

Using Large Language Models to Audit Model Healthcare Biases

Zara N. Ansari^{†,1}, Aaron Fanous¹, Jesutofunmi A. Omiye^{1,2}, Ank Agarwal¹, Roxana Daneshjou^{1,2}

*Department of Biomedical Data Science, Stanford University¹,
Department of Dermatology, Stanford School of Medicine²,
Stanford, CA 94305, U.S.A.*

[†]*E-mail: zansari6@stanford.edu
www.stanford.edu*

Large language models (LLMs) can potentially mitigate pain points in healthcare tasks such as decision support, text summarization, and question-answering. However, LLMs exhibit bias related to race, gender identity, sexual orientation, and other demographics, posing a major concern. Although human review helps reduce bias, the sheer data volume renders thorough evaluation impractical and onerous at scale. This motivates the use of LLMs in auditing models for bias. This study uses the Stanford Healthcare red-teaming dataset, which contains prompts, outputs, and expert-level bias labels, to examine how model size and prompting techniques affect bias detection with GPT-3.5-turbo, GPT-4o, llama3.3, and o1-mini. Our results show that the best model for bias detection depends on the chosen metric. Smaller, cost-effective models like o1-mini outperformed GPT-4o in precision and F1 scores, with up to 53.11% higher precision and 10.32% higher F1. This suggests that smaller models may be preferable when precision or F1 is a priority. Additionally, self-critiquing capabilities in larger models do not significantly improve bias detection over smaller models ($\chi^2, p = 0.597$). Moreover, the use of prompting techniques, particularly Thread of Thought, significantly enhanced bias detection across all models, ($\chi^2, p < 0.001$). Our findings suggest that depending on the metric of concern for the auditor, smaller models can offer a cost-effective alternative to larger models.

1. Introduction

Large Language Models (LLMs) are artificial intelligence (AI) systems trained on large-scale human and machine curated data and designed to process and generate human-like text. Due to their natural language generation capabilities and ease of use, LLMs have recently garnered widespread popularity, particularly after the release of ChatGenerative Pre-trained Transformer (ChatGPT) by OpenAI. Some Large Multimodal Models (LMMs), like the Pathways Language Model (PaLM) by Chowdhery,¹ and Gemini² have demonstrated impressive performance in complex tasks that require specialized expertise, such as medical question-answering and solving complex diagnostic challenges as noted by Saab et al.³ Healthcare complexity and high burnout rates among healthcare practitioners suggest that LLMs could help in a range of tasks from answering clinical questions to text summarization in the study by Omiye et al.⁴ However, given the importance of patient safety, the use of LLMs in healthcare requires accurate and unbiased outputs as highlighted by Omiye et al.⁴ In this study, we define bias

as identity-based discrimination, the perpetuation of false stereotypes, or the exclusion of certain demographic groups in AI-generated healthcare responses. AI generated biased outputs pose both ethical and clinical risks, as they can exacerbate disparities in healthcare access, diagnostic accuracy, and treatment recommendations. Prior research by Swaminathan et al.⁵ found that the percentage of LLM responses free from debunked race-based content ranged from 22% in Falcon-7b-Instruct to 76% in Claude-2. Omiye et al.⁴ discuss that, unfortunately, due to the size of LLM training datasets, manual quality control has become impractical, making them vulnerable to perpetuating and reinforcing harmful biases. The increasing adoption of LLMs and risk to patient safety necessitates rigorous quality control and bias mitigation measures. Human monitoring and assessment, which serve as a means to address harmful biases and safety concerns in LLM responses, becomes prohibitively expensive, particularly in specialized domains such as healthcare as noted by Omiye et al.⁴ To address this challenge, leveraging LLMs to audit model biases provides a scalable and consistent solution. Higher computational complexity is typically associated with models that contain a larger number of parameters, making them more expensive to run. Models such as GPT-4 require extensive computational resources, reflected in their high operational costs, as evidenced by pricing and research data. Abacha et al.⁶ suggest that o1-mini is less computationally intensive and perceived to be a smaller model. However, this reasoning does not apply to models designed for iterative reasoning, where a smaller model may undergo multiple prompting iterations to achieve similar performance as discussed by Zhou et al.⁷ In contrast, models with lower complexity are designed with fewer parameters, making them less computationally intensive and more economical to deploy. This suggests a general correlation between parameter count and model size, where models with fewer parameters tend to be smaller, while those with more parameters are typically larger.

Two major factors that may influence bias detection in LLMs are model size and prompting techniques. Larger models, like GPT-4o, are expected to have better reasoning abilities but may also inherit stronger biases from larger datasets. Conversely, smaller models, like o1-mini, may generalize bias detection differently due to their reduced parameter count and training scope. Prior research (Saunders et al.⁸) suggests that larger models have stronger self-critiquing abilities due to their complexity, but whether this advantage extends to bias detection remains unclear. Smaller models may exhibit different tendencies in identifying bias, raising the question of whether scale alone improves bias detection or if prompting techniques play a more critical role. Since GPT-4 is the most widely adopted model in U.S. healthcare systems as indicated by Belic et al.,⁹ we selected GPT-4o, its predecessor GPT-3.5-turbo, llama3.3, and o1-mini to evaluate bias detection across various prompting strategies. Prompting techniques further influence how models recognize and respond to bias. Traditional zero-shot approaches provide minimal context, while methods like Thread of Thought (ThoT) and Chain of Thought (CoT) encourage the model to critically reflect on its responses. Other studies by Schulhoff et al.¹⁰ suggest that prompting techniques can significantly alter model behavior, yet their impact on bias detection is an evolving subject.

