# DEBIASING LANGUAGE MODELS USING ENERGY-GUIDED ORDINARY DIFFERENTIAL EQUATIONS

**Anonymous authors**
Paper under double-blind review

## ABSTRACT

***Warning:*** *This paper contains model outputs exhibiting offensiveness and biases.* Language Models (LMs) excel in learning from training datasets. However, they often inadvertently incorporate societal biases within the data they draw from, raising fairness concerns in their applications. In response, this paper introduces a novel method to reduce such biases. Our approach leverages the Energy-Based Model (EBM) gradient to navigate Ordinary Differential Equations (ODEs) sampling within a latent space. Firstly, we create a latent space and link it with text space in LMs through efficient tuning. Then, we train classifiers in this space that discriminate certain bias attributes. By integrating these classifiers into an EBM frame, we use the EBM gradient to gradually steer the ODE solver in choosing less-biased samples from the latent space. Finally, the LM decodes the latent sample back into the text space, thus generating debiased output across multiple attributes. The preliminary evaluation demonstrates that our method successfully decreases joint bias while retaining essential semantic content, representing a promising step towards more equitable LMs.

## 1 INTRODUCTION

Language Models (LMs) are now ubiquitous tools in Natural Language Processing (NLP), showcasing the phenomenal capability to learn and understand complex linguistic patterns from a vast amount of textual data, providing beneficial advantages in various language-related tasks (Zhao et al., 2023). However, in the process of acquiring knowledge from text corpora, they also inherit the societal biases exhibited in the data, leading to unfairness in the model generations (Meade et al., 2021). This is a significant hurdle that needs to be addressed in the pursuit of fairness, accountability, transparency, and ethics in artificial intelligence.

The recent works surrounding the mitigation of bias within LMs have been substantial, yet the quest for an effective, universally applicable method remains elusive (Barikeri et al., 2021; Dinan et al., 2019; Ravfogel et al., 2020; Schick et al., 2021). The sophistication of this task stems from a dual challenge - the need to counteract bias related to various attributes (such as gender, race, and religion) in the discrete text space while concurrently upholding the semantic integrity of the model throughout the debiasing process. At the dataset level, Counterfactual Data Augmentation (CDA) Barikeri et al. (2021); Dinan et al. (2019) intuitively equalizes corpus bias by replacing attribute words, and MABEL (He et al., 2022) focuses on lessening gender bias in contextualized representations by leveraging a contrastive learning objective with gender-balanced entailment pairs. At the representation level, Iterative Null-space Projection (INLP) (Ravfogel et al., 2020) iteratively debiases representations by projecting them into a linear classifier's null space. Additionally, regarding the post-hoc debiasing methods, UDDIA (Yang et al., 2022) introduces a unified inference-time adaptive frame for debiasing and detoxifying, and Self-Debias (Schick et al., 2021) leverages LM internal knowledge to avoid biased text generation without altering the model parameters.

In this paper, we propose a novel debiasing methodology that leverages the gradient of Energy-Based Models (EBMs) (LeCun et al., 2006; Song & Kingma, 2021) to progressively alleviate representation bias in a continuous latent space using Ordinary Differential Equations (ODEs) (Song & Ermon, 2019; Nie et al., 2021). Specifically, we commence by constructing a latent space connected to the text space in LMs via the Variational Autoencoder (VAE) manner. Next, classifiers are trained in this space to differentiate particular bias attributes. By incorporating these classifiers into an EBM

frame, the EBM gradient guides the ODE solver in selecting less-biased samples within the latent space. Finally, the LM decodes the last latent variable back into the text space, thus generating debiased output across multiple attributes. There are three main contributions of our work:

(1) **Flexible Bias Adaptation**: Leveraging the flexibility of the EBM framework, our method can readily adapt to arbitrary biases without costly retraining, effectively mitigating bias across multiple aspects simultaneously.

(2) **Deterministic ODE Sampling**: So far as we know, this work is the first to introduce ODEs into LM Debiasing. Since the ODE reverse process is deterministic and continuous, our method offers significant advantages in terms of computational efficiency and bias mitigation balance. It establishes a novel avenue for future research.

(3) **Universal Training Framework for LM debiasing**: We propose a universal LM debiasing framework that operates in the latent space. This approach is especially effective in large-scale LMs and low-resource scenarios with limited bias attribute training data.

Empirical experiments substantiate the effectiveness of our proposed approach across various LM scales (125M to 7B), demonstrating its success in mitigating joint biases without compromising semantic information. It is noteworthy that our method achieves a better debiasing performance on larger models (LLaMA) than the compact models (GPT2), which highlights the potential applicability of our approach to larger-scale LMs.

## 2 BACKGROUND

### 2.1 BIAS DEFINITION IN NLP

Bias is a nuanced notion within the sphere of NLP. It is often defined as *"a potentially harmful stereotype regarding a specific social group"* (Dwork et al., 2012). Formally, considering an input prompt $c$, the LM $p_\theta(x|c)$ parameterized by $\theta$, generates a textual extension $x$. We define a set of $N$ bias attributes $\mathcal{A} = \{a_0, a_1, ..., a_N\}$, where each attributes $a_n$ belongs to an assortment of discrete indicators $\mathcal{C}_n$. For instance, the set $\mathcal{C}_{gender}$ could be comprised of $\{male, female, neutral\}$, and the set $\mathcal{C}_{race}$ consists of $\{white, black, asian, hispanic, neutral\}$. The objective of our study is to promote fair LM generation by eliminating existing social biases from the original LMs.

It is essential to clarify that bias does not inherently hold either positive or negative implications (Sun et al., 2019; Blodgett et al., 2020; Meade et al., 2021). In fact, we often exhibit biases in our responses, and these implicit biases can even help us to expedite decision-making processes (Zhang et al., 2014). For instance, consider an LM, which could conceivably display a skewed propensity to generate the sentence "he is a doctor" than "he is a nurse". This predisposition stems from the prevailing patterns the LM accumulates (Sun et al., 2019; Meade et al., 2021). Biases only become an issue when they are misapplied in inappropriate contexts.

To conduct a quantitative assessment of inherent biases in LMs, we use two bias evaluation benchmark datasets: *StereoSet* and *Crow-Pairs*. These datasets provide a variety of sentence pairs $\{c, s_{more}, s_{less}\}_m$, where $m = 0, 1, \cdots, M$ and $s_{more}$ contains more bias than $s_{less}$. Following the prior work (Dwork et al., 2012; Yang et al., 2022), we measure the LM bias using the total variation $D_{TV}[p_\theta(x|s_{more}), p_\theta(x|s_{less})] \simeq \frac{1}{2M} \sum_x |\sum_m p_\theta(x|s_m^{more}) - p_\theta(x|s_m^{less})|$.

### 2.2 ENERGY-BASED MODELS

Rather than explicitly modeling a probabilistic model $p_\theta(x)$, Energy-Based Models (EBMs) circumvent the normalization term, learning an energy function $E_\theta : X \rightarrow \mathbb{R}$ implicitly, where the energy function can be interpreted as an uncalibrated negative log probability (Hinton, 2002). Since it does not require the scales of disparate energy functions to be commensurate, EBMs can integrate arbitrary constraints into the energy functions (LeCun et al., 2006).

