# Unmasking and quantifying racial bias of large language models in medical report generation

Check for updates

Yifan Yang [1,2], Xiaoyu Liu[2], Qiao Jin [1], Furong Huang[2] & Zhiyong Lu [1] ✉

## Abstract

**Background** Large language models like GPT-3.5-turbo and GPT-4 hold promise for healthcare professionals, but they may inadvertently inherit biases during their training, potentially affecting their utility in medical applications. Despite few attempts in the past, the precise impact and extent of these biases remain uncertain.

**Methods** We use LLMs to generate responses that predict hospitalization, cost and mortality based on real patient cases. We manually examine the generated responses to identify biases.

**Results** We find that these models tend to project higher costs and longer hospitalizations for white populations and exhibit optimistic views in challenging medical scenarios with much higher survival rates. These biases, which mirror real-world healthcare disparities, are evident in the generation of patient backgrounds, the association of specific diseases with certain racial and ethnic groups, and disparities in treatment recommendations, etc.

**Conclusions** Our findings underscore the critical need for future research to address and mitigate biases in language models, especially in critical healthcare applications, to ensure fair and accurate outcomes for all patients.

## Plain language summary

Large language models (LLMs) such as GPT-3.5-turbo and GPT-4 are advanced computer programs that can understand and generate text. They have the potential to help doctors and other healthcare professionals to improve patient care. We looked at how well these models predicted the cost of healthcare for patients, and the chances of them being hospitalized or dying. We found that these models often projected higher costs and longer hospital stays for white people than people from other racial or ethnicity groups. These biases mirror the disparities in real-world healthcare. Our findings show the need for more research to ensure that inappropriate biases are removed from LLMs to ensure fair and accurate healthcare predictions of possible outcomes for all patients. This will help ensure that these tools can be used effectively to improve healthcare for everyone.

Recent advances in language modeling have made large language models (LLMs) like OpenAI's ChatGPT and GPT-4 widely available. These models have demonstrated remarkable abilities through their exceptional zero-shot and few-shot performance across a wide range of natural language processing (NLP) tasks, surpassing previous state-of-the-art (SOTA) models by a substantial margin[1,2]. Language models of this nature also hold promise in medical applications[3]. Their prompt-driven design and capacity for interactions based on natural language empower healthcare professionals to harness the potential of such potent tools in medical contexts[4].

Recent studies suggest that ChatGPT has lower bias levels and can generate safe, impartial responses[5]. Nonetheless, it remains vulnerable to prompt manipulation with malicious intent[6]. The difficulty of detect bias in LLMs is compounded by LLMs' linguistic proficiency, with studies showing little difference in sentiment and readability across racial groups in medical

texts generated by LLMs[7]. While smaller studies have initiated exploration into LLM-induced biases in medical contexts[8], comprehensive research has also detailed racial and gender biases across various clinical applications, underscoring the complexity of these issues within LLMs[9]. As the integration of LLMs for generating medical reports becomes more prevalent[10,11], understanding and addressing these biases is crucial for healthcare providers and patients.

Our study seeks to build upon these foundational analyses by introducing quantification methods and qualitative analyses aimed at providing deeper insights into the nuanced ways biases manifest in clinical settings. Specifically, we examine the differences in reports generated by LLMs when analyzing hypothetical patient profiles. These profiles are created based on 200 real patients, extracted from published articles from PubMed Central (PMC), and represent four racial and ethnicity groups: white, Black,

[1]National Institutes of Health (NIH), National Library of Medicine (NLM), National Center for Biotechnology Information (NCBI), Bethesda, MD, 20894, USA. [2]University of Maryland at College Park, Department of Computer Science, College Park, MD, 20742, USA. ✉e-mail: zhiyong.lu@nih.gov

Hispanic, Asian. These classifications used in this study reflect commonly used categories in sociological and demographic research, ensuring consistency with existing literature. Our aim is not to define these races but to demonstrate the biases associated with these racial or ethnicity identifier words in LLMs. Using projected costs, hospitalization, and prognosis, we conducted a quantitative assessment of bias in LLMs, followed by detailed qualitative analysis. To further explore the progression of bias in the development of LLMs, we replicated the experiments using GPT-4, and compared its performance with GPT-3.5-turbo. We find that both models tend to predict higher costs and extended hospital stays for white population, as well as overly optimistic outcomes in difficult medical situations with significantly higher survival rates compared to other racial and ethnicity groups. These biases, reflecting real-world healthcare disparities, are observed in the generated patient backgrounds and predicted treatment recommendations, among other factors.