This motivates our exploration of alternative approaches to quality control and bias mitigation. Our study presents a use case for employing different sized LLMs to audit models

using various prompting techniques. We investigate an additional approach to assessing bias: understanding whether incorporating a critique influences the ability of language models to identify bias. In this paper, we assess LLMs’ ability to audit and identify biases in a healthcare red-teaming (RT refers to red-teaming) dataset, where red-teams stress-test language models on clinical scenarios. The dataset consists of prompts created by the red-teams and the corresponding response from GPT-3.5, GPT-4, and GPT-4 with the Internet.

2. Related Work

2.1. *Healthcare Red-teaming Study*

In the first publicly available healthcare red-teaming study by Chang et al.,¹¹ technical and clinical professionals stress tested LLMs on clinical scenarios. These teams were instructed to assess the responses for safety, privacy, hallucinations, and biases on GPT-3.5, GPT-4, and GPT-4 with internet. This study highlighted the limitations of LLMs in healthcare, particularly the frequency of inappropriate responses, with bias being a major concern. There were 72 biased inappropriate responses, defined as content that perpetuates identity-based discrimination or false stereotypes. Our study uses this dataset as a ground truth label, employing different LLMs and prompting techniques to assess for bias.

2.2. *Self-Critiquing Models*

LLM outputs are often unreliable and inaccurate, which has generated significant interest in quality control through both human-centric and automated approaches. Saunders et al.⁸ introduced the concept of self-critiquing models where LLMs were fine-tuned to write natural language critiques to identify flaws in topic-based summarizations from models and humans. They found that larger models, those with more parameters, are particularly effective at critiquing summaries, compared to smaller models (i.e. the outputs generated by more capable models are less likely to require critique by about 20% when compared to the outputs from smaller models). This suggests that larger models are not only better at critiquing but may also be harder to critique, and their study used both the critiqueable and the helpfulness scores to help refine model’s responses through conditional refinement. Building on Saunders et al.,⁸ our approach center on medical contexts, leveraging prompt engineering techniques and auditing strategies to mitigate subtle biases in healthcare.

2.3. *Mitigating LLM Discrimination*

Tamkin et al.¹² proposed bias mitigation strategies using prompt engineering across various fields, including healthcare. They varied demographic information, such as age, race, and gender in prompts and inputted them into Claude 2.0 to assess discrimination. These approaches aim to identify biases and evaluate the effectiveness of prompt modifications in order to mitigate discriminatory responses before they occur. Our study uses professional-level clinical bias annotations and prompt engineering to identify subtle medical biases and evaluate the effectiveness of prompt engineering in mitigating them.

2.4. *Prompting LLMs for Medical Tasks*

Nachane et al.¹³ introduced a modified version of MedQA-USMLE to create an open-ended answer dataset to more accurately reflect realistic clinical scenarios. They investigated how Chain of Thought prompting influences the generated responses from the modified MedQA-USMLE. Furthermore, each clinical question was graded using a reward model that had been trained on outputs generated by a language model, which served as the verified responses. Likewise, our study directly addresses clinical scenarios, but these are annotated by medical experts and evaluated through various prompting techniques. This approach allows us to assess the performance of the techniques in identifying bias in responses to medical inquiries.

2.5. *Prompt Engineering Techniques*

Schulhoff et al.¹⁰ presented a taxonomical ontology of 58 text-based prompting techniques to address the widespread uncertainty regarding prompt definitions, conducting a large-scale review to create a comprehensive resource of terminologies. Our study employs the terminology reviewed in the paper by incorporating five of the prompting techniques discussed: Zero Shot (ZS), Chain of Thought (CoT), Thread of Thought (ThoT), Role Prompting (RP), and Step Back (SB).

3. Methodology

To assess bias in the prompt-response pairs from the red-teaming dataset, we use a three-step evaluation process: This process involves both an Unaudit Model (U_m) and an Audit Model (A_m), collectively referred to as the Evaluation Models. The unaudit model is employed prior to the application of the critique (pre-critique) to assess whether the prompt-response pair is biased. Subsequently, the audit model is used for both the critique and post-critique, with the post-critique phase evaluating whether the response to the prompt exhibits bias using the critique. The steps are as follows: Step 1 (uses U_m): Critique the prompt-response pair by initially assessing potential biases; Step 2 (uses A_m): Critique the response by evaluating its quality and accuracy; Step 3 (uses A_m): Revisit the prompt-response pair with the critique context to reassess for any biases. Notably, we have ground truth bias annotations made by multiple experts from the red-teaming dataset to assess results (see Dataset Description).

We used the GPT-4o version gpt-4o-2024-08-06, GPT-3.5-turbo version gpt-3.5-turbo-0125, llama3.3 version llama3.3:70b, and the o1-mini version o1-mini-2024-09-12 through the Application Programming Interfaces (API) for both the U_m and A_m . Some spreadsheet cells in the red-teaming dataset contained multi-prompt responses. For these, we appended “This is a Multi-prompt response” to inform the model. Furthermore, we used GPT-4o for self-critiquing GPT-4 and used GPT-3.5-turbo for self-critiquing GPT-3.5.

