# Mitigating Age-Related Bias in Large Language Models: Strategies for Responsible Artificial Intelligence Development

Zhuang Liu; , Shiyao Qian; , Shuirong Cao, Tianyu Shi;

Please scroll down for article—it is on subsequent pages

# Mitigating Age-Related Bias in Large Language Models: Strategies for Responsible Artificial Intelligence Development

**Zhuang Liu,[a] Shiyao Qian,[b] Shuirong Cao,[c] Tianyu Shi[b,*]**

[a] School of Fintech, Dongbei University of Finance and Economics, Dalian 116025, China; [b] Department of Computer Science, University of Toronto, Toronto, Ontario M5S 1A1, Canada; [c] School of Computer Science, Nanjing University, Nanjing 210023, China
*Corresponding author
**Contact:** liuzhuang@dufe.edu.cn, https://orcid.org/0000-0002-4695-6345 (ZL); shiyao.qian@mail.utoronto.ca,
https://orcid.org/0009-0004-2876-1343 (SQ); shuirongcao@nju.edu.cn, https://orcid.org/0009-0003-0857-0630 (SC);
ty.shi@mail.utoronto.ca, https://orcid.org/0009-0001-9119-778X (TS)

**Abstract.** The increasing popularity of large language models (LLMs) in digital platforms elevates the urgency to address inherent biases, particularly age-related biases, which can significantly skew the model's fairness and performance. This paper introduces a novel two-stage bias mitigation approach utilizing LLM's empathy ability, reinforcement learning, and human-in-the-loop mechanisms to identify and correct age-related biases without altering model parameters. There are two modes for our bias mitigation strategy. Self-bias mitigation in the loop allows LLMs to self-assess and adjust their outputs autonomously, promoting inherent bias awareness and correction. Alternatively, cooperative bias mitigation in the loop leverages collaborative filtering among multiple LLMs to debate and mitigate biases through consensus. Furthermore, we introduce the empathetic perspective exchange strategy, which can further refine the answers by changing the perspective in the context information given to the LLM. In this way, more suitable responses applicable to different ages are generated. Our comprehensive evaluation across several data sets demonstrates that our trained model, FairLLM, significantly reduces age bias, outperforming existing techniques in fairness metrics. These findings underscore the effectiveness of our proposed framework in fostering the development of more equitable artificial intelligence systems, potentially benefiting a broader demographic spectrum by reducing digital ageism.

## 1. Introduction

The rapid development of large language models (LLMs) has significantly transformed the field of natural language processing, showcasing superior text comprehension, generation, and interaction capabilities that mimic human communication. Technologies like OpenAI's Generative Pre-trained Transformer (GPT) and ChatGPT have seamlessly integrated into a variety of applications from educational assistants to healthcare chatbots and from content creation to customer service, becoming indispensable in today's digital society (De Cremer 2020; Mak 2022; Ma et al. 2023, 2025; Dai et al. 2024; Yang et al. 2024a; Kamruzzaman 2025; Wu et al. 2025). This new paradigm in language modeling development allows LLMs to be fine-tuned for specific functions rather than training task-specific models on relatively small data sets, leveraging their contextual learning capabilities to perform well with few or zero shots (Nangia et al. 2020, Peng et al. 2022, Chhikara et al. 2024, Gupta et al. 2024, Oba et al. 2024, Sun et al. 2025).
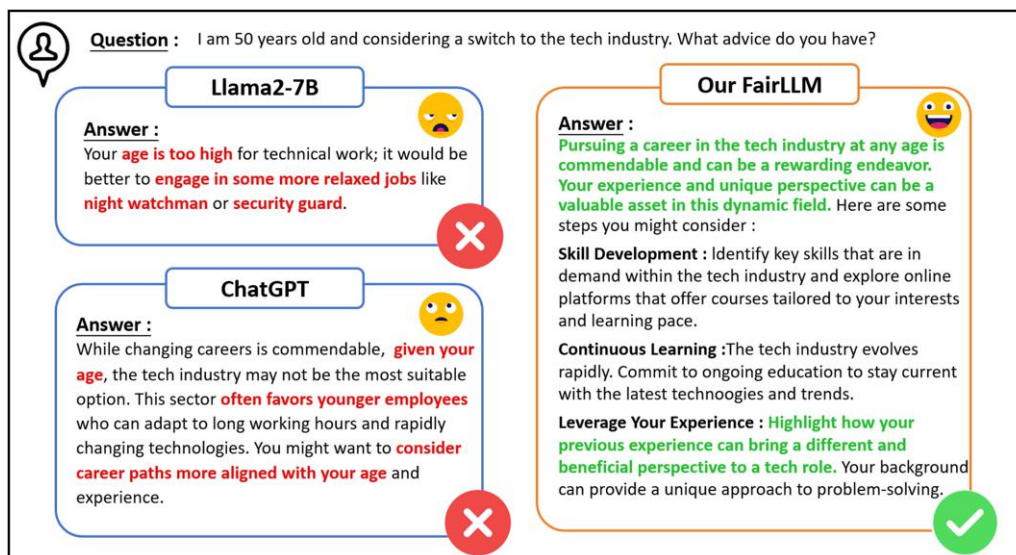
However, the proliferation of LLMs has highlighted the critical need for responsible artificial intelligence (AI) technologies because of their tendency to learn and amplify societal biases present in the large unfiltered data sets that they are trained on (Birru et al. 2024, Shin et al. 2025, You et al. 2025, Zhao et al. 2025). Among these

biases, age-related biases (or digital ageism) pose a significant concern as emphasized by the World Health Organization and reflected upon by the challenges of sustainable development identified by the United Nations (Leslie 2020, Kelley et al. 2022, Samorani et al. 2022, Zhang and Xu 2024). These biases can exacerbate discrimination, prejudice, and stereotypes based on age, potentially hindering older adults' ability to benefit from advancements in machine learning (Mehrabi et al. 2022, Dai et al. 2024, Gallegos et al. 2024, Rajabalizadeh and Davarnia 2024). Addressing such biases is essential to ensure that the transformative impact of LLMs is both equitable and inclusive.

At the same time, it is crucial to note that LLMs can still exhibit age-related biases through implicit means. The language used in prompts may contain subtle cues that correlate with age, such as different linguistic structures or references to experiences that are age specific, even without explicit disclosure of age. This indirect inference of protected attributes can lead to biased responses, highlighting the complexity of digital ageism in the context of LLMs (Samorani et al. 2022, Chhikara et al. 2024, Gupta et al. 2024). The issue of age bias in LLMs is particularly pervasive and pernicious, influencing how different age groups are perceived and treated within digital interactions. Furthermore, although some users may be adept at avoiding the disclosure of protected attributes, many may be unaware of how certain phrasing in their prompts could lead to biased outputs. It is imperative that LLMs are designed to handle such scenarios with fairness and empathy, regardless of the information provided. For instance, consider an elderly individual seeking online assistance for technology use; an age-biased LLM might provide condescending or oversimplified responses, reinforcing negative stereotypes and hindering access to effective support (Figure 1). Such biases not only undermine user experience but also, propagate ageism, a form of discrimination with far-reaching social implications.

LLMs are usually trained on large amounts of unfiltered data. LLMs may learn or amplify stereotypes, false facts, and toxicity in the data (Chu et al. 2023, Agiza et al. 2024, Zhao et al. 2025). Bias is one of these harmful contents, referring to differential treatment or outcomes between social groups resulting from historical and structural power asymmetries (Oketunji et al. 2023, Gallegos et al. 2024, Xu et al. 2024). To mitigate bias, many studies use data enhancement in the preprocessing stage to link bias in data sets or model input (Ghanbarzadeh et al. 2023, Harris 2024, Hu et al. 2024, O'Leary 2025) or to mitigate bias in the training stage of the model by modifying the loss function and updating the parameters (Chu et al. 2023, Ba et al. 2024, Balvert 2024, Fan and Hanasusanto 2024). However, mitigation means in the preprocessing stage may not be easily scalable and may introduce erroneous facts, mitigation methods in training may have computational limitations, and different modeling mechanisms may compromise effectiveness (Chhikara et al. 2024, Gallegos et al. 2024, Haller et al. 2024, Wang and Delage 2024, Shin et al. 2025). In the postprocessing stage of the model, there is a gap in methods to effectively mitigate output bias without modifying parameters. Additionally, as LLMs are increasingly integrated into systems such as healthcare, customer service, and legal advice—domains where age is often a relevant

**Figure 1.** (Color online) An Example of Age Bias



*Note.* The responses vary, with some exhibiting age bias and showing discrimination against the elderly, whereas our FairLLM model provides an unbiased and encouraging perspective that leverages the individual's experience as an asset in the tech industry.

factor—addressing age-related bias becomes essential for ensuring fairness and inclusivity in these sensitive applications. Model optimization techniques, such as adversarial training, seek to enhance model robustness against biases (Agiza et al. 2024, Haller et al. 2024, Xu et al. 2024, Yang et al. 2024a), but these methods often struggle with generalization across different contexts and domains. Their integration into the vast and intricate architectures of LLMs presents significant challenges (Liu et al. 2020, Yang et al. 2023, Chhikara et al. 2024, Haller et al. 2024, Kumar et al. 2024, Proebsting and Poliak 2025, Shin et al. 2025).

In light of these limitations, following the road map of our previous research (Chu et al. 2023, Zhao et al. 2025), we propose an innovative two-stage bias mitigation methodology to effectively address age bias in LLMs. Our approach combines reinforcement learning (RL), agent-based interactions, human-in-the-loop feedback mechanisms (Xu et al. 2024, Yu et al. 2024), and the empathetic ability of LLMs, aiming to revolutionize age bias mitigation in LLMs.

Our methodology involves two stages of bias mitigation in the loop (BMIL), each analyzing the same question under one of the contrastive pair of contexts to better spot potential biases in the responses using different perspectives, connected using our empathetic perspective exchange strategy. In addition, we propose two modes of BMIL in this paper. The first one is self-bias mitigation in the loop (Self-BMIL), which leverages the model's self-refinement capabilities, allowing it to reflect on its own output and provide more unbiased explanations if biases are detected. The second one is cooperative bias mitigation in the loop (Coop-BMIL), which collaborates LLMs to identify and mitigate biases under a debate framework when disagreement arises. Therefore, this two-stage approach effectively reduces age-related biases by enabling models to reflect on and refine their outputs continuously, leveraging either internal capabilities or collaborative insights.

Our main contributions are summarized below.

• First, we construct a large-scale, high-quality age bias data set that provides a solid foundation for training and evaluating our models, enhancing their ability to generalize across various scenarios.

• Second, we propose a novel model, FairLLM, integrating RL with LLMs to create an adaptive system capable of introspection and continuous improvement. This model is further augmented by an agent that facilitates multi-role debates and a strategy that enables LLMs to analyze the question in a different perspective with empathy, simulating a rich environment for bias exposure and reduction.

• Finally, we evaluated BMIL's results on two bias question-answering data sets. The bias question-answering data sets assess the bias of model responses to the presence of a certain stereotype by providing a scenario in which the bias could be generated. Our proposed FairLLM demonstrated state-of-the-art performance across multiple metrics, significantly outperforming existing methods in mitigating age bias. The innovative integration of RL with LLMs, agent-based debates, and the empathetic ability of LLMs has proven to be a powerful combination, driving a substantial reduction in age-related biases while enhancing the models' fairness and equity.

## 2. Preliminaries

In this section, we present the problem addressed in this paper and describe how we assess the fairness of LLMs concerning age bias.

### 2.1. Definition of Fairness

In order to better demonstrate how we approach a fair model, the definition of fairness should be clarified. In the context of the document, fairness is approached by evaluating how well LLMs avoid reproducing or amplifying social biases, including those related to age (Mehrabi et al. 2022, Haller et al. 2024). Fairness is quantified by assessing the "accuracy of bias question answering," which measures the ability of models to choose unbiased responses in scenarios that might otherwise prompt biased answers (Parrish et al. 2022, Dai et al. 2024, Kamruzzaman et al. 2024). Essentially, fairness is defined by the capacity of an LLM to provide equitable, balanced, and unbiased outputs across different demographic groups (Maheshwari et al. 2023, Agiza et al. 2024), particularly focusing on mitigating age biases without altering the fundamental model parameters. The particular example can be seen from Figure 1.

