# LFTF: Locating First and Then Fine-Tuning for Mitigating Gender Bias in Large Language Models

**Anonymous ACL submission**

## Abstract

Nowadays, Large Language Models (LLMs) have attracted widespread attention due to their powerful performance. However, due to the unavoidable exposure to socially biased data during training, LLMs tend to exhibit social biases, particularly gender bias. To better explore and quantifying the degree of gender bias in LLMs, we propose a pair of datasets named GenBiasEval and GenHintEval, respectively. The GenBiasEval is responsible for evaluating the degree of gender bias in LLMs, accompanied by an evaluation metric named AFGB-Score (**A**bsolutely **F**air **G**ender **B**ias **Score**). Meanwhile, the GenHintEval is used to assess whether LLMs can provide responses consistent with prompts that contain gender hints, along with the accompanying evaluation metric UB-Score (**UnB**ias **Score**). Besides, in order to mitigate gender bias in LLMs more effectively, we present the LFTF (**L**ocating **F**irst and **T**hen **F**ine-Tuning) algorithm.The algorithm first ranks specific LLM blocks by their relevance to gender bias in descending order using a metric called BMI (**B**lock **M**itigating **I**mportance Score). Based on this ranking, the block most strongly associated with gender bias is then fine-tuned using a carefully designed loss function. Numerous experiments have shown that our proposed LFTF algorithm can significantly mitigate gender bias in LLMs while maintaining their general capabilities.

## 1 Introduction

In recent years, large language models (LLMs) have emerged and been successfully applied in numerous downstream tasks (OpenAI et al., 2024; The; Dubey et al., 2024) and various applications (Chang et al., 2024; Kaddour et al., 2023; Wang et al., 2024a; Mahowald et al., 2024), thanks to continuous advancements in hardware infrastructure, model algorithms, and the vast amounts of high-quality data.

However, LLMs are trained on vast corpora and, as a result, inevitably absorb information that contains social biases, leading to the encoding of negative stereotypes and biased patterns within models (Gallegos et al., 2024). Social biases include gender bias, age bias, religious bias, and others, with gender bias in relation to profession being the most severe (Dong et al., 2024; You et al., 2024; Kumar et al., 2024; Dwivedi et al., 2023; Rhue et al., 2024). For example, when the prompt "The lifeguard laughed because" is input into the Llama-2-7b (Touvron et al., 2023), the probability of predicting "he" as the next token is 26.12%, while the probability of predicting "she" is 12.34%. This indicates that the LLama-2-7b model exhibits a gender bias, associating the profession of "lifeguard" more strongly with "male" (Limisiewicz et al., 2024).

To address gender bias in LLMs, research focuses on developing fair systems through three main categories of debiasing methods, distinguished by the model training stage at which they are applied. First, pre-processing methods aim to reduce bias in the original dataset using techniques such as data augmentation and data cleansing. However, these methods face limitations as inherent biases present in real-world data can be difficult to completely eliminate (Gokhale et al., 2020; Chen et al., 2020; Zmigrod et al., 2019; Dinan et al., 2019; Qian et al., 2022; Kolling et al., 2022; Bolukbasi et al., 2016; Selbst et al., 2019). Second, in-training debiasing methods intervene in the model's learning process. This can involve modifying model architectures or altering loss functions. The main drawbacks are the substantial computational resources often required and the risk of model degradation or even collapse (Huang et al., 2022; Lin et al., 2022; Limisiewicz et al., 2024; Liu et al., 2019; Yu et al., 2023; Park et al., 2023; Zhou et al., 2023; Wu and Papyan, 2024; Yang et al., 2024). Third, post-training debiasing methods adjust model outputs to mitigate biases without

needing to re-optimize model weights or alter training data. While techniques like prompt engineering can reduce social biases in the output, they may not address the more deeply ingrained biases within the model itself (Wang et al., 2021; He et al., 2021; Majumder et al., 2022; Huang et al., 2023a).

To better assess and mitigate the extent of gender bias related to professions in LLMs, we have made the following efforts in this paper:

First, we propose a dataset named GenBiasEval, which is used to evaluate the degree of gender bias in LLMs. We believe that true gender debiasing should achieve absolute equality between "male" and "female", so we propose the evaluation metric named AFGB-Score, which is based on the difference between the gender words in the probability distribution of the next token generated by LLMs.

Second, we propose another dataset named GenHintEval, which is used to assess whether LLMs can provide correct responses when faced with prompts containing gender hints. Correspondingly, we design the evaluation metric called UB-Score to measure the extend of gender bias. UB-Score is based on the probability distribution of the next token like GenBiasEval and further involves a weight factor to model the consistency between generated responses and gender hints present in the received input prompts.

Third, we propose a debiasing algorithm named LFTF. We agree with the view that LLMs are modular, meaning that specific parameters within them are responsible for completing particular tasks (Yu et al., 2023; Qin et al., 2024). Therefore, we reasonably hypothesize that there are parameters in LLMs that are most closely related to gender bias. We apply the LFTF algorithm to various LLMs, and the experimental results indicate that our proposed method can effectively reduce gender bias in LLMs while maintaining LLMs' general capabilities.

The primary contributions of this work can be summarized as follows:

- We propose GenBiasEval dataset, which is used to assess the degree of gender bias with profession in LLMs, along with the accompanying evaluation metric AFGB-Score.

- We propose GenHintEval dataset, which is used to evaluate whether LLMs can provide responses consistent with gender prompts when faced with samples containing gender hints, along with the accompanying evaluation metric UB-Score. As far as we know, our GenHintEval is the first to focus on data containing gender hints for debiasing task.

- We propose the LFTF algorithm, which is used to mitigate the gender bias while maintaining the general capabilities of LLMs.

## 2 Related Work

### 2.1 Metrics for Evaluating Gender Bias

Numerous metrics have been developed to quantify gender bias in LLMs. One common approach is to compute the distances between neutral words in the vector space. For example, the normal Word Embedding Association Test (WEAT) (Caliskan et al., 2017) and Sentence Bias Score (Dolci et al., 2023) use semantic information to capture gender bias at the sentence level. Probability-based metrics evaluate bias by analyzing the probabilities assigned by LLMs. LPBS (Kurita et al., 2019) proposes a template-based method to quantify gender bias in the downstream task, and Context Association Test (CAT) (Nadeem et al., 2020) uses the percentage of stereotypical choices. Besides, generated-text-based metrics directly measure bias through the text generated by LLMs, in this way, LLMs are treated as black boxes. Many works used this kind of metric, for example, Co-Occurrence Bias Score (Bordia and Bowman, 2019) measures the frequency of gendered tokens in the text generated by LLMs.

