



# Auditing large language models for race & gender disparities: Implications for artificial intelligence-based hiring

## Authors

**Johann D. Gaebler**   
Harvard University,  
Cambridge, MA, USA

**Sharad Goel**   
Harvard University,  
Cambridge, MA, USA

**Aziz Huq**  
University of Chicago,  
Chicago, IL, USA

**Prasanna Tambe**   
University of Pennsylvania,  
Philadelphia, PA, USA

## Corresponding author:

Johann D. Gaebler, Harvard  
University, 1 Oxford Street,  
Cambridge, MA 02138,  
USA.

Email: [jgaebler@fas.harvard.edu](mailto:jgaebler@fas.harvard.edu)

## Keywords

algorithms, artificial  
intelligence, machine  
learning, large language  
models, bias, algorithm  
audits, correspondence  
experiments

Behavioral Science & Policy  
1–10

© Behavioral Science  
and Policy Association 2025  
Article reuse guidelines:  
[sagepub.com/journals-permissions](https://sagepub.com/journals-permissions)  
DOI: 10.1177/23794607251320229  
[journals.sagepub.com/home/bsx](https://journals.sagepub.com/home/bsx)

## Abstract

Rapid advances in artificial intelligence (AI), including large language models (LLMs) with abilities that rival those of human experts on a wide array of tasks, are reshaping how people make important decisions. At the same time, critics worry that LLMs may inadvertently discriminate against some groups. To address these concerns, recent regulations call for auditing the LLMs used in important decisions such as hiring. But neither current regulations nor the scientific literature offers clear guidance on how to conduct these audits. In this article, we propose and investigate one approach for auditing algorithms: *correspondence experiments*, a widely applied tool for detecting bias in human judgments. We applied this method to a range of LLMs instructed to rate job candidates using a novel data set of job applications for K-12 teaching positions in a large American public school district. By altering the application materials to imply that candidates are members of specific demographic groups, we measured the extent to which race and gender influenced the LLMs' ratings of the candidates' suitability. We found moderate race and gender disparities, with the models slightly favoring women and non-White candidates. This pattern persisted across several variations in our experiment. It is unclear what might be driving these disparities, but we hypothesize that they stem from posttraining efforts, which are part of the LLM training process and intended to correct biases in these models. We conclude by discussing the limitations of correspondence experiments for auditing algorithms.

Artificial intelligence (AI) tools increasingly assist employers with many aspects of decision-making. One prominent example is human resources (HR) management, an area in which AI tools have been used to facilitate benefits administration, coaching, development, and applicant screening. Prior to 2022, the global market for AI hiring technology was already growing rapidly. But the advent of large language models

(LLMs)—AI models that are highly adept at understanding and generating text—has dramatically boosted interest in AI hiring tools.

LLMs are a potential boon for any setting where decisions must be made on the basis of a large volume of text. LLMs can be used to do work ranging from grading student essays to evaluating proposals from potential



vendors. In the context of HR, LLMs can ingest entire application dossiers—including résumés, essays, and interview transcripts—to produce seemingly cogent assessments of candidates' qualifications. But using LLMs in this way can be worrisome as well. Even as people put these LLMs to work in hiring, employers and policymakers are racing to establish guidelines for the algorithmic evaluation of candidates.

Bias is a particular concern. Traditional AI tools are “supervised,” meaning users train these programs by inputting labeled data—for instance, previous candidates' application materials labeled with the ratings recruiters gave them. An AI model can learn statistical patterns from these labeled data and then use those patterns to generate predictions about new applicants. But LLMs go through a more complex and opaque training process that often does not involve any prelabeled data. In an unsupervised process called *pretraining*, LLMs examine and find patterns in a huge, unlabeled text corpus or data set. These training data sets may contain the equivalent of tens of billions of books<sup>1</sup> and are drawn largely from data publicly available on the internet. After being pretrained, LLMs are *posttrained* by being provided with a smaller, carefully curated set of data meant to enable them to learn patterns that will improve the accuracy, helpfulness, and safety of their outputs. Once an LLM is trained, users can deliver *prompts* (written questions and instructions or other text), such as a directive to the LLM to rank job candidates on the basis of their written application materials. Compared with traditional AI tools, the LLMs' responses more closely resemble those of human evaluators, who might, when evaluating job candidates, produce candidate ratings that are based on an intuitive understanding of how professional experiences and responses to interview questions relate to a candidate's competency and fit for a job. But for various reasons—including discriminatory content in the pretraining corpus and the complexity of the overall training processes—LLMs might also produce discriminatory or distorted responses that are hard to anticipate.

Employers and policymakers therefore fear that AI could run afoul of employment discrimination laws or otherwise produce unintended, undesirable effects, and they want to find ways to audit these tools to determine whether they are discriminatory. In this article, we demonstrate and examine one potential method for auditing LLMs—*correspondence experiments*—illustrating the approach by applying it to evaluate LLMs that could be used in making hiring decisions. We demonstrate the technique's potential and discuss limitations.

## Background

The ethical and legal implications of using AI tools in high-stakes settings, such as HR, have motivated much academic work.<sup>2,3</sup> And policymakers have become as interested as firms and researchers, introducing a wave of legislation governing the use of algorithms in different contexts, especially hiring. These measures incorporate auditing requirements as an important element, in part because of the belief that informative audits help regulators protect the public from the potential harms of AI tools and that users and firms can use audit results to make informed decisions about which tools to deploy and how to deploy them.

