

---

# Small Changes, Large Consequences: Analyzing the Allocational Fairness of LLMs in Hiring Contexts

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Large language models (LLMs) are increasingly being deployed in high-stakes  
2 settings like hiring, yet their potential for unfair decision-making remains under-  
3 studied in generation and retrieval. In this work, we examine the allocational  
4 fairness of LLM-based hiring systems through two tasks that reflect actual HR  
5 usage (resume summarization and applicant ranking), using a synthetic resume  
6 dataset with demographic perturbations and curated job postings. Our findings  
7 reveal that generated summaries exhibit meaningful differences more frequently  
8 for race than for gender perturbations. Additionally, retrieval models exhibit high  
9 ranking sensitivity to both gender and race perturbations, and can show compara-  
10 ble sensitivity to both demographic and non-demographic changes. Overall, our  
11 results indicate that LLM-based hiring systems, especially in the retrieval stage,  
12 can exhibit notable biases that lead to discriminatory outcomes.

## 13 1 Introduction

14 Large language models (LLMs) are increasingly being adopted in high-stakes domains like hiring  
15 [1], where their decisions can directly shape career opportunities. Responsible deployment requires  
16 anticipating risks such as *allocational harms* (i.e., allocating resources or opportunities unfairly to  
17 different social groups) [2, 3], since automated hiring systems may produce unfair outcomes and  
18 reinforce systemic inequalities. While prior work has extensively examined representational harms  
19 (i.e., representing certain social groups negatively, demeaning them, or erasing their existence) in  
20 LLMs [4, 5, 6, 7, 8], allocational harms—the primary harm at play in high-stakes situations—remain  
21 understudied beyond discriminative systems.

22 The few studies that evaluate allocational harms for LLMs [9, 10, 11, 12] focus on simplified  
23 classification or prediction tasks (e.g., binary hiring decisions, salary estimates), which do not reflect  
24 real-world deployment [13]. These setups risk poor *ecological validity* [14, 15, 16], since harms  
25 must be evaluated in realistic contexts or with predictive proxies. Yet there is limited work on  
26 allocational harms in generative settings without adding a simplification layer, with [17] being a  
27 notable exception, since measuring how generated text might yield disparities is more complex than  
28 analyzing classification predictions.

29 In this work, we examine whether LLMs behave fairly in real-world hiring contexts. We focus on two  
30 core tasks that mirror how real-world usage in hiring workflows [18, 19]: (1) ranking candidates with  
31 respect to a job posting and (2) summarizing resumes, as illustrated in Figure 1. These tasks represent  
32 key stages where automation can influence which candidates are surfaced and considered for a role.  
33 To evaluate fairness, we examine model sensitivity to gender and race perturbations in resumes by

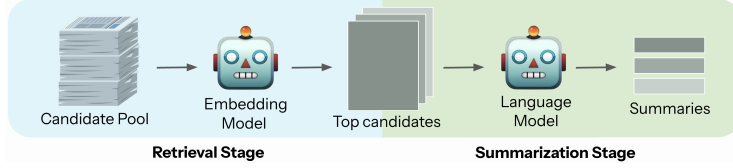


Figure 1: We investigate the fairness of an LLM hiring pipeline with a **retrieval stage** (filters the top- $n$  candidates with respect to a job posting) and a **summarization stage** (generates resume summaries for filtered candidates). We assess fairness issues at each stage separately.

asking: **RQ1**) Do generated summaries differ meaningfully across demographic groups? **RQ2**) How sensitive are model rankings to demographic and non-demographic perturbations in resumes?

To this end, we: (1) construct a new benchmark consisting of a synthetic resume dataset with controlled demographic perturbations (varying names and extracurricular content) and curated job postings, (2) design a holistic evaluation framework with fairness metrics tailored to both generative and retrieval settings, validated by an expert human preference study, (3) conduct a comprehensive fairness analysis of 10 large language models (6 generative, 4 retrieval) based on real-world hiring tasks. We will make all data and code publicly available.

Our findings show that LLM hiring systems display substantial bias, primarily in the retrieval stage. Summaries differ across racial groups in up to 20% of cases versus 3% for gender (**RQ1**), while retrieval is highly brittle, excluding up to 74% of candidates after demographic perturbations (**RQ2**). Models are also highly sensitive to non-demographic changes, indicating fairness concerns arise from general brittleness rather than demographic bias alone (**RQ2**). Overall, even small changes can yield major disparities, raising concerns about the fairness and robustness of LLMs in hiring.

## 2 Methodology and Setup

To study fairness in hiring, we consider an LLM-based pipeline with two components: resume retrieval with respect to a job post (using an embedding model) and resume summarization (using an LLM). This pipeline reflects real-world practices, as informed by interviews with corporations using LLMs for hiring.<sup>1</sup> We focus on summarization first because it is more neglected in research, though in a pipeline it would come after retrieval, as shown in Figure 1 (as summarization would be of retrieved resumes).