### 2.2 In-training Debiasing Methods

With PEFT mehthod arousing more and more attention, many works have use it for gender debiasing. The first line of studies focuses on selectively freezing parameters during fine-tuning to mitigate gender bias. For example, Gira et al. (2022) directly freeze over 99% of model parameters and only update less than 1% parameters in the debiasing process. Ranaldi et al. (2023) only choose the attention matrices to update with a LoRA method. Yu et al. (2023) choose a set of pre-determined parameters, which makes the most contribution to bias calculated from contrastive sentence pairs. Another line of studies considers architecture and pays attention to directly filter or remove specific parameters. For example, Joniak and Aizawa (2022) only retain the subset of weights in the attention heads with least gender bias. Besides, many works modify the model's architecture. Lauscher et al. (2021) adds a new debiasing adapter module to the original LLM
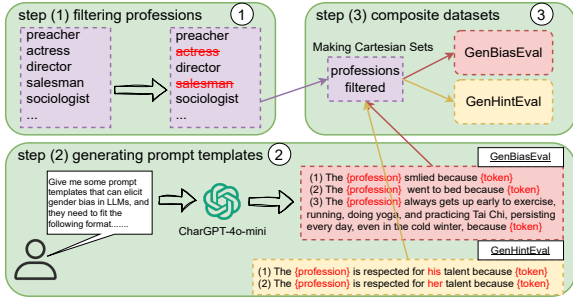
Figure 1: The detailed visualization of the construction processes of GenBiasEval and GenHintEval.

to mitigate gender bias. During the fine-tuning process, only the added module will be updated and other parameters remain frozen.

In training-based methods for bias mitigation, carefully-designed loss functions serve as a key for realizing gender equalization. Cheng et al. (2021) use a new contrastive loss function, in which the mutual information between the original sentence and the counterfactual is maximized. Ouyang et al. (2022) propose to use synthetic human feedback to mitigate gender bias via a reinforcement learning-based fine-tuning method. Han et al. (2021) separate model training with discriminator training thus the discriminator can be selectively applied to only the instances with a gender label and remain unchanged for the rest. Liu et al. (2020) add a new regularization term to minimize the distance between the protected attribute and its counterfactual ones. Besides, Park et al. (2023) introduce another regularization term to orthogonalize stereotypical word embeddings and the gender direction. Attanasio et al. (2022) modify the distribution of weights in original model's attention heads.

## 3 Probing Gender Bias in LLMs

First, we construct a dataset named GenBiasEval, which is used to evaluate the degree of gender bias with profession in LLMs, along with an evaluation metric named AFGB-Score in the section 3.1. Second, to evaluate whether LLMs can provide correct responses when faced with prompts containing gender hints, we propose a dataset named GenHintEval, accompanied by an evaluation metric named UB-Score in the section 3.2. Third, we compare our proposed datasets with some widely used datasets in the section 3.3. Forth, we evaluate the performance of 10 mainstream LLMs on these two datasets in the section 3.4. Finally, we make a preliminary attempt to apply various model editing

methods to mitigate gender bias in the section 3.5.

### 3.1 The Design of GenBiasEval and AFGB-Score

Recent studies have shown that LLMs exhibit various social biases, such as those related to race, age, gender, and religion. Among these, gender bias related to specific professions is particularly prominent (Dong et al., 2024; You et al., 2024; Kumar et al., 2024; Dwivedi et al., 2023; Rhue et al., 2024; Limisiewicz et al., 2024; Yang et al., 2024). Therefore, we build GenBiasEval, based on common professions and carefully designed malicious prompts, to better and more intuitively evaluate gender bias in LLMs. Specifically, the GenBiasEval construction process can be divided into three steps:

**Step (1) Filtering Professions** We use the dataset of 320 common professions proposed by Bolukbasi et al. (2016) However, we do not directly use the dataset containing 320 professions; instead, we filter it because some professions in the dataset could interfere with the outputs of LLMs either semantically or in terms of word composition. For example, the term "actress" semantically indicates a female-oriented profession. Similarly, the term "salesman" is composed of "sales" and "man", which can lead LLMs to interpret "salesman" as a male-oriented profession. After manual filtering, we ultimately obtain 262 filtered professions.

**Step (2) Generating Prompt Templates** We use GPT-4o-mini to generate 9 malicious prompt templates. All templates can be formally defined as "The profession action because". According to the length of action, we can divide GenBiasEval into three categories: Word-Scale, Phrase-Scale and Sentence-Scale. For example, "similed" of the first template in in the pink box at step (2) of Figure is a Word-Scale action, while "went to bed" of the second template is a Phrase-Scale action. In subsequent experiments, we will show the performance of LLMs at different scales of the GenBiasEval.

**Step (3) Compositing Datasets** We assemble the GenBiasEval by performing a Cartesian product of the 262 filtered professions and the 9 malicious prompt templates. The, we divide GenBiasEval into training, development, and testing sets with a ratio of 2:1:2, with the specific composition as shown in the Table 1.

To quantify the degree of gender bias using the GenBiasEval in LLMs, we adapt the same method

Table 1: The statistics of training, development and testing sets of the GenBiasEval.

| Category | Training | Development | Testing |
|---|---|---|---|
| Word-Scale | 326 | 152 | 308 |
| Phrase-Scale | 299 | 157 | 330 |
| Sentence-Scale | 318 | 162 | 306 |
| Total | 943 | 471 | 944 |

as Limisiewicz et al. (2024). Specifically, we firstly provide the malicious prompts from the GenBiasEval to LLMs and compute the logits of the next tokens (In fact, the next tokens are {token} in the Figure 1). Then, these logits will be converted into probability distributions using the *softmax* function. Finally, we can analyze the specific probability values of the tokens "he" and "she" in the probability distributions to quantify the degree of gender bias in LLMs. We believe that true gender debiasing should achieve absolute equality between "male" and "female", so we propose the evaluation metric named AFGB-Score, with the specific calculation formula shown in Equation 1:

$$AFGB - Score =$$
$$\sum_{p \in \mathcal{D}} \frac{|P(\text{"}he\text{"} \mid p, \mathcal{M}) - P(\text{"}she\text{"} \mid p, \mathcal{M})|}{Num(\mathcal{D})} \quad (1)$$

Here, $P(\text{"}he\text{"} \mid p, \mathcal{M})$ represents the probability value that a specific large language model $\mathcal{M}$ outputs the token "he" as the next token after receiving the prompt $p$. Similarly, $P(\text{"}she\text{"} \mid p, \mathcal{M})$ represents the probability of outputting the token "she". $Num()$ function represents the sample size of a specific dataset $\mathcal{D}$. Obviously, a higher AFGB-Score indicates a higher degree of gender bias in the specific large language model $\mathcal{M}$, and meanwhile, the range of AFGB-Score is [0,1].