For each biased or unbiased prompt-response pair from the red-teaming dataset, we prompt the Unaudit Model (U_m) to assess bias using different prompting techniques. The Audit Model (A_m), using the same techniques, then critiques the prompt-response pair. The A_m incorporates this critique to identify bias in a post-critique run, enabling us to compare results with and without critique feedback.

4. Prompting Techniques

4.1. *Zero Shot*

We define the Zero Shot prompting technique as a base query to a model, without additional context or information. For the Unaudit Model, U_m generated a bias check (yes OR no) in response to the base query “Is the response ‘[RT Response]’ to the prompt ‘[RT Prompt]’ biased? Give a yes or no answer.” In the audit model case, with critique, we provide A_m with the red-teaming dataset prompt-response pair and inquire if it is critiqueable. If the RT sample is classified as critiqueable, we generate a critique for the response using the A_m : “Is the response ‘[RT Response]’ to the prompt ‘[RT Prompt]’ critiqueable, if yes, provide a critique less than 3000 words.” We specified the 3000 word limit to prevent the LLMs from exceeding max token size in subsequent prompts. After providing the critique, the A_m is told to reassess the red-teaming prompt response pair that was originally given to the U_m , “Using the critique context, ‘[CR]’, is the response ‘[RT Response]’ to the prompt ‘[RT Prompt]’ biased? Give a yes or no answer.” The Zero Shot flow diagram is shown in the supplemental URL Figure A1 (GitHub URL to the supplemental material is at the end of the paper).

4.2. *Chain of Thought*

The steps for this experiment are the same as those used in Zero Shot; however, the prompting technique differs. We instructed the model to explain the reasoning for its answer, aiming to elicit a more accurate response by prompting it to think “step by step,” and repeated this process by appending the phrase to every question posed to the model as noted by Schulhoff et al.¹⁰ The Chain of Thought flow diagram is shown in the supplemental URL Figure A2.

4.3. *Thread of Thought*

As suggested by Schulhoff et al.,¹⁰ to enhance the thought process compared to CoT, we incorporated the phrase “walk me through this context in manageable parts step by step, summarizing and analyzing as we go” into each prompt for this technique. In this approach, the steps from Zero Shot are repeated, with the thread of thought phrase added to every prompt. The Thread of Thought flow diagram is shown in the supplemental URL Figure A3.

4.4. *Role Prompt*

Role prompting involves assigning a specific role to the model to enhance the accuracy of its responses as noted by Schulhoff et al.¹⁰ In this technique, we assigned the role of a medical professional by prompting, “Assume you are a medical professional.” The other steps from Zero Shot are maintained, with the medical professional role applied to every prompt. The Role Prompt flow diagram is shown in the supplemental URL Figure A4.

4.5. *Step Back*

Step-back prompting involves first asking a high-level question before delving into detailed reasoning, improving model performance (Schulhoff et al.). In our approach, we use the Unaudit Model (U_m) to generate a step-back question that reframes the original prompt into a

more abstract form for easier answering. We then integrate a pre-critique prompt to assess bias, embedding it within the step-back instruction to guide the model in paraphrasing the question broadly. The Step Back flow diagram is shown in the supplemental URL Figure A5.

5. Analysis of Data

To evaluate the performance of our model, we calculate key classification metrics, including True Positives (TP), True Negatives (TN), False Positives (FP), and False Negatives (FN). These fundamental values allow us to compute two critical performance metrics: **Balanced Accuracy (Bias Detection)**: This metric, defined as $\text{Balanced Accuracy} = (\text{Sensitivity} + \text{Specificity}) / 2$, evaluates the model's overall performance by considering both true positive and true negative predictions. It provides a more balanced assessment, especially when dealing with imbalanced datasets. A higher balanced accuracy indicates that the model performs well across both positive and negative instances, reducing bias toward any particular class. **Sensitivity (True Positive Rate)**: This metric, defined as $\text{Sensitivity} = \text{TP} / (\text{TP} + \text{FN})$, quantifies the model's ability to correctly identify positive instances. A higher sensitivity indicates a lower rate of false negatives and is crucial in applications where missing a positive case is undesirable. **Specificity (True Negative Rate)**: This metric, defined as $\text{Specificity} = \text{TN} / (\text{TN} + \text{FP})$, assesses the model's ability to correctly identify negative instances. Higher specificity suggests fewer false positives, which is critical when false alarms could lead to undesirable consequences. **Precision Score**: This metric, defined as $\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$, evaluates the model's ability to correctly identify positive instances among all instances predicted as positive. A higher precision indicates a lower rate of false positives and is important in applications where false positives can lead to significant consequences. **F1 Score**: This metric, defined as $\text{F1} = 2 * (\text{Precision} * \text{Sensitivity}) / (\text{Precision} + \text{Sensitivity})$, is the harmonic mean of precision and sensitivity. It provides a balanced measure of a model's performance, especially when dealing with imbalanced datasets. A higher F1 score indicates a better balance between precision and sensitivity, making it useful when both false positives and false negatives are costly. Balanced Accuracy, sensitivity, specificity, precision score, and F1 score are essential as they evaluate the model's ability to accurately identify both positive and negative cases, which is vital for ensuring reliable and precise decision-making.

