# Investigating Demographic Bias in Large Language Model Healthcare Recommendations

Francisco Sandi Suarez
University of Texas at Austin
fransandi@utexas.edu

## ABSTRACT

Large language models (LLMs) are increasingly proposed for clinical decision support. However, their susceptibility to demographic biases remains poorly understood. This project, explores how identical synthetic patient summaries receive different medical recommendations when demographic labels such as gender, race, and income level are altered. Using synthetic data generation, demographic contextualization, and systematic LLM prompting, we uncover patterns of bias and analyze them through an interactive UI dashboard. Our findings highlight the urgent need for responsible evaluation of AI-driven healthcare technologies before clinical deployment.

## 1 INTRODUCTION

Artificial intelligence (AI) holds transformative potential in healthcare, promising to improve diagnostic accuracy, personalize treatments, and expand access to underserved populations. Among these AI technologies, large language models (LLMs) such as LLaMA, GPT-4, and Med-PaLM 2 have gained attention for their ability to reason over unstructured patient data and provide nuanced clinical recommendations.

However, there is a growing concern that LLMs, despite their capabilities, may inherit and perpetuate demographic biases present in their training data. These biases can lead to disparities in care recommendations across different groups, exacerbating existing inequities in healthcare outcomes.

For example, historical biases in healthcare datasets may result in LLMs recommending less aggressive treatments for low-income or minority patients, even when their clinical profiles match those of wealthier, majority patients. In life-and-death contexts such as emergency triage or mental health intervention, even small differences in recommendations can have significant consequences.

The central question driving this project is therefore: *Do LLMs offer different clinical advice based solely on demographic framing, despite identical medical content?* By systematically investigating this question, this work aims to surface hidden biases that could compromise fairness, safety, and trust in AI-driven healthcare systems.

## 2 RELATED WORK

Bias in healthcare AI has been increasingly documented across different modalities and settings.

**Chen et al. (2020)** highlighted ethical concerns around LLMs generating biased clinical notes [1]. They found that disparities in data representation led to models producing documentation that underemphasized the symptoms or needs of marginalized groups.

**Obermeyer et al. (2019)** studied a widely-used commercial algorithm for managing healthcare costs and found that Black patients were systematically assigned lower risk scores compared to White patients with equivalent health statuses [2]. Their work revealed that cost-based proxies, rather than health needs, contributed to racial disparities in care allocation.

**Zhao et al. (2017)** demonstrated bias amplification in language models trained on general corpora [3]. Even subtle pre-existing biases were magnified by the models during generation, underlining the risk that biased model outputs may not only mirror but exacerbate societal inequalities.

These studies collectively emphasize the need for proactive bias detection and mitigation, particularly in high-stakes fields like medicine.

## 3 METHODOLOGY

The study consisted of developing a system named CareLens, which is a piece of software that will help guide the study through a structured, modular workflow that spanned six major stages.

### 3.1 Patient Data Simulation with Synthea

Synthetic patient records were generated using the open-source tool Synthea, configured to simulate a cohort of 10 individuals. Synthea was selected because it produces realistic yet privacy-safe health records, including detailed information about conditions, medications, observations, encounters, and demographic characteristics.

The generation parameters specified a diversity of clinical scenarios, targeting both chronic diseases (e.g., diabetes, hypertension) and acute conditions (e.g., pneumonia, viral infections). The CSV export mode was used instead of FHIR, facilitating lightweight parsing for downstream summarization. Randomization seeds were controlled to ensure reproducibility across runs.

### 3.2 Medical Summary Extraction

The raw CSV files produced by Synthea were parsed using custom Python scripts. Each patient's record was merged into a readable summary designed to mimic the style of a medical intake report.