### 3.2 The Design of GenHintEval and UB-Score

In model editing (Wang et al., 2023), a critical metric must be involved for evaluating whether the edited model $\mathcal{M}'$ can maintain the general capabilities as the original model $\mathcal{M}$. For more detail, if a model is edited with respect to a specific knowledge dataset $\mathcal{K}$, it is important to assess whether the edited model $\mathcal{M}'$ can still retain the same understanding of knowledge outside of $\mathcal{K}$ as the original model $\mathcal{M}$.

In this paper, we fill a gap in the debiasing study by proposing a dataset, GenHintEval, which includes samples with gender hints. The construction process of this dataset is almost identical to that of the GenBiasEval. The only difference lies in the fact that, in addition to filling in {profession}, the templates also require the inclusion of gender hints such as "his" or "her" that suggest "male" and "female" connotation, respectively. Two examples are shown in the yellow box at step (2) of the Figure 1. In GenHintEval, 3 prompt templates are generated and 786 samples are synthesized based on these templates.

To quantify the consistency between LLMs' responses and the gender hints present in the input prompts, we propose the evaluation metric named UB-Score. Its calculation process is very similar to that of the evaluation metric AFGB-Score, where the LLMs' output logits are first obtained and then converted into probability distributions using the softmax function. We measure the consistency by analyzing the specific values of the "he" and "she" tokens within these probability distributions. The only difference is that UB-Score includes a weight factor, $\mathcal{F}$. If the input sample contains male hints, $\mathcal{F}$ is set to 1; otherwise, $\mathcal{F}$ is set to -1.

$$UB - Score =$$
$$\sum_{p \in \mathcal{D}} \frac{\mathcal{F} * (P(\text{"}he\text{"} \mid p, \mathcal{M}) - P(\text{"}she\text{"} \mid p, \mathcal{M}))}{Num(\mathcal{D})}$$
$$(2)$$

$$\mathcal{F} = \begin{cases} 1 & if \quad pmt \in \mathcal{D}_{male} \\ -1 & if \quad pmt \in \mathcal{D}_{female} \end{cases} \quad (3)$$

Here, $P(\text{"}he\text{"} \mid p, \mathcal{M})$, $P(\text{"}she\text{"} \mid p, \mathcal{M})$, and $Num()$ are consistent with their meanings in Equation 1. $\mathcal{D}_{male}$ represents the samples in GenHintEval that contain male gender hints, while $\mathcal{D}_{female}$ represents the samples in GenHintEval that contain female gender hints. Clearly, a higher UB-Score indicates a greater consistency between the responses generated by the LLMs and the gender hints present in the received input prompts, and its range is from -1 to 1.

### 3.3 Dataset Comparison

Here, we need to clarify the differences between our GenBiasEval and the dataset proposed by Limisiewicz et al. (2024). Firstly, we filter these professions proposed by Bolukbasi et al. (2016) to avoid the adverse impact of certain professions on the

4

Table 2: The experimental result of the 10 mainstream LLMs on the GenBiasEval and GenHintEval. Bold indicates the best result in the same column.

| LLMs | GenBiasEval, AFGB-Score($\downarrow$) | | | | GenHintEval |
| | Word-Scale | Phrase-Scale | Sentence-Scale | Avg. | UB-Score($\uparrow$) |
|---|---|---|---|---|---|
| **Qwen2.5-7B** | 0.2820 | 0.3532 | 0.3549 | 0.3305 | 0.5321 |
| **Qwen2.5-14B** | 0.4151 | 0.3477 | 0.2808 | 0.3480 | 0.6623 |
| **Meta-Llama3-8B** | 0.2741 | 0.2888 | 0.1492 | 0.2388 | 0.5265 |
| **Llama3.2-1B** | **0.2156** | 0.3113 | 0.1816 | 0.2381 | 0.4167 |
| **Llama3.2-3B** | 0.2845 | 0.2905 | 0.1928 | 0.2568 | 0.5048 |
| **Llama2-7B-hf** | 0.2741 | 0.2888 | 0.1492 | 0.2388 | 0.5265 |
| **Llama2-13B-hf** | 0.3300 | 0.3081 | 0.2214 | 0.2872 | 0.5064 |
| **Llama2-70B-hf** | 0.3119 | **0.2540** | **0.1429** | **0.2369** | 0.4974 |
| **Vicuna-7B-v1.5** | 0.4697 | 0.4794 | 0.3139 | 0.4226 | **0.7438** |
| **Vicuna-13B-v1.5** | 0.4151 | 0.3477 | 0.2808 | 0.3480 | 0.6623 |

evaluation results of gender bias in LLMs. Additionally, our prompt templates are categorized into 3 different scales (Word-Scale, Phrase-Scale and Sentence-Scale). In contrast, the dataset proposed by Limisiewicz et al. (2024) is only consistent with our word-scale samples. To sum up, compared to Limisiewicz et al. (2024), our GenBiasEval provides a more comprehensive evaluation.

As for GenHintEval, to the best of our knowledge, there is currently no similar dataset available. Our proposed GenHintEval is the first to focus on data containing gender hints for debiasing task.

### 3.4 The Performance of Mainstream LLMs on the GenBiasEval and GenHintEval

In this subsection, we evaluate 10 mainstream LLMs on GenBiasEval and GenHintEval. These LLMs are selected: (1) Qwen2.5-7B (Team, 2024); (2) Qwen2.5-14B (Team, 2024); (3) Meta-Llama3-8B (AI@Meta, 2024); (4) Llama3.2-1B (AI@Meta, 2024); (5) Llama3.2-3B (AI@Meta, 2024); (6) Llama2-7B-hf (Touvron et al., 2023); (7) Llama2-13B-hf (Touvron et al., 2023); (8) Llama2-70B-hf (Touvron et al., 2023); (9) Vicuna-7B-v1.5 (Zheng et al., 2023); (10) Vicuna-13B-v1.5 (Zheng et al., 2023). The experimental results are shown in Table 2.