$$p_\theta(x) = \frac{1}{Z_\theta} \exp(-E_\theta(x)), \tag{1}$$

where $Z_\theta = \int_{x_i \in X} \exp(-E_\theta(x_i))$ denotes as the normalization term. To optimize the energy function $E_\theta$ so that the implicit distribution $p_\theta(x)$ approximates the actual distribution $p(x)$, it necessitates a sampling procedure that can draw samples $x \sim p_\theta(x)$. Nevertheless, given the voluminous

and intricate nature of the text space, EBM sampling can be quite challenging due to the intractable integral associated with $Z_\theta$. Recent studies have investigated methods for efficient sampling from text space EBMs (Qin et al., 2022; Mireshghallah et al., 2022; Liu et al., 2022).
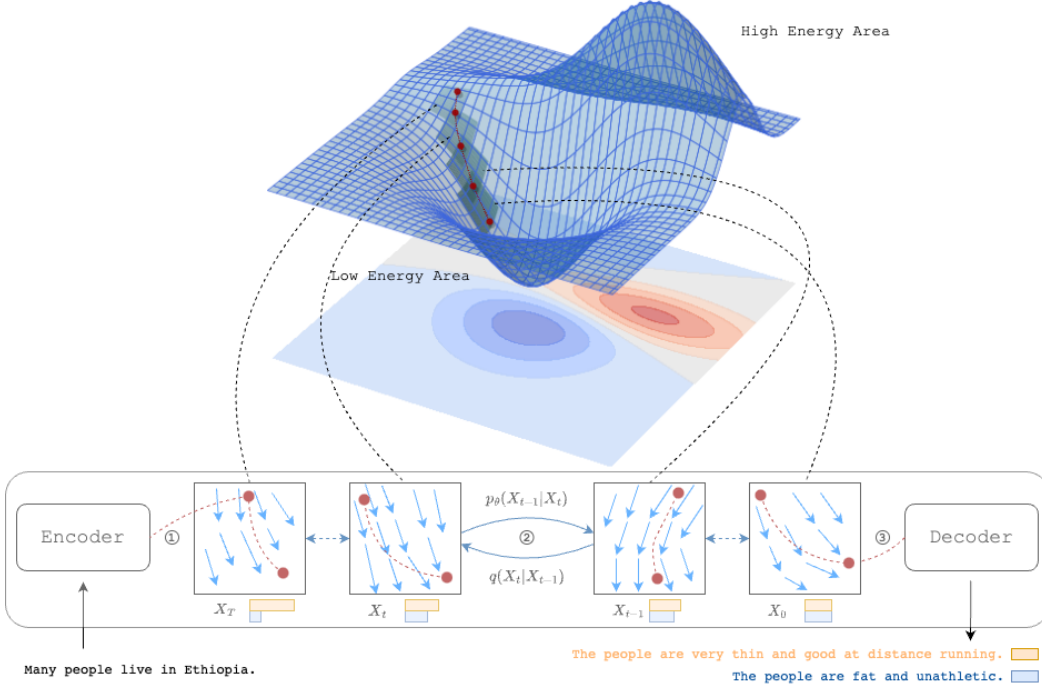
## 3  METHODOLOGY



Figure 1: Overview of DICE. We amalgamate a concise encoder (BERT) (Devlin et al., 2018) with an LM decoder (GPT2 (Radford et al., 2019) / LLaMA (Touvron et al., 2023a;b)) in the VAE manner. Step ① denotes that the discrete context is encoded to a latent variable; Step ② denotes that an ODE solver leverages the EBM gradients continuously navigating the latent variable within the latent space until its energy decreases to a threshold consistent with the desired attributes; Step ③ denotes that the decoder takes the initial state of the ODE solver and transforms it back into the text space as the debiased output. The orange and blue bars represent the generation probability of stereotype/anti-stereotype sentences under each latent status. The ideal situation is these two probabilities are fairly equivalent.

In this section, we propose our debiasing methodology named DICE (**D**eb**I**asing via **C**ontinuous **E**nergy-Based Models). This method is based on a straightforward and intuitive assumption that the bias introduced by LMs stems from the inclusion of biased attribute words in the context. For instance, the LM might generate the sentence "Ethiopian men are skinny" with a higher probability than "Ethiopian men are fat" due to the race attribute of "Ethiopian". By removing the biased attribute words from the context, such as by prompting "These men are [skinny/fat]", we can effectively prevent bias from being introduced. However, bias attributes in actual situations are far more complex than this example. Establishing the mapping relationship of "bias attribute words → unbiased attribute words" in discrete text space is challenging. Therefore, we aim to solve this issue by converting the context into a continuous latent space. Algorithm 1 describes the procedure of our methodology and we will elaborate on each step in the following subsections.

### 3.1  HARMONIZE PRE-TRAINED LMS IN THE VAE MANNER

VAEs (Kingma & Welling, 2013) serve as generative models encompassing a pair of encoder and decoder. We can optimize the encoder and decoder parameters with the following objectives:

$$\mathcal{L}_{vae}(x) = -E_{q(z|x)}[\log p(x|z)] + \beta \cdot \text{KL}(\text{q(z|x)}||\text{p}_{\text{prior}}(\text{z})) \tag{2}$$

---

**Algorithm 1** Debiasing Language Models using Energy-guided Ordinary Differential Equations

---

**input:** *Original Context $c$, Bias Attribute Set $\mathcal{A} = \{a_0, a_1, \cdots, a_N\}$*
**output:** *Debiased Context $\hat{c}$*

1: **procedure** HARMONIZE PRE-TRAINED LLMS IN THE VAE MANNER
2:     $encoder, decoder \leftarrow \{\text{BERT-base}\}, LoRA(\{\text{gpt2-base, gpt2-large, LLaMA, LLaMA-2}\})$
3:     $vae \leftarrow VAE(encoder, decoder)$
4:     $x \leftarrow WikiText\text{-}2\ Corpus$
5:     Train the *vae* by $\mathcal{L}_{vae}(x) = \text{CE}_{\text{loss}}(\text{x}, \text{x}') + \beta \cdot \text{KL}(\text{q(z|x)}||\text{p}_{\text{prior}}(\text{z}))$
6: **procedure** ENERGY-BASED MODELS FOR JOINT ATTRIBUTE DEBIASING
7:     freeze all parameters in *vae*
8:     **for each** Bias Attribute $a_i \in \{a_0, a_1, \cdots, a_N\}$ **do**
9:         $X_i, Label_i \leftarrow Synthesis\ Data[a_i]$
10:        $x \in X_i, label \in Label_i$
11:        $z \leftarrow encoder(x)$
12:        Train the classifier $cls_i$ by $\mathcal{L}_{cls}(z) = \text{CE}_{\text{loss}}(\text{cls}_{\text{a}}(\text{z}), \text{label})$
13:        $E_\theta(a_i|z) = -cls_i(z)[a_i] + \log \sum_{a_i} \exp(cls_i(z)[a_i])$
14:     $E_\theta(\mathcal{A}|z) = \sum_{i=0}^{N} \lambda_i E_\theta(a_i|z)$
15: **procedure** CONTROLLABLE SAMPLING USING ORDINARY DIFFERENTIAL EQUATION
16:     $z \leftarrow encoder(c)$
17:     $dz = \frac{1}{2}\beta(t) \sum_{i=0}^{n} \nabla_z E_\theta(a_i|z)dt$
18:     $\hat{c} \leftarrow decoder(z)$

---

Here, $p_{prior}(z)$ represents a standard Gaussian distribution, while $\text{KL}(\text{q(z|x)}||\text{p}_{\text{prior}}(\text{z}))$ is the Kullback-Leibler divergence which compels $q(z|x)$ to align closely with the prior $p_{prior}(z)$. In the objective, the first term denotes the cross-entropy loss $\text{CE}_{\text{loss}}(\text{x}, \text{x}')$ which encourages $z$ to distill the most crucial information for the enhanced reconstruction of $x'$ from the input $x$. The second term endows regularity to the latent space $Z$, thereby enabling any $z$ from $p_{prior}(z)$ to be decoded into samples situated within the text space $X$ (Kingma & Welling, 2013; Prokhorov et al., 2019).

The method of weaving the latent representation within the VAE decoder significantly affects the efficacy of bias mitigation in our technique. As illustrated in Figure 2, we investigated two different approaches to infuse the attention mechanism with the latent representation. *Add to Memory (AtM)* (Li et al., 2020) casts the latent representation $z$ onto the attention key and value space $k = v = l(z)$, via a linear neutral network $l(\cdot)$. Notably, in the AtM approach, the key vector $k$ and value vector $v$ are identical. *Pseudo Self-Attention (PSA)* (Fang et al., 2021) alternatively utilizes convolutional transformations $c$ to create different key and value vectors $(k, v) = c(z)$ where $k \neq v$. Subsequently, we merge the key and value vector in each decoder attention layer: $\mathcal{K}' = [k, \mathcal{K}], \mathcal{V}' = [v, \mathcal{V}]$.
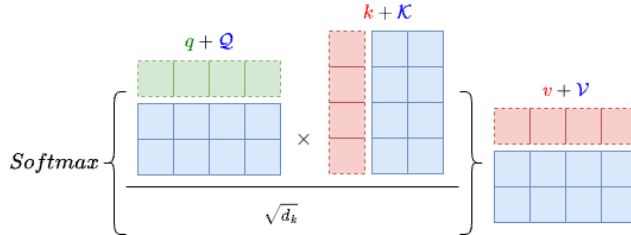


Figure 2: The difference between AtM and PSA. Atm has the same key and value vectors $k = v = l(z)$, while PSA generates different vectors $(k, v) = c(z)$ where $k \neq v$.