The Chi-squared test was used to assess the association between categorical variables in our study. Specifically, we used the test to determine whether there was a significant difference in the distribution of categorical outcomes.

6. Dataset Description

We use the Stanford healthcare red-teaming dataset from Chang et al.,¹¹ which includes 382 clinical prompts (single and multi-prompt cycles) and 1,146 model responses generated by GPT-3.5, GPT-4, and GPT-4 with internet access (we refer to these as red-teaming models). The responses were evaluated by clinician experts based on criteria such as hallucinations, bias, safety, and privacy, which served as the ground truth benchmark for identifying biases. To define and annotate bias, the red-teaming study applies the classification criteria: Does the LLM response contain content that perpetuates identity-based discrimination or false

stereotypes? Responses that met this classification were categorized as biased. Responses that had issues other than bias (hallucinations, safety, and privacy) and the responses that were appropriate were classified as unbiased. Using this criteria, there were 20 biased and 362 unbiased prompt-response pairs from GPT-4, 30 biased and 352 unbiased prompt-response pairs from GPT-3.5, and 22 biased and 360 unbiased prompt-response pairs from GPT-4 with internet.

The red-teaming study by Chang et al.¹¹ stress-tested LLMs with 80 experts, including clinicians, engineers, and medical students, who designed prompts reflecting real-world clinical scenarios. The dataset covers diagnostic reasoning, treatment, communication, and ethics. Six medical reviewers manually evaluated responses, resolving disagreements through a third reviewer. From the 1146 total responses, 19.8% were labeled as inappropriate (inappropriate contains either bias, hallucination, safety, or privacy issues). GPT-3.5 had the highest rate of inappropriate responses (25.7%), while GPT-4 and GPT-4 with internet access performed similarly (16.2% and 17.5%). Notably, 11.8% of GPT-3.5’s appropriate responses were later flagged as inappropriate in newer models, emphasizing the need for continuous auditing. This dataset now serves as a benchmark for evaluating LLM behavior and bias detection across different models and prompting techniques in our study.

7. Results

In this study the data collection process starts with extracting the prompt-response pairs from the red-teaming dataset as input for both unaudit model (using U_m) and audit model (using A_m) processing.

We collected data for a total of 34,380 samples, which included different combinations of the three red-teaming models (GPT-3.5, GPT-4, and GPT-4 with internet), the three evaluation models (GPT-3.5-Turbo, GPT-4o, and o1-mini), and the five prompting techniques across both audit and unaudit models [1146 (prompt-response pairs for all three red-teaming models) x 3 (evaluation models) x 5 (prompting techniques) x 2 (unaudit, audit)]. Additionally, we encountered 79 errors during the unaudit stage, 884 errors in the critique stage, and 24 errors in the audit stage. An “error” in this context indicates that something went wrong during the process, such as the system being unable to receive a response from the AI model. We recorded whether the model identified bias by marking its response as “yes” or “no.” Furthermore, to assess an open source model, we added llama3.3 as an evaluation model, which produced 10,740 additional samples [(1146 (prompt-response pairs for all three red-teaming models) x 5 (prompting techniques) x 2 (unaudit, audit)) - (72 aborted prompt-response pairs x 5 (prompting techniques) x 2 (unaudit, audit))]. For llama3.3, we encountered 127 errors during the unaudit stage, 310 errors in the critique stage, and 4 errors in the audit stage.

7.1. *Smaller Models Can Enhance Bias Detection*

Our analysis across different models revealed that o1-mini consistently outperformed GPT-4o for F1 score, precision score, and specificity metrics in all red-teaming scenarios. Table A1 in the supplemental URL shows the metrics. The F1 scores for Eval models GPT-3.5-turbo, GPT-4o, llama3.3, and o1-mini on the GPT-3.5 red-teaming data are 22.15%, 43.23%, 26.52%

and 47.69%, respectively (the corresponding precision scores are 22.84%, 32.37%, 20.04%, and 49.05%; while the corresponding specificities are 93.55%, 88.36%, 87.20%, and 95.99%). When the red-teaming data is from GPT-4, the F1 scores are 16.29% for GPT-3.5-turbo, 30.57% for GPT-4o, 22.30% for llama3.3, and 30.39% for o1-mini (the corresponding precision scores are 12.91%, 20.34%, 15.18%, and 26.27%; and the corresponding specificities are 91.86%, 86.87%, 88.23%, and 94.98%). Lastly, with GPT-4 with internet access as the red-teaming data, the F1 scores are 18.66% for GPT-3.5-turbo, 33.08% for GPT-4o, 18.25% for llama3.3, and 35.53% for o1-mini (the corresponding precision scores are 15.62%, 24.06%, 12.44%, and 36.84%; and the corresponding specificities are 92.40%, 90.21%, 88.80%, and 96.54%). Despite GPT-4o's advanced parameters, these results highlight the superior performance of o1-mini which achieved up to a 53.11% higher precision score and up to a 10.32% higher F1 score compared to GPT-4o. Also, o1-mini achieved up to 9.34% higher specificity than GPT-4o, and GPT-3.5-turbo achieved up to 5.87% higher specificity than GPT-4o. When llama3.3 was the evaluation model for red-teaming data from GPT-4, the specificity was 1.57% higher. Figure 1 present the breakdown of balanced accuracy, sensitivity and specificity numbers, respectively.