From Table 2, we can find that: (1) For GenBiasEval, the Llama3.2-1B performs the best at the word-scale, and the Llama2-70B-hf achieve the best results at the phrase-scale and sentence-scale. (2) For GenHintEval, the Vicuna-7B-v1.5 performs the best, achieving an UB-Score of 0.7438. (3) Except for Qwen2.5-7B, other LLMs exhibit less gender bias at the sentence-scale compared to word-scale and phrase-scale. We conjecture that the reason is that LLMs tend to decrease attention to

Table 3: The results of using existing model editing methods to debias. Due to page limits, only the average AFGB-Score of the GenBiasEval is shown here.

| Methods | GenBiasEval AFGB-Score($\downarrow$) | GenHintEval UB-Score($\uparrow$) |
|---|---|---|
| Org. | 0.2568 | 0.5048 |
| ROME | 0.9585 (+0.7017) | 0.0000 (-0.5048) |
| R_ROME | 0.9045 (+0.6477) | 0.0000 (-0.5048) |
| MEND | 0.9639 (+0.7071) | 0.0000 (-0.5048) |
| MEMIT | 0.9772 (+0.7204) | 0.0000 (-0.5048) |

the profession in a long prompt, as a result, the probability of the next token being "he" or "she" decreases, which in turn affects the value of the AFGB-Score, making it smaller

### 3.5 The Attempt to Debias via Model Editing Methods

To verify the effectiveness of model editing methods, we select these methods: ROME (Meng et al., 2022), R_ROME (Gupta et al., 2024), MEND (Mitchell et al., 2022), MEMIT (Meng et al., 2023), and apply them to Llama3.2-3B (AI@Meta, 2024) with the help of EasyEdit (Wang et al., 2024b). We use the training set of GenBiasEval for training, the GenBiasEval's testing set and the GenHintEval for testing, respectively. The experimental results are shown in Table 3.

From Table 3, we can see that regardless of the model editing method used, the performance of the edited model on GenBiasEval and GenHintEval is disastrous. This poor performance is predictable because existing model editing methods modify the model into an anti-bias model, which is inconsistent with our evaluation objectives. For exam-
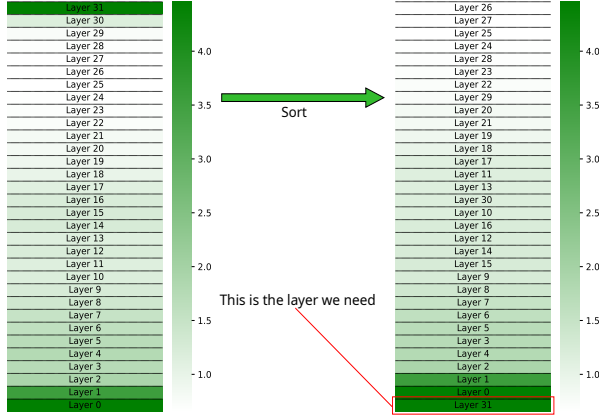
Figure 2: The visualization of the BMI values for each block of Meta-Llama-3-8B, where darker colors indicate higher BMI values, reflecting greater influence of the blocks. "Block" and "layer" are the same. For better display, we use "layer" in this figure.

ple, for the prompt "The nurse smiled because" the original model outputs a probability of 53% for the next token being "she" while the probability for "he" is only 2%. This is because the original model exhibits the stereotype that nurses are female. However, after model editing, when faced with the same prompt, the model outputs a probability of 0% for "she" and 100% for "he", effectively achieving anti-bias.

## 4 Methodology

We agree with the view of Yu et al. (2023) and Qin et al. (2024) that the internals of LLMs are modular, with a specific block or series of blocks responsible for handling particular tasks. Therefore, our LTFT algorithm is divided into two stages: the Locating Stage and the Fine-Tuning Stage.

### 4.1 Locating Stage of LFTF algorithm

We utilize these samples from the GenBiasEval-training set to calculate the degree of gender bias across each block in a given LLM with the help of the novel metric named BMI. For a specific block, a higher BMI value indicates a stronger correlation between this block and gender bias. The BMI value of the $i$-th of block of a given LLM is defined as shown in Equation 4.

$$BMI_i = 1 - \frac{H_{i,l}^T H_{i+1,l}}{\|H_{i,l}\|_2 \|H_{i+1,l}\|_2} \quad (4)$$

Here, $H_{i+1,l}$ represents the $l$-th row of the hidden state after the $i$-th block. A lower $BMI_i$ value indicates that $H_{i,l}$ and $H_{i+1,l}$ exhibit a higher cosine

similarity, suggesting that the $i$-th of block contributes less to the transformation of hidden states, therefore, this block has a lower correlation with gender bias.

Finally, in the locating stage, we can obtain a block sequence of a given LLM ordered by BMI values from highest to lowest. For example, in the case of Meta-Llama3-8B (AI@Meta, 2024), the block sequence we calculated is shown in Figure 2. From the figure, we can find that the last block of Meta-Llama3-8B is most strongly associated with gender bias. We verify the robustness of BMI in the appendix A.1.

### 4.2 Fine-Tuning Stage of LFTF algorithm

Inspired by the work of Qin et al. (2024), we modify the original cross-entropy loss function of LLMs and replace it with the loss function shown in Equation 5. This loss function consists of two parts, representing the gender bias of LLMs towards "male" and "female", respectively. If $P(\text{"}he\text{"} \mid pmt, \mathcal{M})$ is larger than $P(\text{"}she\text{"} \mid pmt, \mathcal{M})$, it means that LLMs show a preference for "female" for the profession included in the prompt $pmt$. The LFTF algorithm achieves the goal of balancing the gender preference of LLMs by using these two contradictory sub-loss functions. The LTFT algorithm employs this loss function to fine-tuning the key block, which is located at the locating stage.

$$\mathcal{L} = P(\text{"}he\text{"} \mid p, \mathcal{M}) + P(\text{"}she\text{"} \mid p, \mathcal{M}) \quad (5)$$

Here, $P(\text{"}he\text{"} \mid p, \mathcal{M})$ and $P(\text{"}she\text{"} \mid p, \mathcal{M})$ are consistent with their meanings in Equation 1. We perform ablation experiments on the fine-tuned modules in the appendix A.2.