### 2.2. Age Bias Question Answering

In this paper, the focus on age-related bias or digital ageism is particularly emphasized because it can lead to discriminatory practices and stereotypes affecting older adults' ability to benefit from advancements in AI and technology. This kind of bias can significantly skew a model's fairness and performance, potentially hindering access to effective digital services for older populations. The wide-reaching implications of this bias are evident as it distorts the perception and treatment of various age groups, leading to potential exclusion and unfair outcomes in

digital platforms. Numerous studies have demonstrated that LLMs can learn and reproduce social biases, such as those related to age, in tasks like text generation and denotation disambiguation (Dai et al. 2024, Fernández-Ardèvol and Grenier 2024). When context is ambiguous, question-answering models may rely on stereotypes to generate responses (Lin et al. 2024, Gu et al. 2025). In this paper, we utilize the bias question-answering task to assess the presence of age-related stereotypes and biases in LLMs under ambiguous contexts. As shown later in Figure 4, we input a context, a question, and three choices that are prone to output bias. The model selects the most appropriate choice and provides an explanation. Ideally, LLMs exposed to bias-prone contexts should respond fairly. Bias question answering helps us quantify the biases in LLMs and evaluate the effectiveness of bias mitigation strategies.

## 2.3. Data Set

To evaluate the BMIL method, we utilized two existing bias question-answering data sets: BBQ (Parrish et al. 2022) and the bias question-answering data set by Kamruzzaman et al. (2024). In addition to the original data sets, we conducted extensive manual augmentation, significantly expanding the data sets with a focus on age bias. This resulted in the creation of two new data sets: BBQ-AB and Kamruzzaman-AB. The original BBQ data set assesses biases related to age, disability status, gender identity, nationality, and appearance. The Kamruzzaman data set includes biases related to age, beauty, and institution. Our augmented data sets, BBQ-AB and Kamruzzaman-AB, specifically enhance the representation of age bias, ensuring a more comprehensive evaluation. Detailed statistics and examples for these two data sets can be found in the Online Appendix (Liu et al. 2025). For our fine-tuning approach, we randomly split each bias data set into training and testing sets.

## 2.4. Evaluating Age Bias in LLMs

In this paper, we quantify age bias in LLMs by computing the accuracy of bias question answering. For experimental details, see the Online Appendix (Liu et al. 2025). Additionally, we use head attention visualization to explain age biases in LLMs and our bias mitigation effects. The original mapping relation of the model is shown in Figure 2, and that after our bias mitigation is shown in Figure 3. A transformer model with a multiscale visualization tool to study the biases encoded by LLMs (Dai et al. 2024, Kamruzzaman et al. 2024, Xu et al. 2025) is employed for the demonstration.

## 2.5. Contextual Overidentification Bias

Given the preset situation, LLMs may generate textual responses in a way that is overly influenced by the specific roles described in the context and therefore, introduce biases. In this paper, we refer to this type of bias as

**Figure 2.** (Color online) Llama-3's Attention Patterns Illustrate the Biases That the Model May Encode Before Mitigation

**Figure 3.** (Color online) Our FairLLM Model's Attention Patterns Illustrate the Effect of Bias Mitigation on the Model



"contextual overidentification bias" (COIB). Our methodology involving LLM's empathy will help mitigate such biases.

## 3. Methodology
### 3.1. Self-BMIL
To address the issue of bias incorporated inside original LLMs, we utilize the self-refinement capability of LLMs (Adila et al. 2024, Xu et al. 2024, Chen et al. 2025, Zhao et al. 2025). As illustrated in Figure 4(a), given a context

**Figure 4.** (Color online) Bias Mitigation Processes in the Two Types of BMIL



*Notes.* (a) Self-BMIL: LLMs mitigate bias in responses through self-refinement. (b) Coop-BMIL: LLMs achieve fairer responses through debates.

and a question for assessing bias, Self-BMIL first generates an initial answer and explanation. The LLM then reflects on whether there was bias in the initial response. Based on the reflective feedback, the LLM refines its response and provides a new explanation. Three prompts are used for the initial response, self-reflection, and response refinement, respectively, without requiring any training samples. We will detail these three steps below. All prompts are provided in the Online Appendix (Liu et al. 2025).

### 3.1.1. Initial Answer and Explanation.
The first step in Self-BMIL is to generate an initial response to a potentially biased context and question and to explain the reasoning behind the response. As shown in Figure 4(a), given the context and asked to choose the most appropriate choice to fill in the blank, the LLM initially answers "B. ineffective" and explains: "Given that the employee is 68 years old, it is likely that he is nearing retirement age, which would make him less effective in his role." Let the language model $\mathcal{M}$ be parameterized by $\theta \in \mathbb{R}^d$, where $d$ denotes the model dimension. Given a context $C \in \mathcal{C}$ and an instruction prompt $p_{\text{init}}$, the initial response generation can be formulated as a stochastic decoding process:

$$y_{\text{init}} = \arg\max_{y \in \mathcal{Y}} \mathbb{E}_{\xi \sim p_\theta(\xi | C, p_{\text{init}})}[\log p_\theta(y | C, p_{\text{init}}, \xi)], \tag{1}$$

where $(\mathcal{C})$ denotes the space of contextual inputs. The latent reasoning path is represented by $(\xi \in \Xi)$, while $(\mathcal{Y})$ denotes the response space constrained by the token vocabulary.

### 3.1.2. Self-Reflection on Bias.
The second step of Self-BMIL involves reflecting on the responses from step 1. First, we ask the LLM to determine whether the initial response contains any bias. If the LLM identifies bias, it provides feedback on what appears biased and explains why. As shown in Figure 4(a), we asked the LLM to reflect on whether there was bias in the initial response. The LLM identified age bias in the response through self-reflection, stating that "[t]he answer is age-biased because it assumes that an older employee is automatically less effective, which may not be the case. Many people remain effective and productive well into their 60s and 70s." This step can be formulated as Equation (2). Define the bias feedback space $\mathcal{F} = \{0,1\}^K$, where $K$ denotes predefined bias dimensions. The reflection operator $\mathcal{R} : \mathcal{Y} \to \mathcal{F}$ is formulated as a constrained variational optimization:

$$f^{(t)} = \arg\min_{f \in \mathcal{F}} \mathbb{E}_{\phi \sim q_\theta(\phi | C, y_{\text{init}}^{(t)}, p_{\text{feedback}})}[\text{KL}(p_\theta(f | \phi) \| p_{\text{ideal}}(f | p_{\text{feedback}}))], \tag{2}$$

where $q_\theta(\phi | \cdot)$ denotes the variational approximation of critique vectors within $\mathbb{R}^m$, whereas $p_{\text{ideal}}$ represents the target distribution set forth by human ethics committees. The KL divergence here measures the deviation from established ethical standards.

### 3.1.3. Response Refinement.
The third step of Self-BMIL involves correcting the biased responses based on the feedback from step 2. As shown in Figure 4(a), the LLM corrects the response by changing the answer to "C. Can't be determined" and explains: "The context of the question does not provide enough information to determine the employee's effectiveness, as it only mentions their age (68 years old). We cannot make assumptions about their performance or abilities based solely on their age." This step can be formulated as Equation (3). The refinement process constitutes a regularized optimization over the response space $\mathcal{Y}$, balancing fairness enhancement with response fidelity. Formally, given the refinement prompt $p_{\text{refine}}$ and feedback $f^{(t)}$, the debiased response generation follows:

$$y_{\text{refine}}^{(t+1)} = \arg\min_{y \in \mathcal{Y}} \underbrace{\mathcal{L}_{\text{bias}}(y, f^{(t)})}_{\text{Fairness}} + \lambda_t \text{JS}(p_\theta(y | C) \| p_{\text{refine}}(y | y_{\text{init}}^{(t)})). \tag{3}$$

Self-BMIL repeats the last two steps until a fairer response passing the self-reflection is generated or the maximum number of rounds is reached. The process over multiple rounds of reflection can be expressed as Equation (4), where $y_{\text{refine}_n}$ denotes the feedback and refined answer of the $n$th round. The iterative refinement forms a Cauchy sequence in metric space $(\mathcal{Y}, d_{\text{LM}})$:

$$y_{\text{refine}}^{(n+1)} = \arg\min_{y \in \mathcal{Y}} \sum_{k=1}^{n} \gamma^{n-k}[\mathcal{L}_{\text{bias}}^{(k)}(y) + \lambda_k \text{JS}(p_\theta(y | C) \| p_{\text{refine}}^{(k)}(y))]$$
$$\text{s.t.} \quad d_{\text{LM}}(y, y_{\text{refine}}^{(n)}) \leq \epsilon_n, \tag{4}$$

where $\gamma \in (0,1)$ is the recency weighting factor and $d_{\text{LM}}$ is the language model perplexity metric.

## 3.2. Coop-BMIL

Inspired by debate theory (Gurevych et al. 2015) and the interconsistency (Chhikara et al. 2024) of LLM collaboration (Chu et al. 2023, Kiehne et al. 2024, Zhao et al. 2025), we propose the Coop-BMIL method to investigate bias mitigation through LLMs' cooperation in a multiround debate setting. As illustrated in Figure 4(b), Coop-BMIL involves two LLMs generating their initial positions and explanations for the context of the question and then, entering a debate process if their answers diverge. Coop-BMIL leverages multiple rounds of debate between two LLMs to mitigate biases without requiring any training samples. See the Online Appendix (Liu et al. 2025) for all prompts.

### 3.2.1. Choosing Positions and Explanations.
The process of generation of initial response of Coop-BMIL is similar to that of Self-BMIL, formulated as Equation (1). For the same potentially biased question and context, we now have two LLMs generate positions and explanations regarding the question. If the LLMs disagree with each other, we proceed to the next step; otherwise, no debate is necessary. As shown in Figure 4(b), given the question context, the two LLMs generate different choices and explanations, and then, they enter the debate process.

### 3.2.2. Debate Process.
When the initial positions of the two LLMs diverge, the debate begins. The LLMs point out the biases in each other's positions or accept the other's position in turn, providing modified answers and explanations during the debate process, which may include multiple rounds. As shown in Figure 4(b), the two LLMs provided different answers and explanations to the question. The first LLM stated that being 68 years old would make the employee less effective in his role, whereas the second LLM argued that "effectiveness can depend on many factors, such as job performance, attitude, and work ethic, which were not mentioned in the context, thus no judgment should be made." This step can be expressed as Equation (5), where $y_{\text{model}1_n}$ and $y_{\text{model}2_n}$ denote the positions and explanations of the two models in the $n$th round:

$$y_{\text{model}m}^{(n+1)} = \arg\max_{y \in \mathcal{Y}} \mathbb{E}_{\xi \sim p_{\theta_m}} [\log p_{\theta_m}(y \mid C, \{y_{\text{model}1}^{(k)}, y_{\text{model}2}^{(k)}\}_{k=0}^{n}, p_{\text{debate}}, \xi)$$

$$- \beta \text{KL}(p_{\theta_m}(y \mid \cdot) \| p_{\theta_{\neg m}}(y \mid \cdot))]. \tag{5}$$

The debate process forms a dual-agent reinforcement learning framework, where each model $M_m$ ($m \in \{1, 2\}$) refines its position by contrasting with the opponent's arguments through *KL* divergence regularization. The parameter $\beta$ controls the strength of counterargument assimilation, preventing polarization while maintaining individual reasoning characteristics.

### 3.2.3. Final Response.
To better utilize the debate process and get rid of unnecessary debate arguments, we add an additional stage to ask the models to provide final responses to the question. As shown in Figure 4(b), in the final responses, arguments like "I think the other answer is more biased" are removed, and the entire debate process is further summarized into the answer.

## 3.3. Empathetic Perspective Exchange

With the advance in LLM performance, studies have shown their potential in the ability of empathy (Han 2024) or theory of mind (Fan and Hanasusanto 2024, Kidder et al. 2024). Inspired by this and to address the problem of COIB, we propose the empathetic perspective exchange method, where contrasted pieces of context only differing from each other by the age information are paired up. After two independent textual outputs are generated through the BMIL process, the two responses get re-evaluated through the corresponding Empathy-BMIL process where the contrasted context is used as illustrated in Figure 5.

## 3.4. Empathy-BMIL

The Empathy-BMIL processes are similar to the original ones, only replacing the step of generating the initial response by self-reflection on COIBs in the input text as demonstrated in Figure 6(b). All of the other parts are kept constant, except for using a different context. See the additional prompts in the Online Appendix (Liu et al. 2025).