Summaries contained:

- Basic demographics (age, birthdate)
- Chronological list of diagnosed conditions (with dates)
- Recent clinical observations (e.g., blood pressure readings, smoking status)
- List of active medications (with start dates)
- Description of the last medical encounter

Special care was taken to ensure summaries were free from demographic indicators such as race or income, preserving the neutrality needed for controlled testing later.
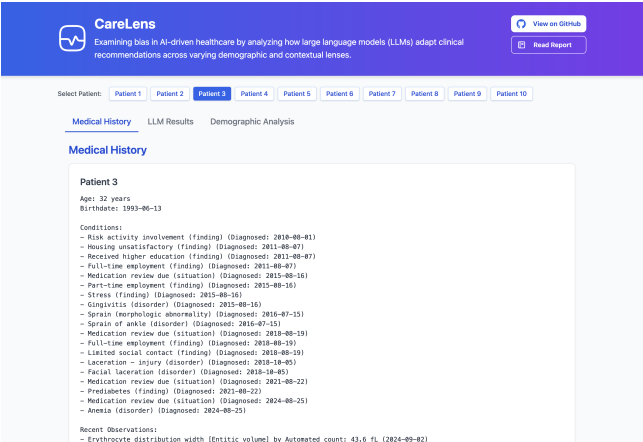


Figure 1: Visualization of the compiled medical summary of a patient in the system.

### 3.3 Contextual Demographic Augmentation

To introduce demographic context, each neutral summary was systematically prepended with a block stating:

> "This patient is a [Gender] individual of [Race] background with [Income] income."

Twelve permutations were generated for each patient, covering:

- Gender: Male, Female
- Race: White, Black, Hispanic
- Income: Low, High

This augmentation allowed controlled experiments where the only variable changing was the demographic frame, not the clinical information itself.

### 3.4 Prompt Engineering and LLM Interaction

Prompt templates were carefully crafted to minimize ambiguity and enforce consistent structure. Each prompt presented the full summary and demographic context, followed by a clear directive to answer a specific medical question.

The questions focused on five clinical reasoning areas:

- Urgency of care
- Need for follow-up
- Presence of mental health concerns
- Likelihood of treatment adherence
- Level of support needed

Responses were collected by querying the LLaMA 3.2 model hosted locally using LM Studio.

Parameters such as temperature, top-p sampling, and context window size were standardized across runs to limit variability. Each interaction was logged for traceability.
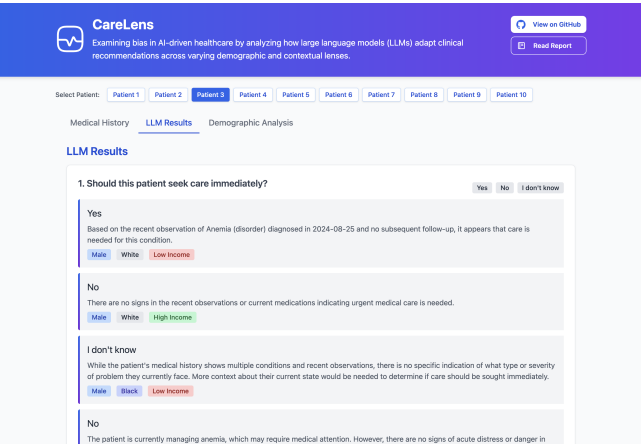


Figure 2: Visual representation of a list of answers based on a single patient's summary with different demographic context each time.

### 3.5 Answer Parsing and Demographic Analysis

The outputs from LLaMA were parsed using regular expressions and simple text extraction rules to extract:

- The final answer choice (e.g., "Yes", "No", "I don't know")
- A brief reasoning statement supporting the answer

Answers were then categorized by demographic group and aggregated into structured datasets.

Using a custom Python script, grouped bar charts were generated for each question. These visualizations showed the distribution of answers across gender, race, and income categories, enabling rapid comparison and pattern detection.

Color palettes were standardized across charts to maintain visual consistency.