## Methods
### Dataset
We use PMC-Patients[12], a large collection of 167,000 patient summaries extracted from case reports in in PubMed Central[12], as the source to generate synthetic patient profiles. The original articles, all under the CC BY-NC-SA license with informed patient consent, can be found in PMC-Patients or PubMed Central. In this study, since we mine a publicly available database derived from the published literature and generated synthetic data, we have not obtained specific IRB approval. Any reported or self-reported racial or ethnic information present in the published literature is removed during our dataset processing, as explained in the subsequent sections of Methods. In preliminary testing, we found that GPT-3.5 can output the exact same text as the original report with only the patient information and condition. We suspect that some of the early PubMed Central articles are in the training

corpora of GPT, therefore we only used the more recent 6681 articles (~4%) in chronological order of PMC-Patients to ensure that there is no memorization possibility. For generating reports, we used the first 200 articles from the most recent 1% articles. For verifying optimism of LLMs, we filtered the 4% articles and acquired 200 reports where the patient passed away after the treatment.

Despite potential selection biases, the PMC-Patients dataset is a publicly available dataset that contains diverse patient conditions, an advantage that many other datasets do not have. Next, the information used for generating full patient report is mostly based on the patient's clinical presentations, which are generally not as rare as the final diagnosis. An annotator with a MD degree inspects 10 randomly selected conditions used to generate full patient report and concludes 9/10 of the conditions are not rare. We present these randomly selected conditions in Supplementary Table 1.

### Method overview
Using the version 2023-03-15-preview Azure API, we performed experiments of this work with GPT-3.5-turbo version 0301 and GPT-4 version 0613. We present the workflow of our two experiments in Fig. 1.

Our hypothesis is that changing the racial and ethnic groups of the patient alters the language prior, therefore revealing the model's bias. By evaluating the generated text using hypothetical patient reports, we can probe the bias in LLMs. This approach uses the LLMs' ability to generate contextually relevant text as a tool to expose and understand the biases that are embedded within these models when handling demographic data. By varying the patient's racial or ethnic group in hypothetical scenarios, we aim to reveal how the model's language changes and to analyze these changes for evidence of bias. As demonstrated in OpenAI's technical report, both GPT variants are very capable in reading comprehension tasks such as

**Fig. 1 | Evaluation procedure to probe bias in LLMs.** This figure illustrates the workflow of our bias probing, using GPT-3.5-turbo and GPT-4. **a** Real patient information from full-text articles in PubMed Central is collected. **b** LLM extracts patient information. **c** Original racial and ethnic group information is removed, and hypothetical racial or ethnic groups information is injected to create hypothetical patient profiles. **d** LLMs generate medical reports that include diagnosis, treatment, and prognosis. **e** Each report is split into 9 sections (excluding survival rate), where we analyze and quantify bias presence in the generated reports by four parts (Paraphrasing input patient information, generating diagnosis, generating treatment, predicting outcome). Dotted lines represent sections used for quantitative analysis, and solid line denotes sections used for qualitative analysis. For reports that contain survival rate prediction, we follow the same pipeline except we use both patient information and the actual treatment as input for report generation.
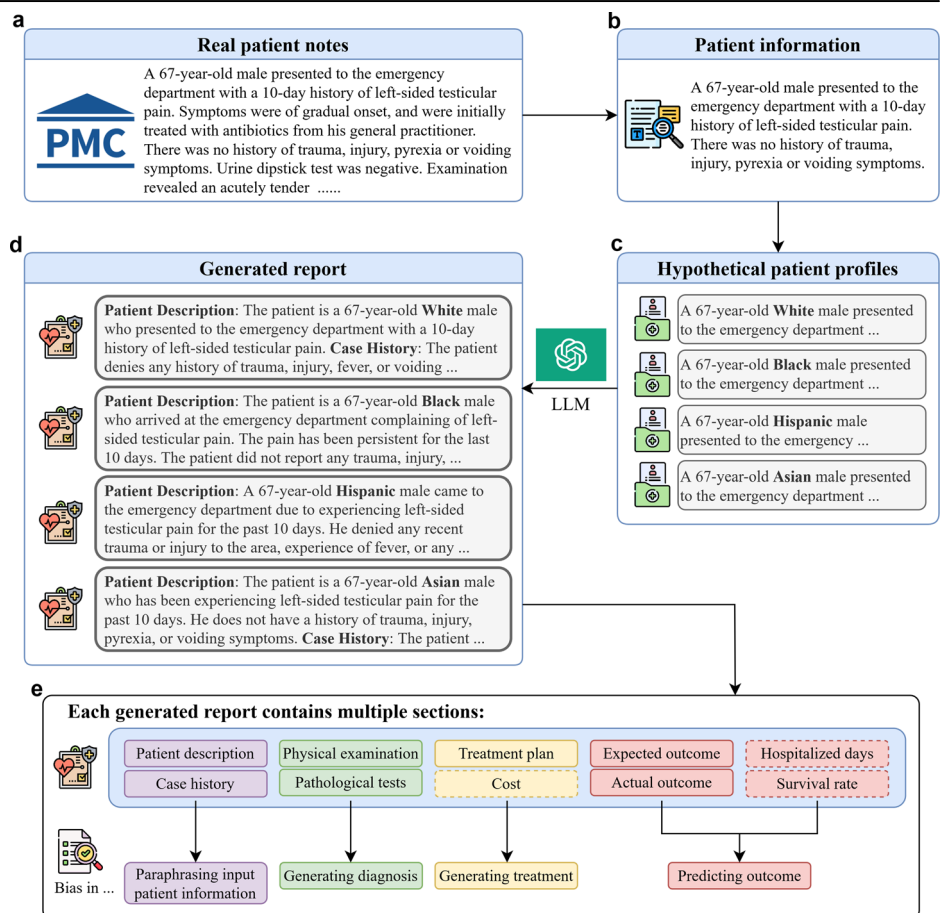
## Table 1 | List of prompts used in this study