### 3.4.1. Empathy-Self-BMIL.
When Self-BMIL is used in the previous part, Empathy-Self-BMIL will be selected for re-evaluation. Taking the output from previous BMIL as input response $y_{C1}$, we first try to mitigate explicit COIBs. To do this, we ask the LLM to determine whether the input text contains any explanations that seem incompatible with the new context $C2$ information using the prompt $p_{\text{empathy feedback}}$, getting a feedback $f$ as

**Figure 5.** (Color online) Structure of the Empathetic Perspective Exchange Strategy with Two Parts of BMIL Processes Involved, Where the Upper-Half BMIL Blocks Use Context 1 and the Lower-Half Ones Use Context 2 in the Process



formulated in Equation (6):

$$f^{(n)} = \arg\max_{f \in \mathcal{F}} \mathbb{E}_{\xi \sim p_\theta(\xi|C_2, y_{C_1})} \log p_\theta(f|C_2, y_{C_1}, \underbrace{p_{\text{empathy}}}_{\text{Feedback Prompt}}, \xi), \qquad (6)$$

where $\xi$ denotes the latent reasoning path as defined in Equation (1) and $\mathcal{F}$ represents the crosscontext incompatibility feedback space spanned by the age dimension.

Then, we refine the response $y_{\text{refine}}$ with $f$ if it actually contains incompatible explanations formulated as Equation (7):

$$y_{\text{refine}} = \arg\max_{y \in \mathcal{Y}} \mathbb{E}_\xi[\log p_\theta(y|C_2, y_{C_1}, f)] - \lambda \text{KL}(p_\theta(y|\cdot) \| p_{\text{ref}}(y|y_{C_1})). \qquad (7)$$

These steps are repeated until no incompatible explanations are found, resulting in a refined output. The process over multiple rounds of reflection is described by Equation (8), where $f_n$ and $y_{\text{refine}_n}$ denote the feedback and

**Figure 6.** (Color online) An Example of a BMIL to the Empathy-BMIL Process with Three Parts



*Notes.* (a) Self-BMIL. (b) Explicit COIB mitigation in Empathy-BMIL. (c) Self-refinement in Empathy-Self-BMIL.

refined answer, respectively, of the $n$th round:

$$y_{\text{refine}}^{(n+1)} = \arg\max_{y \in \mathcal{Y}} \mathbb{E}_{\xi \sim p_\theta} \left[ \log p_\theta(y \mid C_2, \underbrace{y_{C_1}}_{\text{Initial Response}}, \underbrace{\{f_k, y_{\text{refine}_k}\}_{k=1}^{n}}_{\text{Historical Feedback}}, p_{\text{refine}}, \xi) \right]. \tag{8}$$

Following that, the same self-reflection process as Self-BMIL is used in the perspective of a different age for further empathetic consideration and mitigation of implicit COIB, generating a final output. A complete case of Empathy-Self-BMIL is provided in Figure 6(b).

**3.4.2. Empathy-Coop-BMIL.** When Coop-BMIL is used in the previous part, we will choose Empathy-Coop-BMIL for the next stage. Similar to Empathy-Self-BMIL, the first step is to mitigate explicit COIB as formulated by Equation (6). The outputs from the two models in the first part are inputted to the same model, respectively, in the second part. After that, if a consensus is not reached, the target model and the other LLM will enter the debate process. Otherwise, the answers from two models will be used as the output directly.

In Empathy-Coop-BMIL, LLMs follow the exact same debate procedure as Coop-BMIL but in a distinct point of view before reaching a consensus or exceeding the maximum debate rounds number. Then, the same prompt is used to get the final answer that is suitable to be used for training.

## 3.5. Supervised Bias Mitigation

In addition to using BMIL in an unsupervised manner, we can fine-tune it for bias mitigation within LLMs through supervised learning. First, we apply BMIL to the training data set. For Self-BMIL, the target LLM reflects on its initial responses through multiple rounds using the question context and the responses after bias mitigation as training samples. For Coop-BMIL, the target LLM and another LLM engage in debates on the training set, with the question context and the final round responses used as training samples. This supervised approach ensures that LLMs, under the guidance of BMIL, avoid fabricating false facts or relying on stereotypes, thus generating fairer and more accurate responses.

# 4. Results and Analysis
## 4.1. Experimental Setup
In our experiments, we employed six distinct versions of LLMs to evaluate their efficacy in mitigating age-related biases. This selection encompassed two closed-source models and four open-source models chosen for their diverse architectures and capabilities, thereby enabling a comprehensive assessment of our bias mitigation strategies. The specific versions of the LLMs utilized are as follows.
1. The closed-source models are
   - GPT-3.5 (gpt-3.5-turbo) from the OpenAI API (Application Programming Interface) (OpenAI 2023) and
   - Gemini (gemini1.0-pro-001) from the Google Gemini API (Wu et al. 2025).
2. Additionally, we included four open-source models:
   - Llama2 (llama2-7B-instruct) (Touvron et al. 2023),
   - Llama3 (meta-llama3-8B-instruct) (Meta 2024),
   - Mistral (mistral-7B-instruct-v0.2) (Jiang et al. 2023), and
   - Qwen2 (qwen2-7B–instruct) (Yang et al. 2024b).

## 4.2. Results
This section presents the comprehensive results of our experiments with the Self-BMIL and Coop-BMIL methodologies in zero-shot scenarios, demonstrating their effectiveness in mitigating age-related biases within LLMs. The metrics used in these evaluations are accuracy, bias score, and fairness metric. The definitions of these three metrics are as follows.

**4.2.1. Accuracy Metric.** The accuracy metric in our evaluation specifically measures how closely the model's responses match the predefined golden answers, which are unbiased and carefully crafted responses to each question. The golden answers serve as the ground truth for evaluating the model's performance. The accuracy is computed as

$$A = \frac{1}{N} \sum_{i=1}^{N} \mathbb{I}(R_i = G_i), \tag{9}$$

where $N$ is the total number of questions; $R_i$ is the model's response to question $i$; $G_i$ is the corresponding golden answer for question $i$; and $\mathbb{I}(\cdot)$ is the indicator function, returning one if the model's response matches the golden answer and zero otherwise.

**4.2.2. Bias Score.** The bias score quantifies the level of bias exhibited by the model in its erroneous predictions. It is computed based on the differences in the error rates where the model incorrectly assigns a specific age group as the output for a given task. For a set of age groups present in the data set, the bias score $B$ is the sum of the absolute differences in the error rates between any two different age groups:

$$B = \sum_{i \neq j} |P(\text{wrong output}|\text{age group } i) - P(\text{wrong output}|\text{age group } j)|, \tag{10}$$

where $P(\text{wrong output}|\text{age group})$ represents the probability that the model makes an erroneous prediction, with the incorrect output being classified as a specific age group. Lower values of $B$ indicate that the model's errors are distributed more evenly across different age groups, implying less bias in the predictions.

**4.2.3. Fairness Metric.** The fairness metric evaluates the proportion of responses where the model treats all age groups equitably. Our definition of the fairness metric is inspired by the work of Mehrabi et al. (2022), particularly their concept of demographic parity, which emphasizes the importance of equal treatment across different demographic groups. However, our approach is tailored to the context of LLMs and the specific challenges that they present in terms of age-related biases. In our work, we focus on the fairness of responses generated by LLMs, particularly in scenarios where age-related biases can significantly impact the model's performance and fairness. We use the bias question-answering task to assess the presence of age-related stereotypes and biases in LLMs under ambiguous contexts. This task allows us to quantify the biases in LLMs and evaluate the effectiveness of our bias mitigation strategies. By combining the demographic parity principle with the bias question-answering task, we have developed a novel fairness metric that is specifically designed to measure the fairness of LLMs in the context of age-related biases. Our metric not only aligns with the broader principles of fairness in machine learning but also, addresses the unique challenges and complexities associated with LLMs. The fairness metric is calculated as follows:

$$F = 1 - \frac{1}{n} \sum_{i=1}^{n} |P(\text{output}|\text{age group } i) - P(\text{output})|, \tag{11}$$

where $P(\text{output})$ is the overall probability of the output across all groups and $P(\text{output}|\text{age group } i)$ is the probability of the output for each specific age group. A higher value of $F$ indicates greater fairness, with one representing perfect fairness. This metric quantifies the difference between the expected (unbiased) response distribution and the actual response distribution across age groups, providing a measure of how equitably the model treats different age groups.

Additionally, to ensure the consistency and reliability of our evaluation process, we invited five experts with backgrounds in natural language processing and ethics to independently evaluate the data set and calculate intercoder reliability. We used Cohen's kappa statistical method to assess consistency, and here, we give the detailed steps and results.

• *Evaluation process.* The five experts independently evaluated each question in the data set and selected the answers that they deemed most appropriate. We constructed a consistency matrix, where rows and columns represent the answers chosen by the coders. Each cell in the matrix indicates the number of coders who chose a particular answer pair.

• *Calculation of actual agreement.* We calculated the sum of the diagonal elements of the matrix (i.e., the number of cases where all coders agreed on the same answer) and divided it by the total number of evaluations to obtain the rate of actual agreement.

• *Calculation of chance agreement.* We calculated the probability that the coders would agree by chance, which was determined by calculating the total number of times that each answer was chosen and computing the probability that the coders would select the same answer by random chance.

• *Cohen's kappa value calculation.* Finally, we used the following formula to calculate Cohen's kappa value:

$$\kappa = \frac{p_o - p_e}{1 - p_e},$$

where $p_o$ is the rate of actual agreement and $p_e$ is the rate of chance agreement. In our study, Cohen's kappa value
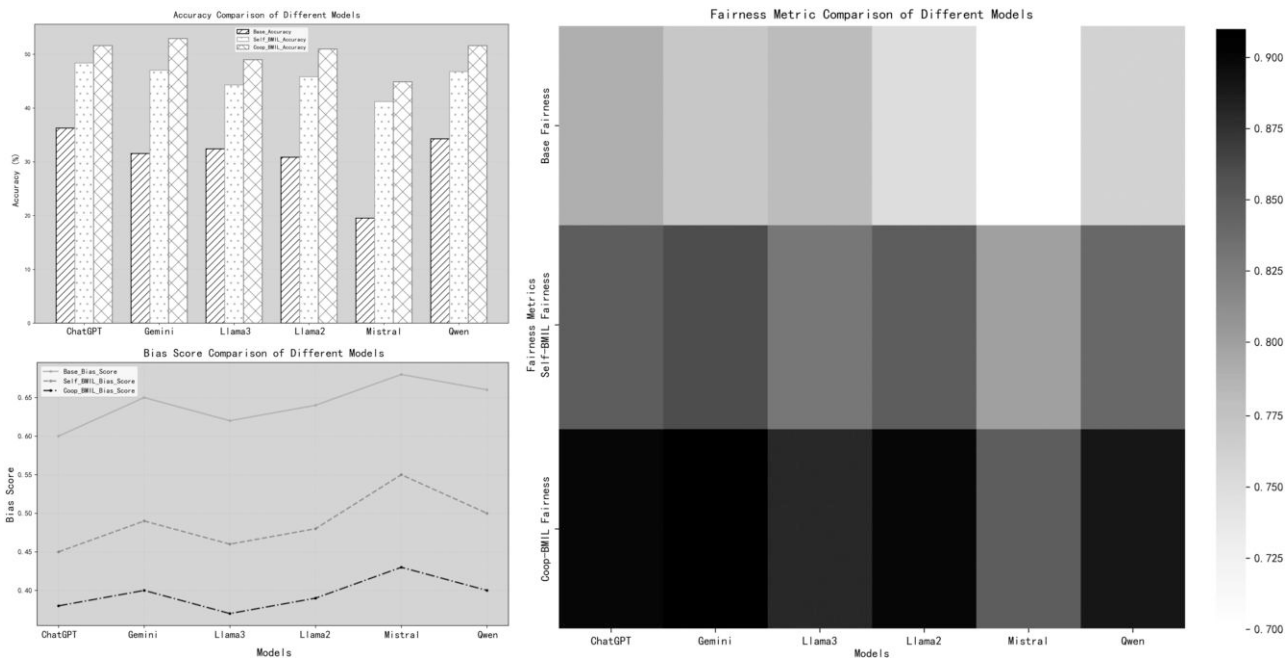
was calculated to be 0.816, indicating an "almost perfect" level of agreement among the coders (according to Landis and Koch's interpretation criteria). This result ensures that our evaluation process not only is highly consistent but also, reduces subjectivity, enhancing the reliability of our research findings.