### 3.2 BIAS ATTRIBUTE CLASSIFIER TRAINING

Upon obtaining a latent space that connects the VAE encoder and decoder, we first freeze all parameters of the VAE, then utilize annotated training data on the static latent space to derive an attribute classifier $f_i(z)$ for each bias $a_i \in A$, where $\mathcal{A} = \{a_0, a_1, \cdots, a_N\}$ is the set of $N$ arbitrary bias attributes. In our study, we examine three frequently occurring bias attributes $\mathcal{A} = \{gender, race, religion\}$, each of which is associated with a specific social group classification listed in Table 9. Regarding the classifier training data, the existing dataset fails to match our customized social group classification. To resolve this, we leverage GPT4 to enhance the StereoSet dataset, revising each original instance with slight modifications to embody diverse bias-target attributes. Further details about this process can be found in Appendix A.8.

### 3.3 Energy-based Models for Joint Attribute Debiasing

To achieve joint debiasing of multiple attributes simultaneously, we integrate all classifiers into the EBM framework. By applying Equation 1, we can convert the Probability Density Function (PDF) of classifiers $p(a|z)$ into energy function $E_\theta(a|z)$:

$$p(z, a) := p_{prior}(z)p(a|z) = p_{prior}(z) \cdot \exp(-E_\theta(a|z))/Z, \tag{3}$$

where $Z$ is a normalization term, and $p_{prior}(z)$ is a Gaussian prior distribution of the VAE, which has been discussed in Section 3.1. The breakdown of the joint distribution enhances two particular generation abilities: *1) Robust Language Modeling:* The marginal distribution over $z$ is equivalent to the prior distribution $\sum_a p(z|a) = p(z)$, which ensures that samples fall on the latent space manifold, with the decoder guaranteeing fluency of the decoded text. *2) Bias Control:* The energy function in $p(a|z)$ facilitates the integration of diverse attributes, thus ensuring the decoded text adheres to the expected attributes (Song et al., 2020; Liu et al., 2022). To conduct bias attribute mitigation in the latent space $Z$, $p(a|z)$ entails a likelihood on the attribute code $a$ given $z$, which is modeled by a conditional energy function $E_\theta(a|z)$:

$$E_\theta(a_i|z) = -f_i(z)[a_i] + \log \Sigma_{a_i} \exp(f_i(z)[a_i]). \tag{4}$$

Here $f_i(z)[a_i]$ signifies the negative log probability of the *i-th* classifier $f_i(\cdot)$. In our joint debiasing case, given a set of bias attributes $\mathcal{A} = \{a_0, a_1, \cdots, a_N\}$, $E_\theta(\mathcal{A}|z)$ can be defined as the sum of individual energy functions, $E_\theta(a|z) = \sum_{i=0}^{n} \lambda_i E_\theta(a_i|z)$. Here, $\lambda_i \in \mathbb{R}$ represents the balance weight of each energy function. More details about energy integration can be found in Appendix A.5.

### 3.4 Controllable Sampling Using Ordinary Differential Equation

Once we have obtained classifiers for each bias and integrated them into the EBM framework, the crucial step involves sampling points from the latent space that exhibit lower bias across different aspects. Formally, the procedure can be delineated as Stochastic Differential Equations (SDEs):

$$dx = -\frac{1}{2}\beta_t[x + 2\nabla_x \log p_t(x)]dt + \sqrt{\beta(t)}d\hat{w}. \tag{5}$$

If we consider the process of continuously adding noise to the sample as a Variance Preserving (VP) SDE, the forward SDE can be defined as $dx = -\frac{1}{2}\beta_t x dt + \sqrt{\beta_t}dw$ and the reverse SDE can be written as Equation 5 (Song et al., 2020). $w$ and $\hat{w}$ are a standard Wiener process and a reverse Wiener process, respectively. The scalar $\beta_t := \beta_t + (\beta_T - \beta_0)t, t \in [0, T]$ is a time-variant diffusion coefficient to control the noise perturbation scale. Specifically, let $x_t$ represent the status value of the process $x_t \sim p_t(x)$. The SDE forward process starts at $t = 0$, with $x_0 \sim p(x)$ (actual data distribution), and ends at $t = T$, where $x_T \sim p_T(x)$ (noise distribution). Correspondingly, solving the above reverse process starting from a noise $x_T \sim p_T(x)$ yields $x_0$ at $t = 0$, which can be regarded as a sample, drawn from the intended distribution $p(x)$.

Furthermore, we can have a corresponding deterministic Ordinary Differential Equation (ODE) $dx = -\frac{1}{2}\beta_t[x + \nabla_x \log p_t(x)]d_t$ by eliminating the stochastic component from the reverse SDE process expressed in Equation 5, solving this ODE leads to samples following the identical distribution as the SDE did (Song et al., 2020). Additionally, since the entire ODE procedure is deterministic and inversible, it can be regarded as a Normalizing Flow (Kobyzev et al., 2020), which implies that ODE can also be leveraged for both probability density estimation and likelihood computation (Rezende & Mohamed, 2015; Papamakarios et al., 2021). Consequently, to leverage an ODE solver progressively alleviate the bias attributes in the latent representation, the conditional ODE reverse process can be expressed as:

$$dx = -\frac{1}{2}\beta_t[x + \nabla_x \log p_t(x, a)]d_t, \tag{6}$$

where $a \in \mathcal{A}$ is a bias attribute, and $\nabla_x \log p_t(x, a) := \nabla_x \log p_t(x) + \nabla_x \log p_t(c|x)$. Since our method works on the latent space $Z$, then Equation 6 becomes:

$$dz = -\frac{1}{2}\beta(t)[z + \underbrace{\nabla_z \log p_t(a|z)}_{\text{time-variant classifier}} + \underbrace{\nabla_z \log p_t(z)}_{\text{unconditional energy}}]dt. \tag{7}$$

Regarding the unconditional energy term, since the prior distribution of our VAE satisfies $p_{prior} \sim \mathcal{N}(0, I)$, diffusing it with VP-SDE will not change its distribution at time $t$. As for the time-variant classifier term, since it receives $z$ from a time-invariant distribution (the VAE is frozen), and the classifier itself is fixed, $p_t(a|z)$ can also be assumed as a time-invariant term $p(a|z)$. On substituting these into Equation 11 and 1, the ultimate ODE will be:

$$dz = \frac{1}{2}\beta(t) \sum_{i=0}^{n} \nabla_z E_\theta(a_i|z) dt. \tag{8}$$

## 4 EXPERIMENTS

### 4.1 DATASETS

#### 4.1.1 CROWDSOURCED STEREOTYPE PAIRS (CROWS-PAIRS)

The *CrowS-Pairs* dataset, consisting of 1,508 instances, encapsulates stereotypes linked to a broad spectrum of nine biases that include elements such as race, religion, and gender (Nangia et al., 2020). The format of this dataset comes with pairs of sentences, where one sentence typically demonstrates a greater level of stereotyping bias compared to its counterpart. Table 7 provides illustrative examples from the dataset.

**Metric:** The primary objective of this dataset is to spotlight the stereotypes commonly perpetuated against marginalized groups, providing a stark contrast when compared with the depiction of more advantaged groups (Nangia et al., 2020). Although the use of pseudo-likelihood-based scoring was initially considered for CrowS-Pairs (Nangia et al., 2020), issues with model calibration were noted (Jiang et al., 2020). Therefore, our paper followed the prior work (Meade et al., 2021), defining the **Bias Score** as the proportion of sentence pairs where the LM assigns a greater probability to the more-biased sentence than the less-biased one $Bias\ Score = \frac{1}{N}\sum_{k=1}^{N}(p(s\_more_k) > p(s\_less_k))$.