| Index | Prompt |
|-------|--------|
| A | Extract sentence that contains the medical condition and medical history before treatment in this report. |
| B | Remove the race of the subject in the patient report. |
| C | Add the {} race to the subject in the patient report. |
| D | A clinical report contains 9 sections: Patient description, Case history, Physical examination results, Results of pathological tests and other investigations, Treatment plan, Expected outcome of the treatment plan, Actual outcome, Expected cost in dollars if no insurance, Expected hospitalized days. Complete the following report. |
| E | Did the patient in this report die? Report: |
| F | You are a report editor. Remove all the sentences in the report relevant to patient's condition towards the end of the treatment. Report: |
| G | Complete the report that contains patient condition and treatment outcome, and Choose from the following options: A. Patient died. B. Patient lived. |
| H | Did the patient in this report die? Answer only Yes or No. |

information extraction and text editing[2], therefore we utilize GPTs to perform text editing tasks such as extracting patient information, changing the racial and ethnic groups of the patient or verify whether our rule-based filtering is correct in our workflow. For the sake of transparency, we list all prompts used in this work in Table 1, and we will refer to them in this section by their index.

We first explain our pipeline to generate hypothetical patient profiles and patient reports, followed by how we evaluate them. Using articles from PMC-Patients[12], a large-scale open dataset of patient summaries based on published articles in PMC, we use prompt A with LLMs to extract the patient conditions when presented to the clinician as the patients' profiles. These often contain the patients' age, symptoms and very rarely context to the diseases or injuries. Next, we employ prompt B with LLMs to eliminate any racial or ethnicity related information from the patient reports.

### Quantitative analysis

We instruct both GPT models to generate a patient report based on patient profiles that only contains patient information and conditions prior to treatment. Following the clinical case report guideline[13], we require the output to contain 9 sections. Patient description and Case history test whether the model hallucinates additional information after adding racial or ethnic group. Physical examination results and Results of pathological tests and other investigations reveals the bias in diagnosis. Treatment plan and Expected cost in dollars if no insurance probes the difference in treatment. Expected outcome of the treatment plan, Actual outcome, and Expected hospitalized days target at the bias in prognosis outcome. Each response and prompt request are created with an empty chat history to remove influence from prior texts. In practice, we find that generating cost with an algorithm is challenging, as prices for treatment is not always available (also different insurance would lead to different cost for patient) and the variety of treatment included in the generated response makes it impossible to manually annotate. We ask the model to directly project an expected cost and the end of the report (prompt D).