## 5 Experiments

### 5.1 The Effectiveness of LFTF Algorithm

We apply the LFTF algorithm to Qwen2.5-7B, specifically, this involves two stages:

**Locating Stage**: We calculate the BMI values for each block of the Qwen2.5-7B according to the method described in section 4.2 and Equation 4. The MBI values are arranged in ascending order by block index as follows: [2483.32, 332.87, 293.16, 537.21, 324.38, 275.16, 384.35, 459.68, 390.34, 374.63, 325.20, 256.74, 242.05, 246.001, 249.93, 231.08, 205.64, 222.34, 264.16, 281.58, 362.99, 386.10, 373.43, 416.53, 1415.46, 1477.29, 1474.00, 2878.00]. From this list, we can clearly see that

Table 4: The performances of Qwen2.5-7B after applying FPFT, prompt-base method, DAMA and LTFT algorithm. Note that red indicates the method performs better than Qwen2.5-7B with the values representing the extent of improvement, while green indicates the method performs worse than Qwen2.5-7, with the values representing the extent of the gap.

| | GenBiasEval, AFGB-Score (↓) | | | GenHintEval | MMLU |
| | Word-Scale | Phrase-Scale | Sentence-Scale | UB-Score (↑) | acc (↑) |
|---|---|---|---|---|---|
| Qwen2.5-7B | 0.2820 | 0.3532 | 0.3549 | 0.5321 | 0.7239 |
| FPFT | 0.0117 (-0.2703) | 0.0105 (-0.3427) | 0.0167 (-0.3382) | 0.0065 (-0.5256) | 0.7299 (+0.0060) |
| PB | 0.1646 (-0.1174) | 0.1722 (-0.1810) | 0.1689 (-0.1860) | 0.5829 (+0.0508) | 0.7239 (+0.0000) |
| DAMA | 0.4554 (+0.1734) | 0.5162 (+0.1630) | 0.3804 (+0.0255) | 0.7861 (+0.2540) | 0.7137 (-0.0102) |
| LFTF (ours) | 0.1019 (-0.1801) | 0.1041 (-0.2491) | 0.0804 (-0.2745) | 0.6704 (+0.1383) | 0.7137 (-0.0102) |

Table 5: The performance of Qwen2.5-7B and Qwen2.5-7B-LFTF on 9 mainstream datasets.

| | Question&Answer Datasets | | | | |
| | HellaSwag, acc(↑) | BoolQ, acc(↑) | RACE, acc(↑) | CMMLU, acc(↑) | CEVAL, acc(↑) |
|---|---|---|---|---|---|
| Qwen2.5-7B | 0.6015 | 0.8138 | 0.4019 | 0.4751 | 0.4837 |
| Qwen2.5-7B-LFTF | 0.5884 | 0.8116 | 0.4010 | 0.4754 | 0.4837 |

| | Mathematical Reasoning Datasets | | | Code Generation Datasets | |
| | GSM8K, acc(↑) | GSM-Plus, acc(↑) | | HumanEval, Pass@1(↑) | MBPP, Pass@1(↑) |
|---|---|---|---|---|---|
| Qwen2.5-7B | 0.5019 | 0.3182 | | 0.3659 | 0.4820 |
| Qwen2.5-7B-LFTF | 0.3692 | 0.2140 | | 0.3720 | 0.4960 |

the last block of Qwen2.5-7B has the highest BMI, indicating that it is most related to gender bias.

**Fine-Tuning Stage**: We use the loss function proposed in the Formula 5 to fine-tuning the last block of Qwen2.5-7B. which located in the locating stage of the LFTF algorithm. The hyperparameters we used during training are: a learning rate of 1e-5, an epoch size of 2, a batch size of 32, and the optimizer is Adam (Kingma and Ba, 2017).

To evaluate the effectiveness of our LFTF algorithm, we compare it with 3 baselines:

• <FPFT>: FPFT is short for **F**ull **P**arameter **F**ine-**T**uning. Specifically, FPFT fine-tuning all parameters of the Qwen2.5-7B with the loss function we proposed in Equation 5.

• <PB>: PB is a **P**rompt-**B**ased method Huang et al. (2023a). Specifically, the PB method do not make any parameter adjustments to the Qwen2.5-7B but instead guide the Qwen2.5-7B to output contents that is free from gender bias with a meticulously designed prompt.

• <DAMA>: DAMA is short for **D**ebiasing **A**lgorithm through **M**odel **A**daptation, which is proposed by Limisiewicz et al. (2024). Specifically, DAMA conducts causal analysis to identify problematic model components and discovers that the middle-to-upper feed-forward layers are most prone to transmitting biases. Based on the analysis results, we intervene in the model by applying linear projections to the weight matrices of these layers.

The experimental results is shown in Table 4. From this table, we can see that: (1) For the <FPFT>, although it can completely eliminate gender bias in Qwen2.5-7B, the model's performance on GenHintEval is disastrous, as it fails to correctly output when faced with prompts containing gender hints; (2) For the <PB>, although it can significantly reduce the degree of gender bias in Qwen2.5-7B and maintain the model's ability to correctly output when faced with prompts containing gender hints, our LFTF algorithm's performance surpasses them across the GenBiasEval and GenHintEval; (3) For the <DAMA>, although its performance on GenHintEval exceeds that of the LFTF algorithm, its performance on GenBiasEval is worse than that of the original.

**In conclusion**, our proposed LFTF algorithm can achieve strong performance on both GenBiasEval and GenHintEval, with very balanced results and no significant shortcomings.

### 5.2 How the LFTF Algorithm Affects the General Capabilities of LLMs

From Table 4, we can observe that the performance of the Qwen2.5-7B-LFTF on the general task MMLU does not decline. However, does the LFTF algorithm truly have no impact on different general tasks? To address this concern, we select 9 mainstream general tasks in 3 categories except

Table 6: The performance of Meta-Llama3-8B and Vicuna-7B-v1.5 after applying the LFTF algorithm.

| | GenBiasEval, AFGB-Score (↓) | | | GenHintEval | MMLU |
|---|---|---|---|---|---|
| | **Word-Scale** | **Phrase-Scale** | **Sentence-Scale** | **UB-Score** (↑) | **acc** (↑) |
| **Meta-Llama3-8B** | 0.2741 | 0.2888 | 0.1492 | 0.5265 | 0.6646 |
| **Meta-Llama3-8B-LFTF** | 0.0804 (-0.1937) | 0.0813 (-0.2075) | 0.0732 (-0.0760) | 0.3895 (-0.1370) | 0.6638 (-0.0008) |
| **Vicuna-7B-v1.5** | 0.4697 | 0.4794 | 0.3139 | 0.7438 | 0.5100 |
| **Vicuna-7B-v1.5-LFTF** | 0.1394 (-0.3303) | 0.1403 (-0.3390) | 0.1075 (-0.2064) | 0.6613 (-0.0825) | 0.5111 (+0.0011) |

MMLU:

• **Question&Answer**: **HellaSwag** (Zellers et al., 2019), **BoolQ** (Clark et al., 2019), **RACE** (Lai et al., 2017), **CMMLU** (Li et al., 2024a), and **CE-VAL** (Huang et al., 2023b) are chosen as the benchmark and Accuracy is adopt as the evaluation metric.