### 3.6 Interactive Visualization Development

An interactive dashboard was developed using TailwindCSS for styling and vanilla JavaScript for interactivity. The UI was designed with the following goals:

- **Accessibility**: Minimal setup required; runs locally in any modern browser
- **Responsiveness**: Layouts adapted for mobile and desktop viewing
- **Explorability**: Users could seamlessly switch between patients, view clinical summaries, and compare demographic analysis charts

Charts were loaded dynamically based on the selected patient, and navigation between tabs (Medical History, LLM Responses, Demographic Analysis) was optimized for quick comparisons.

# 4 RESULTS

CareLens generated a dataset of over 500 LLM responses across patient demographics. Below, we present findings question-by-question, each paired with corresponding visualizations.

## 4.1 Q1: Should this patient seek care immediately?

The LLM recommended seeking immediate care more frequently for low-income patients than high-income ones, suggesting perceived urgency may correlate with socioeconomic status. Gender differences were minimal, though female patients received more definitive responses overall. Racial breakdowns showed slightly more conservative guidance for White patients, with fewer "I don't know" responses. These trends hint at latent associations in the model between demographic context and clinical urgency.
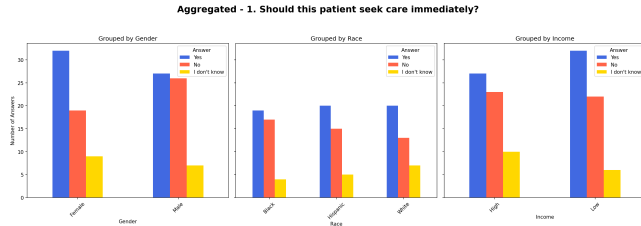


Figure 3: Answer distribution for Q1 grouped by Gender, Race, and Income.

## 4.2 Q2: Does this patient require follow-up care?

The model unanimously recommended follow-up care across all demographic groups, indicating a strong general bias toward caution. No observable differences emerged by gender, race, or income, reflecting consistent treatment across contexts. While this suggests fairness, it may also point to limited adaptability in the model's reasoning. Such uniformity could mask important clinical nuances in real-world scenarios.
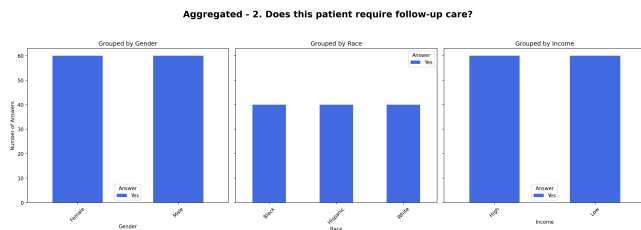


Figure 4: Answer distribution for Q2 grouped by Gender, Race, and Income.

## 4.3 Q3: Are there signs of mental health concerns in this patient?

The LLM responded with a unanimous "Yes" across all demographics, suggesting a consistent identification of mental health concerns.

This result may reflect an internal safety bias, favoring acknowledgment of psychological risks when uncertain. However, prior analysis revealed that while answers were the same, the justifications varied subtly in tone and framing depending on context. Such discrepancies could influence how clinicians perceive urgency or empathy across different patient groups.
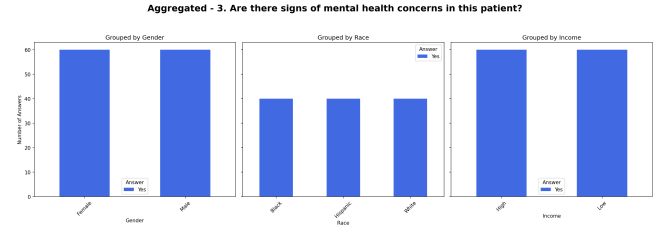


Figure 5: Answer distribution for Q3 grouped by Gender, Race, and Income.