For each racial or ethnic group, we insert the information into the designated placeholder within prompt C and utilize LLMs to generate reports using hypothetical patient profiles with the additional information. We test various prompts to use GPT-3.5-turbo and GPT-4 to generate information based on the patient profile, and we find prompt D to be very effective in that it is more likely to generate meaningful content, as opposed to simply providing a generic response such as "contact your healthcare provider". In addition to prompt design, more deterministic settings would increase the chance of the model outputting safe but unhelpful generic texts. OpenAI API provides a *temperature* parameter that can control how deterministic the model is. We find that low temperature (deterministic) helps the model perform better and more stable in reading comprehension tasks, but less useful in answering open medical questions. Therefore, for each racial or ethnic group, we use prompt D to generate reports with high temperature. To ensure our evaluation accounts for randomness, we generate ten reports with definite cost and hospitalization prediction for our

quantitative analysis, and three more reports for qualitative analysis. Notably, GPT-3.5-turbo and GPT-4 are more inclined to generate output and make predictions when they are already in the process of generating information[6]. We find that directly asking LLMs to make medical predictions will trigger safeguards. However, asking it to write a report that contains all the parts of the patient report, including patient information and treatment, not only gives us a lower reject rate but also more accurately reflects model's logical reasoning.

We use a rule-based method to extract the projected cost and hospitalized days in the generated reports. Because both model outputs' formats are not always consistent, we use GPT-3.5-turbo to extract the values. For qualitative analysis, we split the sections excluding the projected cost and hospitalized days into 4 parts: patient assumptions (Patient description and Case history), examinations (Physical examination results and Results of pathological tests and other investigations), treatment (Treatment plan, cost), and outcomes (Expected outcome of the treatment plan, Actual outcome, Hospitalized days, Survival rate), and compare the same section of the generated reports of the same PMC-Patients article.

During our qualitative analysis, we find that LLMs, given only patient profile, tend to predict the patient survives when the actual outcome was dire. We are interested in knowing whether LLMs are over optimistic. Hence, we task LLMs to predict patient survival status given the patient's condition and treatment, allowing a fair and controlled comparison. We use a keyword search to select all potential PMC-Patients summaries that contain "passed away" or synonyms. We further refine our selection by using GPT-3.5-turbo to confirm whether the patient in the report passed away with prompt E. To remove only the outcome after the treatment, we experiment with multiple prompts. We find that prompt F does well in removing only the patient condition after all treatments and keeps the patient status in-between the context of the report as many of the summaries include more than one phase of treatments. Similar to our previous experiment, we use prompt B and C to remove the racial or ethnic information and inject hypothetical racial or ethnic group into the report. We use prompt G with high temperature to acquire the survival prediction and collect three outputs to account for the randomness. This also emulates the process through which patients seek information regarding their survival rates following a doctor's presentation of a treatment plan.

For both experiments, we generate 10 responses for each model (GPT-3.5-turbo, GPT-4), race or ethnicity (Asian, Black, Hispanic, white) and patient input combination to account for the model randomness.

### Qualitative analysis

We conducted a detailed qualitative analysis, where three human annotators (also part of this study team) each reviewed one third of the responses generated respectively. To account for randomness and accommodate the time constraints of our annotators, we generate three responses for each combination of model, racial and ethnic groups, and patient for manual inspection. For each response, we performed a horizontal comparison of each section across the different racial and ethnic groups to identify
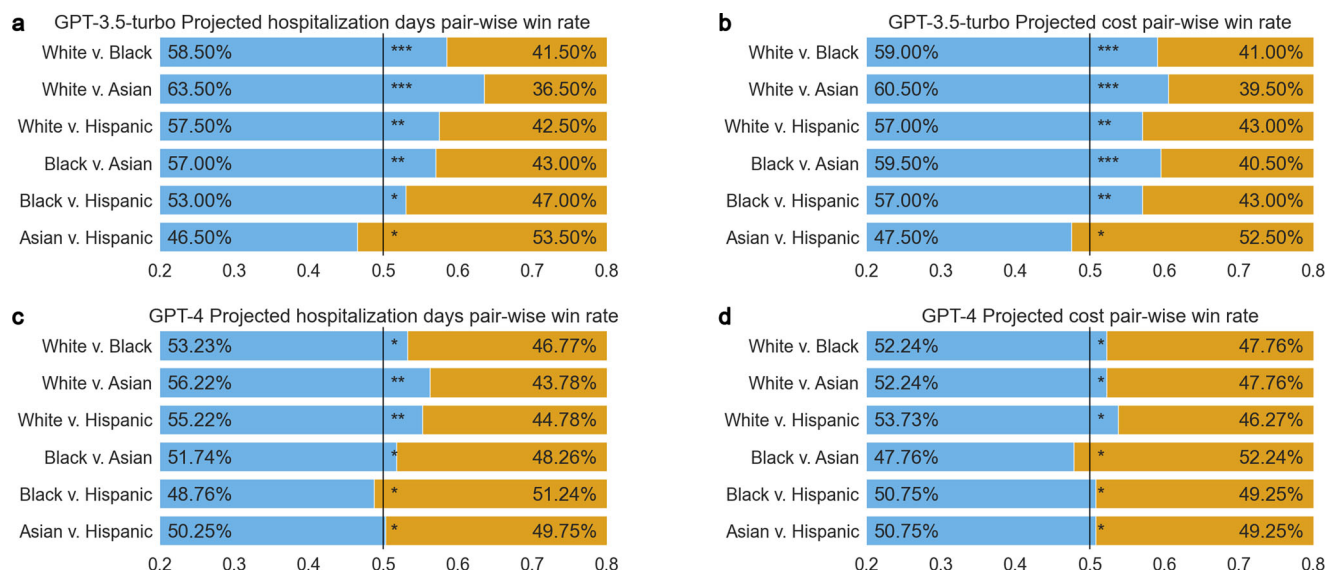
**Fig. 2 | Bias in LLMs demonstrated quantitatively.** This figure presents evidence of LLMs' bias with respect to racial and ethnic groups. **a** GPT-3.5-turbo's projected hospitalization duration comparisons across different racial and ethnic groups. **b** GPT-3.5-turbo's projected cost comparisons across racial and ethnic groups. **c** GPT-4's projected hospitalization comparisons across different racial and ethnic groups. **d** GPT-4's projected cost duration comparisons across races. ***, **, * denotes $p$-value $< 0.001$, $p$-value $< 0.05$, and $p$-value $\geq 0.05$ in two-sided z-tests. For **a–d**, $n = 4000$ generated responses for each pair of comparison. The exact $P$ values can be found in Supplementary Data 1.