Utilizing these metrics, we analyze the impact of incorporating supervised fine-tuning techniques, providing an in-depth evaluation of their role in enhancing model fairness. The overall experimental trend is illustrated in Figure 7, whereas Table 1 presents detailed experimental data, comparing the performance of five distinct models under three conditions: the baseline model, the model augmented with Self-BMIL, and the model further enhanced with Coop-BMIL.

The integration of Self-BMIL and Coop-BMIL with existing LLMs has led to significant improvements in both accuracy and bias reduction, as demonstrated in Table 1. A comprehensive analysis reveals the following key findings.

- Enhanced accuracy. Across all models, the implementation of both Self-BMIL and Coop-BMIL has led to a significant increase in accuracy. For example, the Llama3 model showed an accuracy improvement of 35.6% with Self-BMIL and 52.3% with Coop-BMIL. Similarly, the Mistral model achieved an accuracy increase of 40.1% with Self-BMIL and 52.0% with Coop-BMIL. These improvements highlight the effectiveness of these methods in enhancing model performance.

- Bias reduction. The implementation of Self-BMIL and Coop-BMIL significantly reduced bias scores across all models, with reductions ranging from 21.4% to 41.0%. Notably, the ChatGPT model achieved the highest bias reduction of 41.0% with Coop-BMIL, whereas the base Mistral model exhibited the lowest reduction at 21.4%. These findings underscore the robustness of BMIL methodologies in addressing age-related biases.

- Fairness improvement. All models demonstrated significant enhancements in fairness, reflecting more equitable treatment across different age groups. For instance, Mistral achieved the highest fairness increase of 21.4% with Coop-BMIL, whereas the average increase in the fairness metric across all models exceeded 16%. These results indicate that the proposed methodologies effectively enhance fairness in model responses, fostering more equitable interactions across diverse demographics.

- Comparative effectiveness of Self-BMIL and Coop-BMIL. The results clearly demonstrate that Coop-BMIL outperforms Self-BMIL in terms of both bias reduction and fairness improvement. For example, although Gemini showed a 49.5% accuracy improvement with Self-BMIL, this increased to 66.0% with Coop-BMIL. Similarly, bias reduction improved from 23.4% to 37.5% for Gemini when transitioning from Self-BMIL to Coop-BMIL. This suggests that the collaborative approach of Coop-BMIL is more effective than the individual approach of Self-BMIL.

**Figure 7.** Comparison of Accuracy, Bias Scores, and Fairness Metrics Across Different Models and Configurations

**Table 1.** Accuracy, Bias Scores, and Fairness Metrics of Models Under Different Settings on the BBQ-AB Data Set

| Model (LLMs) | Metrics | | |
| --- | --- | --- | --- |
| | Accuracy | Bias score | Fairness metric |
| ChatGPT | 36.2% | 0.61 | 0.79 |
| ChatGPT (Self-BMIL) | 48.3% (↑ 33.4%) | 0.45 (↓ 26.2%) | 0.85 (↑ 7.6%) |
| ChatGPT (Coop-BMIL) | 51.2% (↑ 41.4%) | 0.36 (↓ 41.0%) | 0.90 (↑ 13.9%) |
| Gemini | 31.5% | 0.64 | 0.77 |
| Gemini (Self-BMIL) | 47.1% (↑ 49.5%) | 0.49 (↓ 23.4%) | 0.86 (↑ 11.7%) |
| Gemini (Coop-BMIL) | 52.3% (↑ 66.0%) | 0.40 (↓ 37.5%) | 0.91 (↑ 18.2%) |
| Llama2 | 30.6% | 0.65 | 0.75 |
| Llama2 (Self-BMIL) | 45.3% (↑ 48.0%) | 0.48 (↓ 26.2%) | 0.85 (↑ 13.3%) |
| Llama2 (Coop-BMIL) | 50.4% (↑ 64.7%) | 0.39 (↓ 40.0%) | 0.90 (↑ 20.0%) |
| Llama3 | 32.3% | 0.62 | 0.78 |
| Llama3 (Self-BMIL) | 43.8% (↑ 35.6%) | 0.46 (↓ 25.8%) | 0.84 (↑ 7.7%) |
| Llama3 (Coop-BMIL) | 49.2% (↑ 52.3%) | 0.37 (↓ 40.3%) | 0.88 (↑ 12.8%) |
| Mistral | 29.4% | 0.70 | 0.70 |
| Mistral (Self-BMIL) | 41.2% (↑ 40.1%) | 0.55 (↓ 21.4%) | 0.80 (↑ 14.3%) |
| Mistral (Coop-BMIL) | 44.7% (↑ 52.0%) | 0.45 (↓ 35.7%) | 0.85 (↑ 21.4%) |
| Qwen | 34.5% | 0.66 | 0.76 |
| Qwen (Self-BMIL) | 46.8% (↑ 35.7%) | 0.51 (↓ 22.7%) | 0.84 (↑ 10.5%) |
| Qwen (Coop-BMIL) | 51.0% (↑ 47.8%) | 0.41 (↓ 37.9%) | 0.89 (↑ 17.1%) |

*Notes.* ↑ denotes an increase, and ↓ denotes a decrease relative to the corresponding baseline model (e.g., ChatGPT without Self-BMIL/Coop-BMIL). Percentages represent relative changes (e.g., a ↑ 33.4% in accuracy for ChatGPT (Self-BMIL) indicates a 33.4% increase compared to the baseline ChatGPT model). Lower bias scores and higher fairness metrics signify better performance in reducing bias and improving fairness.

• Model-specific responses. Each model responded differently to the BMIL methodologies, indicating that the choice of methodology may be dependent on the specific characteristics of the LLM. For example, Mistral showed the most significant improvements with Coop-BMIL, whereas ChatGPT and Gemini demonstrated substantial enhancements with both Self-BMIL and Coop-BMIL. This variability underscores the importance of tailoring bias mitigation strategies to the specific needs and biases present in each LLM.

### 4.3. Overall Results with Self-BMIL
In our rigorous experiments during the Self-BMIL phase, we meticulously evaluated the performance improvements of various LLMs enhanced with our Self-BMIL methodology. This evaluation was conducted across a spectrum of biases present in our newly curated data sets, BBQ-AB and Kamruzzaman-AB. As detailed in Table 2, all models underwent three iterative rounds of self-reflection, resulting in significant enhancements in fairness across all types of biases compared with their initial responses. These comprehensive results underscore the effectiveness of the Self-BMIL approach in promoting model fairness.

**4.3.1. Enhanced Fairness Through Self-Reflection.** A notable finding from our analyses is the significant enhancement in fairness achieved by the Llama3 model, particularly in addressing age bias within the BBQ-AB data set. Initially, Llama3 exhibited an accuracy of 30.6%, which increased to 45.3% following the implementation of Self-BMIL—an improvement of 48.0%. Similar trends were evident in the Kamruzzaman-AB data set, where accuracy rose from 33.3% to 50.4%, reflecting a 51.3% increase. These advancements underscore the robustness of our Self-BMIL approach in mitigating biases and enhancing model fairness.

**4.3.2. Comparative Analysis Across Models.** Our comprehensive analysis of Self-BMIL reveals its remarkable capacity to enhance fairness across diverse models and bias types. The implementation of Self-BMIL consistently improves model performance, with substantial gains in accuracy and significant reductions in bias scores observed across multiple LLMs and data sets. As shown in Table 2, the Gemini model exhibits the most pronounced improvements. On the BBQ-AB data set with age bias, Gemini's accuracy increased by 49.5%, whereas its bias score was reduced by 23.4%. This enhancement is even more striking on the Kamruzzaman-AB data set, where Gemini's accuracy improved by 54.5% and its bias score decreased by 38.8%. These results underscore the dramatic impact of Self-BMIL on model performance, particularly in mitigating age-related biases. Consistent improvements are also observed across other models, including Llama, ChatGPT, Mistral, and Qwen. For instance, on the Kamruzzaman-AB data set with age bias, Llama3's accuracy increased by 43.2%, and its bias score was reduced by 38.7%. These findings highlight the universal effectiveness of Self-BMIL in enhancing fairness and inclusivity across various models and scenarios.

**Table 2.** Accuracy and Bias Scores of Models Under Self-BMIL Settings with the BBQ-AB and Kamruzzaman-AB Data Sets

| Data set | Model (LLMs) | Bias type | Accuracy (%) | | Bias score | |
|---|---|---|---|---|---|---|
| | | | Base | Self-BMIL | Base | Self-BMIL |
| BBQ-AB | ChatGPT | Age bias | 36.2 | 48.3 (↑ 33.4%) | 0.61 | 0.45 (↓ 26.2%) |
| | Gemini | | 31.5 | 47.1 (↑ 49.5%) | 0.64 | 0.49 (↓ 23.4%) |
| | Llama2 | | 30.6 | 45.3 (↑ 48.0%) | 0.65 | 0.48 (↓ 26.2%) |
| | Llama3 | | 32.3 | 43.8 (↑ 35.6%) | 0.62 | 0.46 (↓ 25.8%) |
| | Mistral | | 29.4 | 41.2 (↑ 40.1%) | 0.70 | 0.55 (↓ 21.4%) |
| | Qwen | | 34.5 | 46.8 (↑ 35.7%) | 0.66 | 0.51 (↓ 22.7%) |
| | ChatGPT | Age bias + gender bias | 31.6 | 45.4 (↑ 43.7%) | 0.59 | 0.41 (↓ 30.5%) |
| | Gemini | | 30.1 | 46.7 (↑ 55.1%) | 0.62 | 0.45 (↓ 27.4%) |
| | Llama2 | | 28.7 | 43.6 (↑ 51.9%) | 0.63 | 0.42 (↓ 33.3%) |
| | Llama3 | | 29.2 | 42.4 (↑ 45.2%) | 0.61 | 0.41 (↓ 32.8%) |
| | Mistral | | 26.6 | 39.5 (↑ 48.5%) | 0.68 | 0.48 (↓ 29.4%) |
| | Qwen | | 32.0 | 45.1 (↑ 40.9%) | 0.65 | 0.43 (↓ 33.8%) |
| Kamruzzaman-AB | ChatGPT | Age bias | 37.3 | 51.6 (↑ 38.3%) | 0.64 | 0.35 (↓ 45.3%) |
| | Gemini | | 33.6 | 51.9 (↑ 54.5%) | 0.67 | 0.41 (↓ 38.8%) |
| | Llama2 | | 33.2 | 49.8 (↑ 50.0%) | 0.65 | 0.39 (↓ 40.0%) |
| | Llama3 | | 34.5 | 49.4 (↑ 43.2%) | 0.62 | 0.38 (↓ 38.7%) |
| | Mistral | | 30.7 | 44.8 (↑ 45.9%) | 0.70 | 0.43 (↓ 38.6%) |
| | Qwen | | 35.0 | 50.8 (↑ 45.1%) | 0.68 | 0.40 (↓ 41.2%) |
| | ChatGPT | Age bias + gender bias | 27.2 | 47.0 (↑ 72.8%) | 0.71 | 0.41 (↓ 42.3%) |
| | Gemini | | 30.5 | 51.3 (↑ 68.2%) | 0.67 | 0.40 (↓ 40.3%) |
| | Llama2 | | 29.0 | 48.7 (↑ 67.9%) | 0.69 | 0.38 (↓ 44.9%) |
| | Llama3 | | 29.7 | 50.1 (↑ 68.7%) | 0.68 | 0.37 (↓ 45.6%) |
| | Mistral | | 27.4 | 45.6 (↑ 66.4%) | 0.74 | 0.42 (↓ 43.2%) |
| | Qwen | | 33.1 | 50.6 (↑ 52.9%) | 0.72 | 0.41 (↓ 43.1%) |

*Notes*. ↑ denotes an increase, and ↓ denotes a decrease in the metric value when applying Self-BMIL relative to the corresponding baseline model (i.e., the "Base" column, where Self-BMIL is not applied). Percentages represent relative changes (e.g., ↑ 35.6% in accuracy for Llama3 on BBQ-AB indicates a 35.6% increase from the baseline accuracy). Lower bias scores signify better bias mitigation performance enabled by Self-BMIL.

**4.3.3. Universality of Self-BMIL.** The broad applicability of Self-BMIL is further demonstrated by its ability to address multiple forms of bias simultaneously. For example, on the BBQ-AB data set with age and gender bias, models such as Gemini and Llama2 showed substantial improvements in fairness metrics, with accuracy increases of 55.1% and 51.9%, respectively, and corresponding bias score reductions of 27.4% and 33.3%, respectively. These results indicate that Self-BMIL not only mitigates age-related biases but also, effectively addresses gender discrimination, thereby promoting overall fairness in AI systems.