For example, on October 30, 2023, President Joe Biden issued an executive order imposing mandates on federal agencies' use of AI and calling for regulatory efforts by agencies with authority over private AI uses.<sup>4</sup> The order twice mentions audits as tools to advance fair public decision-making and to ensure AI safety. Meanwhile, the Digital Services Act (DSA) in Europe, which entered into force in February 2024, requires what it terms “very large online platforms” to conduct audits to promote transparency and accountability.<sup>5</sup> And in March 2024, the European Parliament adopted the more general AI Act, which imposes a variety of requirements that potentially involve audits. Most notably, it compels users of *high-risk* systems—for instance, systems used in critical infrastructure, biometric applications, law enforcement, and certain other high-stakes domains—to create quality management systems that may include a kind of audit.<sup>6</sup>

The most developed audit mandate for AI HR tools was introduced in the United States under a New York City ordinance, Local Law (LL) 144, effective July 2023.<sup>7</sup> LL 144 requires a *bias audit* when employers use “any computational process, derived from machine learning, statistical modeling, data analytics, or artificial intelligence” to classify or recommend persons for employment. In these bias audits, independent third parties must calculate and publicly report an *adverse impact ratio*. This is defined by LL 144 as the rate at which individuals in a race or gender group are hired or move forward in the hiring process relative to people in the most frequently selected race or gender group. LL 144 imposes no legal obligations when a disparity is identified—in other words, no particular course of action is required. Employers also have broad discretion to determine whether they are covered by the measure. In the first 6 months after LL 144's entry into force, only 19 audits linked to the law were published.<sup>8</sup>

The adverse impact ratio is a relatively common auditing tool in both policymaking and scientific research—required

not only by LL 144 but also by traditional hiring regulations such as the U.S. Equal Employment Opportunity Commission's four-fifths rule (which holds that selecting candidates from a protected or minority group less than 80% as often as members of the group with the highest selection rate can be considered evidence of discrimination). But the adverse impact ratio is known to be imperfect: It cannot reveal whether disproportionate selection of one kind of candidate over another indicates a true psychological bias or if it simply reflects average differences in the qualifications of applicants from different groups.

Beyond the adverse impact ratio, relatively few tools exist for identifying potentially biased decision-making by LLMs. To date, no scientific consensus exists on how best to audit algorithms for bias, although researchers have proposed a wide variety of algorithmic fairness metrics.<sup>9–26</sup> Here, we explore using correspondence experiments to audit LLMs for race and gender bias in high-stakes decision settings such as HR, in part because in behavioral science research, correspondence experiments have proven valuable for identifying discrimination in hiring decisions.

In correspondence experiments (also known as *audit studies*), researchers assume that two otherwise identical individuals from different demographic groups should receive similar treatment and that divergence is evidence of improper discrimination.<sup>27</sup> These studies are typically set in contexts where decision-makers learn about individuals exclusively through written documents (for example, an initial screening of job applicants). Researchers then experimentally manipulate elements of those materials that could suggest an individual's race and gender. (For simplicity, throughout the balance of this article, we use the term *race* to refer to race or ethnicity.) In the study described here, we mainly assessed the effect of the name applicants use on their résumé, a practice in line with most research on this topic in the social sciences.<sup>28</sup> Names are the strongest signal of race typically perceived by recruiters in the United States. In other contexts or in countries with different hiring conventions, auditors might manipulate other application elements, such as the applicant photos included with résumés in Germany, Japan, China, and many other countries. (See note A for more information on the rationale for our decision.)

For at least 50 years, social scientists and government agencies have used correspondence experiments to study discrimination in hiring,<sup>29</sup> housing,<sup>30</sup> prosecutorial charging decisions,<sup>31</sup> and other domains.<sup>32–34</sup> More recently, correspondence experiments have been proposed to similarly identify evidence of bias in the algorithms used in

AI.<sup>35–45</sup> For example, Latanya Sweeney found that Google searches of Black-sounding names were more likely to generate advertisements suggesting the named individual had an arrest record than were comparable searches of White-sounding names.<sup>46</sup>

In our exploration of the potential value of using correspondence experiments to audit AI-based hiring decisions, we found overall that correspondence experiments are useful for identifying race and gender bias, and we provide a workable strategy for carrying out mandated algorithm audits. However, as we explain in the Discussion section, we also identified some key conceptual and technical limitations of this approach for auditing algorithms.

## Empirical Analysis & Results

### LLMs for Candidate Evaluation

To evaluate the use of correspondence experiments for auditing hiring algorithms, we first gathered a novel corpus from 1,373 applications to K-12 teaching positions in a large public school district in Texas by filing a public records request. (To our knowledge, this school district does not use algorithms to evaluate applicants.) Application materials included the applicants' résumés as well as transcripts that we produced from self-recorded videos from the applicants. In these videos, applicants answered written questions about previous teaching experience, teaching style, hypothetical classroom situations, and other job-related subjects. Ultimately, we restricted our analysis to the 801 applicants who had provided both a résumé and video responses. These applicants represented a diverse pool, of which 67% were women, 2% were Asian, 45% were Black, 10% were Hispanic, and 38% were White. (We use "White" throughout to mean non-Hispanic White.)

For each applicant, we provided the LLM with (a) a description of the requirements for the teaching position based on a job posting from the school district; (b) the applicant's résumé; (c) a written transcript of the applicant's self-recorded responses to interview questions; (d) a request for the model to summarize the applicant's qualifications in prose; and (e) a request for the model to provide numerical evaluations, on a scale from 1 to 5, on several measures, including the applicant's experience, professionalism, and fit, as well as the model's overall hiring recommendation, which was expressed as a score on a scale ranging from 1 (*definitely do not hire*) to 5 (*definitely hire*). (See the Supplemental Material for more information.) Although our primary statistical analysis focused on the models' overall numerical hiring recommendation for the applicants, we

requested additional information from the models—including written summaries and other scores—to improve the quality of the LLM results, a common strategy when using these models.