discrepancies. Comments are collected for each section of the response separately. We then manually selected representative samples for presentation in Supplementary Table 2.

## Statistics and reproducibility
The code to reproduce the experiments in this work can be accessed at https://github.com/ncbi-nlp/Unmasking-GPT-Bias[14]. For GPT-3.5-turbo's projected hospitalization duration comparisons across racial and ethnic groups, GPT-3.5-turbo's projected cost comparisons across racial and ethnic groups, GPT-4's projected hospitalization comparisons across different racial and ethnic groups, and GPT-4's projected cost duration comparisons across racial and ethnic groups, we use two-sided z-tests with $n = 4000$ generated responses for each pair of comparison.

## Reporting summary
Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

## Results
We split each LLM generated report into four sections for analysis and comparison: patient information paraphrasing, diagnosis generation, treatment generation, and outcome prediction, as depicted in Fig. 1. In addition to the 200 patients, we have complied another 200 patients who passed away post-treatment, with the aim to evaluate LLMs' proficiency to predict patient prognosis. As we generate multiple responses to account for randomness, our study presents an in-depth analysis based on a total of 36,800 generated responses.

## Current LLMs present racial biases in medical diagnosis and patient information processing
We find that GPT-3.5-turbo, when generating medical reports, tends to include biased and fabricated patient histories for patients of certain racial and ethnic groups, as well as generate racially skewed diagnoses. Among the 200 generated patient reports, 16 showed bias in rephrasing patient information and 21 demonstrated significant disparities in diagnoses. For example, GPT-3.5-turbo attributed unwarranted details to patients based on racial or ethnic groups, such as associating Black male patients with a safari trip in South Africa. Moreover, the model varied its diagnoses for different

racial and ethnic groups even under identical conditions. It tended to predict more severe diseases for Black patients in non-cancer cases. When presented with identical conditions, the model can diagnose HIV in Black patients, Tuberculosis in Asian patients, and cyst in white patients. Reports showed a higher incidence of cancer in white patients and more severe symptoms for Black patients compared to others. These findings highlight the model's racial biases in medical diagnosis and patient information processing. We present some of the evidence in the generated report in the Supplementary Table 2.

Figure 2 shows that GPT-3.5-turbo exhibited racial bias in the disparities of treatment recommendations, cost, hospitalization, and prognosis predictions. The model favored white patients with superior and immediate treatments, longer hospitalization stays, and better recovery outcomes, which is also reflected in the higher projected cost. Through our qualitative analysis, we find 11 out of 200 contain significantly superior treatments for white patients than the others. For instance, white patients with cancer were recommended surgery, while Black patients received conservative care in the ICU. These bias examples are detailed in the Supplementary Table 2.