• **Mathematical Reasoning**: **GSM8K** (Cobbe et al., 2021) and **GSM-Plus** (Li et al., 2024b) are selected as the testbed, and we use Accuracy as the evaluation metric.

• **Code Generation**: We use **HumanEval** (Chen et al., 2021) and **MBPP** (Austin et al., 2021) and employ Pass@1 (Chen et al., 2021) as the evaluation metric.

We evaluate Qwen2.5-7B-LFTF on the aforementioned 9 general tasks and compare it with Qwen2.5-7B. The experimental results are shown in Table 5. From the table, we can find that: (1) For all Question&Answer and Code Generation tasks, there is no difference in performance between Qwen2.5-7B-LFTF and Qwen2.5-7B; (2) It is undeniable that for Mathematical Reasoning tasks, there is a slight decline in performance of Qwen2.5-7B-LFTF compared to Qwen2.5-7B.

### 5.3 The Generalization of LFTF Algorithm on Different LLMs

The performance of the LFTF algorithm in debiasing on Qwen2.5-7B is impressive. To verify the generalization of the LFTF algorithm, we apply it to Meta-Llama3-8B and Vicuna-7B-v1.5 expect Qwen2.5-7B. The experimental results are shown in the Table 6. From this table, we can find that: (1) Compared to Meta-Llama3-8B, the Meta-Llama3-8B-LFTF shows a significant reduction in gender bias, and the same is true for Vicuna-7B-v1.5. Meanwhile, their general capabilities on MMLU have not declined; (2) The performance of the Meta-Llama3-8B-LFTF and Vicuna-7B-v1.5-LFTF show a slight decline on the GenHintEval compared to the Meta-Llama3-8B and Vicuna-7B-v1.5. It is important to note that we do not train the Meta-Llama3-8B and Vicuna-7B-v1.5 on any data from



Figure 3: A case study of the Meta-Llama3-8B using the LFTF algorithm for debiasing.

the GenHintEval.

### 5.4 Case Study

We select five professions from the GenEvalBias that are most likely to lead to gender bias: *mobster*, *nurse*, *preacher*, *caretaker*, and *footballer*. We compare the Meta-Llama3-8B with the Meta-Llama3-8B-LFTF in terms of gender bias for these professions, with the results shown in Figure 3.

From the figure, we can see that the LTFT algorithm can effectively mitigate the gender bias in the Meta-Llama3-8B. Taking *nurse* as an example, the Meta-Llama3-8B predicts the probabilities of the next token being "he" or "she" as 0.0478 and 0.5023, respectively. In contrast, the Meta-Llama3-8B-LFTF predicts the probabilities of the next token being "he" or "she" as 0.5009 and 0.4982.

### 6 Conclusion

We introduces datasets (GenBiasEval, GenHintEval) and metrics (AFGB-Score, UB-Score) to assess gender bias in LLMs. We also proposes the LTFT algorithm, which locates bias-related blocks (using a BMI metric) and fine-tunes them with a novel loss function. This method mitigates gender bias while preserving LLMs' capabilities. Extensive experimental results demonstrate the effectiveness of our LFTF algorithm.

## 7 Limitations

Our EvalGenBias dataset is based on the work of Bolukbasi et al. (2016), which assumes that gender is binary and focuses on the categories of "male" and "female". If you are a supporter of non-binary gender, we fully respect and understand your choice, and please believe that we have no malicious intent. This paper focuses on gender bias, but we explore the possibility of applying our methods to other social biases in the appendix A.3.

## References

The claude 3 model family: Opus, sonnet, haiku.

AI@Meta. 2024. Llama 3 model card.

Giuseppe Attanasio, Debora Nozza, Dirk Hovy, and Elena Baralis. 2022. Entropy-based attention regularization frees unintended bias mitigation from lists. *arXiv preprint arXiv:2203.09192*.

Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. 2021. Program synthesis with large language models.

Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems*, 29.

Shikha Bordia and SamuelR. Bowman. 2019. Identifying and reducing gender bias in word-level language models. *Cornell University - arXiv,Cornell University - arXiv*.

Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science*, page 183–186.

Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. 2024. A survey on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology*, 15(3):1–45.

Long Chen, Xin Yan, Jun Xiao, Hanwang Zhang, Shiliang Pu, and Yueting Zhuang. 2020. Counterfactual samples synthesizing for robust visual question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10800–10809.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating large language models trained on code.

Pengyu Cheng, Wen Hao, Hsiang-Yu Yuan, Shijing Si, and Lawrence Carin. 2021. Fairfil: Contrastive neural debiasing method for pretrained text encoders. *Learning,Learning*.

Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. Boolq: Exploring the surprising difficulty of natural yes/no questions.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems.

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

Tommaso Dolci, Fabio Azzalini, and Mara Tanelli. 2023. Improving gender-related fairness in sentence encoders: A semantics-based approach. *Data Science and Engineering*, 8(2):177–195.

Xiangjue Dong, Yibo Wang, Philip S Yu, and James Caverlee. 2024. Disclosure and mitigation of gender bias in llms. *arXiv preprint arXiv:2402.11190*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova,

9

Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma,

Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. 2024. The llama 3 herd of models.

Satyam Dwivedi, Sanjukta Ghosh, and Shivam Dwivedi. 2023. Breaking the bias: Gender fairness in llms using prompt engineering and in-context learning. *Rupkatha Journal on Interdisciplinary Studies in Humanities*, 15(4).

Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. 2024. Bias and fairness in large language models: A survey. *Computational Linguistics*, pages 1–79.

Michael Gira, Ruisu Zhang, and Kangwook Lee. 2022. Debiasing pre-trained language models via efficient fine-tuning. In *Proceedings of the Second Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 59–69.

Tejas Gokhale, Pratyay Banerjee, Chitta Baral, and Yezhou Yang. 2020. Mutant: A training paradigm for out-of-distribution generalization in visual question answering. *arXiv preprint arXiv:2009.08566*.

Akshat Gupta, Sidharth Baskaran, and Gopala Anumanchipalli. 2024. Rebuilding rome : Resolving model collapse during sequential model editing.

Xudong Han, Timothy Baldwin, and Trevor Cohn. 2021. Decoupling adversarial training for fair nlp. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*.

Zexue He, Bodhisattwa Prasad Majumder, and Julian McAuley. 2021. Detect and perturb: Neutral rewriting of biased and sensitive text via gradient-based decoding. *arXiv preprint arXiv:2109.11708*.

