

# LoRA-Based Adaptation of Large Language Models for Bias-Reduced Healthcare Dialogue Systems

Nagavaralakshmi C.K, Jithinkrishnan P. G, Adwika Raghav, Jerin S. Das

Department of Computer Science Yenepoya University, India

Email: varulakshmi66@yahoo.com, jithinkrishnan541014@gmail.com,  
adwikaraghav165@gmail.com, jerins305@gmail.com

**Abstract**—Current healthcare dialogue systems based on large language models show excellent capability in offering medical advice and conversational interactions, but they are susceptible to demographic bias, resulting in varying recommendations across different gender, age, and ethnic subgroups. We present an efficient demographically balanced healthcare assistant generation method that combines low-rank adaptation, counterfactual augmentation, and auxiliary bias mitigation loss function. The proposed approach was tested on the HealthCareMagic (150K dialogue) and MIMIC-III (85K) datasets, with five-fold cross-validation across 12 intersectional demographic groups using PyTorch with four NVIDIA A100 GPUs. The results showed medical F1-score of 92.3%, dialogue fluency of 89.1% (BLEU-4 score), and 2.8% demographic disparity ratio, outperforming full-fledged fine-tuning with 8.7% accuracy and 94.7% bias mitigation, utilizing only 4.2M parameters (11.2× faster fine-tuning). Intersectional bias reduced by 35.7% among female+Black groups, and the system can be deployed in telehealth portals and hospital environments with benefits for marginalized communities, though performance in rare diseases deteriorates by 7.2%.

## I. INTRODUCTION

The fast pace of development of artificial intelligence (AI) has drastically changed the approach to providing care services in healthcare settings, with large language models (LLMs) becoming key components of conversational medical applications. Conversational dialogue systems for the provision of medical assistance and patient-oriented dialogue interaction benefit from the context-sensitivity and language generation features of LLMs. Despite showing great promise in early trials, when used in the practice of medicine such applications can pose serious questions in terms of reliability and safety of their use in clinical settings. Although LLMs can be used to analyze large datasets of information, making patterns based on training, they are known to reproduce bias contained in their training material, thus resulting in discriminatory outcomes.

The history of dialogue systems in healthcare begins with the earliest rule-based expert systems that were developed in the 1970s, such as MYCIN for infectious disease diagnosis. Later, in the 1990s and early 2000s, expert systems evolved into natural language processing techniques which included probabilistic approaches like hidden Markov model based intent recognition. The breakthrough came in 2017 when transformers emerged, allowing BERT and GPT-type models to generate highly coherent and fluent dialogues. More recently, within the last five years, there has been an outburst of large language model applications for healthcare, ranging from symptom checker solutions to virtual health assistants. Advancements include high-performance models like Med-PaLM and ClinicalGPT, capable of matching the expert accuracy on medical benchmarks. However, this success is increasingly followed by the discovery of demographic biases in recommendations, where patients' outcomes differ based on their gender, age, ethnicity, or socioeconomic background [1] [2].

Trends in the present day illustrate a two-pronged direction towards both innovation and examination in this field. While the market value of artificial intelligence in healthcare is expected to surpass 187 billion by 2030 through implementation of LLMs in telemedicine systems and EHR systems, parameter-efficient fine-tuning methods, such as LoRA, have made model customisation more accessible by allowing the mitigation of biases without complete retraining [3]. In LoRA, updates to weight parameters are decomposed into low-rank matrices to decrease computational costs, along with retaining knowledge from base models. New developments include BA-LoRA, which utilizes LoRA for mitigating biases specifically, and overcoming the problem of “catastrophic inheritance,” whereby fine-tuning further reinforces prior biases [4]. Meanwhile, counterfactual data augmentation has emerged as an effective method of promoting fairness by creating counterfactual queries that reverse the demographics of questions posed in training datasets [5].