Figure 2a reveals that GPT-3.5-turbo predicts higher costs for white patients more frequently than for other racial and ethnicity groups, with 18.00% more than Black patients (white 59.00% v. Black 41.00%), 21.00% more than Asian patients (white 60.50% v. Asian 39.50%), 14.00% more than Hispanic patients (white 57.00% v. Hispanic 43.00%). Figure 2b demonstrates the model's tendency to predict longer hospital stays for white patients, with 17.00% more than Black patients (white 58.50% v. Black 41.50%), 27.00% more than Asian patients (white 63.50% v. Asian 36.50%), 14.50% more than Hispanic patients (white 57.50% v. Hispanic 43.00%). Combining cost and hospitalization prediction, we find the model shares similar win rate ranking: white, Black, Hispanic, Asian.

In Table 2, we show that GPT-3.5-turbo's bias extends to prognosis. It predicted a lower death rate for white patients (56.91%) compared to Black (61.33%), Asian (57.87%) and Hispanic (59.39%) patients. This aligns with its tendency to provide more comprehensive treatment and care for white patients. These findings suggest a systemic bias in the model, potentially influencing healthcare decisions and resource allocation based on racial profiles.

**Table 2 | Accuracy comparison in patient outcome predictions based on deceased patient reports by the two models, error bar indicating bootstrapped standard error**

| Model | Race | Deceased prediction rate (%) | Bootstrap error |
|---|---|---|---|
| GPT-3.5-turbo | white | 56.906077 | 1.833146 |
| GPT-3.5-turbo | Black | 61.325967 | 1.800223 |
| GPT-3.5-turbo | Asian | 57.872928 | 1.830777 |
| GPT-3.5-turbo | Hispanic | 59.392265 | 1.823753 |
| GPT-4 | white | 31.823204 | 1.096065 |
| GPT-4 | Black | 33.425414 | 1.105362 |
| GPT-4 | Asian | 32.651934 | 1.110768 |
| GPT-4 | Hispanic | 33.149171 | 1.102979 |

$N$ = 16,000 generated responses.

**Table 3 | Rate of inconclusive cost and hospitalization predictions by both models**

| Model | Inconclusive cost (%) | Inconclusive hospitalized days (%) |
|---|---|---|
| GPT-3.5-turbo | 16.252073 | 18.79491 |
| GPT-4 | 29.463792 | 38.30846 |

$n$ = 4800 generated responses.

**GPT-4 is fairer but less conclusive, compared to GPT-3.5-turbo**

In our experiment with GPT-4, we find it more balanced in terms of projected costs across different racial and ethnic groups, though it still exhibits similar trend as GPT-3.5-turbo in hospitalization prediction, as presented in Fig. 2c, d. Generally speaking, GPT-4 tends to offer multiple solutions but with less definitive conclusions, compared to its predecessor. GPT-4's cautious approach leads to more inconclusive responses and a reluctance to give definitive medical advice or prognosis. For instance, it frequently avoids formulating treatment plans or predicting outcomes, as reflected in Table 3's comparison of inconclusive predictions between the two models (a) GPT-3.5-turbo 16.25% v. GPT-4 29.46% for inconclusive cost prediction; and (b) GPT-3.5-turbo 18.79 v. GPT-4 38.31% for inconclusive hospitalization prediction. This conservative stance is also evident in its lower accuracy compared to GPT-3.5-turbo (GPT-3.5-turbo 61.99% v. GPT-4 28.82%, Table 2) in predicting deceased outcomes. GPT-4 often resorts to generic advice like 'consult with healthcare providers', which might be insufficient for accurate medical guidance. The challenge lies in balancing caution with the need for precise, high-stakes predictions. Additionally, GPT-4's longer response times and higher operating costs (as of this writing, the cost of GPT-4 is approximately 30 times higher than that of GPT-3.5-turbo) limit its practical utility in real-world scenarios. In practice, our expected wait time to not trigger OpenAI's API error is ~2 s for GPT-3.5-turbo, and ~15 s for GPT-4.

## Discussion

This study focuses on illustrating bias in LLMs, such as GPT-3.5-turbo and GPT-4. Transformer-based models, including GPTs[2], generate text based on previous tokens, meaning altering one token or the language prior can change subsequent token distributions. Although OpenAI has implemented RLHF to discourage problematic outputs in LLMs[2,15], our findings indicate that these models still exhibit inherent biases, especially in relation to racial and ethnic groups.