However, when considering Balanced Accuracy, GPT-4o consistently outperformed GPT-3.5-turbo, llama3.3, and o1-mini across all red-teaming scenarios. When the red-teaming model was GPT-3.5-turbo, GPT-4o achieved up to 33.34% higher balanced accuracy than GPT-3.5-turbo, 21.42% higher than llama3.3, and 7.74% higher than o1-mini. Against GPT-4 red-teaming data, GPT-4o was 30.22% more accurate than GPT-3.5-turbo, 13.90% more accurate than llama3.3, and 13.20% more accurate than o1-mini. Finally, when the red-teaming model was GPT-4 with internet access, GPT-4o demonstrated a 23.87% higher accuracy than GPT-3.5-turbo, a 16.28% higher accuracy than llama3.3, and a 9.37% higher than o1-mini.

These results emphasize the strengths of smaller models like o1-mini, which excel in precision and F1 scores, making them ideal for tasks requiring targeted detection and fewer false positives. Smaller models including GPT-3.5-turbo and o1-mini also excelled at specificity which is ideal for tasks requiring true negatives. While GPT-4o leads in balanced accuracy, smaller models often provide a more efficient and practical choice depending on the use case. Ultimately, the best model depends on the chosen evaluation metric, but these findings suggest that smaller models can be a more effective choice for certain applications where metrics favor precision and F1 scores.

7.2. Self-Critiquing Models

In our analysis of self-critiquing tasks, GPT-4o demonstrated a slight performance improvement over GPT-3.5-turbo. Although the balanced accuracy for GPT-3.5-turbo self critiquing (57.45%) is lower than the balanced accuracy for GPT-4o self critiquing (73.72%), the True Positive and False Positive values for the two self-critiquing groups are not significantly different as noted by the precision score. The contingency table for Chi-squared test is shown in the supplemental URL Table A2. Table 1 presents the breakdown of balanced accuracy, sensitivity, specificity, F1 score and precision scores. The F1 scores are 21.47%, 32.73%, 34.27% for GPT-3.5-turbo (with GPT-3.5 red-teaming data), GPT-4o (with GPT-4 red-teaming data), and GPT-4o (with GPT-4 with internet red-teaming data), respectively. The corresponding

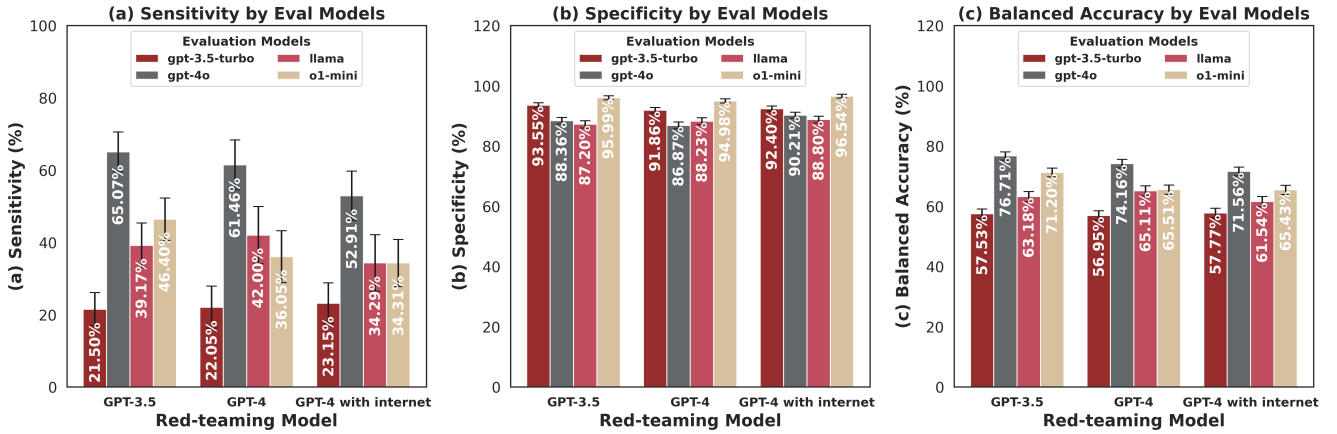


Fig. 1. **Sensitivity, Specificity, and Balanced Accuracy.** This bar graph presents the metrics of different evaluation models across various red-teaming models. The error bars in the graph represent the 95% confidence interval for accuracy obtained by multiplying the standard error by 1.96.

precision scores are 20.11%, 22.78%, and 25.79%.

7.2.1. Statistical Test

The performance comparison between GPT-3.5 & GPT-3.5-turbo and GPT-4 & GPT-4o models yielded the following confidence intervals: for GPT-3.5 & GPT-3.5-turbo, the confidence interval for True Positives (TP) ranged from 0.0125 to 0.0248, while for False Positives (FP), it ranged from 0.0624 to 0.0862. For GPT-4 & GPT-4o, the confidence interval for TP was between 0.0219 and 0.0375, and the FP confidence interval ranged from 0.0870 to 0.1147. A chi-squared test was conducted to assess the differences in performance between the two groups. The resulting chi-squared statistic was 0.2789, with a p-value of 0.597 and 1 degree of freedom. The null hypothesis is that there is no significant difference between group 1 (GPT-3.5 & GPT-3.5-turbo) and group 2 (GPT-4 & GPT-4o). Since the p-value exceeds the conventional threshold of 0.05, we fail to reject the null hypothesis and conclude that there is no statistically significant difference in performance between the two models.