We audited 11 LLMs: OpenAI’s GPT-3.5, GPT-4, GPT-4o, and GPT-4o Mini models;<sup>47–50</sup> Mistral’s Mistral 7B and Mixtral 8x7B models;<sup>51,52</sup> and Anthropic’s Claude Instant, Claude 2, Claude 3 Haiku, Claude 3 Sonnet, and Claude 3.5 Sonnet models.<sup>53–55</sup> Experts considered the OpenAI and Anthropic models to be the best of their kind available at the time of our evaluation. Mistral’s models are popular open-source competitors. We did not formally assess how well each LLM’s candidate ratings aligned with candidate qualifications, but an informal inspection suggested that the highest rated candidates generally had more experience and gave more polished responses to interview questions than did those receiving lower ratings. This rough assessment of the LLM ratings’ validity—and the LLMs’ relative ease of use—makes it likely that a pipeline like the one we implemented here will soon be used by employers to screen applicants, if one has not been launched already.

### Assessing Adverse Impact Ratios

As we have previously noted, each LLM rated candidates’ overall suitability for the job on a scale of 1 to 5, where a higher number indicates a stronger positive recommendation. OpenAI’s GPT-3.5, one of the most popular and widely used LLMs at the time of our experiments, gave 20% of candidates an overall score of 5, 51% a 4, 24% a 3, 4% a 2, and 1% a 1.

We then turned to the issue of whether each LLM rated candidates similarly across demographic groups. We first looked at the adverse impact ratio, given its prominence in past research and the New York law. That is, we determined the rate at which individuals in one race or gender category were positively selected relative to those in another category. Per the Equal Employment Opportunity Commission’s four-fifths rule, a ratio of 80% or lower is particularly concerning and likely warrants some response, although federal law does not require a particular course of action.<sup>56</sup>

The results of this adverse impact ratio analysis for OpenAI’s GPT-3.5 are shown in Figure 1; results across all models were similar. (See Figures S7–S10 in the Supplemental Material for results from other models.) When comparing the proportion of applicants across groups who received a 5, the highest threshold for recommendation, we found that female applicants received positive assessments more often than male applicants did and Black and Hispanic applicants received positive assessments more often than

White applicants did. At a recommendation threshold of 4, the pattern flips for race—with White applicants receiving a positive assessment more often than Black and Hispanic applicants did—and we found near parity for gender. Finally, at a threshold of 3, we found near parity across both gender and race groups.

This simple analysis suggests that GPT-3.5 and the other LLMs we studied might be favoring certain demographic groups. However, even with the comparatively large minority applicant pools in our corpus, we can estimate adverse impact ratios only imprecisely. In fact, of nine adverse impact ratios calculated at this stage, only three were statistically significant, suggesting this approach found only slight evidence of meaningful disparities.

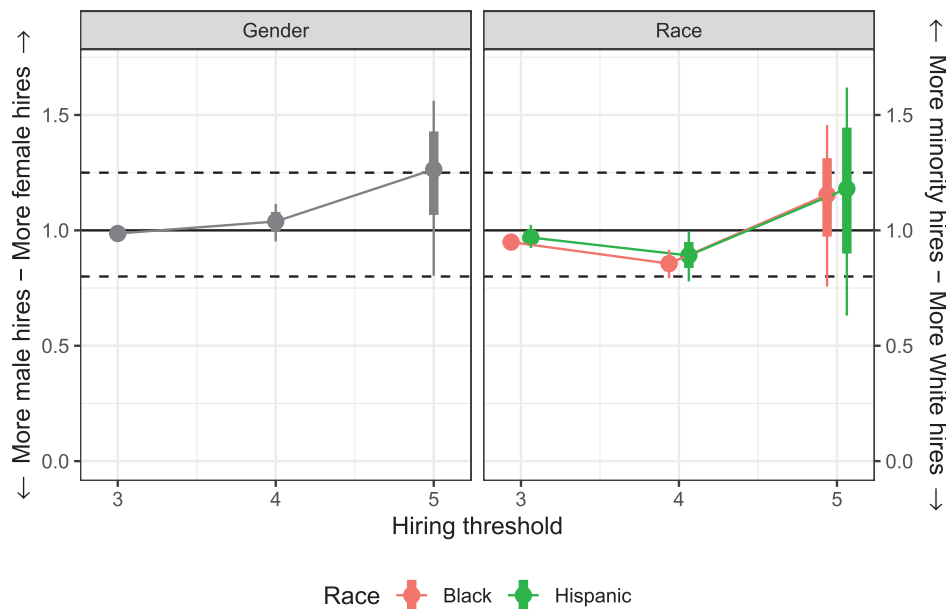
Without further evidence, we cannot definitively say whether these disparities are due to algorithmic bias or group-specific differences in the applicant pool. Put another way, the women and racial minorities in our applicant pool might simply be more qualified than the other applicants for these positions, which would explain why the model ratings they received were higher than the ratings received by the other applicants. This makes it difficult to draw clear conclusions, which indicates that an adverse impact ratio analysis alone is too limited for auditing LLMs.

### Assessing Correspondence Experiments

Correspondence experiments make it possible to differentiate between algorithmic bias and differences in candidate qualifications. We began by manipulating the real application materials to create synthetic application dossiers that differed only in details that strongly signaled an applicant’s race or gender. For each real applicant, we generated eight synthetic applications, corresponding to a particular race (Asian, Black, Hispanic, or White) and gender (female or male). We then replaced the applicant’s actual name throughout the materials with one that strongly signaled membership in that group. We similarly changed any mention of the applicant’s pronouns in the materials to match the assigned group, as well as other indicia of race or gender. (See the Supplemental Material for a detailed description of how we generated these synthetic applications.)