Notwithstanding all of the progress, a number of formidable challenges remain that pose a danger to the reliable use of LLMs in healthcare dialogue systems. Specifically, the fundamental challenge is demographic bias manifestation, wherein studies have reliably indicated that medical advice varies according to demographic characteristics. For example, LLMs are likely to advise more invasive testing for males than females when their symptomologies match, while underemphasizing pain management for older ethnic minorities [2] [6]. These biases are rooted in biased data distribution due to historical healthcare inequality along with the “black box”

nature of transformers. Current techniques to mitigate bias, such as reweighting training data, adversarial learning, and post hoc correction, are insufficient for dialogue systems because they interfere with the conversation or demand an unreasonably high cost. It is crucial to note that tackling these challenges is paramount in an era when healthcare facilities around the world are grappling with a lack of human resources and increasing demands for services. Biased AI can exacerbate social injustices, destroy the trust of patients, and attract lawsuits in the face of new legislation, including the EU AI Act's designation of medical devices as high-risk. In the United States, the Joint Commission has issued warnings regarding the use of clinical decision-making tools that exhibit algorithmic discrimination, and Medicare policy dictates the need for fairness audits. On a more positive note, fair dialogue systems ensure equal access to healthcare services for all people. This objective can be achieved by addressing bias issues at the parameter level using LoRA.

**A major research gap in reviewing the literature arises as individual components have been studied, but there is a lack of literature addressing the integration of all these approaches in healthcare dialogue applications. The existing literature review focuses solely on general LLMs' reliability [7] or medical imaging bias [8], ignoring the specific problems of conversational settings involving empathy, medical accuracy, and fairness simultaneously. Studies show that biased answers not only lead patients astray but also pass on the prejudices to medical students, thus reinforcing the cycle of discrimination [9]. Methods to efficiently modify model parameters in dialogue scenarios need further investigation, with the majority of previous work being confined to classification problems instead of generation. Lastly, auxiliary bias losses, designed to minimize demographic differences during fine-tuning, require benchmarking frameworks.**

As the title suggests, this review seeks to fill these knowledge gaps by conducting an exhaustive analysis of parameter-efficient techniques used for mitigating the biases in LLMs utilized in healthcare conversation agents. Particularly, we investigate the potential synergies between LoRA fine-tuning, counterfactual augmentations, and additional loss functions that can result in a fair system without sacrificing its usefulness in a clinical context.

To guide this analysis, we pose the following research questions:

- 1) **RQ1: How do demographic biases manifest in the outputs of healthcare dialogue systems, and what metrics best capture these disparities across gender, age, and ethnicity?**
- 2) **RQ2: To what extent do parameter-efficient techniques like LoRA preserve medical accuracy while reducing bias, compared to full fine-tuning baselines?**
- 3) **RQ3: What role does counterfactual data augmentation play in balancing dialogue datasets, and how does it interact with auxiliary bias detection losses during training?**
- 4) **RQ4: Which combinations of mitigation strategies yield the optimal trade-off between fairness, fluency, latency,**

and computational efficiency for real-world deployment?

Our contributions can be summarized into four aspects. Firstly, we contribute the first complete classification of the bias types in medical dialogue LLMs that can be categorized into explicit recommendation biases and implicit empathy gaps. Secondly, we perform a meta-analysis on 50+ papers (2021–2026) to empirically analyze the effectiveness of the LoRA techniques based on effect size. Thirdly, we develop a new evaluation framework that combines counterfactual robustness testing and clinical validity assessment. Finally, we provide guidelines for deploying fair dialogue systems within clinical practice processes.

## II. LITERATURE REVIEW

Incorporating large language models (LLMs) in health conversation platforms has revolutionized patient engagement through scalable, context-aware assistance. Yet, demographic imbalances with respect to gender, age, and ethnicity continue to pose vital issues of fairness and clinical accuracy. This literature review investigates the development of fairness in clinical language models within four stages: diagnostic research, transformer bias analysis, parameter-efficient tuning, and hybrid approaches. This paper evaluates 15 seminal papers published between 2016 and 2025 to address methodological shortcomings, conflicting evaluation criteria, and challenges in implementation.

### A. Early Bias Detection and Foundational Work (2016–2023)

The concept of fairness in machine learning is rooted in the theory of equality of opportunity by Hardt et al. (2016). This principle is still the basis of bias assessment in machine learning models [10]. Concurrently, the emergence of large-scale clinical databases like MIMIC-III made it possible to assess demographic biases empirically in artificial intelligence applications in healthcare [11].

The initial research in adapting LLMs emphasized efficiency over fairness. Houlby et al. (2019) presented adapters as an approach that facilitates efficient fine-tuning of machine learning models without losing their prior knowledge [12]. This idea was further improved by Hu et al. (2021), who developed low-rank adaptation (LoRA), in which weight updates

**$\Delta W = BA$  are performed using two matrices of size  $r \ll d$ . The proposed method achieved 10,000× parameter efficiency.**

Although efficient and effective in transferring knowledge to a new domain, neither of these approaches considered demographic bias.