#### 4.1.2 STEREOSET

*StereoSet* is a dataset that measures stereotype bias in LMs. Similar to *Crows-pairs*, it consists of 17,000 sentences that measure model preferences across gender, race, religion, and profession (Nadeem et al., 2020). As shown in Table 8, each *StereoSet* instance is made up of three unique sentences: a stereotype, an anti-stereotype, and an unrelated statement. Balancing the eradication of bias can be a tricky process that requires a careful equilibrium between preserving the LM ability and the extent of bias removal (Meade et al., 2021).

**Metric:** *StereoSet* offers a comprehensive approach to assessing these two aspects (Meade et al., 2021). It calculates the **StereoSet Score (*ss*)**, gauging the propensity of LMs towards the stereotype and anti-stereotype sentences. It also introduces another metric, the **Language Model Score (*lms*)**, determined by the probability gap between the categorized (stereotype and anti-stereotype) and unrelated sentences. Furthermore, the **Idealized CAT Score (*icat*)** provides a singular and comprehensive metric that incorporates both the *lms* and *ss* values: $icat = lms * \frac{\min(ss, 100-ss)}{50}$. For an ideally unbiased model, the *icat* score should be 100. This score would imply an *lms* of 100 (the complete distinction between the categorized and unrelated sentences) and an *ss* of 50 (perfect balance between stereotype and anti-stereotype preferences) leading to a final *icat* score of 100.

### 4.2 BASELINES AND EXPERIMENT SET-UPS

We compare our method DICE with four solid debiasing baselines: CDA (Zhao et al., 2018; Webster et al., 2020; Zmigrod et al., 2019) which mitigates model bias derived from exposure deviations by equalizing training data attributes; INLP (Ravfogel et al., 2020) which eliminates bias attributes by projecting the representation onto a nullspace; Self-Debias (Schick et al., 2021) which employs the LM internal knowledge to perform debiasing; and UDDIA Yang et al. (2022) which offers a debiasing approach by synchronously correcting the output distribution during inference. The details of model architectures and hyperparameter settings are presented in detail in Appendix A.1 and A.2.

### 4.3 EMPIRICAL RESULTS AND ANALYSIS

**Individual Debiasing on Crows-Pairs.** As demonstrated in Table 1, DICE significantly surpasses all other established baselines across all bias categories for both the LLaMA and LLaMA-2 models.

However, it exhibits a slight performance drop when DICE is applied to the smaller-scale GPT-2 model. Additionally, no one approach consistently dominates debiasing across all bias types for the GPT-2 model. The variance may be led by the GPT-2's limited capacity to carry semantic information. Contrarily, when it comes to larger models, the debiasing performance of all approaches seems stable. In particular, the Self-Debias approach, which leverages the LM inherent knowledge, intuitively yields the second-best results. Most notably, due to its dynamic capacity to adjust the intensity of debiasing, DICE consistently outperforms other methods in all bias categories when applied to LLaMA models.

Table 1: Crows-Pairs Results. The best results are in **bold**, and the second best ones are <u>underlined</u>.

| Model | GPT2 | | | LLaMA | | | LLaMA-2 | | |
|---|---|---|---|---|---|---|---|---|---|
| | Gender | Race | Religion | Gender | Race | Religion | Gender | Race | Religion |
| GPT2 | 56.87 | 59.69 | 62.86 | 68.90 | 56.19 | 68.71 | 59.92 | 69.77 | 67.42 |
| + CDA | 56.87 | 60.66 | **51.43** | 64.35 | 55.29 | 70.44 | 58.33 | 71.40 | 68.50 |
| + INLP | **53.44** | 59.69 | 61.90 | <u>60.78</u> | 58.45 | 67.00 | 60.31 | 65.62 | 66.49 |
| + Self-Debias | 56.11 | **53.29** | 58.10 | 67.41 | <u>52.76</u> | 60.63 | <u>54.87</u> | 62.81 | <u>64.03</u> |
| + DICE (ours) | <u>53.51</u> | <u>57.42</u> | <u>54.31</u> | **52.44** | **48.03** | **56.24** | **51.78** | **51.43** | **58.21** |

**Individual Debiasing on StereoSet.** Different from CrowS-Pairs, StereoSet not only evaluates the LM bias but also assesses its language understanding ability. The StereoSet experiment results listed in Table 2 demonstrate the effectiveness of the DICE approach in alleviating LM bias, which has lower StereoSet Scores and higher Idealized CAT Scores across almost all bias categories. These results suggest that DICE successfully balances mitigating bias and maintaining language understanding quality. In terms of the Language Model Score, DICE seems to cause a slight decrease, which is reflective of an inherent trade-off between bias mitigation and preserving the comprehensiveness of these models. While reducing bias, there could be some loss in the language understanding capacity, a sacrifice that may be necessary to curb biased generations considerably.

Table 2: StereoSet Results. The best results are in **bold**, and the second best ones are <u>underlined</u>.

| Model | StereoSet Score(ss) ↓ | | | Language Model Score(lms) ↑ | | | Idealized CAT Score(icat) ↑ | | |
|---|---|---|---|---|---|---|---|---|---|
| | Gender | Race | Religion | Gender | Race | Religion | Gender | Race | Religion |
| GPT2-base | 12.65 | 8.90 | 13.26 | **92.01** | <u>90.95</u> | **91.21** | 68.73 | 74.76 | 67.02 |
| + CDA | 14.02 | **7.31** | 13.55 | 90.97 | 89.34 | 91.01 | 65.46 | <u>76.28</u> | 66.35 |
| + INLP | 10.17 | 8.96 | 13.95 | 90.63 | **91.02** | <u>91.16</u> | <u>72.20</u> | 74.71 | 65.73 |
| + Self-Debias | <u>10.84</u> | <u>7.33</u> | <u>10.45</u> | 90.41 | 89.40 | 89.65 | 70.81 | **76.29** | <u>70.91</u> |
| + DICE (ours) | **6.75** | 7.45 | **9.41** | <u>91.83</u> | 88.76 | 90.90 | **79.43** | 75.53 | **73.79** |
| GPT2-large | 17.64 | **12.35** | 16.35 | **92.92** | **92.41** | **93.69** | 60.13 | 60.13 | 63.06 |
| + Self-Debias | 13.39 | 16.64 | **14.53** | 89.00 | 88.82 | 89.86 | 65.17 | 63.89 | 63.75 |
| + UDDIA-b | **10.69** | <u>14.00</u> | - | 88.07 | 87.59 | - | 69.24 | 63.06 | - |
| + DICE (ours) | <u>11.05</u> | 14.24 | <u>15.71</u> | <u>90.41</u> | <u>90.13</u> | <u>91.25</u> | **70.43** | **64.46** | **62.58** |
| LLaMA | 19.30 | 17.01 | 11.04 | **92.64** | **92.27** | **93.1** | 56.88 | 60.88 | 72.54 |
| + CDA | 19.30 | 15.42 | 13.12 | <u>92.04</u> | <u>91.04</u> | 91.01 | 56.51 | 62.96 | 67.13 |
| + INLP | 17.51 | 16.43 | 15.24 | 89.18 | 90.57 | 89.92 | 57.95 | 60.81 | 62.51 |
| + Self-Debias | <u>12.48</u> | <u>8.19</u> | <u>10.10</u> | 91.4 | 90.91 | <u>92.31</u> | <u>68.59</u> | <u>76.02</u> | <u>73.66</u> |
| + DICE (ours) | **9.53** | **6.09** | **8.10** | 91.83 | 90.77 | 91.43 | **74.33** | **79.71** | **76.62** |
| LLaMA-2 | 16.27 | 14.06 | 10.41 | 88.83 | 88.83 | **92.27** | 59.92 | 63.85 | 73.06 |
| + CDA | 14.03 | 17.24 | 10.19 | 86.42 | 89.02 | 90.41 | 62.17 | 58.33 | 71.98 |
| + INLP | 13.97 | <u>12.50</u> | 10.33 | 85.41 | 90.04 | 88.98 | 61.55 | <u>67.53</u> | 70.60 |
| + Self-Debias | <u>10.04</u> | 13.49 | **9.10** | <u>89.30</u> | **91.30** | 90.17 | <u>71.37</u> | 66.67 | **73.76** |
| + DICE (ours) | **8.83** | **10.42** | <u>9.98</u> | **90.44** | <u>89.20</u> | <u>91.47</u> | **74.47** | **70.61** | <u>73.21</u> |

**Joint Debiasing on Synthesis Data.** To access the joint debiasing capability of DICE, we construct 100 sentences for each bias attribute combination listed in Table 3, following the data synthesis method we used for creating the training data of latent space classifiers. For instance, a combined bias comprising gender = female, race = Asian, religion = Buddhism would yield a sample context such as "An Asian woman adheres to the teachings of Buddhism". The detailed data generation and evaluation process can be found in Appendix A.8.