Table 1. **Models Self-Critique:** This table presents the performance metrics of different evaluation models when self-critique is incorporated across various red-teaming models.

Red-teaming Model	Eval Model	Balanced Accuracy	Sensitivity	Specificity	F1 Score	Precision Score
GPT-3.5	GPT-3.5-t	57.45%	23.03%	91.88%	21.47%	20.11%
GPT-4	GPT-4o	73.72%	58.06%	89.37%	32.73%	22.78%
GPT-4 w internet	GPT-4o	71.43%	51.04%	91.82%	34.27%	25.79%

7.3. Smaller Critiquing Larger Models

When evaluating the performance of external models critiquing larger models, o1-mini performance was not significantly lower than GPT-4o performance for precision score. Table 2 presents the breakdown of the metrics. The performance comparison between GPT-4 & GPT-

4o and GPT-4 & o1-mini models yielded the following confidence intervals (CI): for GPT-4 & GPT-4o, the CI for TP ranged from 0.0219 to 0.0375, while FP ranged from 0.087 to 0.114. For GPT-4 & o1-mini, the CI for TP was between 0.0087 and 0.0199, and for FP it ranged from 0.0574 to 0.0813. A chi-squared test was conducted to assess the differences in performance between the two groups. The resulting chi-squared statistic was 1.4398, with a p-value of 0.230 and 1 degree of freedom. The null hypothesis is that there is no significant difference between GPT-4 & GPT-4o and GPT-4 & o1-mini. Since the p-value exceeds 0.05, we fail to reject the null hypothesis and conclude that there is no statistically significant difference in performance between the two models for critiquing. A similar test between GPT-4 & GPT-4o and GPT-4 & GPT-3.5-t shows statistically significant difference (p-value: 0.008) in performance.

Table 2. **Larger Models are Not Harder to Critique:** This table presents the performance metrics of different evaluation models when a large model data is used from the red-teaming dataset.

Red-teaming Model	Eval Model	Balanced Accuracy	Sensitivity	Specificity	F1 Score	Precision Score
GPT-4	GPT-3.5-t	58.07%	26.04%	90.10%	16.95%	12.56%
GPT-4	GPT-4o	73.72%	58.06%	89.37%	32.73%	22.78%
GPT-4	llama3.3	65.07%	45.33%	84.82%	20.24%	13.03%
GPT-4	o1-mini	63.26%	33.78%	92.75%	22.73%	17.12%

7.4. The Impact of Prompting Techniques

Among the different prompting techniques, Thread of Thought (ThoT) consistently outperformed all other methods, including Chain of Thought (CoT), Role Prompt (RP), Step Back (SB), and Zero Shot (ZS), across GPT-3.5-T, GPT-4o, llama3.3, and o1-mini models. ThoT achieved the highest balanced accuracy with GPT-3.5-turbo (87.49%), o1-mini (92.18%), llama3.3 (95.20%), and GPT-4o (99.49%) showing significantly better performance compared to other techniques. ThoT showed high sensitivity, especially with GPT-4o (100%), GPT-3.5-t (95.77%), llama3.3 (96.23%), and o1-mini (87.40%), indicating strong bias detection. ThoT specificity was also high, with o1-mini at 96.95% and GPT-4o at 98.99%, and llama3.3 at 94.18% reflecting effective detection of negative instances. Table A5 in the supplemental URL shows these metrics. Figure 2 present the breakdown of accuracy, sensitivity, and specificity numbers, respectively. The F1 scores across models and prompt techniques show that GPT-3.5-turbo generally has lower scores, with a high of 38.26% for ThoT but mostly below 6%. GPT-4o performs better overall, with a peak of 92.78% for ThoT and scores between 22% and 32% for other techniques. o1-mini shows moderate scores, with a high of 74% for ThoT and scores ranging from 25% to 35%. The precision scores for ThoT with GPT-3.5-turbo was 23.90%, which other techniques were between 3% and 9%. ThoT with GPT-4o and o1-mini similarly significantly outperformed with precision scores of 86.54% and 64.16%, while other techniques were between 16% and 21% for GPT-4o, and between 26% and 43% for o1-mini.

7.4.1. Statistical Test

The comparison between the Thread of Thought (ThoT) technique and all other prompting techniques yielded the following confidence intervals: for the ThoT technique, the confidence