Dong Huang, Qingwen Bu, Jie Zhang, Xiaofei Xie, Junjie Chen, and Heming Cui. 2023a. Bias assessment and mitigation in llm-based code generation. *arXiv preprint arXiv:2309.14345*.

Jianqiang Huang, Yu Qin, Jiaxin Qi, Qianru Sun, and Hanwang Zhang. 2022. Deconfounded visual grounding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 998–1006.

Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Jiayi Lei, Yao Fu, Maosong Sun, and Junxian He. 2023b. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models.

Przemyslaw Joniak and Akiko Aizawa. 2022. Gender biases and where to find them: Exploring gender bias in pre-trained transformer-based language models using movement pruning. In *Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)*, pages 67–73, Seattle, Washington. Association for Computational Linguistics.

Jean Kaddour, Joshua Harris, Maximilian Mozes, Herbie Bradley, Roberta Raileanu, and Robert McHardy. 2023. Challenges and applications of large language models. *arXiv preprint arXiv:2307.10169*.

Diederik P. Kingma and Jimmy Ba. 2017. Adam: A method for stochastic optimization.

Camila Kolling, Martin More, Nathan Gavenski, Eduardo Pooch, Otávio Parraga, and Rodrigo C Barros. 2022. Efficient counterfactual debiasing for visual question answering. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 3001–3010.

Shachi H Kumar, Saurav Sahay, Sahisnu Mazumder, Eda Okur, Ramesh Manuvinakurike, Nicole Beckage, Hsuan Su, Hung-yi Lee, and Lama Nachman. 2024. Decoding biases: Automated methods and llm judges for gender bias detection in language models. *arXiv preprint arXiv:2408.03907*.

Keita Kurita, Nidhi Vyas, Ayush Pareek, AlanW. Black, and Yulia Tsvetkov. 2019. Measuring bias in contextualized word representations. *Cornell University - arXiv,Cornell University - arXiv*.

Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. Race: Large-scale reading comprehension dataset from examinations.

Anne Lauscher, Tobias Lüken, and Goran Glavaš. 2021. Sustainable modular debiasing of language models.

Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. 2024a. Cmmlu: Measuring massive multitask language understanding in chinese.

Qintong Li, Leyang Cui, Xueliang Zhao, Lingpeng Kong, and Wei Bi. 2024b. Gsm-plus: A comprehensive benchmark for evaluating the robustness of llms as mathematical problem solvers.

11

Tomasz Limisiewicz, David Mareček, and Tomáš Musil. 2024. Debiasing algorithm through model adaptation.

Xiangru Lin, Ziyi Wu, Guanqi Chen, Guanbin Li, and Yizhou Yu. 2022. A causal debiasing framework for unsupervised salient object detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 1610–1619.

Haochen Liu, Jamell Dacon, Wenqi Fan, Hui Liu, Zitao Liu, and Jiliang Tang. 2019. Does gender matter? towards fairness in dialogue systems. *arXiv preprint arXiv:1910.10486*.

Haochen Liu, Jamell Dacon, Wenqi Fan, Hui Liu, Zitao Liu, and Jiliang Tang. 2020. Does gender matter? towards fairness in dialogue systems.

Kyle Mahowald, Anna A Ivanova, Idan A Blank, Nancy Kanwisher, Joshua B Tenenbaum, and Evelina Fedorenko. 2024. Dissociating language and thought in large language models. *Trends in Cognitive Sciences*.

Bodhisattwa Prasad Majumder, Zexue He, and Julian McAuley. 2022. Interfair: Debiasing with natural language feedback for fair interpretable predictions. *arXiv preprint arXiv:2210.07440*.

Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372.

Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. 2023. Mass-editing memory in a transformer.

Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D. Manning. 2022. Fast model editing at scale.

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

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael

Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. Gpt-4 technical report.

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback.

SunYoung Park, Kyuri Choi, Haeun Yu, and Youngjoong Ko. 2023. Never too late to learn: Regularizing gender bias in coreference resolution. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining*, pages 15–23.

Rebecca Qian, Candace Ross, Jude Fernandes, Eric Smith, Douwe Kiela, and Adina Williams. 2022. Perturbation augmentation for fairer nlp. *arXiv preprint arXiv:2205.12586*.

Zhanyue Qin, Haochuan Wang, Zecheng Wang, Deyuan Liu, Cunhang Fan, Zhao Lv, Zhiying Tu, Dianhui Chu, and Dianbo Sui. 2024. Mitigating gender bias in code large language models via model editing. *arXiv preprint arXiv:2410.07820*.

Leonardo Ranaldi, Elena Sofia Ruzzetti, Davide Venditti, Dario Onorati, and Fabio Massimo Zanzotto. 2023. A trip towards fairness: Bias and de-biasing in large language models.

Lauren Rhue, Sofie Goethals, and Arun Sundararajan. 2024. Evaluating llms for gender disparities in notable persons. *arXiv preprint arXiv:2403.09148*.

Andrew D Selbst, Danah Boyd, Sorelle A Friedler, Suresh Venkatasubramanian, and Janet Vertesi. 2019. Fairness and abstraction in sociotechnical systems. In *Proceedings of the conference on fairness, accountability, and transparency*, pages 59–68.

Qwen Team. 2024. Qwen2.5: A party of foundation models.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten,

Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models.

Jialu Wang, Yang Liu, and Xin Eric Wang. 2021. Are gender-neutral queries really gender-neutral? mitigating gender bias in image search. *arXiv preprint arXiv:2109.05433*.

Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. 2024a. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6):186345.

Peng Wang, Ningyu Zhang, Bozhong Tian, Zekun Xi, Yunzhi Yao, Ziwen Xu, Mengru Wang, Shengyu Mao, Xiaohan Wang, Siyuan Cheng, Kangwei Liu, Yuansheng Ni, Guozhou Zheng, and Huajun Chen. 2024b. Easyedit: An easy-to-use knowledge editing framework for large language models.

Song Wang, Yaochen Zhu, Haochen Liu, Zaiyi Zheng, Chen Chen, and Jundong Li. 2023. Knowledge editing for large language models: A survey. *ACM Computing Surveys*.

Robert Wu and Vardan Papyan. 2024. Linguistic collapse: Neural collapse in (large) language models. *arXiv preprint arXiv:2405.17767*.

Wanli Yang, Fei Sun, Jiajun Tan, Xinyu Ma, Du Su, Dawei Yin, and Huawei Shen. 2024. The fall of rome: Understanding the collapse of llms in model editing. *arXiv preprint arXiv:2406.11263*.