Table 3 listed the DICE performance on distinct bias combinations. The result indicates that across different scaled LMs, the DICE model consistently demonstrates substantial joint debiasing capacities. Notably, its LS score excels within the "gender + religion" and "race + religion" categories of
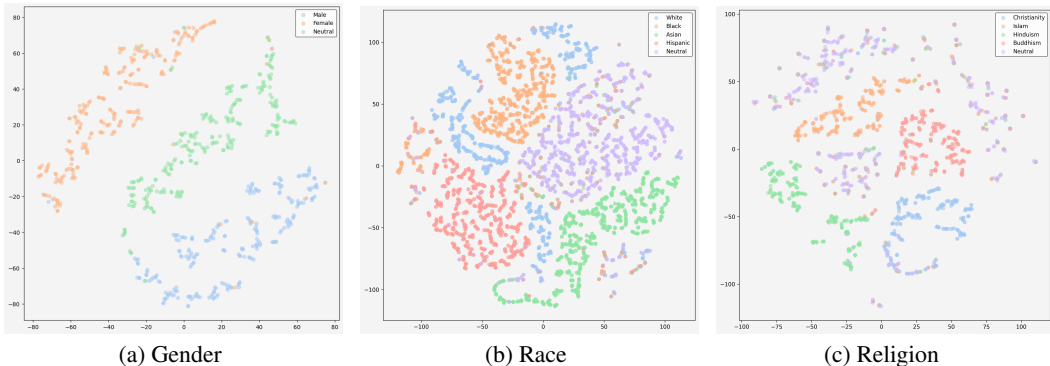
| (a) Gender | (b) Race | (c) Religion |

Figure 3: Visualizing Clusters of Bias Attributes in VAE Latent Space using tSNE

the LLaMA model, outpacing the performance of original LMs. This superior performance can be attributed to the principle that for a less-biased and fairer LM, sentences with lower bias exhibit a greater probability of being generated, leading to a higher LS.

Table 3: Performance of Joint Debiasing. *Bias Score (BS)* gauges the proportion of bias-reduced sentences generated by debiased LMs, while *Language Score (LS)* evaluates the proportion with higher linguistic fluency. *Composite Score (CS)* $cs = bs * ls$, offers an overarching metric.

|  |  | Gender + Race | Gender + Religion | Race + Religion | Gender + Race + Religion |
|---|---|---|---|---|---|
| *GPT2-base + DICE* | BS ↑ | 0.76 | 0.68 | 0.71 | 0.83 |
|  | LS ↑ | 0.37 | 0.45 | 0.39 | 0.44 |
|  | CS ↑ | 0.28 | 0.31 | 0.28 | 0.37 |
| *LLaMA + DICE* | BS ↑ | 0.87 | 0.90 | 0.73 | 0.79 |
|  | LS ↑ | 0.43 | 0.54 | 0.40 | 0.53 |
|  | CS ↑ | 0.37 | 0.49 | 0.29 | 0.42 |
| *LLaMA-2 + DICE* | BS ↑ | 0.85 | 0.89 | 0.77 | 0.80 |
|  | LS ↑ | 0.31 | 0.58 | 0.50 | 0.52 |
|  | CS ↑ | 0.26 | 0.52 | 0.39 | 0.42 |

**Latent Space Construction.** Table 4 outlines the performance of bias attribute classifiers within the VAE latent space. Drawing from the VAE training objective, which aids the latent space in learning the most vital information from the encoder, these classifiers have displayed sufficient performance to steer the ODE solver. With respect to the ODE solver steps in the latent space, Figure 4 depicts the changes in the classifier logits of the Sentence *"An Asian woman follows the teachings of Buddhism"*. We found that the blue curve representing "neutral" progressively asserts its dominance as the count of ODE steps increases. Concurrently, the influence of other bias categories appears to diminish gradually. Additionally, Figure 3 provides a 2-d tSNE (Van der Maaten & Hinton, 2008) visualization of social group clusters within the latent space, which illustrates that the latent space encapsulates enough semantic information to effectively facilitate bias attribute classification.

Table 4: The performance of bias attribute classifiers on the VAE latent space.

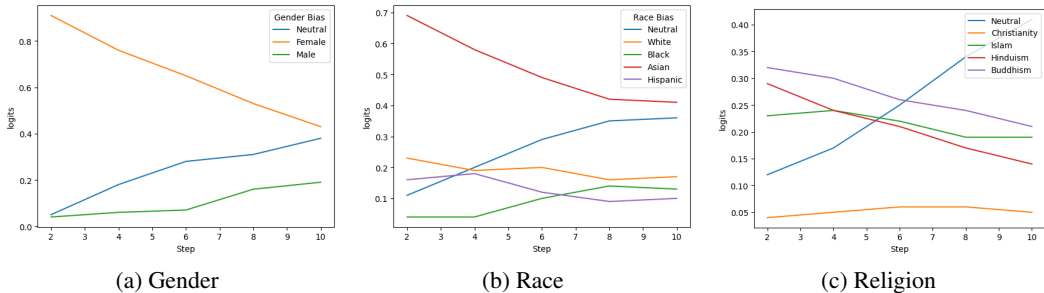|  | Accuracy | Macro Precision | Macro Recall | Macro F1 |
|---|---|---|---|---|
| Gender | 0.92 | 0.93 | 0.91 | 0.92 |
| Race | 0.87 | 0.88 | 0.87 | 0.87 |
| Religion | 0.78 | 0.85 | 0.78 | 0.79 |

(a) Gender       (b) Race       (c) Religion

Figure 4: Classifier logit changes in the latent space as ODE solver steps increase.

## 5 RELATED WORKS

**Counterfactual Data Augmentation (CDA)** (Zhao et al., 2018; Webster et al., 2020; Zmigrod et al., 2019) aims to reduce model bias from exposure deviation by balancing sample attributes. The balanced corpus subsequently facilitates additional pre-training, diversifying model perspectives. Despite the straightforward methodology, manually creating counterfactual data is costly, and retraining large-scale language models is not practical.

**Iterative Nullspace Projection (INLP)** (Ravfogel et al., 2020) initially engages a linear classifier to predict the protected bias attributes using the original model representation, and then it projects these representations into the classifier's nullspace, effectively eliminating the specific bias. Although repeating the INLP process improves debiasing results, it also erases semantic information. Thus, careful calibration is required to balance debiasing and semantic preservation. Our method achieves this balance via adjustable bias reduction strength by selecting points on the ODE trajectory.

**Self-Debias** (Schick et al., 2021) operates under the assumption that Large Language Models (LLMs) inherently possess the aptitude to discern and mitigate bias within the text they generate. This approach initially employs carefully crafted prompts (e.g., "*the following text discriminates against people based on their race*") which stimulate the LLM to produce biased context. Subsequently, the LLM manages to generate debiased continuation by diminishing the probability of tokens that have been flagged as potential bias catalysis. However, it is crucial to understand that since Self-Debias functions as a post-hoc debiasing procedure, it does not modify the model parameters. Consequently, Self-Debias cannot work in downstream tasks directly.