We presented the synthetic application materials to the LLMs, instructing the models to report the race and gender of the synthetic applicants. We found generally high agreement between the race and gender we intended to convey and the model’s perception of these attributes, which in most cases exceeded 90% but with the precise level of agreement varying from case to case (see Figures S1 and S2 in the Supplemental Material). This level of agreement is

Figure 1. GPT-3.5's adverse impact ratios highlighting disparities that may be the result of algorithmic bias while nevertheless being inconclusive on their own



*Note.* We calculated adverse impact ratios for the outputs of the large language model GPT-3.5 at a variety of hiring thresholds. The y-axis shows the ratio of female to male (left) and Black and Hispanic to White (right) applicants who would have been hired if one hired all the individuals to whom the model gave a rating of at least 3, 4, or 5, respectively. (A rating of 5 denotes the strongest endorsement for hiring.) At the lowest threshold, we observed near parity across both race and gender; however, at higher thresholds, we found some evidence of disparities in hiring rates across demographic groups, although the estimates are imprecise. Dots denote the mean adverse impact ratio, and the thick and thin bars indicate 70% and 95% confidence intervals, respectively. The dashed line indicates the standard set by the U.S. Equal Employment Opportunity Commission's four-fifths rule; hiring ratios exceeding these bounds are considered particularly concerning. However, it is not possible to know on the basis of this analysis alone whether the adverse impact ratio reflects differences in applicant quality or bias in the model's selection. (The confidence intervals were calculated using 1,000 pivotal bootstrap resamples.)

comparable to what has been found in studies that manipulate human perceptions of race by altering names in résumés.<sup>29</sup>

Finally, we asked our LLMs to provide hiring recommendations for the synthetic candidates. The results of these correspondence experiments are shown in Figure 2. (Results for Mixtral 8x7B should be interpreted with caution, because the model often failed to follow the instructions, producing responses without ratings; see the Supplemental Material for details.) For each model, we created comparisons of groups of candidates relative to a historically favored reference group. To assess gender bias, we made men a reference group and looked at the estimated difference in average score for women versus men in each model. In the case of race, we made White people the reference group and looked for the difference in average scores for all other groups in comparison with the average

score of these White candidates. (See the Supplemental Material for further methodological details.)

Across models, we found that the LLMs rated synthetic female candidates moderately higher than they did the synthetic male candidates. Models also generally rated synthetic Black, Hispanic, and Asian candidates moderately higher than they did synthetic White candidates, although we found more variation between models, with Mistral's models exhibiting smaller disparities.

Our correspondence experiments suggest that the models' perceptions of race and gender influence algorithmic candidate assessments, at least to some degree. At each hiring threshold—the minimum rating (3, 4, or 5) at which we imagine hiring candidates—race and gender disparities were generally a few percentage points (see Figures S14 and S15 in the Supplemental Material). These disparities were



Figure 2. Correspondence experiments showing that the large language models tested modestly favor women & non-White candidates



*Note.* For each large language model, the colored dots indicate the difference in mean applicant rating, reported in estimated population standard deviations, with 70% and 95% confidence intervals indicated by thick and thin bars, respectively. Positive values indicate that the model rated female or racial minority applicants higher than it rated male (left panels) or White (right panels) applicants, respectively. The dashed reference line shows the value corresponding to no average difference in ratings between groups. Most large language models showed the same pattern, with the exception of Mixtral 8x7B, a model that frequently failed to follow user instructions. (The confidence intervals are clustered at the level of the real application dossier used to generate the synthetic application.)

modest but within the typical range identified in recent studies of human recruiters (for an example, see Reference 27). In other words, the LLMs showed a level of bias smaller than but still comparable to what people tasked with identifying the best job applicants in similar experiments often show.

### Sensitivity to Prompt Variation & Context

An LLM's approach to a given task—such as rating job applicants—can be modified and potentially improved by changing the directions the LLM is given. So we repeated our analysis with several variants of the initial prompt—that is, the set of directions we gave an LLM for evaluating our applicants and their suitability for the job. (See the Supplemental Material for more details on the prompt variants investigated.) For simplicity, we ran these robustness tests on only one model, OpenAI's GPT-3.5, which exhibited

approximately average disparities in our primary analysis. Regardless of the prompt used, we found the same general patterns persisted across all variants: a slight bias in favor of women over men and in favor of Hispanic, Black, and Asian applicants over White applicants. (See Figures S3 and S4 in the Supplemental Material.)

We then reran our primary analysis while inputting only an applicant's résumé into the model, excluding interview transcripts. This variation on our experiment allowed us to gauge a scenario that frequently occurs in real-world settings, when employers must make decisions based on résumés alone. As with the robustness tests, we conducted this analysis on only GPT-3.5, and we again found that women and racial minority applicants received moderately higher scores than did men and White applicants, respectively.

Many LLMs can access virtually encyclopedic knowledge and incorporate this information into their evaluations. For example, in our analysis, the American school district to which our candidates applied is especially racially diverse—a fact known by the LLMs we examined. If a model took this information into account, the diversity might influence its ratings in some way. We therefore ran another analysis on GPT-3.5 in which we replaced all mentions of the district (and the city and state) with the name of a predominately White school district in West Virginia. Once again, we found disparities mirroring our primary results.