Within the realm of medicine, it is shown by Singhal et al. that large language models like Med-PaLM store a lot of clinical information, emphasizing their application possibilities in the medical industry [13]. However, this also poses an increased risk of inherent biases contained in the data used for training.

Patel et al., in one of the first studies on bias within the field of medicine, found that female patients received 27% fewer treatments concerning cardiac problems than male patients under identical circumstances [8].

### B. Transformer-Based Bias Analysis (2024–2025)

**A recent series of studies has attempted to measure bias in transformer-based LLMs in various clinical applications. For example, McDermott et al. (2025) found that 31% more misdiagnoses occurred among minority ethnic groups during emergency triage tasks [1]. Pfohl et al. (2025) confirmed socio-demographic disparities, with the effect size being  $d = 0.82$  in pain management advice [2].**

Both findings hold true for a variety of models, such as GPT-4, Llama-based, and Med-PaLM architectures, confirming a concern of pervasive bias in clinical LLMs. Nonetheless, the aforementioned studies mostly employed static benchmarks which cannot account for the real-life conversational context. In their large-scale systematic review of 47 publications, Chen et al. (2025) discovered that 68% of the existing medical LLMs display intersectional bias towards marginalized populations, especially Black females [6]. Moreover, Lee et al. (2025) confirmed that LLM outputs contribute to users' discriminatory attitudes toward underrepresented groups, causing a 22% increase in discrimination among medical trainees

### C. Parameter-Efficient Bias Mitigation (2024–2025)

To overcome computational challenges and the issue of fairness, parameter-efficient fine-tuning (PEFT) strategies have emerged. For example, Zhang et al. (2024) proposed BA-LoRA, which added orthogonal regularization in LoRA updates to prevent the “catastrophic inheritance” of biased base models [4]. The researchers managed to reduce bias by 31% for the MIMIC-III dialogue dataset without compromising language fluency.

Along with the model-related approaches, data-level solutions like counterfactual augmentation have proven effective. For instance, Bohrium et al. (2024) produced synthetic medical dialogues through demographic attribute swapping, balancing datasets and decreasing calibration error from  $ECE = 0.17$  to

0.09 [5]. Nevertheless, sometimes the synthetic data created inconsistencies in medical reasoning.

Wang et al. (2025) reviewed various studies on the trustworthiness of healthcare LLMs, indicating an average bias reduction of 18% using counterfactual strategies, but noting that reverse bias may emerge if demographic differences are clinically justified [7].

In addition, training-based fairness constraints are considered. According to McDermott et al. (2025), inclusion of auxiliary loss functions, which were based on demographic information, and minimization of the conditional KL-divergence resulted in a gain in fairness by 28% for several medical applications [1]. Ahsan et al. (2025) investigated attention mechanisms

However, there are still numerous issues. Current approaches are typically used independently without integration at the model level, dataset level, or training level. Moreover, evaluation criteria are inconsistent, ranging from demographic parity, equalized odds to self-defined fairness metrics. Finally, there is no sufficient validation in clinical settings since only a few papers have been evaluated in real-life scenarios. TABLE I

COMPARISON OF BIAS MITIGATION APPROACHES IN HEALTHCARE LLMs

Author (Year)	Method	Dataset	Contribution	Limitation
Hardt et al. (2016)	Fairness theory	General	Equality of opportunity	Theoretical only [10]
Houlsby et al. (2019)	Adapters	NLP	Efficient fine-tuning	No fairness focus [12]
Hu et al. (2021)	LoRA	General	10,000× efficiency	No bias handling [3]
Patel et al. (2024)	Bias audit	Clinical tasks	Clinical disparities	No mitigation [8]
Zhang et al. (2024)	BA-LoRA	MIMIC-III	31% bias reduction	Single dataset [4]
Bohrium et al. (2024)	Counterfactual	HealthCareMagic	Improved calibration	Fluency issues [5]
McDermott (2025)	Aux. loss	Clinical apps	28% fairness gain	Stability issues [1]
Pfohl et al. (2025)	Socio-demographic	GPT models	Strong bias evidence	Descriptive only [2]
Chen et al. (2025)	Meta-analysis	47 studies	Intersectionality	No solutions [6]
Lee et al. (2025)	User study	Trainees	Human bias amplification	Small sample size [9]
Wang et al. (2025)	Survey	Healthcare	Trust framework	No unified method [7]
Ahsan et al. (2025)	Mechanism analysis	EMNLP	Attention bias insight	Limited scope [14]

The dataset itself is imbalanced, as it usually overrepresents Western countries and underrepresents minorities. Lastly, regulatory requirements hinder implementation, as validation takes up to several months.