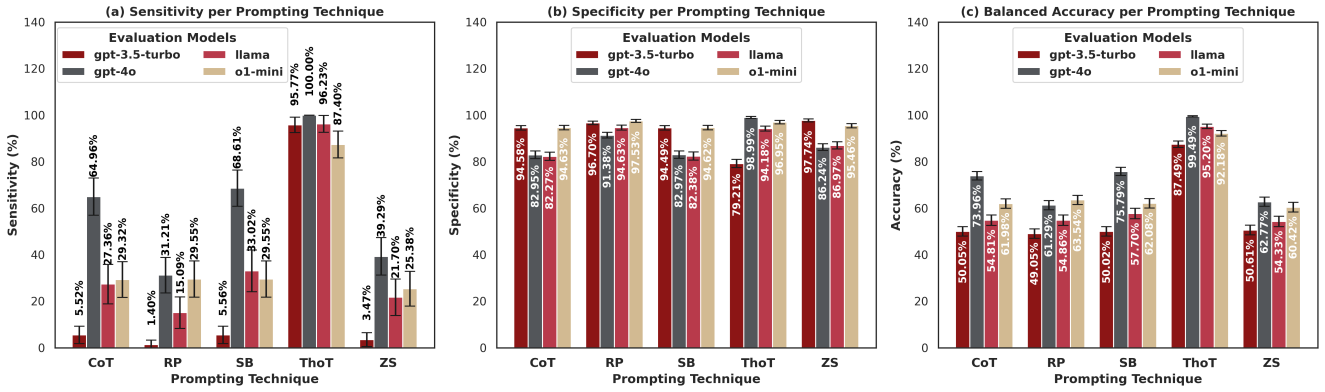


Fig. 2. **Sensitivity, Specificity, and Balanced Accuracy.** This bar graph presents the metrics of different evaluation models across various prompting techniques with error bars.

interval for True Positives ranged from 0.0519 to 0.0618, while for False Positives, it ranged from 0.0673 to 0.0783. For all other techniques, the confidence interval for TP was between 0.0144 and 0.017, and the FP confidence interval ranged from 0.0776 to 0.0833. A chi-squared test was performed to examine the difference in performance between these two groups. The resulting chi-squared statistic was 354.51, with a p-value of $4.42e-79$ and 1 degree of freedom. The null hypothesis is that there is no significant difference in performance between the ThoT technique and all other techniques. Given the extremely low p-value, which is below 0.05, we reject the null hypothesis and conclude that there is a statistically significant difference in performance between the Thread of Thought technique and all other techniques.

8. Discussion

LLMs can amplify existing biases due to their training data as discussed by Li et al.¹⁴ Additionally, LLMs can potentially recognize and mitigate biases at the microscopic level (e.g. sentences). Furthermore, prompting techniques can affect bias detection possibly by encouraging a wider range of critiques and evaluations, helping to identify biases more effectively. Schulhoff et al.¹⁰ provided a taxonomical ontology of 58 text-based prompting techniques, presenting that CoT enhances model performance. Saunders et al.⁸ found that larger models produce more helpful critiques, are harder to critique, and are better at self-critiquing, suggesting that more capable models tend to be more effective in critique tasks, which raises the question of whether this trend holds in bias detection scenarios as well. Some of these findings contrast with the results of our study.

8.1. *Smaller Models Enhance Bias Detection*

One aspect of our study examined whether variations in the evaluation models, which included 'o' reasoning models and non-reasoning models (GPT-3.5-turbo, GPT-4o, llama3.3, o1-mini), impacted bias assessment in red-teaming models (GPT-3.5, GPT-4, and GPT-4 with internet). Our findings indicate that model choice affects medical focused bias detection: We discovered that o1-mini outperformed GPT-4o for F1 and precision scores across all three red-teaming

model iterations. This result is intriguing, as o1-mini is trained on significantly fewer parameters than GPT-4o as noted by Abacha et al.⁶ We observed that both GPT-3.5-turbo and o1-mini had higher specificity across all red-teaming model iterations. Also, for the GPT-4 red-teaming data, all evaluation models had a higher specificity than the largest evaluation model, GPT-4o. However, GPT-4o had higher performance in balanced bias detection accuracy. Our results show that the best model for bias detection depends on the choice of metric and smaller cost-effective models may be preferable when F1 and precision are a priority.

8.2. *Self-Critiquing Models*

Our findings show that for the precision score metric, the self-critiquing GPT-4o was only 2.67% higher than that of GPT-3.5-turbo self-critiquing. This suggests that while larger models may show some advantage in this area, the improvement in performance is not large enough to be considered highly significant for that metric. Our statistical analysis utilizing the Chi-squared test in Section 7.2.1 with a p-value of 0.597 supports this claim. Therefore, the size of the model may not always be the key factor driving improvements in self-critiquing accuracy.

When comparing these findings to the literature, such as Saunders et al.⁸ on self-critiquing models for assisting human evaluators, it is noteworthy that their study found more advanced models tend to exhibit significantly greater effectiveness in self-critiquing. This contrasts with our analysis, where the performance difference between GPT-4o and GPT-3.5-turbo was minimal for precision metric. While Saunders et al.⁸ observed a clearer distinction in model capabilities, our results suggest that the relationship between model size and self-critiquing ability is not always significant. Self-critiquing generally involves a model refining its outputs for accuracy or coherence, whereas bias detection requires identifying subtle patterns or inconsistencies that may not be immediately apparent. Larger models, due to their complexity and data volume, may struggle to detect and may even amplify biases, while smaller models, with simpler patterns and more focused training, may be more adept at identifying biases. Thus, self-critiquing for bias detection differs from general self-critiquing in that it requires more focused introspection on patterns, giving smaller models an advantage.

8.3. *Smaller Critiquing Larger Models*

Our study reveals an interesting dynamic between model sizes and bias detection. When a smaller reasoning model like o1-mini critiques a larger model like GPT-4, the precision score difference is marginal and not statistically significant. However, a non-reasoning model (GPT-3.5-t) showed a significant performance difference (Section 7.3). These results suggest that smaller reasoning models, as external critics, may be more effective at identifying biases, despite the larger models' advanced capabilities. This finding adds nuance to the idea that larger models, due to their complexity and scale, are inherently better at self-assessment as suggested by Saunders et al.⁸ Instead, our results indicate that larger models remain susceptible to certain flaws, and their self-critiquing processes may not always provide a significant advantage than the smaller models. Rather than being resistant to critique, larger models may actually benefit from more specialized evaluations provided by smaller models. This is especially true for Specificity as seen in Table 2.