## Discussion

Our results demonstrate that correspondence experiments can reveal race and gender disparities in LLM outputs. Unlike the commonly used adverse impact ratio analysis, correspondence experiments revealed not only differences in qualifications across groups in the candidate pool but also the effect of race and gender on LLM outputs. These experiments therefore offer policymakers a potentially useful tool for auditing algorithms. The patterns we observed were substantive and robust, persisting across several variations of our study, including changes to the instructions provided to the model and the specific application materials we inputted.

It bears emphasizing that the specific pattern we observed—with models favoring female applicants over male applicants and favoring Black, Hispanic, and Asian applicants over White applicants—may not generalize across contexts in which people use LLMs. Indeed, although some recent studies that focused on hiring decisions also found disparities similar to those we found in our work,<sup>36,43</sup> others have reported disparities in the opposite direction.<sup>35,38–40,42,44</sup> Such contrasting results are not surprising given the complex and often inscrutable ways in which LLMs are trained. In particular, recall that developers typically fine-tune models in a final alignment or posttraining phase in part to avoid mirroring overt discrimination in the training data. But this step may leave traces of bias in ways that are hard to predict. For example, efforts to mitigate discriminatory associations the model learned in pretraining might overshoot the mark, causing a distortion in the other direction.

Although correspondence experiments may be a useful tool for auditing algorithms, they also have notable limitations. First, all experiments—correspondence or otherwise—are inherently limited in their ability to manipulate race or gender while leaving other characteristics of people untouched. In the case of our study, names are imperfect

signals of race and gender, and other aspects of the application dossier (such as history of employment in a male- or female-dominated field) can further attenuate that signal. Indeed, in our own pool of applicants, we find that LLM-like statistical models can nearly always infer an applicant's race and gender from unaltered application materials as well as from masked application materials (that is, application materials with names and pronouns removed; see Figure S6 in the Supplemental Material). Nevertheless, as already noted, the models still identified the intended race and gender of our synthetic candidates more than 90% of the time, indicating that our name manipulation largely worked as designed.

Further, names may reveal information beyond applicants' race and gender. Changing an individual's name can also change an LLM's perception of their age, their socioeconomic status, and other characteristics aside from race. This effect means we cannot be certain that we have accounted for all the ways in which LLMs rated our applicants or the precise degree to which race and gender alone could have influenced these decisions. These confounds limit our statistical conclusions.

Conceptually, it is challenging to rigorously define what it would even mean to manipulate an LLM's perception of an individual applicant's race in isolation from other factors. Considerable literature has explored the question of whether it is possible to discuss and perceive race in isolation, given the ways in which ideas around racial identity so often couple with perceptions of social status, wealth, religion, neighborhood of residence, and other variables (for examples, see references 57 and 61). For that reason, some scholars conceive of race as a “bundle of sticks,” lumping all of those elements together, along with skin color and ancestral origin, when determining an individual's racial identity.<sup>57,62</sup> Similar considerations apply to gender.<sup>59</sup> Which of these factors can or should be manipulated in an audit of race or gender bias is a difficult question.

Finally, even if models are not directly influenced by an individual's race or gender in the way we test here, biases against other aspects of a candidate's application could still produce undesirable outcomes. For instance, a model might, in theory, prioritize the applications of individuals who attended private school, regardless of their race or other qualifications. That pattern could unjustly disadvantage qualified minority applicants.<sup>63</sup>

Ultimately, LLMs need to be audited for each specific use, and the conclusions drawn could vary on the basis of both the task and the pool of individuals evaluated. Context- and

scenario-specific correspondence experiments can serve as tools for evaluating LLMs in many situations. For example, LLMs are already being used to determine credit risk, and policymakers keen to craft legislation and guidance for this use would benefit from the data generated by correspondence experiment audits of credit risk LLM outputs. The calls to audit algorithms are likely to increase as LLMs become even more capable and widespread and are used in ever more varied decision-making contexts. Correspondence experiments, despite their important limitations, represent one promising method of auditing algorithms for race and gender bias. We hope our work aids in the ongoing regulatory efforts to ensure that LLMs yield equitable outcomes.

### Author Note

We are grateful to Avi Bagchi and Marissa Gerchick for research assistance. Data and replication code are available at <https://github.com/jgaeb/llm-audit>.

### Declaration of Conflicting Interests


The author(s) declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.

### Funding

The author(s) disclosed receipt of the following financial support for the research, authorship, and/or publication of this article: The

authors received financial support from the Wharton AI & Analytics Initiative. Aziz Huq acknowledges support from the Frank J. Cicero Fund.

### ORCID iDs

Johann D. Gaebler  <https://orcid.org/0000-0003-3340-9542>

Sharad Goel  <https://orcid.org/0000-0002-6103-9318>

Prasanna Tambe  <https://orcid.org/0000-0002-5308-1057>

### Supplemental Material

Supplemental material for this article is available at <https://doi.org/10.1177/23794607251320229>.

### Note

- A. In general, an auditor's main concern when designing a correspondence experiment is for manipulated applications to maintain consistency. If they do not—for instance, if the auditor pairs a male name with an application that uses female pronouns to describe the applicant or pairs a picture of a White applicant with a résumé that lists attendance at a historically Black college or university—differences in recruiters' behavior could reflect their reactions to the incongruity rather than their perception of the manipulated demographics. In our case, any manipulation in an experiment being designed to measure the effect of gender would need to alter names, because in the United States, names are typically very strong gender signals; name manipulation in this country is both necessary and sufficient for altering an applicant's apparent gender (and, to a lesser extent, race).