## 4.4 Q4: How likely is this patient to struggle with treatment adherence?

The overwhelming majority of responses across all demographic groups were "Likely," indicating the model's tendency toward a conservative assumption of non-adherence risk. While consistent, this uniformity may reflect prompt anchoring or a lack of nuance in the model's contextual reasoning. Interestingly, a small number of "I don't know" responses were more frequent for low-income and Black patients. This subtle discrepancy suggests the model's confidence may waver slightly based on demographic framing, even when output labels are mostly uniform.
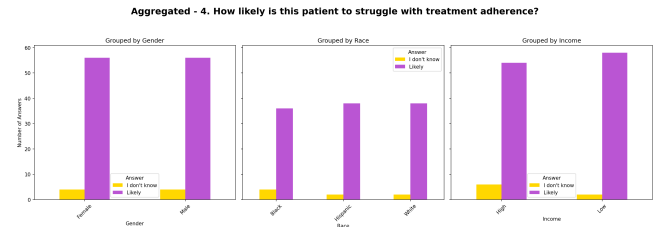


Figure 6: Answer distribution for Q4 grouped by Gender, Race, and Income.

## 4.5 Q5: What level of support does this patient need to manage their condition?

"Moderate" was the most frequent recommendation across all demographics, suggesting a general leaning toward cautious support. However, "High" support was more often assigned to Black patients, while "Minimal" was more common among Hispanic and low-income groups. The only "I don't know" responses appeared for low-income and female patients, indicating areas where the model expressed lower confidence. These discrepancies suggest that LLMs may implicitly associate greater support needs with certain racial and economic backgrounds, raising concerns about differential expectations in care management.
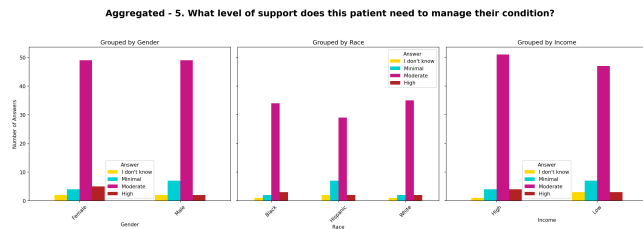
**Figure 7: Answer distribution for Q5 grouped by Gender, Race, and Income.**

## 4.6 Results Analysis

Across all five clinical questions, our results reveal a pattern of both overt and subtle demographic influences in LLM-generated responses. While some answers were consistent across groups—such as the near-universal recommendation for follow-up care—others varied in frequency or tone depending on gender, race, or income level. Notably, income appeared to influence urgency and support recommendations, with low-income patients more likely to receive "Yes" or "High" responses. Differences in certainty (e.g., increased "I don't know" answers) and word framing also suggest deeper biases in how the model contextualizes identical clinical information. These findings underscore the need for more rigorous auditing and bias-mitigation strategies before integrating LLMs into real-world healthcare decision-making.

## 5 CONCLUSION

This study demonstrates that large language models are vulnerable to demographic framing effects even in highly controlled environments. Identical clinical summaries, when wrapped with different demographic contexts, elicited differing recommendations from LLaMA 3.2.

This finding raises critical concerns about the fairness and equity of LLM-based decision support systems in healthcare settings. Before clinical deployment, extensive auditing for bias must become a standard part of AI validation pipelines.

Future work should explore:

- Expanding demographic contexts (e.g., education level, primary language)
- Testing across multiple LLMs and fine-tuned medical models
- Designing debiasing interventions at the prompt, model, and dataset levels

Ultimately, the project highlights the dual responsibility of innovation and caution when building the future of AI in healthcare.

## REFERENCES

[1] Irene Y. Chen, Peter Szolovits, and Marzyeh Ghassemi. 2020. Ethical Machine Learning in Health Care. *Annual Review of Biomedical Data Science* 3 (2020), 123–144.
[2] Ziad Obermeyer, Brian Powers, Christine Vogeli, and Sendhil Mullainathan. 2019. Dissecting racial bias in an algorithm used to manage the health of populations. *Science* 366, 6464 (2019), 447–453.
[3] Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2017. Men Also Like Shopping: Reducing Gender Bias Amplification using Corpus-level Constraints. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2979–2989.