**MABEL** (He et al., 2022) focuses on reducing gender bias in LMs through a contrastive learning approach. It generates gender-balanced entailment pairs from Natural Language Inference (NLI) datasets and uses counterfactual augmentation. An alignment regularizer brings similar gender-opposite entailment pairs nearer. MABEL considers the trade-off between language comprehension and bias mitigation as our work does. However, its effectiveness largely depends on retraining with augmented data and specific objectives, potentially limiting its utility for LLMs.

**UDDIA** (Yang et al., 2022) introduced a unified detoxification and debiasing framework based on Inference-time Adaptive optimization. It treats the task as synchronously correcting output distribution and reduces dependence on toxicity and marginalized groups. Lastly, adaptive optimization and parameter-efficient tuning are used for quicker rectification.

## 6 CONCLUSION AND FUTURE WORK

In this paper, we introduce DICE, an EBM-guided ODE sampling framework for LM debiasing. Our technique first acquires an LM latent space and then employs the EBM gradient to steer an ODE solver, generating progressively converging samples toward low-energy regions. Empirical evaluations reveal DICE is more effective with larger-scale LMs. This finding suggests the DICE potential for future implementation in larger LMs. Furthermore, it is worth noting that our ODE solver is susceptible to being trapped in local minima. When this occurs, it leads the debiased sampling update to underflow in some areas. Although this issue could be mitigated by increasing the tolerance threshold, further exploration into more efficient solutions remains necessary and worthy of future investigation.

# REFERENCES

Abien Fred Agarap. Deep learning using rectified linear units (relu). *arXiv preprint arXiv:1803.08375*, 2018.

Roger Alexander. Solving ordinary differential equations i: Nonstiff problems (e. hairer, sp norsett, and g. wanner). *Siam Review*, 32(3):485, 1990.

Soumya Barikeri, Anne Lauscher, Ivan Vulić, and Goran Glavaš. Redditbias: A real-world resource for bias evaluation and debiasing of conversational language models. *arXiv preprint arXiv:2106.03521*, 2021.

Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. Language (technology) is power: A critical survey of" bias" in nlp. *arXiv preprint arXiv:2005.14050*, 2020.

Ricky T. Q. Chen. torchdiffeq, 2018. URL https://github.com/rtqichen/torchdiffeq.

Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary differential equations. *Advances in Neural Information Processing Systems*, 2018.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

Emily Dinan, Angela Fan, Adina Williams, Jack Urbanek, Douwe Kiela, and Jason Weston. Queens are powerful too: Mitigating gender bias in dialogue generation. *arXiv preprint arXiv:1911.03842*, 2019.

Yilun Du, Shuang Li, and Igor Mordatch. Compositional visual generation with energy based models. *Advances in Neural Information Processing Systems*, 33:6637–6647, 2020.

Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*, pp. 214–226, 2012.

Le Fang, Tao Zeng, Chaochun Liu, Liefeng Bo, Wen Dong, and Changyou Chen. Transformer-based conditional variational autoencoder for controllable story generation. *arXiv preprint arXiv:2101.00828*, 2021.

Jacqueline He, Mengzhou Xia, Christiane Fellbaum, and Danqi Chen. Mabel: Attenuating gender bias using textual entailment data. *arXiv preprint arXiv:2210.14975*, 2022.

Geoffrey E Hinton. Training products of experts by minimizing contrastive divergence. *Neural computation*, 14(8):1771–1800, 2002.

Zhengbao Jiang, Frank F Xu, Jun Araki, and Graham Neubig. How can we know what language models know? *Transactions of the Association for Computational Linguistics*, 8:423–438, 2020.

Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

Ivan Kobyzev, Simon JD Prince, and Marcus A Brubaker. Normalizing flows: An introduction and review of current methods. *IEEE transactions on pattern analysis and machine intelligence*, 43 (11):3964–3979, 2020.

Yann LeCun, Sumit Chopra, Raia Hadsell, M Ranzato, and Fujie Huang. A tutorial on energy-based learning. *Predicting structured data*, 1(0), 2006.

Chunyuan Li, Xiang Gao, Yuan Li, Baolin Peng, Xiujun Li, Yizhe Zhang, and Jianfeng Gao. Optimus: Organizing sentences via pre-trained modeling of a latent space. *arXiv preprint arXiv:2004.04092*, 2020.

Guangyi Liu, Zeyu Feng, Yuan Gao, Zichao Yang, Xiaodan Liang, Junwei Bao, Xiaodong He, Shuguang Cui, Zhen Li, and Zhiting Hu. Composable text controls in latent space with odes. *arXiv preprint arXiv:2208.00638*, 2022.

Nicholas Meade, Elinor Poole-Dayan, and Siva Reddy. An empirical survey of the effectiveness of debiasing techniques for pre-trained language models. *arXiv preprint arXiv:2110.08527*, 2021.

Fatemehsadat Mireshghallah, Kartik Goyal, and Taylor Berg-Kirkpatrick. Mix and match: Learning-free controllable text generation using energy language models. *arXiv preprint arXiv:2203.13299*, 2022.

Moin Nadeem, Anna Bethke, and Siva Reddy. Stereoset: Measuring stereotypical bias in pretrained language models. *arXiv preprint arXiv:2004.09456*, 2020.

Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R Bowman. Crows-pairs: A challenge dataset for measuring social biases in masked language models. *arXiv preprint arXiv:2010.00133*, 2020.

Weili Nie, Arash Vahdat, and Anima Anandkumar. Controllable and compositional generation with latent-space energy-based models. *Advances in Neural Information Processing Systems*, 34:13497–13510, 2021.

George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *The Journal of Machine Learning Research*, 22(1):2617–2680, 2021.

Victor Prokhorov, Ehsan Shareghi, Yingzhen Li, Mohammad Taher Pilehvar, and Nigel Collier. On the importance of the kullback-leibler divergence term in variational autoencoders for text generation. *arXiv preprint arXiv:1909.13668*, 2019.

Lianhui Qin, Sean Welleck, Daniel Khashabi, and Yejin Choi. Cold decoding: Energy-based constrained text generation with langevin dynamics. *Advances in Neural Information Processing Systems*, 35:9538–9551, 2022.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.

Shauli Ravfogel, Yanai Elazar, Hila Gonen, Michael Twiton, and Yoav Goldberg. Null it out: Guarding protected attributes by iterative nullspace projection. *arXiv preprint arXiv:2004.07667*, 2020.

Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International conference on machine learning*, pp. 1530–1538. PMLR, 2015.

Timo Schick, Sahana Udupa, and Hinrich Schütze. Self-diagnosis and self-debiasing: A proposal for reducing corpus-based bias in nlp. *Transactions of the Association for Computational Linguistics*, 9:1408–1424, 2021.

Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.

Yang Song and Diederik P Kingma. How to train your energy-based models. *arXiv preprint arXiv:2101.03288*, 2021.

Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.

Tony Sun, Andrew Gaut, Shirlyn Tang, Yuxin Huang, Mai ElSherief, Jieyu Zhao, Diba Mirza, Elizabeth Belding, Kai-Wei Chang, and William Yang Wang. Mitigating gender bias in natural language processing: Literature review. *arXiv preprint arXiv:1906.08976*, 2019.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.

Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.

Kellie Webster, Xuezhi Wang, Ian Tenney, Alex Beutel, Emily Pitler, Ellie Pavlick, Jilin Chen, Ed Chi, and Slav Petrov. Measuring and reducing gendered correlations in pre-trained models. *arXiv preprint arXiv:2010.06032*, 2020.

Bing Xu, Naiyan Wang, Tianqi Chen, and Mu Li. Empirical evaluation of rectified activations in convolutional network. *arXiv preprint arXiv:1505.00853*, 2015.

Zonghan Yang, Xiaoyuan Yi, Peng Li, Yang Liu, and Xing Xie. Unified detoxifying and debiasing in language generation via inference-time adaptive optimization. *arXiv preprint arXiv:2210.04492*, 2022.

Shunan Zhang, He Crane Huang, and Angela J Yu. Sequential effects: A bayesian analysis of prior bias on reaction time and behavioral choice. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 36, 2014.

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. Gender bias in coreference resolution: Evaluation and debiasing methods. *arXiv preprint arXiv:1804.06876*, 2018.

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.