### References

1. Meta. (2024, April 18). Introducing Meta Llama 3: The most capable openly available LLM to date. *AI at Meta Blog*. <https://ai.meta.com/blog/meta-llama-3/>
2. Tambe, P., Cappelli, P., & Yakubovich, V. (2019). Artificial intelligence in human resources management: Challenges and a path forward. *California Management Review*, 61(4), 15–42. <https://doi.org/10.1177/0008125619867910>
3. Raghavan, M., Barocas, S., Kleinberg, J., & Levy, K. (2020). Mitigating bias in algorithmic hiring: Evaluating claims and practices. In *FAT\* '20: Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency* (pp. 469–481). Association for Computing Machinery. <https://doi.org/10.1145/3351095.3372828>
4. Exec. Order No. 14110, 88 C. F. R. 75191. (2023). <https://www.federalregister.gov/documents/2023/11/01/2023-24283/safe-secure-and-trustworthy-development-and-use-of-artificial-intelligence>
5. European Commission. (2024). *The Digital Services Act package*. European Union. <https://digital-strategy.ec.europa.eu/en/policies/digital-services-act-package>
6. Directorate-General for Communications Networks, Content and Technology. (2024). *Proposal for a regulation of the European Parliament and of the council: Laying down harmonised rules on artificial intelligence (Artificial Intelligence Act) and amending certain union legislative acts* (Document 52021PC0206). European Union, European Commission. <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=celex%3A52021PC0206>
7. N.Y.C. Local Law No. 144, N.Y.C. Admin. Code §20-870 (2021).
8. Groves, L., Metcalf, J., Kennedy, A., Vecchione, B., & Strait, A. (2024). *Auditing work: Exploring the New York City algorithmic bias audit regime*. arXiv. <https://doi.org/10.48550/arXiv.2402.08101>
9. Dwork, C., Hardt, M., Pitassi, T., Reingold, O., & Zemel, R. (2012). Fairness through awareness. In *ITCS '12: Proceedings of the 3rd Innovations in Theoretical Computer Science Conference* (pp. 214–226). Association for Computing Machinery. <https://doi.org/10.1145/2090236.2090255>
10. Hardt, M., Price, E., & Srebro, N. (2016). Equality of opportunity in supervised learning. In D. D. Lee & U. von Luxburg (Eds.), *NIPS'16: Proceedings of the 30th International Conference on Neural Information Processing Systems* (pp. 3315–3323). Curran Associates. <https://dl.acm.org/doi/abs/10.5555/3157382.3157469>
11. Kleinberg, J., Mullainathan, S., & Raghavan, M. (2017). Inherent trade-offs in the fair determination of risk scores. In C. H. Papadimitriou (Ed.), *8th Innovations in Theoretical Computer Science Conference* (Vol. 67, pp. 43:1–43:23). Dagstuhl Publishing. <https://doi.org/10.4230/LIPIcs.ITCS.2017.43>
12. Corbett-Davies, S., Pierson, E., Feller, A., Goel, S., & Huq, A. (2017). Algorithmic decision making and the cost of fairness. In *KDD '17: Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 797–806). Association for Computing Machinery. <https://doi.org/10.1145/3097983.3098095>
13. Corbett-Davies, S., Gaebler, J. D., Nilforoshan, H., Shroff, R., & Goel, S. (2023). The measure and mismeasure of fairness. *The Journal of Machine Learning Research*, 24(1), 14730–14846. <https://dl.acm.org/doi/10.5555/3648699.3649011>