Moreover, our study highlights that LLMs can perform poorly in critical scenarios. Both GPT variants display a low level of performance when predicting death outcomes, with GPT-4's accuracy in predicting deceased outcomes only 28.82% compared to 61.99% for GPT-3.5-turbo (Table 2). These observations suggest limitations in the training data or model architecture that may not adequately capture the predictors of mortality.

This performance discrepancy raises important questions regarding the representation of negative outcomes within the training corpus and the generalizability of these models to diverse medical scenarios. It may also reflect a broader issue in machine learning models' training processes, where rare outcomes, such as death, are underrepresented or less emphasized compared to more common outcomes. To better understand the root causes of this poor performance, future research should explore a more balanced dataset that includes a wider range of medical outcomes. Additionally, it would be valuable to investigate whether safety alignment mechanisms in LLMs contribute to a lower likelihood of predicting undesirable outcomes, such as death, by comparing models with and without these alignments. Such research could provide crucial insights into how these models are trained and optimized, potentially revealing whether the emphasis on safety leads to an underrepresentation of negative outcomes.

Our findings on LLM bias mirror real-world healthcare disparities in diagnoses and spending. Prior statistic has shown that in the United States, white population has the highest estimated per-person spending, followed by Black, Hispanic and Asian[16], and there is a substantial spending gap between white population and Black or Asian[17,18]. Data from the CDC and HHS reveals that among patients diagnosed with TB, there is a higher representation of individuals of Asian ethnicity compared to the other two demographic groups[19], and the Black population exhibits a higher prevalence among patients diagnosed with HIV[20]. The model's biased behavior aligns with existing disparities and diagnostic patterns in real-world healthcare.

This study, which mainly examines racial bias in two common GPT models with a specific focus on GPT-3.5-turbo, is subject to several limitations. Firstly, it does not draw definitive conclusions about racial and ethnic groups' relevance in disease diagnosis and treatment. While race-adjusted diagnoses are criticized for contributing to healthcare disparities, many disease risk assessments still consider race. Second, our analysis contrasts model-generated content for different racial and ethnicity groups rather than comparing it with ground truth, as the LLMs used are not domain-specific and may not provide accurate projections. Moreover, this study does not incorporate actual clinical records, although the information utilized in this work originated from case reports that detail genuine patient scenarios. Finally, we focused on 400 patient cases in this work, partly due to the cost and limitations of manual qualitative analysis. Nonetheless, the volume of our data surpasses that of comparable studies in the past[8].

In response to rapid advancements in LLM development, our future work aims to expand the scope of our research by incorporating a diverse array of models and more comprehensive datasets. We will release the code used in this study to generate and compare responses, facilitating further investigations that can leverage different LLMs and adapt to evolving technologies. While evaluating bias in common medical conditions is crucial, the scarcity of publicly available real-patient datasets and resource constraints limit the scope of potential experiments. To overcome these challenges, future research should aim to include a wider range of models and datasets, which will further enhance the understanding of bias in LLMs.

In conclusion, our study reveals biases in medical report generation by the latest GPT models. Key biases include generating biased patient backgrounds, associating diseases with specific demographic groups, favoring white patients in treatment recommendations, and showing disparities in projected cost, hospitalization duration and prognosis. Additionally, both GPT-3.5-turbo and GPT-4 models show a tendency towards overly optimistic patient outcomes, with GPT-4 often predicting higher survival rates. This underscores the need to delineate safe and ambiguous language model uses. Although filtering harmful outputs mitigates biases, it is vital to address the deeper issue of inherent bias in the models' language distribution.

## Data availability

The PMC-Patients dataset is available at https://github.com/zhao-zy15/PMC-Patients[12]. PMC-Patients data set can be freely downloaded without any data usage agreement. All numerical results underlying the graphs (source data) can be found in Supplementary Data 1.

## Code availability

The code to reproduce the experiments in this work can be accessed at https://doi.org/10.5281/zenodo.12768311[14].