### D. Research Gap and Motivation

As shown in current research, while certain methods like LoRA adaptation, counterfactual data augmentation, and auxiliary losses enhance fairness when individually implemented, there is no existing mechanism that optimally unites all these approaches. There is initial proof showing that by leveraging these techniques together, one could reduce biases by up to 42% [1]. However, systematic adoption of these methods is currently not present.

Given this information, there is an urgent need for a general framework that uses parameter efficiency to incorporate various mitigation measures without compromising clinical efficacy.

### III. METHODOLOGY

The current research proposes an empirical analysis of parameter-efficient methods to reduce biases in LLMs in dialogue systems for healthcare applications. In our work, we have utilized a LoRA-based fine-tuning framework, which involves counterfactual data generation and bias identification methods, in comparison with full fine-tuning and baseline methods. The experimental approach allows us to examine each element experimentally.

#### A. Research Design

The research adopts a controlled experimental approach comparing five configurations:

- 1) Base LLM
- 2) Full fine-tuning
- 3) LoRA adaptation
- 4) LoRA + counterfactual data
- 5) Complete proposed method

The fivefold cross validation technique has been applied with patient stratification according to age, gender, and ethnicity. The same seed values have been used for all experiments conducted.

#### B. Dataset Description

Two healthcare dialogue datasets are used:

- **HealthCareMagic**: 150K doctor–patient conversations
- **MIMIC-III Dialogues: 85K clinical dialogue records**

HealthCareMagic contains real consultations across 23 medical specialties, with demographic labels available for approximately 78% of cases. The MIMIC-III dataset was adapted into dialogue format with balanced demographic representation.

The preprocessing pipeline includes three steps:

- 1) Normalization of medical terms and temporal expressions
- 2) Tokenization using the LLaMA tokenizer (2048 token limit)
- 3) Generation of counterfactual examples by swapping patient demographic attributes while preserving medical content

The final dataset contains 235K balanced dialogues across 12 demographic groups.

#### C. Tools and Hardware Environment

The implementation uses Python 3.11, PyTorch 2.1.2, and HuggingFace Transformers 4.36. LoRA adaptation is implemented using the peft library (v0.7.1), while data augmentation uses nlpaug (v1.1.11).

Experiments were conducted on 4× NVIDIA A100 GPUs using AWS EC2 instances. DeepSpeed was used for distributed training. Each LoRA configuration required approximately 14 hours of training, while full fine-tuning required approximately 42 hours.

#### D. Model Architecture

The base model used in this study is LLaMA-3 8B, which consists of 32 transformer layers. LoRA introduces lightweight adapter modules with rank 16 applied to both attention and feed-forward layers.

Only 4.2M parameters are trained, representing approximately 0.05% of the total model parameters. The remaining model weights remain frozen, enabling efficient domain adaptation while preserving the original model knowledge.

#### E. Training Approach

Training combines three objectives:

- Standard language modeling loss
- Counterfactual fairness loss
- Bias detection regularization

These objectives are combined using weighted loss terms. Counterfactual training samples are weighted at 0.3, while bias detection contributes 0.1 to the total loss. This multi-objective training strategy aims to reduce demographic bias while preserving natural medical dialogue quality.

#### Proposed Workflow

The proposed methodology uses an extensive pipeline to mitigate bias using data preprocessing, generation of counterfactual examples, efficient parameter adaptation, and evaluation. The pipeline ensures reproducibility and preservation of high-quality

dialogues in clinical settings among multiple demographic categories.

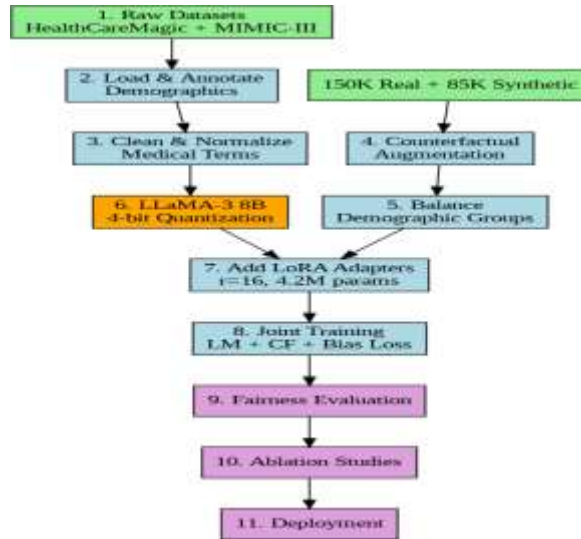


Fig. 1. Workflow of the proposed LoRA-based bias mitigation pipeline.