14. Chouldechova, A. (2017). Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big Data*, 5(2), 153–163. <https://doi.org/10.1089/big.2016.0047>
15. Kilbertus, N., Rojas-Carulla, M., Parascandolo, G., Hardt, M., Janzing, D., & Schölkopf, B. (2017). Avoiding discrimination through causal reasoning. In U. von Luxburg, I. Guyon, S. Bengio, H. Wallach, & R. Fergus (Eds.), *NIPS'17: Proceedings of the 31st International Conference on Neural Information Processing Systems* (pp. 656–666). Curran Associates. <https://dl.acm.org/doi/10.5555/3294771.3294834>
16. Kusner, M., Loftus, J., Russell, C., & Silva, R. (2017). Counterfactual fairness. In U. von Luxburg, I. Guyon, S. Bengio, H. Wallach, & R. Fergus (Eds.), *NIPS'17: Proceedings of the 31st International Conference on Neural Information Processing Systems* (pp. 4069–4076). Curran Associates. <https://dl.acm.org/doi/10.5555/3294996.3295162>
17. Zafar, M. B., Valera, I., Rodriguez, M. G., & Gummadi, K. P. (2017). Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment. In *WWW '17: Proceedings of the 26th International Conference on World Wide Web* (pp. 1171–1180). International World Wide Web Conferences Steering Committee. <https://doi.org/10.1145/3038912.3052660>
18. Loftus, J. R., Russell, C., Kusner, M. J., & Silva, R. (2018). *Causal reasoning for algorithmic fairness*. arXiv. <https://doi.org/10.48550/arXiv.1805.05859>
19. Nabi, R., & Shpitser, I. (2018). Fair inference on outcomes. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1), 1931–1940. <https://doi.org/10.1609/aaai.v32i1.11553>
20. Chiappa, S. (2019). Path-specific counterfactual fairness. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01), 7801–7808. <https://doi.org/10.1609/aaai.v33i01.33017801>
21. Wang, Y., Sridhar, D., & Blei, D. M. (2019). *Equal opportunity and affirmative action via counterfactual predictions* (Version 2). arXiv. <https://doi.org/10.48550/arXiv.1905.10870>
22. Chouldechova, A., & Roth, A. (2020). A snapshot of the frontiers of fairness in machine learning. *Communications of the ACM*, 63(5), 82–89. <https://doi.org/10.1145/3376898>
23. Coston, A., Mishler, A., Kennedy, E. H., & Chouldechova, A. (2020). Counterfactual risk assessments, evaluation, and fairness. In *FAT\* '20: Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency* (pp. 582–593). Association for Computing Machinery. <https://doi.org/10.1145/3351095.3372851>
24. Imai, K., & Jiang, Z. (2022). *Principal fairness for human and algorithmic decision-making* (Version 5). arXiv. <https://doi.org/10.48550/arXiv.2005.10400>
25. Berk, R., Heidari, H., Jabbari, S., Kearns, M., & Roth, A. (2021). Fairness in criminal justice risk assessments: The state of the art. *Sociological Methods & Research*, 50(1), 3–44. <https://doi.org/10.1177/0049124118782533>
26. Raghavan, M. (2024). What should we do when our ideas of fairness conflict? *Communications of the ACM*, 67(1), 88–97. <https://doi.org/10.1145/3587930>
27. Lippens, L., Vermeiren, S., & Baert, S. (2023). The state of hiring discrimination: A meta-analysis of (almost) all recent correspondence experiments. *European Economic Review*, 151, Article 104315. <https://doi.org/10.1016/j.eurocorev.2022.104315>
28. Gaddis, S. M. (2019). Understanding the “how” and “why” aspects of racial–ethnic discrimination: A multimethod approach to audit studies. *Sociology of Race and Ethnicity*, 5(4), 443–455. <https://doi.org/10.1177/2332649219870183>
29. Bertrand, M., & Mullainathan, S. (2004). Are Emily and Greg more employable than Lakisha and Jamal? A field experiment on labor market discrimination. *American Economic Review*, 94(4), 991–1013. <https://doi.org/10.1257/0002828042002561>
30. Wienk, R. E., Reid, C. E., Simonson, J. C., & Eggers, F. J. (1979). *Measuring racial discrimination in American housing markets: The Housing Market Practices Survey*. U.S. Department of Housing and Urban Development, Office of Policy Development and Research, Division of Evaluation. <https://www.huduser.gov/portal/publications/Measuring-Racial-Discrimination-in-American-1979.html>
31. Chohlas-Wood, A., Nudell, J., Yao, K., Lin, Z. (J.), Nyarko, J., & Goel, S. (2021). Blind justice: Algorithmically masking race in charging decisions. In *AIES '21: Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society* (pp. 35–45). Association for Computing Machinery. <https://doi.org/10.1145/3461702.3462524>
32. Ayres, I., Banaji, M., & Jolls, C. (2015). Race effects on eBay. *The RAND Journal of Economics*, 46(4), 891–917. <https://doi.org/10.1111/1756-2171.12115>
33. Gaddis, S. M., & Ghoshal, R. (2020). Searching for a roommate: A correspondence audit examining racial/ethnic and immigrant discrimination among millennials. *Socius: Sociological Research for a Dynamic World*, 6. <https://doi.org/10.1177/2378023120972287>
34. Lyons-Padilla, S., Markus, H. R., Monk, A., Radhakrishna, S., Shah, R., Dodson, N. A. “D.” IV, & Eberhardt, J. L. (2019). Race influences professional investors’ financial judgments. *Proceedings of the National Academy of Sciences, USA*, 116(35), 17225–17230. <https://doi.org/10.1073/pnas.1822052116>
35. Armstrong, L., Liu, A., MacNeil, S., & Metaxa, D. (2024). The silicon ceiling: Auditing GPT’s race and gender biases in hiring. In *EAAMO '24: Proceedings of the fourth ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization*, Article 2. Association for Computing Machinery. <https://doi.org/10.1145/3689904.3694699>
36. An, J., Huang, D., Lin, C., & Tai, M. (2024). *Measuring gender and racial biases in large language models*. arXiv. <https://doi.org/10.48550/arXiv.2403.15281>
37. Bandy, J. (2021). Problematic machine behavior: A systematic literature review of algorithm audits. *Proceedings of the ACM on Human–Computer Interaction*, 5(CSCW1), Article 74. <https://doi.org/10.1145/3449148>
38. Eloundou, T., Beutel, A., Robinson, D. G., Gu-Lemberg, K., Brakman, A.-L., Mishkin, P., Shah, M., Heidecke, J., Weng, L., & Kalai, A. T. (2024). *First-person fairness in chatbots*. arXiv. <https://doi.org/10.48550/arXiv.2410.19803>
39. Salinas, A., Haim, A., & Nyarko, J. (2025). *What’s in a name? Auditing large language models for race and gender bias* (Version 3). arXiv. <https://doi.org/10.48550/arXiv.2402.14875>
40. Lippens, L. (2024). Computer says “no”: Exploring systemic bias in ChatGPT using an audit approach. *Computers in Human Behavior: Artificial Humans*, 2(1), Article 100054. <https://doi.org/10.1016/j.chbah.2024.100054>
41. Metaxa, D., Park, J. S., Robertson, R. E., Karahalios, K., Wilson, C., Hancock, J., & Sandvig, C. (2021). Auditing algorithms: Understanding algorithmic systems from the outside in. *Foundations and Trends® in Human–Computer Interaction*, 14(4), 272–344. <https://doi.org/10.1561/11000000083>
42. Nghiem, H., Prindle, J., Zhao, J., & Daumé, H. III. (2024). “You gotta be a doctor, Lin”: An investigation of name-based bias of large language models in employment recommendations (Version 2). arXiv. <https://doi.org/10.48550/arXiv.2406.12232>
43. Tamkin, A., Askeel, A., Lovitt, L., Durmus, E., Joseph, N., Kravec, S., Nguyen, K., Kaplan, J., & Ganguli, D. (2023). *Evaluating and mitigating discrimination in language model decisions*. arXiv. <https://doi.org/10.48550/arXiv.2312.03689>
44. Veldanda, A. K., Grob, F., Thakur, S., Pearce, H., Tan, B., Karri, R., & Garg, S. (2023). *Are Emily and Greg still more employable than Lakisha and Jamal? Investigating algorithmic hiring bias in the era of ChatGPT*. arXiv. <https://doi.org/10.48550/arXiv.2310.05135>
45. Sandvig, C., Hamilton, K., Karahalios, K., & Langbort, C. (2014). Auditing algorithms: Research methods for detecting discrimination on internet platforms. *Data and Discrimination: Converting Critical Concerns Into Productive Inquiry*, 22, 4349–4357.
46. Sweeney, L. (2013). Discrimination in online ad delivery: Google ads, Black names and White names, racial discrimination, and click advertising. *Queue*, 11(3), 10–29. <https://doi.org/10.1145/2460276.2460278>
47. Brown, T. B., Mann, B., Ryder, N., Subbiah, N., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., ... Amodei, D. (2020). Language models are few-shot learners. In H. Larochelle, M. Ranzato,