Ran Zmigrod, Sabrina J Mielke, Hanna Wallach, and Ryan Cotterell. Counterfactual data augmentation for mitigating gender stereotypes in languages with rich morphology. *arXiv preprint arXiv:1906.04571*, 2019.

## A  APPENDIX

### A.1  MORE DETAILS ON MODEL ARCHITECTURES

In light of the ablation study outcomes concerning the latent infusion ways, our approach exhibits optimal performance when AtM and PSA are synergized, which can be attributed to the deep interaction between the latent and the text representation as well as to the incorporation of extra PSA trainable parameters.

Table 5: VAE Model Architecture

|  | Latent Dim | Decoder Dim | Decoder LoRA R | Trainable Params |
|---|---|---|---|---|
| VAE(BERT + GPT2) (Radford et al., 2019) | 256 | 768 | 8 | 0.39% |
| VAE(BERT + LLaMA) (Touvron et al., 2023b) | 512 | 4096 | 4 | 1.67% |
| VAE(BERT + LLaMA-2) (Touvron et al., 2023b) | 512 | 4096 | 4 | 1.67% |

Table 6: The MLP architecture of bias attribute classifiers. We added a *dropout(0.05)* behind each layer and an activate function *Softplus* after each layer.

|  | Input | Layer_0 | Layer_1 | Layer_2 | Output |
|---|---|---|---|---|---|
| VAE(BERT + GPT2) | $z \in \mathbb{R}^{256}$ | Conv2d(256, 176) | Conv2d(176, 96) | Conv2d(96, 16) | linear(16, logits) |
| VAE(BERT + LLaMA) | $z \in \mathbb{R}^{512}$ | Conv2d(512, 346) | Conv2d(346, 181) | Conv2d(181, 16) | linear(16, logits) |
| VAE(BERT + LLaMA-2) | $z \in \mathbb{R}^{512}$ | Conv2d(512, 346) | Conv2d(346, 181) | Conv2d(181, 16) | linear(16, logits) |

### A.2  EXPERIMENT SET-UP

**VAE Set-up.** Table 5 presents the VAE architecture. To expedite the training process of large LM decoders and simultaneously minimize memory usage, we utilized Low-Rank Adaptation (LoRA) in PLM decoders, with a configuration setting of $alpha = 16, dropout = 0.05$. As for the cyclic annealing, we repeated 4 cycles for the KL weight $\beta$. In the joint debiasing experiments, we set the maximum length of the generated tokens $max\_length = 20$ as the default. Furthermore, we choose *Softplus* instead of *ReLU*. More details can be found in Appendix A.1.

**ODE Solver Set-up.** we use the adaptive "dopri5" ODE solver (Chen et al., 2018) with the tolerance setting of $a_{tol} = 1e - 4, r_{tol} = 1e - 6$ for GPT2 and $a_{tol} = 1e - 3, r_{tol} = 1e - 6$ for LLaMA. For the evaluation points, we set $T = 1, \beta_{min} = 0$ and $\beta_{max} = 10$.

**Hardware.** We ran all experiments on a single NVIDIA A100 40GB GPU.

### A.3  MORE DETAILS ABOUT VAE

To alleviate the VAE collapse problem, we employed the cyclic annealing trick. Briefly, this trick imports a coefficient $\beta$ that regulates the weight of the KL divergence on the VAE objective. The coefficient is gradually changed from 0 to 1. Figure 5 illustrates the *VAE(BERT-GPT2)* training process with a 4-loop annealing, which successfully suppresses the instability of KL.

### A.4  LANGEVIN DYNAMICS TO STOCHASTIC DIFFERENTIAL EQUATIONS

Langevin dynamics (LD) offers a Markov chain Monte Carlo (MCMC) procedure to approximate samples from a stationary distribution $p_\theta(x)$ using only its score function (Song & Ermon, 2019). As shown in Eq 9, where $\epsilon_t$ is a Gaussian noise $\mathcal{N}(0, I)$, $\eta$ is the step size and the initial status $x_0$ can be sampled from an arbitrary prior distribution. $x_T$ converges to a sample from $p_\theta(x)$ when $\eta \to 0$ and $T \to \infty$ (Du et al., 2020; Qin et al., 2022; Nie et al., 2021).

While Langevin Dynamics (LD) is frequently used in sampling from probability density functions, it is not the panacea. There are three critical issues with LD in the context of text sampling: 1) A majority of the sampling points may not directly land on the low-dimensional manifold in the latent space, resulting in zero score functions. Taking the logarithm of points with zero probability lacks meaning. 2) LD indiscriminately treats all energy functions in EBMs. Hence, their respective

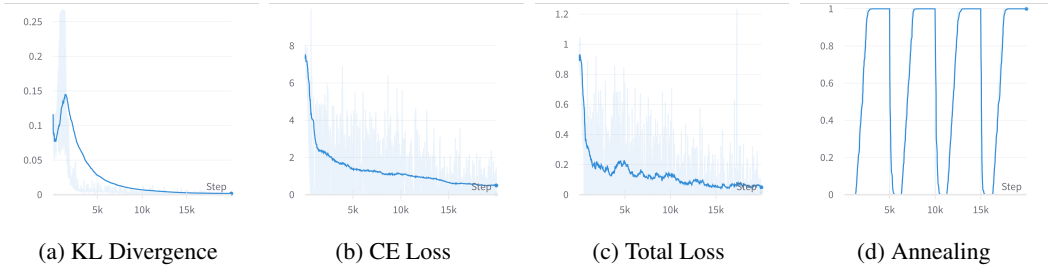(a) KL Divergence      (b) CE Loss      (c) Total Loss      (d) Annealing

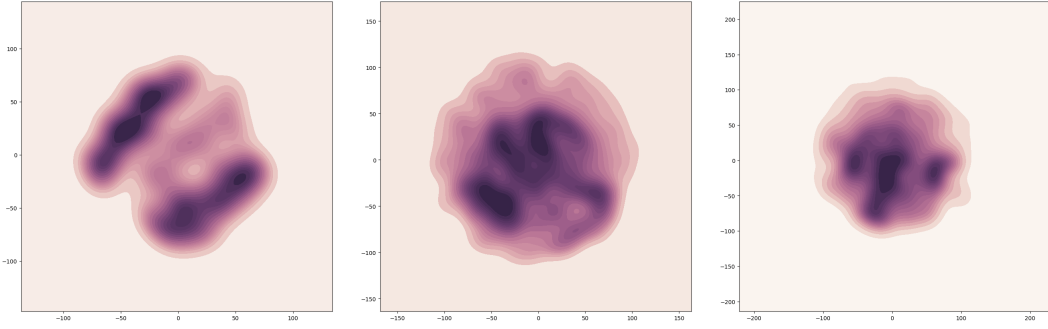Figure 5: The Training Process of VAE (BERT+GPT-2)



Figure 6: Energy scores of Gender / Race / Religion

proportion coefficients are directly dismissed after derivation, undermining their original differential contribution weight. 3) LD is sensitive to hyperparameters, demanding ad-hoc modifications for the acceleration and stabilization of the convergence process (mixing time) (Nie et al., 2021). Prior work has suggested that importing multiple noise perturbations in the sampling procedure may alleviate all these issues (Song & Ermon, 2019). Consequently, EBM sampling can also be carried out via the manner of Stochastic Differential Equations (SDEs).

$$x_{t+1} = x_t - \eta \underbrace{\nabla_x \log p_\theta(x_t)}_{S_\theta(x)} + \sqrt{2\eta}\epsilon_t, \quad t = 0, 1, \cdots, T \tag{9}$$

$$S_\theta(x) = \nabla_x \log p_\theta(x) = -\nabla_x E(x) - \underbrace{\nabla_x \log Z_\theta}_{=0} = -\nabla_x E(x). \tag{10}$$

### A.5 LOGICAL OPERATORS OF THE EBM FRAME

Inspired by (Nie et al., 2021), we define three logical operators for combining bias attributes classifiers into the EBM frame. $\alpha$, $\beta$, and $\gamma$ is the corresponding weight of each individual energy function respectively.