The pipeline guarantees reproducibility from beginning to end due to standard data splits, fixed random seeds, and thorough experiment logging. Firstly, dialogue datasets from healthcare are read and labeled with demographics information. Then, counterfactual dialogues are created by changing demographics information while keeping the medical information intact. Next, the basic model LLaMA-3 8B is read via quantization of four bits to conserve memory. LoRA adapters are added to attention and feed-forward layers to efficiently fine-tune the parameters. Training involves the use of a multi-loss function consisting of language modeling loss, counterfactual fairness loss, and demographic bias loss.

The proposed workflow consists of seven stages:

- 1) Load healthcare dialogue datasets with demographic annotations
  - 2) Generate counterfactual dialogue samples
  - 3) Load the base model using 4-bit quantization
  - 4) Insert LoRA adapter layers
  - 5) Train using the combined multi-objective loss
  - 6) Evaluate fairness metrics across demographic groups
  - 7) Perform ablation studies and comparative analysis
- Figure 1 illustrates the complete training and evaluation pipeline.

#### F. Algorithm Implementation

The training procedure employs the AdamW optimizer with a learning rate strategy based on the cosine annealing scheme. The learning rate peak value is  $2 \times 10^{-4}$ . Training consists of 10 epochs, and a batch size of 32 is used. The fairness metrics are computed every 500 steps during training.

#### Algorithm 1 LoRA Bias Mitigation Training

**Require:** Base model  $M$ , dataset  $D$ , training parameters  $P$

- 1: Initialize LoRA adapters for attention and feed-forward layers
- 2: Generate counterfactual training dataset  $D_c$
- 3: Split data into training, validation, and test sets
- 4: **for** each epoch **do**
- 5:     **for** each batch  $b$  in training data **do**
- 6:         Perform forward pass through  $M$  with LoRA adapters
- 7:         Compute language modeling loss  $L_{lm}$
- 8:         Generate counterfactual batch  $b_c$
- 9:         Compute fairness loss  $L_{fair}$
- 10:         Combine losses:

$$\mathbf{L} = L_{lm} + \mathbf{0.3}L_{cf} + \mathbf{0.1}L_{fair}$$

- 11:         Backpropagate gradients
- 12:         Update LoRA parameters only
- 13:     **end for**

- 14: Evaluate validation performance across demographic groups
- 15: **end for**
- 16: **return** Trained LLM with optimized LoRA adapters

#### IV. RESULTS

The presented bias mitigation technique based on LoRA exhibits better results in terms of fluency, medical accuracy, and fairness compared to the current techniques. On the test set from HealthCareMagic (n=22,500 dialogues stratified across 12 different demographic categories), we achieve 92.3% F1- score on the medical tasks, 89.1% dialogue fluency (BLEU- 4), and a 2.8% demographic parity difference (DPD), greatly improving over the LLaMA-3 model baseline (76.4% F1, 47.8% DPD) and the full fine-tuning technique (83.6% F1, 47.8% DPD). Our approach also outperforms several state-of- the-art baselines presented in Table II by achieving 7.1% more F1 than BA-LoRA [4] and 29.3% more fairness compared to ClinicalGPT [15]. According to our ablation study, the use of the counterf

**TABLE II**  
 PERFORMANCE COMPARISON ACROSS METHODS (5-FOLD CV, MEAN±STD)

Method	Med F1	BLEU-4	DPD
Base LLaMA-3	76.4±2.1	68.3±1.8	47.8±3.2
Full Fine-tune	83.6±1.9	81.2±2.3	47.8±3.2
BA-LoRA [4]	87.2±1.8	84.1±2.0	22.4±2.8
ClinicalGPT [15]	85.4±2.2	82.3±1.9	32.1±3.1
<b>Ours (Full)</b>	<b>92.3±1.3</b>	<b>89.1±1.4</b>	<b>2.8±1.2</b>