- R. Hadsell, M. F. Balcan, & H. Lin (Eds.), *NIPS '20: Proceedings of the 34th International Conference on Neural Information Processing Systems* (pp. 1877–1901). Curran Associates. <https://dl.acm.org/doi/abs/10.5555/3495724.3495883>
48. OpenAI, Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., Avila, R., Babuschkin, I., Balaji, S., Balcom, V., Baltescu, P., Bao, H., Bavarian, M., Belgum, J., . . . Zoph, B. (2024). *GPT-4 technical report* (Version 6). arXiv. <https://doi.org/10.48550/arXiv.2303.08774>
  49. OpenAI. (2024, July 18). *GPT-4o mini: Advancing cost-efficient intelligence*. <https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/>
  50. OpenAI. (2024, May 13). *Hello GPT-4o*. <https://openai.com/index/hello-gpt-4o/>
  51. Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., de las Casas, D., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., Lavaud, L. R., Lachaux, M.-A., Stock, P., Le Scao, T., Lavril, T., Wang, T., Lacroix, T., & El Sayed, W. (2023). *Mistral 7B*. arXiv. <https://doi.org/10.48550/arXiv.2310.06825>
  52. Jiang, A. Q., Sablayrolles, A., Roux, A., Mensch, A., Savary, B., Bamford, C., Chaplot, D. S., de las Casas, D., Hanna, E. B., Bressand, F., Lengyel, G., Bour, G., Lample, G., Lavaud, L. R., Saulnier, L., Lachaux, M.-A., Stock, P., Subramanian, S., Yang, S., . . . El Sayed, W. (2024). *Mixtral of experts*. arXiv. <https://doi.org/10.48550/arXiv.2401.04088>
  53. Anthropic. (n.d.). *Model card and evaluations for Claude models*. Retrieved February 15, 2025, from <https://www-cdn.anthropic.com/bd2a28d2535bfb0494cc8e2a3bfl35d2e7523226/Model-Card-Claude-2.pdf>
  54. Anthropic. (n.d.). *Claude 3.5 Sonnet model card addendum*. Retrieved February 15, 2025, from [https://www-cdn.anthropic.com/fed9cc193a14b84131812372d8d5857f8f304c52/Model\\_Card\\_Claude\\_3\\_Addendum.pdf](https://www-cdn.anthropic.com/fed9cc193a14b84131812372d8d5857f8f304c52/Model_Card_Claude_3_Addendum.pdf)
  55. Anthropic. (n.d.). *The Claude 3 model family: Opus, Sonnet, Haiku*. Retrieved February 15, 2025, from [https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model\\_Card\\_Claude\\_3.pdf](https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf)
  56. Tobia, K. (2017). Disparate statistics. *The Yale Law Journal*, 126(8), 2382–2420. <http://hdl.handle.net/20.500.13051/10305>
  57. Sen, M., & Wasow, O. (2016). Race as a bundle of sticks: Designs that estimate effects of seemingly immutable characteristics. *Annual Review of Political Science*, 19, 499–522. <https://doi.org/10.1146/annurev-polisci-032015-010015>
  58. Gaebler, J., Cai, W., Basse, G., Shroff, R., Goel, S., & Hill, J. (2022). A causal framework for observational studies of discrimination. *Statistics and Public Policy*, 9(1), 26–48. <https://doi.org/10.1080/2330443X.2021.2024778>
  59. Hu, L., & Kohler-Hausmann, I. (2020). *What's sex got to do with fair machine learning?* (Version 2). arXiv. <https://doi.org/10.48550/arXiv.2006.01770>
  60. Greiner, D. J., & Rubin, D. B. (2011). Causal effects of perceived immutable characteristics. *The Review of Economics and Statistics*, 93(3), 775–785. [https://doi.org/10.1162/REST\\_a\\_00110](https://doi.org/10.1162/REST_a_00110)
  61. Holland, P. W. (1986). Statistics and causal inference. *Journal of the American Statistical Association*, 81(396), 945–960. <https://doi.org/10.2307/2289064>
  62. Kang, S. K., DeCelles, K. A., Tilcsik, A., & Jun, S. (2016). Whiteness résumés: Race and self-presentation in the labor market. *Administrative Science Quarterly*, 61(3), 469–502. <https://doi.org/10.1177/0001839216639577>
  63. Jung, J., Corbett-Davies, S., Gaebler, J. D., Shroff, R., & Goel, S. (2024). *Mitigating included- and omitted-variable bias in estimates of disparate impact* (Version 4). arXiv. <https://doi.org/10.48550/arXiv.1809.05651>