Since perturbed resumes are highly similar to original resumes by design, we expect resulting summaries and rankings for original and perturbed resumes to also be similar. Based on this idea, we propose two metrics to assess allocational fairness: *invariance violations* for summarization and *exclusion* for retrieval. These metrics quantify systematic differences in generated resume summaries and meaningful changes in resume rankings with respect to job postings, respectively. While [20] also studies hiring fairness in a retrieval setting, their approach does not explicitly capture how perturbing a resume impacts resume screening outcomes.

### 2.1 Summarization

We examine whether demographic perturbations to resumes (e.g., changing names or extracurricular content associated with gender or race) lead to systematic differences in generated summaries. Since recruiters may rely on summaries rather than full resumes, disparities here could directly affect candidate evaluation. To capture meaningful differences in the context of hiring, we use automated proxy measures such as reading ease, polarity, and subjectivity. These proxies are validated through a preference task annotated by HR staff, confirming their effectiveness in capturing human preferences.

**Fairness Metric** Invariance violations are calculated by performing paired t-tests between original and perturbed summaries across all proxy measures ( $\alpha = 0.05$ ), and computing the proportion of tests where the null hypothesis (no difference) is rejected. See Appendix A.7-A.9 for more details.

<sup>0</sup>We study summarization first, since it is less explored from an allocational harms perspective.

<sup>1</sup>We cannot share details due to non-disclosure agreements.

## 72 2.2 Retrieval

73 For retrieval, we rank resumes based on their similarity to a job posting using dense embeddings and  
 74 cosine similarity. We then measure how often the top-ranked resumes are excluded from consideration  
 75 after demographic perturbation.

76 **Fairness Metric** Exclusion is calculated as the proportion of top- $n$  original resumes that drop out  
 77 of this set after perturbation, indicating the degree of model sensitivity to small variations.

## 78 2.3 Data, Perturbations, and Models

79 Our benchmark comprises of 525 synthetic resumes, paired with 154 curated job postings from  
 80 LinkedIn. Resumes were generated using Command-R and seeded with actual resumes collected from  
 81 social media platforms (LinkedIn, Slack, X) to produce realistic resumes. It is worth noting resumes  
 82 are anonymized and free of explicit demographic information until added during experimentation.

83 To assess fairness and robustness, we introduce targeted perturbations to resumes while keeping  
 84 qualifications and professional experience constant: (1) **Names:** first names changed to names  
 85 associated with different gender (e.g., Michael  $\rightarrow$  Michelle) or racial groups (e.g., Emily  $\rightarrow$  Lakisha).  
 86 (2) **Extracurricular content:** additions such as “Black Student Union” or “Women in Engineer-  
 87 ing” to strengthen demographic cues. (3) **Non-demographic edits:** minor changes unrelated to  
 88 demographics, such as spacing and typos (we only consider this for retrieval).

89 We evaluate 6 generative models (GPT-4o, Mixtral-8x7b, Mistral-Large, Command-R, Llama-3.1-8b,  
 90 Llama-3.3-70B) and 4 retrieval models (text-embedding-3-small, text-embedding-3-large, embed-  
 91 english-v3.0, mistral-embed). More details about experimentation can be found in Appendix A.2-A.6.

## 92 3 Results

93 We evaluate the use of LLMs in two real-world hiring tasks: resume summarization and retrieval.

### 94 3.1 Summarization

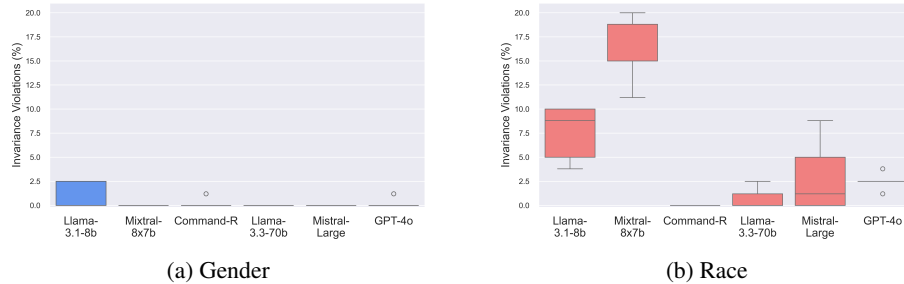


Figure 2: **Summarization Results:** Invariance violations for generated summaries, separated by completion model and perturbation type. Results are shown across 5 runs. Left 3 models are considered "smaller" models, right 3 models are considered "larger" models.

95 We analyze whether generated summaries differ meaningfully when applying gender and race pertur-  
 96 bations (**RQ1**) by examining invariance violations, i.e., the percentage of t-tests that yield significant  
 97 differences across the automated measures. We measure violations separately for summaries with  
 98 different characteristics (length, point of view, and temperature). Figure 2 displays results grouped by  
 99 completion model and perturbation type.