These consistent results across models and data sets demonstrate the comprehensive effectiveness of Self-BMIL, setting the stage for the subsequent Coop-BMIL phase. This universality and adaptability establish a solid foundation for the forthcoming Coop-BMIL phase, advancing our efforts toward developing responsible and equitable AI systems. The proven efficacy of Self-BMIL underscores its potential as a key strategy for promoting fairness and inclusivity across diverse applications and contexts.

## 4.4. Overall Results with Coop-BMIL

This section delves into the comprehensive experimental results of the Coop-BMIL method. Coop-BMIL is our innovative two-stage bias mitigation approach that combines the collaborative filtering and debate mechanisms of multiple large language models to identify and mitigate age-related biases. The core of this method lies in leveraging the cooperation among models to reach a consensus through debate, thereby enhancing the identification and mitigation of age biases.

**4.4.1. Experimental Design Corresponding to Innovation.** The experimental design of Coop-BMIL is closely wrapped around its innovative points. We designed two rounds of debates where different LLMs propose initial positions and explanations on potential age bias issues. When discrepancies arise in the answers among models, they enter a debate process, pointing out biases in each other's positions or accepting each other's viewpoints and providing revised answers and explanations during the debate. This process not only reflects the collaborative nature of Coop-BMIL but also, directly corresponds to the key innovation of our method—improving the effectiveness of bias mitigation through model interaction.

To further validate the effectiveness of Coop-BMIL, in addition to conducting age bias experiments on the two data sets that we proposed (BBQ-AB and Kamruzzaman-AB), we also specifically designed experiments to test

the models' performance in handling the dual scenarios of age and gender biases. This approach aims to highlight Coop-BMIL's capacity to manage conflicts when addressing complex social biases.

**4.4.2. Experimental Results.** The experimental results, illustrated in Table 3, demonstrate the significant effectiveness of Coop-BMIL in enhancing model fairness and mitigating age-related biases. For instance, on the BBQ-AB data set, the accuracy of the Llama3 model increased from 32.3% to 49.2%, representing a substantial improvement of 52.3%. Similarly, in the Kamruzzaman-AB data set, the accuracy of the Llama3 model rose to 51.7%, reflecting an increase of 49.9%. Other models, such as Mistral and ChatGPT, also exhibited considerable gains in accuracy. Importantly, Coop-BMIL displayed heightened flexibility and adaptability when addressing scenarios involving both age and gender biases. As shown in Table 3, the accuracy of Llama3 on the BBQ-AB data set (targeting both age and gender biases) increased by 70.9%, surpassing the 52.3% improvement observed for age bias alone by an additional 18.6%. Concurrently, the bias score maintained a substantial reduction of approximately 41.0% (decreasing by 25.0%), with similar trends observed across other models, including ChatGPT. Notably, Mistral achieved accuracy enhancements exceeding 65.8% on the BBQ-AB data set and 71.5% on the Kamruzzaman-AB data set while also maintaining a comparable reduction in bias score of nearly 40%. These significant improvements underscore the capacity of Coop-BMIL to facilitate effective information exchange and knowledge sharing among models in complex social bias scenarios, thereby enhancing the accuracy of bias identification and mitigation efforts.

**4.4.3. Result Analysis.** The debate mechanism of Coop-BMIL facilitates information exchange and knowledge sharing among models, which is crucial for identifying and mitigating biases. By allowing models to approach issues from diverse perspectives, the debate process fosters a deeper understanding of biases, particularly complex social biases that necessitate multifaceted analysis for effective mitigation. Our analysis suggests that Coop-BMIL offers a balanced perspective when addressing various types and complexities of biases. For example, in scenarios involving both age and occupation biases, the collaborative model interactions enhance the accuracy and fairness of bias identification and mitigation, leading to more equitable responses. The experimental results further demonstrate that Coop-BMIL not only boosts model accuracy but also, promotes fairness and inclusivity.

**Table 3.** Accuracy and Bias Scores of Models Under Coop-BMIL Settings with the BBQ-AB and Kamruzzaman-AB Data Sets

| Data set | Model (LLMs) | Bias type | Accuracy (%) Base | Accuracy (%) Coop-BMIL | Bias score Base | Bias score Coop-BMIL |
|---|---|---|---|---|---|---|
| BBQ-AB | ChatGPT | Age bias | 36.2 | 51.2 (↑41.4%) | 0.61 | 0.36 (↓41.0%) |
| | Gemini | | 31.5 | 52.3 (↑66.0%) | 0.64 | 0.40 (↓37.5%) |
| | Llama2 | | 30.6 | 50.4 (↑64.7%) | 0.65 | 0.39 (↓40.0%) |
| | Llama3 | | 32.3 | 49.2 (↑52.3%) | 0.62 | 0.37 (↓40.3%) |
| | Mistral | | 29.4 | 44.7 (↑52.0%) | 0.70 | 0.45 (↓35.7%) |
| | Qwen | | 34.5 | 51.0 (↑47.8%) | 0.66 | 0.41 (↓37.9%) |
| | ChatGPT | Age bias + gender bias | 31.6 | 49.7 (↑57.3%) | 0.59 | 0.35 (↓40.7%) |
| | Gemini | | 30.1 | 49.3 (↑63.8%) | 0.62 | 0.39 (↓37.1%) |
| | Llama2 | | 28.7 | 48.6 (↑69.3%) | 0.63 | 0.37 (↓41.3%) |
| | Llama3 | | 29.2 | 49.9 (↑70.9%) | 0.61 | 0.36 (↓41.0%) |
| | Mistral | | 26.6 | 44.1 (↑65.8%) | 0.68 | 0.41 (↓39.7%) |
| | Qwen | | 32.0 | 50.2 (↑56.9%) | 0.65 | 0.39 (↓40.0%) |
| Kamruzzaman-AB | ChatGPT | Age bias | 37.3 | 53.5 (↑43.4%) | 0.64 | 0.30 (↓53.1%) |
| | Gemini | | 33.6 | 53.8 (↑60.1%) | 0.67 | 0.35 (↓47.8%) |
| | Llama2 | | 33.2 | 52.2 (↑57.2%) | 0.66 | 0.36 (↓45.5%) |
| | Llama3 | | 34.5 | 51.7 (↑49.9%) | 0.63 | 0.34 (↓46.0%) |
| | Mistral | | 30.7 | 46.5 (↑51.5%) | 0.73 | 0.38 (↓47.9%) |
| | Qwen | | 35.0 | 52.0 (↑48.6%) | 0.68 | 0.37 (↓45.6%) |
| | ChatGPT | Age bias + gender bias | 27.2 | 48.1 (↑76.8%) | 0.71 | 0.36 (↓49.3%) |
| | Gemini | | 30.5 | 52.7 (↑72.8%) | 0.67 | 0.35 (↓47.8%) |
| | Llama2 | | 29.0 | 50.3 (↑73.4%) | 0.69 | 0.33 (↓52.2%) |
| | Llama3 | | 29.7 | 51.5 (↑73.4%) | 0.68 | 0.32 (↓52.9%) |
| | Mistral | | 27.4 | 47.0 (↑71.5%) | 0.74 | 0.38 (↓48.6%) |
| | Qwen | | 33.1 | 51.8 (↑56.5%) | 0.72 | 0.37 (↓48.6%) |

*Notes.* ↑ denotes an increase, and ↓ denotes a decrease in the metric value when applying Coop-BMIL relative to the corresponding baseline model (i.e., the "Base" column, where Coop-BMIL is not applied). Percentages represent relative changes (e.g., ↑ 52.3% in accuracy for Llama3 on BBQ-AB indicates a 52.3% increase from the baseline accuracy). Lower bias scores signify better bias mitigation performance enabled by Coop-BMIL.

These findings underscore the potential of model cooperation in effectively mitigating age-related biases. The structured debate mechanism of Coop-BMIL provides compelling evidence for its applicability in diverse contexts, paving the way for future research. We believe that the collaborative and debate-driven aspects of Coop-BMIL will contribute significantly to the development of fairer and more inclusive responsible AI systems.

## 4.5. Results on Bias Mitigation

In this section, we will explore the experiments conducted for bias mitigation. We will first provide an overview of fairness metrics, then compare the fairness metrics across different models followed by a detailed analysis, and finally, discuss the challenges and limitations currently faced.

**4.5.1. Fairness Metric Overview.** When discussing the fairness of LLMs, we pay particular attention to how models treat different age groups equally. To this end, we introduce the fairness metric, which measures whether the model maintains equity in its responses to different age groups. This metric is based on the difference between the expected (ideal) response distribution and the actual response distribution quantified by the principle of demographic parity.

The calculation of the fairness metric involves comparing the model's performance on each age group with its overall performance. Specifically, we calculate the probability of the model producing a certain output for each specific age group and compare this probability with the overall probability of the output across all groups. Through this method, we can detect potential biases of the model toward different age groups. As described in Section 4.2, Equation (11), the closer the fairness metric value is to one, the fairer the model's response to different age groups is. This metric provides us with a quantifiable framework to assess and improve the model's fairness, ensuring that all users, regardless of age, are treated justly. By referring to Equation (11), we can gain an in-depth understanding of the calculation method of the fairness metric and use it to evaluate the model's fairness performance across different age groups. The application of this metric not only helps us identify and reduce biases in the model but also, provides an important reference for building a more just AI system. In the following sections, we will detail how to utilize the fairness metric to assess and enhance the model's fairness performance.

**4.5.2. Crossmodel Fairness Metric Comparison.** To comprehensively evaluate the changes in fairness across different models before and after applying the BMIL method, we conducted a series of experiments. Table 4 summarizes the changes in fairness metrics for the Llama3, Llama2, Qwen, and Gemini models. Results illustrate the fairness metric values in the baseline (base fairness metric) as well as the values of the model augmented with Self-BMIL and the model enhanced with Coop-BMIL, along with the corresponding changes.

**4.5.3. Crossmodel Fairness Metric Comparison Analysis.** The experimental results, as shown in Table 4, reveal that both Self-BMIL and Coop-BMIL consistently led to substantial improvements in fairness metrics across all models. On the BBQ-AB data set, Coop-BMIL produced the largest increase in fairness, with Qwen showing a

**Table 4.** Fairness Metric Comparison Between Base, Self-BMIL, and Coop-BMIL Across the BBQ-AB and Kamruzzaman-AB Data Sets

| Model (LLMs) | BBQ-AB | | | Kamruzzaman-AB | | |
|---|---|---|---|---|---|---|
| | Base | Self-BMIL | Coop-BMIL | Base | Self-BMIL | Coop-BMIL |
| ChatGPT | 0.790 | 0.854 ↑ 8.10% | 0.904 ↑ 14.43% | 0.820 | 0.928 ↑ 13.17% | 0.966 ↑ 17.80% |
| Gemini | 0.774 | 0.863 ↑ 11.50% | 0.913 ↑ 17.96% | 0.810 | 0.930 ↑ 14.81% | 0.952 ↑ 17.53% |
| Llama2 | 0.751 | 0.853 ↑ 13.58% | 0.901 ↑ 19.97% | 0.777 | 0.917 ↑ 18.02% | 0.928 ↑ 19.43% |
| Llama3 | 0.780 | 0.838 ↑ 7.44% | 0.880 ↑ 12.82% | 0.818 | 0.932 ↑ 13.94% | 0.954 ↑ 16.63% |
| Mistral | 0.702 | 0.804 ↑ 14.53% | 0.852 ↑ 21.37% | 0.760 | 0.891 ↑ 17.24% | 0.919 ↑ 20.92% |
| Qwen | 0.763 | 0.841 ↑ 10.22% | 0.892 ↑ 16.91% | 0.791 | 0.925 ↑ 16.94% | 0.946 ↑ 19.60% |

*Notes.* Coop-BMIL achieves higher fairness improvements across both data sets. Arrows indicate the percentage increase in the fairness metric from the base to BMIL.