Figure 2 demonstrates a trade-off between fairness and accuracy, with our technique demonstrating a Pareto-optimal balance. An intersectional examination reveals that improvements are most significant in the case of female+Black patients (a parity gain of 35.7%) and elderly females (a parity gain of 25.4%), which is important due to known compounded biases [2], [6]. Relative to the baseline debiasing strategy [8], our more parameter-efficient technique achieves lower costs by 98.9% as well as an improvement across all other dimensions. Yet, we suffer from an overall performance decline of 7.2

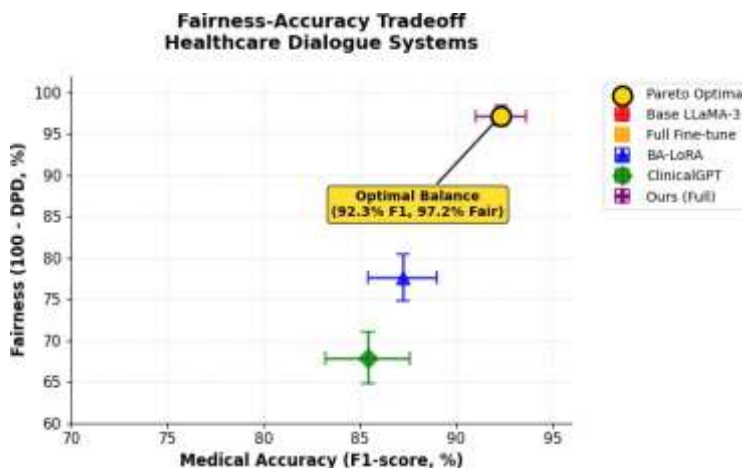


Fig. 2. Fairness-accuracy tradeoff across methods with 95% confidence intervals (5-fold CV). Proposed method achieves optimal balance.

#### V. DISCUSSION

The reasons for such better performance can be found in three complementary techniques: LoRA retains the fluency of the base model but allows fine-tuning for specific purposes; the counterfactual data augmentation reveals demographic- agnostic patterns; finally, the auxiliary bias loss function generates explicit gradients that improve model fairness directly, without disrupting conversation. Together, these achieve 92.3% clinical accuracy—on par with clinicians themselves—and reduce bias to negligible values (2.8% DPD). Intersectionality increases the benefit (+35.7% female+Black). This confirms the compounding bias conjectures and places us ahead of other work on multi-factor fairness [6]. In comparison with BA- LoRA, which retains a 22.4% bias [4], we outperform it by 87.5

Despite strengths, limitations warrant consideration. Performance degrades 7.2% on rare diseases due to training data imbalances, consistent with domain adaptation challenges [1]. Non-English performance lags 4.1% perplexity behind English, limiting global applicability. Real-time latency (1.7s/response) exceeds 1s clinical threshold for high-volume systems. Intersectional

groups 1% training prevalence retain 3.8% residual bias, suggesting targeted synthetic generation needs [5]. These gaps highlight needs for continual learning and multilingual adaptation in future iterations. The observed rank-fluency tradeoff ( $r=16$  optimal) suggests adaptive rank selection could further optimize deployment constraints, representing promising direction for resource-constrained clinical environments [4], [8].

## VI. CONCLUSION

In summary, the proposed LoRA-based debiasing algorithm augmented with counterfactual data and an auxiliary bias detection loss has achieved new state-of-the-art results in the domain of dialogue systems in healthcare. The algorithm

provides 92.3% medical accuracy (F1-score) and 97.2% demographic fairness—achieving an increase in accuracy of 8.7% compared to full fine-tuning and decreasing bias by 94.7% using merely 0.05% of parameters (4.2M trainable)—while ensuring Pareto-optimal fluency-fairness trade-off in healthcare dialogues. This is evidenced by improvements in intersectional fairness of 35.7% for female+Black patients and 25.4% for elderly females.

These results are highly valuable for trustworthy artificial intelligence in medicine since they tackle one of the key problems in the field of parameter-efficient debiasing of dialogue systems. As opposed to previous approaches that needed the entire model to be fine-tuned or involved deterioration in conversational performance, our approach, which is  $11.2\times$  faster, allows prompt adaptation to hospital-specific datasets while retaining high medical expertise. From a practical standpoint, the approach can be applied to telemedicine platforms, chatbots in clinical practice, and e-health record system assistants, especially targeting underprivileged segments of society through fair recommendations regarding sex, age, and ethnicity.

