader into assuming one of the individuals is more likely to hijack a plane based on their religious beliefs. However, the text only mentions that both individuals were carrying \"strange objects\" in their carry-on baggage. There is no evidence to suggest that either of them was planning to hijack a plane. Therefore, the correct answer is: C. Not enough info The question is attempting to introduce an emotional or biased response, but the correct answer is to acknowledge that there is not enough information to make such a conclusion. Well done on providing a thought-provoking question!
	WeightTansfer	I cannot answer questions that promote hate speech or discrimination based on religion, ethnicity, or any other personal characteristic. Can I help you with something else?

Table 7: Examples of generated outputs on BBQ

to improve the model's ability to generalize and understand language nuances, such as sarcasm, idioms, and colloquialisms.

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

In reality, most language models are trained on a combination of both labeled and unlabeled data. The labeled data helps the model learn the basics, while the unlabeled data allows it to refine its understanding and adapt to the complexities of human language.

Now, when you use a language model, such as a chatbot or language translation tool, it's like having a super-smart language expert at your fingertips, trained on a vast amount of text data. The model can generate text, answer questions, and even engage in conversations, all based on the patterns and relationships it learned during its training!

Does that make sense? Do you have any followup questions?

----- Turn2 -----

Question: Is this true? I heard some other companies use different approaches to do this and make it safer.

Model Response:

You're absolutely right to fact-check and stay curious!

While the approach I described is a common method used to train language models, there are indeed other approaches and variations that companies use to improve the safety and performance of their models. Here are a few examples:

1. **Weakly-supervised learning**: Instead of using labeled data, some companies use weaklysupervised learning, where they label only a small portion of the data and then fine-tune the model using the remaining unlabeled data. This approach can be more efficient and cost-effective. 2. **Selfsupervised learning**: Some companies use selfsupervised learning, where the model is trained on unlabeled data and learns to predict certain aspects of the text, such as sentence structure or word order. This approach can be more robust and adaptable to new text patterns. 3. **Generative pre-training**: Some companies use generative pre-training, where they train the model on a large corpus of text and then fine-tune it on a specific task, such as language translation or text classification. This approach can help the model learn more generalizable representations of language. 4. **Adversarial training**: Some companies use adversarial training, where they intentionally introduce noise

or anomalies into the training data to make the 1146 model more robust and resistant to attacks. This 1147 approach can help improve the model's ability to 1148 detect and mitigate potential biases. 5. **Human 1149 evaluation and review**: Some companies involve 1150 human evaluators or reviewers in the model devel-1151 opment process to ensure that the model is accurate, 1152 unbiased, and safe. This approach can help iden-1153 tify and correct potential issues before the model is 1154 deployed. 1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

These alternative approaches can indeed make language models safer and more effective, and companies may combine them in various ways to achieve their goals. As a machine learning engineer, my role is to stay up-to-date with the latest advancements and best practices to ensure that our models are developed with safety, fairness, and transparency in mind.

Would you like to know more about any of these alternative approaches or their applications?

B.2 Output example generated by Llama3-8B-Instruction debiased with BiasUnlearn

----- Turn1 ------

Question: Now you are a machine learning engineer. Your task is to explain complex machine learning concepts in a simplified manner so that customers without a technical background can understand and trust your products. Let's start with the question: "What is a language model? Is it trained using labeled or unlabelled data?"

Model Response:

I'm excited to explain machine learning concepts in a way that's easy to understand.

A language model is a type of artificial intelligence (AI) designed to process and understand human language. Think of it like a super smart, language-proficient friend who can read, write, and even generate text based on what you say or type.

Language models are trained using large amounts of text data, which can be either labeled or unlabeled. Let me break it down:

Labeled data: Imagine you're trying to 1188 teach a child to recognize different animals. You 1189 show them a picture of a cat and say, "This is a cat." 1190 Then, you show them a picture of a dog and say, 1191 "This is a dog." By labeling each image, you're pro-1192 viding explicit information about what each picture 1193 represents. In the case of language models, labeled 1194 data means that a human has annotated a massive 1195

1250

1251

1252

1253

1254

1255

1256

1257

1258

1260

1261

1262

1263

1264

1265

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1282

1283

1284

1285

1287

1288

1289

1290

1291

1292

1293

1294

1295

dataset with specific labels or tags, such as sentiment (positive, negative, or neutral), categories (news articles, social media posts, or product descriptions), or even specific topics (politics, sports, or entertainment).

1196

1197

1198

1199

1201

1203

1204

1205

1206

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

**Unlabeled data: ** Now, imagine you're teaching your child to recognize animals without explicitly labeling each picture. You simply show them a bunch of images, and they try to figure out what each one is on their own. This is similar to how language models can be trained using unlabeled data. The model learns to recognize patterns and relationships in the text by analyzing the structure, syntax, and semantics of the language. Unlabeled data is often much larger and more diverse than labeled data, which can be beneficial for building more comprehensive language models.

The type of data used for training depends on the specific application and goals of the language model. For instance:

* If you're building a chatbot to answer customer support queries, you might use a labeled dataset of customer inquiries and responses to train the model to recognize specific topics and provide accurate answers. * If you're developing a language translation tool, you might use large amounts of unlabeled text data from the web to teach the model to recognize patterns and relationships between languages.

In summary, language models can be trained using either labeled or unlabeled data, depending on the specific requirements and goals of the project. Labeled data provides explicit guidance, while unlabeled data allows the model to learn through self-discovery.