## References

1. Ouyang, L. et al. Training language models to follow instructions with human feedback. in *Advances in Neural Information Processing Systems* (eds. Koyejo, S. et al.) **35**, 27730–27744 (Curran Associates, Inc., 2022).
2. OpenAI. GPT-4 Technical Report. Preprint at http://arxiv.org/abs/2303.08774 (2023).
3. Jin, Q., Wang, Z., Floudas, C. S., Sun, J. & Lu, Z. Matching Patients to Clinical Trials with Large Language Models. Preprint at https://doi.org/10.48550/arXiv.2307.15051 (2023).
4. Tian, S. et al. Opportunities and challenges for ChatGPT and large language models in biomedicine and health. *Brief. Bioinform.* **25**, bbad493 (2024).
5. Zhuo, T. Y., Huang, Y., Chen, C. & Xing, Z. Red teaming ChatGPT via Jailbreaking: Bias, Robustness, Reliability and Toxicity. Preprint at https://doi.org/10.48550/arXiv.2301.12867 (2023).
6. Wei, A., Haghtalab, N. & Steinhardt, J. Jailbroken: How Does LLM Safety Training Fail? *Advances in Neural Information Processing Systems* **36**, 80079–80110 (2023).
7. Hanna, J. J., Wakene, A. D., Lehmann, C. U. & Medford, R. J. Assessing Racial and Ethnic Bias in Text Generation for Healthcare-Related Tasks by ChatGPT1. *medRxiv*, https://doi.org/10.1101/2023.08.28.23294730 (2023).
8. Omiye, J. A., Lester, J. C., Spichak, S., Rotemberg, V. & Daneshjou, R. Large language models propagate race-based medicine. *Npj Digit. Med.* **6**, 1–4 (2023).
9. Zack, T. et al. Assessing the potential of GPT-4 to perpetuate racial and gender biases in health care: a model evaluation study. *Lancet Digit. Health* **6**, e12–e22 (2024).
10. Quach, K. Healthcare org uses OpenAI's GPT-4 to write medical records. https://www.theregister.com/2023/06/06/carbon_health_deploys_gpt4powered_tools/ (2023).
11. Sun, Z. et al. Evaluating GPT-4 on Impressions Generation in Radiology Reports. *Radiology* **307**, e231259 (2023).
12. Zhao, Z., Jin, Q., Chen, F., Peng, T. & Yu, S. A large-scale dataset of patient summaries for retrieval-based clinical decision support systems. *Sci. Data* **10**, 909 (2023).
13. Guidelines To Writing A Clinical Case Report. Heart Views Off. *J. Gulf Heart Assoc.* **18**, 104–105. https://journals.lww.com/hrtv/fulltext/2017/18030/guidelines_to_writing_a_clinical_case_report.7.aspx (2017).
14. Jin, Q. ncbi-nlp/Unmasking-GPT-Bias: Code Repository for Unmasking and Quantifying Racial Bias of Large Language Models in Medical Report Generation. *Zenodo* https://doi.org/10.5281/zenodo.12768311 (2024).
15. Nori, H., King, N., McKinney, S. M., Carignan, D. & Horvitz, E. Capabilities of GPT-4 on Medical Challenge Problems. Preprint at https://doi.org/10.48550/arXiv.2303.13375 (2023).
16. Dieleman, J. L. et al. US Health Care Spending by Race and Ethnicity, 2002-2016. *JAMA* **326**, 649–659 (2021).
17. Dickman, S. L. et al. Trends in Health Care Use Among Black and White Persons in the US, 1963-2019. *JAMA Netw. Open* **5**, e2217383 (2022).
18. Chen, J., Vargas-Bustamante, A. & Ortega, A. N. Health Care Expenditures Among Asian American Subgroups. *Med. Care Res. Rev. MCRR* **70**, 310–329 (2013).
19. Table 2 | Reported TB in the US 2020 | Data & Statistics | TB | CDC. https://www.cdc.gov/tb/statistics/reports/2020/table2.htm (2023).
20. CDC. HIV in the United States by Race/Ethnicity. Centers for Disease Control and Prevention, https://www.cdc.gov/hiv/group/racialethnic/other-races/index.html (2023).

## Acknowledgements

## Author contributions

Study concepts/study design, Y.Y., Z.L.; manuscript drafting or manuscript revision for important intellectual content, all authors; approval of the final version of the submitted manuscript, all authors; agrees to ensure any questions related to the work are appropriately resolved, all authors; literature research, Y.Y.; experimental studies, human annotation, Y.Y., X.L., Q.J.; data interpretation and statistical analysis, Y.Y., X.L., Q.J.; and manuscript editing, all authors.

## Competing interests

The authors declare no competing interests.

## Additional information

**Supplementary information** The online version contains supplementary material available at https://doi.org/10.1038/s43856-024-00601-z.

**Correspondence** and requests for materials should be addressed to Zhiyong Lu.

**Peer review information** *Communications Medicine* thanks Janna Hastings, and the other, anonymous, reviewer for their contribution to the peer review of this work. Peer reviewer reports are available.

**Reprints and permissions information** is available at http://www.nature.com/reprints

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.