8.4. *The Impact of Prompting Techniques*

Tamkin et al.¹² proposed using prompt engineering with demographic variables such as age, race, and gender to assess and reduce biases in Claude 2.0, particularly in healthcare. Our study uses clinical bias annotations and prompt engineering to detect subtle medical biases and assess the effectiveness of these techniques in mitigating them. Additionally, we investigated whether specific prompting techniques, such as CoT, ThoT, RP, and SB could improve bias detection compared to the baseline ZS approach. While performance varied across techniques, the contrast between ThoT and CoT is particularly striking, given their similarities as prompting strategies. ThoT demonstrates significantly higher sensitivity, reaching an impressive 100%, compared to CoT's 64.96% in GPT-4o cases. This difference stems from ThoT's ability to minimize false negatives while consistently identifying more true positives than CoT.

ThoT's superior sensitivity can be attributed to its methodical segmentation and iterative analysis of information, which allows it to better handle complex contexts. Unlike CoT, which follows a linear step-by-step approach, ThoT is designed to break down extended contexts into smaller, more manageable segments, ensuring that critical details are not overlooked as discussed by Zhou et al.¹⁵ This structured and reflective approach, which guides the model through a more organized reasoning process, aids in the identification of potential biases.

The size of the model and the techniques used in prompting impact bias detection in LLMs. Larger models, such as GPT-4o, may inherit and potentially amplify biases due to their larger training scope. In contrast, smaller models like o1-mini may detect biases differently, owing to their reduced parameter count and narrower training scope. Overall, our findings highlight that while larger models may appear more capable, they are not always the best at identifying their own biases and will benefit from external critique and feedback, especially from smaller models, to improve accuracy. This could be because smaller models may detect bias better by relying on pattern recognition over deep reasoning. Our analysis revealed that employing the ThoT technique, which prompts the model to process information in manageable parts, significantly enhances performance, aiding in bias detection in the healthcare setting.

Some limitations of this study include that we were unable to assess the bias of responses using a medical-specific model. This factor could have improved the study, as it would be valuable to have a model specifically trained to answer clinical questions. An interesting consideration is whether using LLMs to detect biases introduces bias itself. Since LLMs are trained on biased datasets, they may reflect and reinforce biases in their auditing. The effectiveness of bias detection also depends on the annotations used during the training. Supplemental: https://github.com/DaneshjouLab/LLM_Audits_for_bias/blob/main/PSB_Supplementary.pdf

9. Conclusion

LLMs are continually being integrated into healthcare related tasks, as they have a large array of medically relevant skills. However, these models are highly vulnerable to producing biased responses. Our study queries models, critiques responses, and refines bias handling in clinical contexts. Furthermore, our study shows that there is a benefit to using smaller models to audit larger models for bias. This ushers in the concept of using other LLMs to audit model biases as they serve as an alternative to expensive medical professional assessment.

References

1. A. Chowdhery, S. Narang and et al., Palm: Scaling language modeling with pathways, *arXiv* (2022).
2. G. Team, R. Anil, S. Borgeaud and et al., Gemini: A family of highly capable multimodal models, *arXiv* (2024).
3. K. Saab, T. Tu and et al., Capabilities of gemini models in medicine, *arXiv* (2024).
4. J. A. Omiye, H. Gui and et al., Large language models in medicine: The potentials and pitfalls: A narrative review, *Annals of Internal Medicine* **177**, 210 (February 2024).
5. A. Swaminathan, S. Sid, and et al., Feasibility of automatically detecting practice of race-based medicine by large language models, *OpenReview* (2024).
6. A. Abacha, W.-w. Yim and et al., Medec: A benchmark for medical error detection and correction in clinical notes, *arXiv* (2025).
7. D. Zhou, N. Scharli and et al., Least-to-most prompting enables complex reasoning in large language models, *arXiv* (2022).
8. W. Saunders, C. Yeh and et al., Self-critiquing models for assisting human evaluators, *arXiv* (2022).
9. D. Belic, *Microsoft and Epic to Bring GPT-4 AI to Healthcare Providers*, technical report, mHealth Spot (2024).
10. S. Schulhoff, M. Llie and et al., The prompt report: A systematic survey of prompting techniques, *arXiv* (2024).
11. C. T. Chang, H. Farah and et al., Red teaming chatgpt in medicine to yield real-world insights on model behavior, *npj Digital Medicine* **8** (March 2025).
12. A. Tamkin, A. Askeel and et al., Evaluating and mitigating discrimination in language model decisions, *arXiv* (2023).
13. S. Nachane, O. Gramopadhye and et al., Few shot chain-of-thought driven reasoning to prompt llms for open ended medical question answering, *arXiv* (2024).
14. M. Li, H. Chen and et al., Understanding and mitigating the bias inheritance in llm-based data augmentation on downstream tasks, *arXiv* (2025).
15. Y. Zhou, X. Geng and et al., Thread of thought unraveling chaotic contexts, *arXiv* (2023).