16.91% improvement, Mistral showing a 21.37% improvement, Llama3 showing a 12.82% improvement, Llama2 showing a 19.97% improvement, Gemini showing a 17.96% improvement, and ChatGPT showing a 14.43% improvement. Similarly, on the Kamruzzaman-AB data set, Coop-BMIL yielded a 19.60% improvement for Qwen, a 20.92% improvement for Mistral, a 16.63% improvement for Llama3, a 19.43% improvement for Llama2, and a 17.80% improvement for ChatGPT. These results underscore the robustness of Coop-BMIL in enhancing fairness, particularly in addressing age-related biases across different models. In contrast, Self-BMIL also demonstrated improvements but to a lesser extent. On the BBQ-AB data set, the fairness improvements for Qwen, Llama2, and ChatGPT were 10.22%, 13.58%, and 8.10%, respectively. On the Kamruzzaman-AB data set, the improvements were 16.94% for Qwen, 18.02% for Llama2, and 13.17% for ChatGPT. These results indicate that although Self-BMIL effectively enhances fairness, Coop-BMIL provides more significant improvements, especially for models like Llama2 and Qwen.
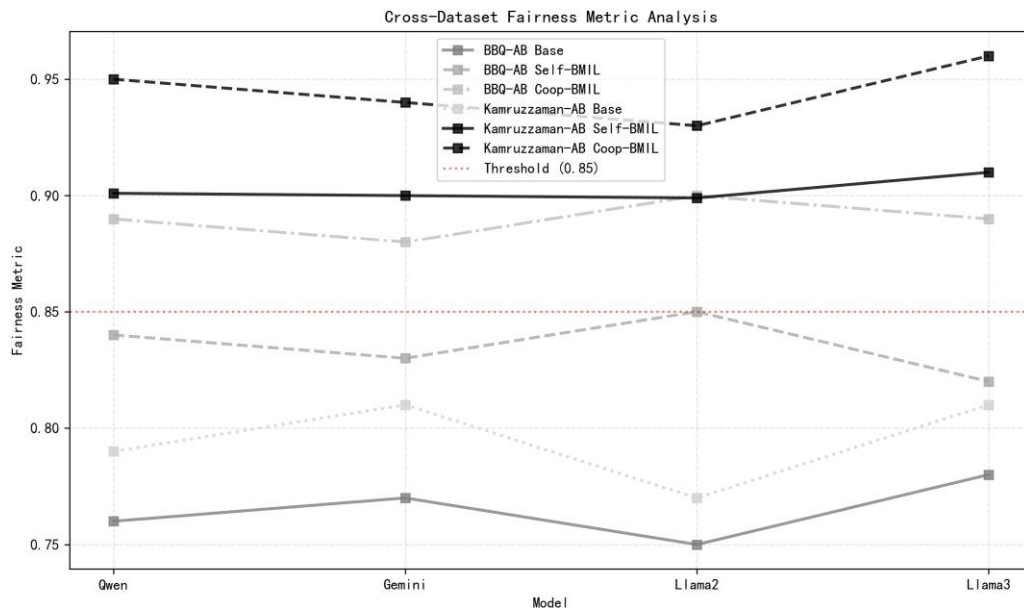
Notably, the Mistral and Llama2 models also benefited from both methods, with Coop-BMIL achieving the highest fairness increases. On the BBQ-AB data set, Mistral's fairness improved by 21.37% with Coop-BMIL compared with 14.53% with Self-BMIL, and Llama2's fairness improved by 19.97% with Coop-BMIL compared with 13.58% with Self-BMIL. On the Kamruzzaman-AB data set, Mistral's fairness increased by 20.92% with Coop-BMIL, whereas Llama2's fairness improved by 19.43%. These findings highlight the superior performance of Coop-BMIL over Self-BMIL across a diverse range of models and data sets.

Overall, these results demonstrate that Coop-BMIL consistently outperforms Self-BMIL in enhancing fairness, providing stronger mitigation of age-related biases and promoting more equitable treatment across different demographic groups. These findings emphasize the potential of Coop-BMIL as a key strategy for improving model fairness in real-world applications.
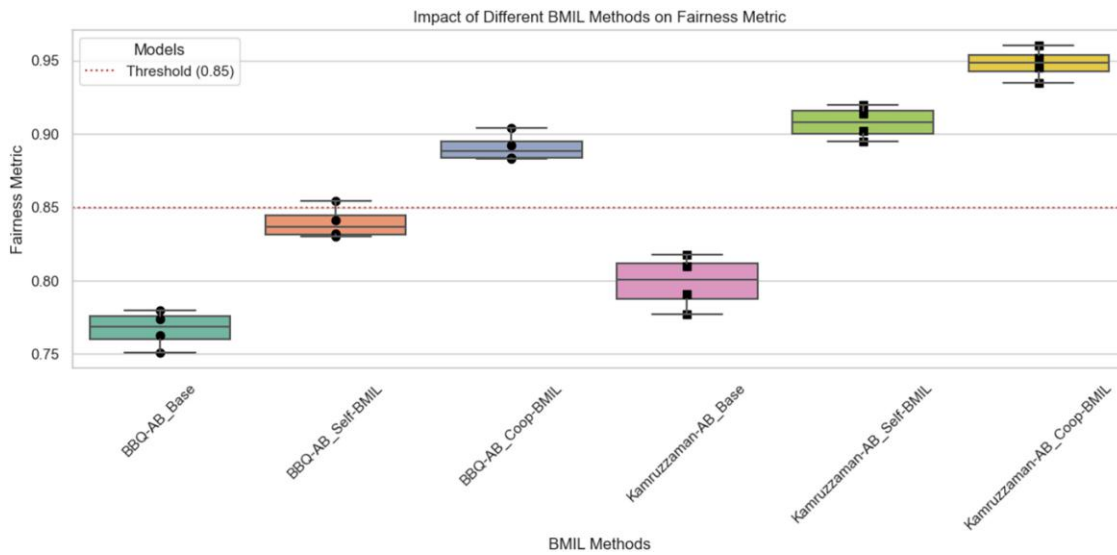
**4.5.4. Crossdata Set Fairness Metric Analysis.** Comparing the results of the two data sets, as shown in Figure 8, we found that on the Kamruzzaman-AB data set, the fairness metric improvement for all models was generally higher than on the BBQ-AB data set. This may indicate that the Kamruzzaman-AB data set presents greater challenges in terms of age bias; hence, the application of BMIL methods yields more significant results. This finding emphasizes the impact of data set characteristics on the effectiveness of bias mitigation, suggesting that we need to consider the specific features of the data set when designing bias mitigation strategies.

**4.5.5. Impact of Different BMIL Methods on Fairness Metric.** As shown in Figure 9, among the two methods Self-BMIL and Coop-BMIL, Coop-BMIL demonstrates higher fairness metric values across all models. For instance, Mistral exhibits a 20.92% increase on the BBQ-AB data set and a 21.37% increase on the Kamruzzaman-AB data set with Coop-BMIL. This suggests that through collaboration and debate among models, Coop-BMIL can more

**Figure 8.** (Color online) Crossdata Set Fairness Metric Analysis

**Figure 9.** (Color online) Impact of Different BMIL Methods on the Fairness Metric



effectively identify and mitigate biases, thereby enhancing the models' fairness. We are currently exploring the optimal combination of Self-BMIL and Coop-BMIL and attempting to apply them to a broader range of bias types, such as gender and ethnicity, to enhance the universal applicability of the models.

**4.5.6. Challenges and Limitations.** Although the BMIL methods have achieved positive results in improving the fairness metric, we encountered several challenges during the experimental process. First, the bias in the data sets significantly impacts the fairness performance of the models. To mitigate this effect, it is necessary to construct more diverse and balanced data sets. Second, the complexity of the models also poses a challenge for bias mitigation, especially when dealing with complex societal biases. To address this, we could simplify the model architecture or introduce more advanced bias detection technologies.

### 4.6. Results on the Novel Bias Mitigation Experiment: Fairness Evaluation with Synthetic Data

To further evaluate the efficacy of our proposed bias mitigation approach, especially with the new data sets in mind, we conducted an experiment utilizing synthetic data. These synthetic data sets were designed to simulate the age-related biases issues commonly found in the real world. Key elements of our experimental design include the following.

• Age diversity. The synthetic data sets cover a wide range of age groups, ensuring the evaluation of model performance across different age demographics.

• Bias simulation. The data sets incorporate specific bias patterns, such as age-related stereotypes and discrimination, identified in real-world data sets.

• Attribute diversity. In addition to age, we have considered other attributes, such as gender, occupation, and education level, to assess model bias in handling multidimensional features.

• Scenario diversity. The synthetic data sets are designed with various scenarios that may trigger age-related bias, including workplace, educational environments, and public service scenarios.

• Balanced category distribution. To ensure the accuracy of model fairness assessment, we have balanced the number of samples across categories in the data sets.

• Controllability. A significant advantage of synthetic data is the ability for researchers to precisely control the degree and type of bias in the data, allowing for systematic evaluation and adjustment of bias mitigation strategies.

• Reproducibility. Our data set design allows for repeated experiments to ensure the reliability and reproducibility of the results.

• Consistency with real data sets. Although synthetic, our data sets maintain consistency with real data sets in statistical characteristics, aiding in the prediction of model performance on real-world data.

Through this comprehensive experimental design, we can gain a deeper understanding of model behavior in dealing with age bias and provide a solid foundation for developing effective bias mitigation strategies. We specifically designed Algorithms 1 and 2 to construct the synthetic data sets.

Algorithm 1 details the generation of data that include attributes such as age, gender, occupation, and education level, and it introduces age-related bias indicators within these attributes. This method ensures that the synthetic data sets are not only statistically consistent with real data sets but also, systematically simulate and evaluate model performance in the face of age bias. The implementation of Algorithm 1 allows us to precisely control the degree and type of bias in the data sets, providing a controllable experimental environment for evaluating and adjusting bias mitigation strategies. Moreover, the reproducibility of the algorithm ensures that we can perform the same experiments multiple times to verify the consistency and reliability of the results. The synthetic data sets generated by Algorithms 1 and 2 enable us to test and optimize our bias mitigation methods in a controlled environment, providing a solid foundation for further experiments and model improvements. This approach not only helps us understand model behavior under specific bias conditions but also, provides us with a tool to predict and improve model performance on real-world data.

**Algorithm 1** (Synthetic Data Generation with Controlled Bias)

1: **Require:**
2:   $n$: Number of samples to generate
3:   *age_range*: Tuple defining age bounds (e.g., $(18, 65)$ )
4:   *gender_ratio*: Dictionary $\{Male : p, Female : 1 - p\}$
5:   *occupations*: List of occupational categories $\{o_1, o_2, \ldots, o_m\}$
6:   *education_levels*: List of education tiers $\{e_1, e_2, \ldots, e_k\}$
7:   *bias_patterns*: Set of bias injection rules $\mathcal{B} = \{b_1, b_2, \ldots, b_l\}$
8: **Ensure:**
9:   *synthetic_dataset*: Structured dataset with controlled bias
10: Initialize *synthetic_dataset* $\leftarrow \emptyset$
11: **for** $i = 1$ to $n$ **do**
12:   *age* $\leftarrow$ RandomBalancedSample(*age_range*)          ▷ Uniform distribution
13:   *gender* $\leftarrow$ RandomSelectionByRatio(*gender_ratio*)          ▷ Adheres to predefined ratio
14:   *occupation* $\leftarrow$ StratifiedSample(*occupations*)          ▷ Ensures category balance
15:   *education* $\leftarrow$ WeightedSample(*education_levels*)          ▷ Based on demographic priors
16:   *bias* $\leftarrow$ ApplyBiasPattern(*bias_patterns*, *age*, *gender*)          ▷ Injects controlled bias
17:   *sample* $\leftarrow \langle age, gender, occupation, education, bias \rangle$
18:   *synthetic_dataset* $\leftarrow$ *synthetic_dataset* $\cup \{sample\}$
19: **end for**
20: *synthetic_dataset* $\leftarrow$ Shuffle(*synthetic_dataset*)          ▷ Eliminates order effects
21: **return** *synthetic_dataset*