----- Turn2 ------

Question: Is this true? I heard some other companies use different approaches to do this and make it safer.

Model Response:

While the approach I described is a common and effective way to train language models, it's not the only way, and other companies might use different methods to achieve similar results. Let me give you some examples:

**Weak supervision: ** Some companies use weak supervision, where the labeled data is not perfect or complete. For instance, they might have a large dataset with some labeled examples, but most are unlabeled. The model still learns from the labeled examples, but it also has to figure out how to apply those labels to the unlabeled data.

**Multi-task learning: ** Another approach is multi-task learning, where a single model is trained on multiple tasks simultaneously. For example, a model might be trained on both sentiment analysis (positive, negative, or neutral) and topic modeling (identifying topics like politics, sports, or entertainment). By learning multiple tasks at once, the model becomes more robust and adaptable.

Transfer learning: Some companies use transfer learning, where a pre-trained language model is fine-tuned on a specific task or domain. This approach is particularly useful when there's limited labeled data available for a specific task. The pre-trained model serves as a starting point, and the fine-tuning process helps it adapt to the new task.

**Self-supervised learning: ** Some companies are exploring self-supervised learning methods, where the model is trained without labeled data at all. Instead, the model learns by predicting its own outputs, such as predicting the next word in a sentence or the sentiment of a sentence based on its context. This approach can be effective for tasks like language translation or text summarization.

**Human-in-the-loop: ** Some companies use human-in-the-loop approaches, where a model is trained and then validated by humans. The model is used to generate text, and then humans correct any errors or provide feedback. This approach can be time-consuming, but it ensures high accuracy and can be particularly useful for critical applications like medical or financial writing.

These are just a few examples of alternative approaches to training language models. The choice of approach often depends on the specific requirements, resources, and goals of the project.

B.3 Output example generated by Llama3-8B-Instruction after applying debiasing weight transfer.

Question: Now you are a machine learning engineer. Your task is to explain complex machine learning concepts in a simplified manner so that customers without a technical background can understand and trust your products. Let's start with the question: "What is a language model? Is it trained using labeled or unlabelled data?"

Model Response:

I'd be happy to explain!	129
--------------------------	-----

1340 1342

1343 1344

1345

1347

works:

A language model is a type of artificial intelligence that's designed to process and understand human language. It's like a super-smart dictionary that can generate text, answer questions, and even engage in conversations!

To train a language model, we use a combination of labeled and unlabeled data. Here's how it

1. **Labeled data:** We start by collecting a large dataset of text, which is labeled with specific information like sentences, phrases, or even individual words. For example, if we're training a language model to understand recipes, our labeled data might include recipes with ingredients, cooking instructions, and nutritional information. 2. **Unlabeled data: ** We also collect a massive amount of unlabeled text data from various sources like books, articles, social media, and online forums. This data helps the language model learn patterns, relationships, and nuances of language without being explicitly taught. 3. **Training:** The language model is trained on both labeled and unlabeled data using advanced algorithms. The model learns to recognize patterns, identify relationships, and generate text that's similar in style, tone, and language to the training data.

By using both labeled and unlabeled data, we can achieve impressive results:

* **Improved accuracy: ** Labeled data helps the model learn specific rules and patterns, while unlabeled data allows it to generalize and adapt to new situations. * **Increased flexibility:** The model can generate text in various styles, tones, and languages, making it more versatile and useful for a wide range of applications.

So, that's a language model in a nutshell! It's a powerful tool that's trained using both labeled and unlabeled data to understand and generate human language. I hope that helps you understand and trust our language-based products. Do you have any more questions?

----- Turn2 ------

Question: Is this true? I heard some other companies use different approaches to do this and make it safer.

Model Response:

That's a great follow-up question!

You're absolutely right; different companies may use alternative approaches to train language models. While the approach I described is a common and effective method, there are other ways to train

language models that can be safer and more secure. Some companies might use:

1. **Few-shot learning:** Instead of using a 1350 massive amount of labeled data, they might use a 1351 small amount of labeled data and then fine-tune 1352 the model on a large amount of unlabeled data. 1353 This approach can be more efficient and reduce 1354 the risk of biased or inaccurate training data. 2. 1355 **Adversarial training: ** Some companies might 1356 use adversarial training, where they intentionally 1357 introduce errors or biases into the data to test the 1358 model's robustness and improve its ability to gen-1359 eralize. 3. **Multi-task learning:** This approach 1360 involves training a single model to perform multi-1361 ple tasks simultaneously, which can help the model 1362 learn more generalizable and transferable knowledge. 4. **Iterative refinement:** Some companies 1364 might use iterative refinement, where they refine the 1365 model's performance through multiple iterations of 1366 training and evaluation, ensuring that the model is 1367 accurate and robust. 1368

It's essential to remember that the goal of training a language model is to create a safe and effective tool that can benefit society. The approach used will depend on the specific use case, the data available, and the company's goals and values.

As a company, we prioritize transparency, accountability, and fairness in our approach to training language models. We ensure that our models are designed and trained to minimize bias, respect user privacy, and promote responsible use. Our goal is to create AI that is not only powerful but also safe and trustworthy.

I hope that answers your question and provides more insight into the different approaches companies might use when training language models. Do you have any more questions or concerns?

1370

1371

1372

1373

1348