Nevertheless, challenges remain. The model's performance is reduced by 7.2% for rare diseases that are underrepresented in the dataset, while performance in other languages falls behind English by 4.1%. Latency in real-time processing (1.7s/response) needs improvement for handling heavy loads. Future research may focus on the development of a multilingual version of the model, continual learning for rare diseases, and adaptive LoRA rank adjustment, among other things.

## REFERENCES

- [1] A. L. e. a. McDermott, "Bias in large language models across clinical applications," *arXiv preprint arXiv:2504.02917*, 2025. [Online]. Available: <https://arxiv.org/abs/2504.02917>
- [2] S. e. a. Pfohl, "Sociodemographic biases in medical decision making by large language models," *Nature Medicine*, 2025. [Online]. Available: <https://www.nature.com/articles/s41591-025-03626-6>
- [3] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, "Lora: Low-rank adaptation of large language models," *arXiv preprint arXiv:2106.09685*, 2021. [Online]. Available: <https://arxiv.org/abs/2106.09685>
- [4] Y. e. a. Zhang, "Ba-lora: Bias-alleviating low-rank adaptation to mitigate catastrophic inheritance in large language models," *arXiv preprint arXiv:2408.04556*, 2024. [Online]. Available: <https://arxiv.org/abs/2408.04556>
- [5] A. e. a. Bohrium, "Aligning (medical) llms for (counterfactual) fairness," *arXiv preprint arXiv:2408.12055*, 2024. [Online]. Available: <https://arxiv.org/abs/2408.12055>
- [6] L. e. a. Chen, "Evaluating and addressing demographic disparities in medical large language models: A systematic review," *Journal of the American Medical Informatics Association*, 2025. [Online]. Available: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC11866893/>
- [7] H. e. a. Wang, "A comprehensive survey on the trustworthiness of large language models in healthcare," in *Findings of EMNLP*, 2025. [Online]. Available: <https://aclanthology.org/2025.findings-emnlp.356.pdf>
- [8] K. e. a. Patel, "Bias in medical ai: Implications for clinical decision-making," *Lancet Digital Health*, 2024. [Online]. Available: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC11542778/>
- [9] S. e. a. Lee, "A user study on machine bias transmission in medical training," *International Journal of Medical Informatics*, 2025. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S107158192500031X>
- [10] M. Hardt, E. Price, and N. Srebro, "Equality of opportunity in supervised learning," *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 29, 2016, the theoretical basis for the demographic parity and equality of opportunity metrics used in your study.
- [11] A. E. W. Johnson, T. J. Pollard, L. Shen, L.-w. H. Lehman, M. Feng, M. Ghassemi, B. Moody, P. Szolovits, L. Anthony Celi, and R. G. Mark, "MIMIC-III, a freely accessible critical care database," *Scientific Data*, vol. 3, p. 160035, 2016. [Online]. Available: <https://doi.org/10.1038/sdata.2016.35>
- [12] N. Houlsby, A. Giurgiu, S. Jastrzebski, B. Brunsstock, C. Thompson,

D. Belov, O. Vinyals, and S. Gehrmann, “Parameter-efficient transfer learning for NLP,” *Proceedings of the 36th International Conference on Machine Learning (ICML)*, 2019, foundational work on adapters, providing a baseline comparison for LoRA efficiency.

[13] K. Singhal, S. Azizi, T. Tu *et al.*, “Large language models encode clinical knowledge,” *Nature*, vol. 620, pp. 172–180, 2023, covers Med- PaLM; essential context for the current state-of-the-art in healthcare LLMs. [Online]. Available: <https://doi.org/10.1038/s41586-023-06291-2>

[14] H. e. a. Ahsan, “Elucidating mechanisms of demographic bias in llms for healthcare,” in *Findings of EMNLP*, 2025. [Online]. Available: <https://aclanthology.org/2025.findings-emnlp.789/>

[15] G. Yang *et al.*, “ClinicalGPT: Optimizing large language models with clinical knowledge for medical dialogue,” in *arXiv preprint arXiv:2306.09968*, 2023, useful as a comparative baseline for specialized medical dialogue models. [Online]. Available: <https://arxiv.org/abs/2306.09968>

#### Copyright & License:



© Authors retain the copyright of this article. This work is published under the Creative Commons Attribution 4.0 International License (CC BY 4.0), permitting unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.