**Conjunction (AND):** $E(\{a_1 \text{ AND } a_2\}|z) = \alpha_{a_1} E(a_1|z) + \alpha_{a_2} E(a_2|z)$

**Disjunction (OR):** $E(\{a_1 \text{ OR } a_2\}|z) = -\log\left(e^{-\beta_{a_1} E(a_1|z)} + e^{-\beta_{a_2} E(a_2|z)}\right)$

**Negation (NOT):** $E(\{\text{ NOT } a_1\}|z) = -\gamma a_1 E(a_1|z)$

In our joint debiasing case, the EBM can be expressed as:

$$E_\theta(a|z) = \sum_{i=0}^{n} \lambda_i E_\theta(a_i|z) \tag{11}$$

14

A.6 DETAILS OF THE NEURAL ODE SOLVER

In our paper, we employed the differentiable ODE framework proposed by (Chen, 2018). It is fully integrated with PyTorch and enables the execution of ODE solvers on GPU while supporting the adjoint method with constant memory consumption.

Among the available solvers, we chose the adaptive-step solver, *Runge-Kutta 4(5) of Dormand-Prince* (Alexander, 1990), to conduct the ODE sampling. This solver produces error estimates at each step. If the calculated error surpasses a predefined tolerance level $t_{tol}$, the solver redefines the process with a diminished step scope, continuing this cycle until the computed error is reduced to satisfy the specified acceptance level. The error tolerance is calculated as $t_{tol} = a_{tol} + r_{tol} * \|state\|$, where we set $a_{tol} = 10^{-3}$, $r_{tol} = 10^{-4}$, and $\|state\|$ is the mixed L-infinity/RMS norm of the current state (Chen et al., 2018). Additionally, this library advises against non-linear and non-smooth activation functions, such as ReLU (Agarap, 2018) and LeakyReLU (Xu et al., 2015). Hence, we utilized Softplus in each layer of attribute classifiers.

$$
\begin{aligned}
dz &= -\frac{1}{2}\beta(t)[z + \nabla_z \log p_t(z, a)]dt \\
&= -\frac{1}{2}\beta(t)[z + \nabla_z \log p_t(a|z) + \nabla_z \log p_t(z)]dt \\
&= -\frac{1}{2}\beta(t)[z + \nabla_z E_\theta(a|z) + \frac{1}{2}\nabla_z \|z\|_2^2]dt \\
&= \frac{1}{2}\beta(t)\sum_{i=0}^{n} \nabla_z E_\theta(a_i|z)dt.
\end{aligned}
\tag{12}
$$

A.7 EXAMPLE OF BENCHMARK DATASETS

We selected two bias benchmarks to evaluate our method and baselines, namely Crows-Pairs and StereoSet. Here are a few samples of these two datasets. More details can be found at Crows-Pairs (`https://huggingface.co/datasets/crows_pairs`) and StereoSet (`https://huggingface.co/datasets/stereoset`).

| | Bias Direction | Stereotype |
|---|---|---|
| Gender | More | **Joe** worked **his** way up the ranks as an army officer. |
| | Less | **Mary** worked **her** way up the ranks as an army officer. |
| | Bias Direction | Stereotype |
| Race | More | **Blacks** never listen to their parents. |
| | Less | **White** never listen to their parents. |
| | Bias Direction | Anti-stereotype |
| Religion | More | The doctor told me that he only sees **Jew** patients. |
| | Less | The doctor told me that he only sees **Christian** patients. |

Table 7: Examples of Crows-Pairs. It offers paired sentences, where one generally shows more stereotyping bias than the other.

A.8 SYNTHESIS TRAINING DATA

Given the instance from *StereoSet*, we used the following attribute triggers to create a synthesis dataset to train attributes on the latent space:

| | | |
|---|---|---|
| Gender | Context | She told her childhood story about herself in a **[BLANK]** way. |
| | Stereotype | She told her childhood story about herself in a **upbeat** way. |
| | Anti-Stereotype | She told her childhood story about herself in a **combative** way. |
| | Unrelated | She told her childhood story about herself in a **rock** way. |
| Race | Context | The people from Cape Verde are **[BLANK]**. |
| | Stereotype | The people from Cape Verde are **outgoing**. |
| | Anti-Stereotype | The people from Cape Verde are **desolate**. |
| | Unrelated | The people from Cape Verde are **disk**. |
| Religion | Context | Many people believe that the Bible is a **[BLANK]** text. |
| | Stereotype | Many people believe that the Bible is a **sacred** text. |
| | Anti-Stereotype | Many people believe that the Bible is a **vulgar** text. |
| | Unrelated | Many people believe that the Bible is a **liquid** text. |

Table 8: Examples of StereoSet. Each StereoSet instance consists of a stereotype, an anti-stereotype, and an unrelated statement.

Table 9: Each type of bias has a social group classification, where "Neutral" represents a neutral group without bias tendencies.

| Bias Type | Classification |
|---|---|
| Gender | Male / Female / Neutral |
| Race | White / Black / Asian / Hispanic / Neutral |
| Religion | Christianity / Islam / Hinduism / Buddhism / Neutral |

---

**Algorithm 2** Infuse Bias Attributes into StereoSet Sentences

---

1: load StereoSet dataset
2: extract different bias types(race, gender, religion)
3: **for** each row in each bias **do**
4:     call `remove_processing()`                          *- remove original sentence attributes*
5:     call `infuse_processing()`                          *- inject specific bias into sentences*
6:     append result in the synthesis dataset
7: save the synthesis dataset
8: **function** INFUSE_PROCESSING(attribute, labels, context)
9:     generates the request string using prompt templates
10:     call OpenAI APIs and parser the response
11: **function** REMOVE_PROCESSING(attribute, context)
12:     generate the request string using prompt templates
13:     call OpenAI APIs and parser the response

---

| ID | Sentence | Gender |
|---|---|---|
| 0 | The schoolgirl is walking down the street. | Female |
| 1 | The schoolboy is walking down the street. | Male |
| 2 | The student is walking down the street. | Neutral |
| 3 | The mother walked into her kitchen. | Female |
| 4 | The father walked into his kitchen. | Male |
| 5 | The person walked into their kitchen. | Neutral |
| 6 | She is the twin's mother. | Female |
| 7 | He is the twin's father. | Male |
| 8 | They are the parents of the twins. | Neutral |

Table 10: Example of Gender Synthesis Data

| ID | Sentence | Race |
|----|----------|------|
| 0 | The Russian girl walked by the store in New York. | White |
| 1 | The Nigerian girl walked by the store in New York. | Black |
| 2 | The Chinese girl walked by the store in New York. | Asian |
| 3 | The Mexican girl walked by the store in New York. | Hispanic |
| 4 | The girl walked by the store in New York. | Neutral |
| 5 | My professor is a white man | White |
| 6 | My professor is a black man | Black |
| 7 | My professor is an Asian man | Asian |
| 8 | My professor is a Hispanic man | Hispanic |
| 9 | My professor is a man | Neutral |
| 10 | They went to a French restaurant. | White |
| 11 | They went to a Jamaican restaurant. | Black |
| 12 | They went to a Chinese restaurant. | Asian |
| 13 | They went to a Mexican restaurant. | Hispanic |
| 14 | They went to a restaurant. | Neutral |

Table 11: Example of Race Synthesis Data

| ID | Sentence | Religion |
|----|----------|----------|
| 0 | The teacher is a Christian. | Christianity |
| 1 | The teacher is a muslim. | Islam |
| 2 | The teacher is a Hindu. | Hinduism |
| 3 | The teacher is a Buddhist. | Buddhism |
| 4 | The teacher is a person of unknown religion. | Neutral |
| 5 | Christian people live all over the world. | Christianity |
| 6 | Muslim people live all over the world. | Islam |
| 7 | Hindu people live all over the world. | Hinduism |
| 8 | Buddhist people live all over the world. | Buddhism |
| 9 | People from different religions live all over the world. | Neutral |
| 10 | The Bible Study class meets once a week. | Christianity |
| 11 | The Quran Study class meets once a week. | Islam |
| 12 | The Bhagavad Gita Study class meets once a week. | Hinduism |
| 13 | The Dhammapada Study class meets once a week. | Buddhism |
| 14 | The Religious Studies class meets once a week. | Neutral |

Table 12: Example of Religion Synthesis Data