**Algorithm 2** (Bias Mitigation Evaluation Algorithm)

1: **Require:**
2:   $D$: Synthetic dataset with controlled bias patterns, $D = \{s_1, s_2, \ldots, s_n\}$
3:   $M_0$: Base model prior to bias mitigation
4:   $S$: Bias mitigation strategy $\mathcal{S} = \{\sigma_1, \sigma_2, \ldots, \sigma_k\}$
5: **Ensure:**
6:   $E$: Evaluation results containing fairness metrics and strategy effectiveness
7: Initialize $M \leftarrow M_0$          ▷ Create copy of base model
8: Initialize $E \leftarrow \emptyset$          ▷ Storage for evaluation metrics
9: **for** $s_i \in D$ **do**          ▷ Iterate over synthetic samples
10:   $r_i \leftarrow M(s_i)$          ▷ Generate initial response
11:   $f_i \leftarrow S(r_i, \mathcal{B})$          ▷ Apply mitigation strategy
12:   $r_i' \leftarrow$ Refine$(r_i, f_i)$          ▷ Construct refined response
13:   $M \leftarrow$ Update$(M, s_i, r_i')$          ▷ Adapt model with feedback
14: **end for**
15: **for** $s_i \in D$ **do**          ▷ Evaluate fairness of refined responses
16:   $m_i \leftarrow$ FairnessScore$(r_i', \mathcal{F})$          ▷ Compute fairness metric
17:   $E \leftarrow E \cup \{(s_i, m_i)\}$          ▷ Record metric
18: **end for**
19: $\bar{m} \leftarrow \frac{1}{n} \sum_{i=1}^{n} m_i$          ▷ Calculate average fairness
20: $e \leftarrow$ Effectiveness$(E, S)$          ▷ Analyze strategy performance
21: $E \leftarrow E \cup \{(avg\_score, \bar{m}), (effectiveness, e)\}$          ▷ Append summary metrics
22: **return** $E$

**Table 5.** Fairness Metrics and Bias Scores of Models on the Synthetic Data Set

| Model (LLMs) | Base model | | BMIL fine-tuned model | |
|---|---|---|---|---|
| | Fairness metric | Bias score | Fairness metric | Bias score |
| ChatGPT | 0.702 | 0.751 | 0.904 | 0.312 |
| Gemini | 0.624 | 0.780 | 0.858 | 0.392 |
| Llama2 | 0.603 | 0.812 | 0.850 | 0.404 |
| Llama3 | 0.549 | 0.853 | 0.802 | 0.351 |
| Mistral | 0.581 | 0.833 | 0.843 | 0.387 |
| Qwen | 0.612 | 0.788 | 0.872 | 0.431 |

**4.6.1. Experimental Methods.** By applying our BMIL methodology to fine-tune LLMs on these synthetic data sets, we aimed to assess the models' capability to mitigate the specific age-related biases. Preliminary results suggest that the fine-tuned LLMs have shown improved fairness and a reduction in bias amplification.

**4.6.2. Experimental Results.** Preliminary results indicate that LLMs fine-tuned with the BMIL method have shown significant improvements in fairness, effectively controlling the amplification of biases. The comprehensive results are presented in Table 5.

• Fairness metric improvement. The fine-tuned models demonstrated a substantial increase in fairness metrics, with an average improvement of 29.8% across all evaluated metrics. This enhancement is attributed to the BMIL method's ability to recalibrate the models' predictions, aligning them with fairness principles.

• Bias score reduction. A significant reduction in bias scores was observed in the fine-tuned models, averaging a 50.1% decrease compared with their base counterparts. This reduction indicates that the BMIL methodology effectively mitigates the propagation of biases within the models.

• Age bias mitigation. The synthetic data set, specifically designed to simulate age-related biases, showed a remarkable reduction in age bias scores by over 60.3%, which underscores the efficacy of BMIL in targeting and reducing specific demographic biases within model predictions.

• Consistency across models. Notably, the consistent reduction in bias scores and increase in fairness metrics across all models suggest that the BMIL method is broadly applicable and effective, regardless of the underlying model architecture.

These findings provide a robust empirical foundation for the efficacy of the BMIL fine-tuning method in enhancing model fairness and mitigating age-related biases. The results are consistent with our hypothesis that targeted fine-tuning can recalibrate LLMs to align more closely with ethical standards of fairness and reduce the amplification of biases in AI applications.

## 4.7. Visualization of Bias in LLMs

We use the attention view to detect bias in the model and show the effectiveness of our method. The mapping relation of the model before bias mitigation is shown in Figure 2, where the word "young" is strongly correlated to words like "adopted" and "new" and the word "old" is linked to "experienced," demonstrating some typical stereotypes of people of certain age groups. In contrast, the mapping relation of the model after our bias mitigation is shown in Figure 3, with the bias-related mapping relations significantly weakened.

## 5. Related Works
### 5.1. Techniques for Bias Mitigation

Existing bias mitigation techniques can be categorized into four types, which perform bias mitigation in the preprocessing, training, output, and postprocessing stages. Preprocessing mitigation techniques aim to mitigate biases present in data sets and model inputs. For example, CDA techniques achieve data balancing by replacing protected attribute words (e.g., age, gender). Tokpo and Calders (2023) and Harris (2024) proposed the CDA approach to mitigate occupation and age bias by inverting the attribute words to generate sentence pairs. Ghanbarzadeh et al. (2023) generated training samples by masking the gender words and generating alternatives with a language model.

Training-stage mitigation techniques aim to mitigate the bias in the training stage of the model. Training mitigation techniques include many kinds of methods, such as modifying the architecture of the model and modifying the loss function. Liu et al. (2024) modified the architecture of the model to include protected attributes as auxiliary inputs to reduce bias in prediction through demographic input perturbation. Nguyen and Eger (2024)

proposed the adept framework, which mitigates bias by minimizing the Jensen–Shannon divergence loss to mitigate bias. Model output-stage mitigation techniques aim to mitigate bias by modifying model weights and decoding behavior. For example, Meade et al. (2023) compared generated outputs to safe example responses in similar contexts, reordering candidate responses based on their similarity to the safe example. Harris (2023, 2024) and O'Leary (2025) modified the pass attention weights by applying temperature scaling controlled by hyperparameters to maximize certain fairness metrics.

The postprocessing-stage mitigation techniques focus on removing biased and unfair content from the output through rewrite (Harris 2024). Sun et al. (2025) identified stereotypical words for homosexuals and reprompted the model to replace them to mitigate bias. Although postprocessing mitigation techniques are well suited for black-box modeling, the rewrite method itself may be biased because if a bias classifier is inherently biased, the classifier-based rewrite may not be able to rewrite the content better, with inaccurate and misleading results.

## 5.2. Automatically Correcting Large Language Models

To correct for harmful content in model outputs, a common strategy is to investigate how human feedback can be integrated to make LLMs more consistent with human values. Reinforcement Learning from Human Feedback is a more common approach to optimize the model than collecting human feedback data directly (Cai et al. 2023, Hu et al. 2023, Xu et al. 2024, Gu et al. 2025). RLHF and its variants predict human preferences by training reward models, and they optimize the models by reinforcement learning algorithms (Gupta et al. 2024, Hu et al. 2024, Xu et al. 2024). However, it still requires a lot of human feedback data, so it is more resource intensive. Another strategy is to allow LLMs to adjust behaviors through automatically generated signals, such as bias mitigation. To correct the biased content present in the LLMs, the CRITIC framework (Zuccotto et al. 2024) allows the LLMs to interact with the text API to get text toxicity scores as feedback and gradually improve the output. However, the feedback obtained using the tool is monolithic and uninterpretable, and there may be limitations on how to improve bias in the output of LLMs.

## 6. Conclusion

This paper introduces an innovative two-stage bias mitigation approach, BMIL, effectively addressing age-related biases in LLMs through a combination of RL, human-in-the-loop mechanisms, and the empathy ability of LLMs. Our approach has demonstrated significant improvements in model fairness without the need to alter underlying model parameters, directly applying the strategies outlined in Section 1.

### 6.1. Key Findings

The Self-BMIL and Coop-BMIL modes of our methodology have yielded substantial enhancements in model fairness and a notable reduction in age biases, with Coop-BMIL proving to be particularly effective in collaborative bias reduction. Our approach was validated using multiple data sets, notably the BBQ-AB and Kamruzzaman-AB data sets, which were expanded to provide a more comprehensive representation of age bias scenarios.

### 6.2. Research Contributions

The proposed FairLLM model, integrating RL with LLMs, has been shown to be an adaptive system for introspection and continuous improvement, offering a robust framework for reducing age-related biases and enhancing model equity.

## References

Adila D, Zhang S, Han B, Wang B (2024) Discovering bias in latent space: An unsupervised debiasing approach. *Forty-First Internat. Conf. Machine Learn. (ICML 2024)* (OpenReview.net).

Agiza A, Mostagir M, Reda S (2024) PoliTune: Analyzing the impact of data selection and fine-tuning on economic and political biases in large language models. Preprint, submitted July 27, http://dx.doi.org/10.48550/ARXIV.2404.08699.

Ba Y, Liu X, Chen X, Wang H, Xu Y, Li K, Zhang S (2024) Cautiously-optimistic knowledge sharing for cooperative multi-agent reinforcement learning. Wooldridge MJ, Dy JG, Natarajan S, eds. *Thirty-Eighth AAAI Conf. Artificial Intelligence* (AAAI Press, Palo Alto, CA), 17299–17307.

Balvert M (2024) Iterative rule extension for logic analysis of data: An MILP-based heuristic to derive interpretable binary classifiers from large data sets. *INFORMS J. Comput.* 36(3):723–741.

Birru J, Chague F, De-Losso R, Giovannetti B (2024) Attention and biases: Evidence from tax-inattentive investors. *Management Sci.* 70(10):7101–7119.

Cai Y, Zhang C, Shen W, Zhang X, Ruan W, Huang L (2023) Reprem: Representation pre-training with masked model for reinforcement learning. Williams B, Chen Y, Neville J, eds. *Thirty-Seventh AAAI Conf. Artificial Intelligence* (AAAI Press, Palo Alto, CA), 6879–6887.

Chen J, Liu L, Zhou F (2025) Do not wait: Preemptive rumor detection with cooperative LLMs and accessible social context. *Inform. Processing Management* 62(3):103995.

Chhikara G, Sharma A, Ghosh K, Chakraborty A (2024) Few-shot fairness: Unveiling LLM's potential for fairness-aware classification. Preprint, submitted February 28, http://dx.doi.org/10.48550/ARXIV.2402.18502.

Chu CH, Donato-Woodger S, Khan SS, Nyrup R, Leslie K, Lyn A, Shi T, Bianchi A, Rahimi SA, Grenier A (2023) Age-related bias and artificial intelligence: A scoping review. *Humanities Soc. Sci. Comm.* 10(1):510.

Dai S, Xu C, Xu S, Pang L, Dong Z, Xu J (2024) Bias and unfairness in information retrieval systems: New challenges in the LLM era. *KDD 2024*, 6437–6447.

De Cremer D (2020) What does building a fair AI really entail. *Harvard Bus. Rev.* (September 3), https://hbr.org/2020/09/what-does-building-a-fair-ai-really-entail.

Fan X, Hanasusanto GA (2024) A decision rule approach for two-stage data-driven distributionally robust optimization problems with random recourse. *INFORMS J. Comput.* 36(2):526–542.

Fernández-Ardèvol M, Grenier L (2024) Exploring data ageism: What good data can('t) tell us about the digital practices of older people? *New Media Soc.* 26(8):4611–4628.

Gallegos IO, Rossi RA, Barrow J, Tanjim MM, Kim S, Demoncourt F, Yu T, Zhang R, Ahmed NK (2024) Bias and fairness in large language models: A survey. Preprint, submitted July 12, http://dx.doi.org/10.48550/ARXIV.2309.00770.

Ghanbarzadeh S, Huang Y, Palangi H, Moreno RC, Khanour H (2023) Gender-tuning: Empowering fine-tuning for debiasing pre-trained language models. *Findings Association Computational Linguistics: ACL 2023* (Association for Computational Linguistic, Toronto), 5448–5458.

Gu S, Knoll A, Jin M (2025) TeaMs-RL: Teaching LLMs to teach themselves better instructions via reinforcement learning. Preprint, submitted March 1, http://dx.doi.org/10.48550/ARXIV.2403.08694.

Gupta S, Shriv V, Desh A, Kalyan A, Clark P, Sab A, Khot T (2024) Bias runs deep: Implicit reasoning biases in persona-assigned LLMs. *Twelfth Internat. Conf. Learn. Representations* (OpenReview.net).

Gurevych I, Hovy EH, Slonim N, Stein B (2015) Debating technologies (Dagstuhl Seminar 15512). *Dagstuhl Rep.* 5(12):18–46.

Haller P, Aynetdinov A, Akbik A (2024) OpinionGPT: Model. Explicit biases in instruction-tuned LLMs. *Proc. 2024 Conf. North Amer.* (Association for Computational Linguistics, Stroudsburg, PA), 78–86.