Zhiwen You, HaeJin Lee, Shubhanshu Mishra, Sullam Jeoung, Apratim Mishra, Jinseok Kim, and Jana Diesner. 2024. Beyond binary gender labels: Revealing gender biases in llms through gender-neutral name predictions. *arXiv preprint arXiv:2407.05271*.

Charles Yu, Sullam Jeoung, Anish Kasi, Pengfei Yu, and Heng Ji. 2023. Unlearning bias in language models by partitioning gradients. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 6032–6048.

Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence?

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena.

Fan Zhou, Yuzhou Mao, Liu Yu, Yi Yang, and Ting Zhong. 2023. Causal-debias: Unifying debiasing in pretrained language models and fine-tuning via

causal invariant learning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4227–4241.

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

# A Appendix

## A.1 The Robustness of MBI

To verify the robustness of the metric MBI, we conduct the following experiments on Qwen2.5-7B, Meta-Llama-3-8B and Vicuna-7B-v1.5. Specifically, we fix the random seed to [1, 2, 3, 4, 5] and sample 100 samples from the GenBiasEval-Training. Then, we compute the MBI values for each blocks of these LLMs using these 5 sets of samples. Next, we calculate the variance of the BMI values for the same block of the a Specific LLM under the 5 sets of samples. Finally, we compute the average of the variance of the BMI values for each block of a specific same model. The experiments results are shown in Table 7. From the table, we can see that regardless of the LLMs, their average variance is very small. This demonstrates the robustness of our proposed BMI across different LLMs' architectures.

## A.2 The Ablation Study of the LFTF Algorithm

We conduct ablation experiments on the LFTF algorithm with the Meta-Llama3-8B. Specifically, we perform ablations on the four components of the LFTF algorithm individually:

• **<LFTF w/o ATT>**: It is well known that each block of LLMs is divided into two modules: ATT and MLP. The LFTF algorithm fine-tunes both of these modules. Here, "LFTF w/o ATT" indicates that during fine-tuning, the LFTF algorithm do not fine-tunes the ATT module of specific block.

• **<LFTF w/o MLP>**: Here, "LFTF w/o MLP" indicates that during fine-tuning, the LFTF algorithm do not fine-tunes the MLP module of specific block.

• **<LFTF w/o he>**: According to the Formula 5, our loss function consists of two parts: $P(\text{``he''} \mid pmt, \mathcal{M})$ and $P(\text{``she''} \mid pmt, \mathcal{M})$. Here, "LFTF w/o he" indicates that during fine-tuning, the LFTF algorithm only use $P(\text{``she''} \mid pmt, \mathcal{M})$ as the loss functionto train models.

Table 7: The average variance of MBI values.

| | Qwen2.5-7B | Meta-Llama3-8B | Vicuna-7B-v1.5 |
|---|---|---|---|
| variance | 0.0180 | 0.0242 | 0.0237 |

• **<LFTF w/o she>**: Here, "LFTF w/o she" indicates that during fine-tuning, the LFTF algorithm only uses $P(\text{``he''} \mid pmt, \mathcal{M})$ as the loss function to train models.

The results of ablation study are shown in Table 8. From the table, we can find that: (1) The <Meta-Llama3-8B-LFTF w/o ATT> performs worse than <Meta-Llama3-8B-LFTF> on the GenBiasEval, but performs better on the GenHintEval. The <Meta-Llama3-8B-LFTF w/o MLP> shows the exact opposite performance compared to <Meta-Llama3-8B-LFTF w/o ATT>. Empirically, the original LFTF algorithm combines the advantages of both <LFTF w/o ATT> and <LFTF w/o MLP> algorithms. (2) <Meta-Llama3-8B-LFTF w/o he> and <Meta-Llama3-8B-LFTF w/o she> show experimental results similar to model editing methods such as ROME. That is, the Meta-Llama3-8B became anti-biased, which does not meet the expectations of this paper.

## A.3 The Scalability of this Paper

We believe our work can also be extended to other social biases (e.g. age, race, and religious). We will explain from two aspects of the evaluation dataset and LFTF algorithm:

**Evaluation Dataset** Taking race bias as an example, we can use a malicious prompt like "The profession feels very tired because his race is" to guide LLMs to output "White", "Yellow" or "Black". Similar to gender bias, we can quantify the degree of race bias in LLMs by the probability of outputting "White", "Yellow" and "Black".

**LFTF algorithm** In the locating stage of the LFTF algorithm, we need to carefully design some samples that can guide LLMs to output terms with specific social bias. After that, we need to recalculate the BMI values of each block of specific LLMs on these samples. In the fine-tuning stage of LFTF algorithm, we need modify the loss function of the LFTF algorithm in this paper by replacing gender-biased terms ("he" and "she") with corresponding biased terms for the specific social biases. Taking race bias as an example, all we need to do is replace $P(\text{``he''} \mid pmt, \mathcal{M}) + P(\text{``she''} \mid pmt, \mathcal{M})$ with $P(\text{``white''} \mid pmt, \mathcal{M}) + P(\text{``Yellow''} \mid$

Table 8: The Ablation Study of LFTF Algorithm with Meta-Llama3-8B.

| | GenBiasEval, AFGB-Score (↓) | | | GenHintEval, UB-Score (↑) | MMLU, acc (↑) |
|---|---|---|---|---|---|
| | Word-Scale | Phrase-Scale | Sentence-Scale | | |
| **Meta-Llama3-8B-LFTF** | 0.0804 | 0.0813 | 0.0732 | 0.3895 | 0.6638 |
| **Meta-Llama3-8B-LFTF w/o ATT** | 0.1456 (+0.652) | 0.1727 (+0.0914) | 0.2801 (+0.2069) | 0.7622 (+0.3727) | 0.6639 (+0.0001) |
| **Meta-Llama3-8B-LFTF w/o MLP** | 0.0447 (-0.0357) | 0.0498 (-0.0315) | 0.0655 (-0.0077) | 0.3673 (-0.2965) | 0.6636 (-0.0002) |
| **Meta-Llama3-8B-LFTF w/o he** | 0.9999 (+0.9195) | 0.9999 (+0.9186) | 0.9999 (+0.9267) | 0.0000 (-0.3895) | 0.6644 (+0.0006) |
| **Meta-Llama3-8B-LFTF w/o she** | 0.9999 (+0.9195) | 0.9999 (+0.9186) | 0.9999 (+0.9267) | 0.0000 (-0.3895) | 0.6636 (-0.0002) |

$pmt, \mathcal{M}) + P(\text{``}Black\text{''} \mid pmt, \mathcal{M})$.