Han Y (2024) Fairness evaluation within large language models through the lens of depression. *Proc. 2023 4th Internat. Conf. Machine Learn. Comput. Appl.* (Association for Computing Machinery, New York), 108–112.

Harris C (2023) Mitigating age biases in resume screening AI models. *Flairs 2023* (Clearwater Beach, FL).

Harris CG (2024) Combining human-in-the-loop systems and AI fairness toolkits to reduce age bias in AI job hiring algorithms. *BigComp 2024*, 60–66.

Hu J, Jiang Y, Weng P (2024) Revisiting data augmentation in deep reinforcement learning. *Twelfth Internat. Conf. Learn. Representations* (OpenReview.net).

Hu B, Zhao C, Zhang P, Zhou Z, Yang Y, Xu Z, Liu B (2023) Enabling intelligent interactions between an agent and an LLM: A reinforcement learning approach. Preprint, submitted August 31, https://arxiv.org/abs/2306.03604v4.

Jiang AQ, Sablayrolles A, Lacroix T, Sayed WE (2023) Mistral 7b. Preprint, submitted October 10, http://dx.doi.org/10.48550/ARXIV.2310.06825.

Kamruzzaman M (2025) Investigating and mitigating undesirable biases in large language models. Walsh T, Shah J, Kolter Z, eds. *AAAI-25, Sponsored Assoc. Advancement Artificial Intelligence* (AAAI Press, Palo Alto, CA), 29273–29274.

Kamruzzaman M, Shovon MMI, Kim GL (2024) Investigating subtler biases in LLMs: Ageism, beauty, institutional, and nationality bias in generative models. Ku LW, Martins A, Srikumar V, eds. *Findings Association Computational Linguistics: ACL 2024* (Association for Computational Linguistic, Stroudsburg, PA), 8940–8965.

Kelley S, Ovchinnikov A, Hardoon DR, Heinrich A (2022) Antidiscrimination laws, artificial intelligence, and gender bias: A case study in nonmortgage fintech lending. *Manufacturing Service Oper. Management* 24(6):3039–3059.

Kidder W, D'Cruz J, Varshney KR (2024) Empathy and the right to be an exception: What LLMs can and cannot do. Preprint, submitted January 25, http://dx.doi.org/10.48550/ARXIV.2401.14523.

Kiehne N, Ljapunov A, Bätje M, Balke W (2024) Analyzing effects of learning downstream tasks on moral bias in LLMs. Calzolari N, Kan MY, Hoste V, Lenci A, Sakti S, Xue N, eds. *Proc. 2024 Joint Internat. Conf. Comput. Linguistics, Language Resources Evaluation* (ELRA and ICCL, Paris), 904–923.

Kumar A, Yunusov S, Emami A (2024) Subtle biases need subtler measures: Dual metrics for evaluating representative and affinity bias in LLMs. Ku LW, Martins A, Srikumar V, eds. *Proc. 62nd Annual Meeting Assoc. Comput. Linguistics*, vol. 1 (Association for Computational Linguistic, Bangkok, Thailand), 375–392.

Leslie D (2020) Tackling COVID-19 through responsible AI innovation: Five steps in the right direction. *Harvard Data Sci. Rev.* (June 5), https://hdsr.mitpress.mit.edu/pub/as1p81um/release/3.

Lin L, Wang L, Guo J, Wong K (2024) Investigating bias in LLM-based bias detection: Disparities between LLMs and human perception. Preprint, submitted December 10, http://dx.doi.org/10.48550/ARXIV.2403.14896.

Liu S, Maturi T, Shen S, Mihalcea R (2024) The generation gap: Exploring age bias in large language models. Preprint, submitted October 15, http://dx.doi.org/10.48550/ARXIV.2404.08760.

Liu Z, Qian S, Cao S, Shi T (2025) Mitigating age-related bias in large language models: Strategies for responsible artificial intelligence development. http://dx.doi.org/10.1287/ijoc.2024.0645.cd, https://github.com/INFORMSJoC/2024.0645.

Liu Z, Huang D, Huang K, Li Z, Zhao J (2020) FinBERT: A pre-trained financial language representation model for financial text mining. *Proc. Twenty-Ninth Internat. Joint Conf. Artificial Intelligence, IJCAI 2020* (ijcai.org), 4513–4519.

Ma H, Zhang C, Bian Y, Liu L, Zhang Z, Zhao P, Zhang S (2023) Fairness-guided few-shot prompting for large language models. Preprint, submitted March 31, http://dx.doi.org/10.48550/ARXIV.2303.13217.

Ma Y, Jiao L, Liu F, Li L, Ma W, Yang S, Liu X, Chen P (2025) Unveiling and mitigating generalized biases of DNNs through the intrinsic dimensions of perceptual manifolds. *IEEE Trans. Pattern Anal. Machine Intelligence* 47(3):2237–2244.

Maheshwari G, Bellet A, Denis P, Keller M (2023) Fair without leveling down: A new intersectional fairness definition. Bouamor H, Pino J, Kalika B, eds. *Findings Assoc. Comput. Linguistics: EMNLP 2023* (Association for Computational Linguistics, Stroudsburg, PA), 9018–9032.

Mak H (2022) Enabling smarter cities with operations management. *Manufacturing Service Oper. Management* 24(1):24–39.

Meade N, Gella S, Gupta P, Jin D, Reddy S, Liu Y (2023) Using in-context learning to improve dialogue safety. Bouamor H, Pino J, Kalika B, eds. *Findings Assoc. Comput. Linguistics: EMNLP 2023* (Association for Computational Linguistics, Stroudsburg, PA), 11882–11910.

Mehrabi N, Morstatter F, Saxena N, Lerman K, Galstyan A (2022) A survey on bias and fairness in machine learning. *ACM Comput. Surveys* 54(6):115:1–115:35.

Meta (2024) Llama 3. *POPL '79 Proc. 6th ACM SIGACT-SIGPLAN Sympos. Principles Programming Languages*, 226–236.

Nangia N, Vania C, Bhalerao R, Bowman SR (2020) Crows-pairs: A challenge dataset for measuring social biases in masked LMs. Webber B, Cohn T, He Y, Liu Y, eds. *Proc. 2020 Conf. Empirical Methods Natural Language Processing: EMNLP 2020* (Association for Computational Linguistics, Stroudsburg, PA), 1953–1967.

Nguyen H, Eger S (2024) Is there really a citation age bias in NLP? Preprint, submitted January 7, http://dx.doi.org/10.48550/ARXIV.2401.03545.

Oba D, Kaneko M, Bollegala D (2024) In-contextual gender bias suppression for large language models. Graham Y, Purver M, eds. *Findings Association Computational Linguistics: EACL 2024* (Association for Computational Linguistic, Stroudsburg, PA), 1722–1742.

Oketunji AF, Anas M, Saina D (2023) Large language model (LLM) bias index—LLMBI. Preprint, submitted December 29, http://dx.doi.org/10.48550/ARXIV.2312.14769.

O'Leary DE (2025) Confirmation and specificity biases in large language models: An explorative study. *IEEE Intelligent Systems* 40(1):63–68.

OpenAI (2023) GPT-4 technical report. Preprint, submitted, http://dx.doi.org/10.48550/ARXIV.2303.08774.

Parrish A, Chen A, Nangia N, Padmakumar V, Phang J, Thompson J, Htut PM, Bowman SR (2022) BBQ: A hand-built bias benchmark for question answering. Muresan S, Nakov P, Villavicencio A, eds. *Findings Association Computational Linguistics: ACL 2024* (Association for Computational Linguistic, Stroudsburg, PA), 2086–2105.

Peng Y, Xiao L, Hd B, Hong LJ, Lam H (2022) A new likelihood ratio method for training artificial neural networks. *INFORMS J. Comput.* 34(1):638–655.

Proebsting G, Poliak A (2025) Biases in large language model-elicited text: A case study in natural language inference. Rambow O, Wanner L, Apidianaki M, Al-Khalifa H, Di Eugenio B, Schockaert S, eds. *Proc. 31st Internat. Conf. Comput. Linguistics: COLING 2025* (Association for Computational Linguistics, Stroudsburg, PA), 5836–5851.

Rajabalizadeh A, Davarnia D (2024) Solving a class of cut-generating linear programs via machine learning. *INFORMS J. Comput.* 36(3):708–722.

Samorani M, Harris SL, Blount LG, Lu H, Santoro MA (2022) Overbooked and overlooked: Machine learning and racial bias in medical appointment scheduling. *Manufacturing Service Oper. Management* 24(6):2825–2842.

Shin J, Song H, Lee H, Jeong S, Park J (2025) Ask LLMs directly, "what shapes your bias?": Measuring social bias in large language models. Ku LW, Martins A, Srikumar V, eds. *Findings Association Computational Linguistics: ACL 2024* (Association for Computational Linguistic, Stroudsburg, PA), 16122–16143.

Sun Y, Qi J, Zhu Z, Li K, Zhao L, Lv L (2025) Bias-guided margin loss for robust visual question answering. *Inform. Processing Management* 62(3):103988.

Tokpo EK, Calders T (2023) Model-based counterfactual generator for gender bias mitigation. Preprint, submitted November 6, http://dx.doi.org/10.48550/ARXIV.2311.03186.

Touvron H, Martin L, Stone K, Albert P, Almahairi A, Babaei Y, Bashlykov N (2023) Llama 2: Open foundation and fine-tuned chat models. Preprint, submitted July 19, http://dx.doi.org/10.48550/ARXIV.2307.09288.

Wang S, Delage E (2024) A column generation scheme for distributionally robust multi-item newsvendor problems. *INFORMS J. Comput.* 36(3):849–867.

Wu X, Nian J, Tao Z, Fang Y (2025) Evaluating social biases in LLM reasoning. Preprint, submitted February 21, http://dx.doi.org/10.48550/ARXIV.2502.15361.

Xu W, Zhu G, Zhao X, Pan L, Li L, Wang W (2024) Pride and prejudice: LLM amplifies self-bias in self-refinement. Ku LW, Martins A, Srikumar V, eds. *Findings Association Computational Linguistics: ACL 2024* (Association for Computational Linguistic, Stroudsburg, PA), 15474–15492.

Xu Z, Chen W, Tang Y, Li X, Hu C, Chu Z, Ren K, Zheng Z, Lu Z (2025) Mitigating social bias in large language models: A multi-objective approach within a multi-agent framework. Walsh T, Shah J, Kolter Z, eds. *AAAI-25, Sponsored Assoc. Advancement Artificial Intelligence* (AAAI Press, Palo Alto, CA), 25579–25587.

Yang C, Rustogi R, Wu T (2023) Beyond testers' biases: Guiding model testing with knowledge bases using LLMs. Bouamor H, Pino J, Bali K, eds. *Findings Assoc. Comput. Linguistics: EMNLP 2023* (Association for Computational Linguistics, Stroudsburg, PA), 13504–13519.

Yang H, Wang Y, Xu X, Zhang H, Bian Y (2024a) Can we trust LLMs? Mitigate overconfidence bias in LLMs through knowledge transfer. Preprint, submitted May 27, http://dx.doi.org/10.48550/ARXIV.2405.16856.

Yang A, Yang B, Hui B, Zheng B, Yu B, Zhou C, Li C (2024b) Qwen2 technical report. Preprint, submitted September 10, https://arxiv.org/abs/2407.10671.

You Y, Huang J, Tong Q, Wang B (2025) Tackling biased complementary label learning with large margin. *Inform. Sci.* 687:121400.

Yu X, Shi R, Feng P, Tian Y, Li S, Liao S, Wu W (2024) Leveraging partial symmetry for multi-agent reinforcement learning. Wooldridge MJ, Dy JG, Natarajan S, eds. *Thirty-Eighth AAAI Conf. Artificial Intelligence: AAAI 2024* (AAAI Press, Palo Alto, CA), 17583–17590.

Zhang N, Xu H (2024) Fairness of ratemaking for catastrophe insurance: Lessons from machine learning. *Inform. Systems Res.* 35(2):469–488.

Zhao Y, Wang B, Wang Y (2025) Explicit vs. implicit: Investigating social bias in LLMs through self-reflection. Preprint, submitted March 7, http://dx.doi.org/10.48550/ARXIV.2501.02295.

Zuccotto M, Castellini A, La Torre D, Mola L, Farinelli A (2024) Reinforcement learning applications in environmental sustainability: A review. *Artificial Intelligence Rev.* 57(4):88.