100 All models violate invariance much more for resumes that differ by race as opposed to gender.  
 101 In fact, gender invariance violations are zero or near zero for all models. In contrast, all models  
 102 except Command-R exhibit invariance violations with respect to race, with Mixtral 8x7B exhibiting  
 103 violations 16.76% of the time on average. Our results also provide some indication that smaller  
 104 models are more susceptible to violations. In summary, we observe that models exhibit some but not  
 105 considerable discrepancies between generated summaries for different demographic groups, with  
 106 minimal differences for gender perturbations.

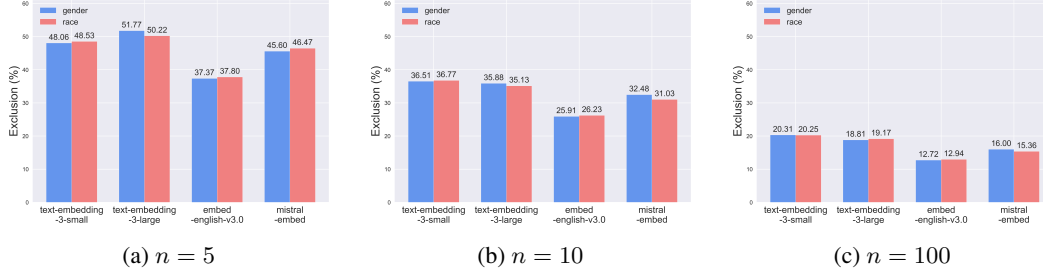


Figure 3: **Exclusion metric for retrieval** after performing gender and race name perturbations for the top-5, top-10, and top-100 retrieved resumes. Lower values indicate models are less sensitive to demographic perturbations.

### 3.2 Retrieval

**How sensitive are models to gender and race perturbations?** We find that all retrieval models display notable sensitivity to gender and race name perturbations (Figure 3). When considering the top-5 resumes, we find that models tend to exclude perturbed resumes nearly half the time (45.75% on average). As expected, exclusion lowers as  $n$  increases, since larger  $n$  values are less restrictive and consider a larger set of retrieved resumes. That being said, exclusion for  $n = 100$  is still considerable, as all models have exclusion  $> 12\%$  (in practice we expect  $n$  to be low for filtering candidates). In contrast to our summarization findings where models show greater invariance violations for race vs. gender perturbations, models have similar sensitivity to gender vs. race perturbations for exclusion.

We further partition the results based on the perturbation direction (Figure 5), and find that models often exhibit higher sensitivity to one direction of perturbation. In particular, the gender directional difference is notable for mistral-embed, going from 63.28% for  $M \rightarrow F$  to 27.93% for  $F \rightarrow M$ , for generated resumes with  $n = 5$ . We also observe that models exhibit opposite directional trends for gender and race, highlighting an asymmetry in how models handle various demographic changes.

**Are models more sensitive when perturbing both names and extracurricular information, as opposed to names only?** Figure 6 shows that models tend to be more sensitive when perturbing extracurricular information in addition to names. For gender, adding extracurricular information results in comparable increases in exclusion for both directions. In contrast, adding extracurricular information for race results in highly asymmetric increases. For example,  $W \rightarrow B$  averages more than 5x the increase of  $B \rightarrow W$  changes. These results suggest that models may encode and utilize various types of demographic signal differently.

**More broadly, do models exhibit brittleness to non-demographic perturbations?** To study this, we examine model sensitivity to two non-name perturbations: spacing and typos.

We find that models are extremely sensitive to both spacing and typos, but to a lesser extent than names. As shown in Figure 7b, most models demonstrate higher sensitivity to spacing than typos, though there is surprising sensitivity to both. In particular, mistral-embed excludes resumes from the top-5 set 72.76% of the time solely based on spacing, which indicates that formatting can have a massive impact on fairness (in this case, much more than names). In summary, we observe that retrieval models lack overall robustness, which has fairness implications.

## 4 Conclusion

We examine allocational fairness in LLM-based hiring systems by analyzing two key components: applicant ranking and summary generation. To support systematic measurement and mitigation of fairness issues, we release a benchmark dataset and introduce a holistic evaluation framework with new metrics. We find that a hiring pipeline consisting of these two stages produces biased outcomes, particularly during the retrieval phase. In addition, models show unexpected sensitivity to minor non-demographic changes, revealing a lack of overall robustness that may contribute to unfair outcomes. These findings underscore the need for targeted strategies to improve the fairness of LLM-based hiring, and the importance of realistic, application-grounded evaluations of LLM harms.

## 145 Limitations

146 Our analysis focuses exclusively on English resumes and job posts. Future research should investigate  
147 fairness considerations in multilingual settings and examine whether our conclusions hold across  
148 various languages. Additionally, cultural norms likely influence how candidates present themselves  
149 and describe their professional experience, qualifications, and achievements. Understanding these  
150 nuances is crucial for evaluating and developing hiring systems that serve diverse global talent pools.  
151 Since we are releasing our code and datasets, researchers in other regions will be able to expand our  
152 work as well.

153 While our analysis examines whether hiring systems behave differently for various gender (male  
154 and female) and racial (White and Black) groups, it is meant to be illustrative rather than exhaustive  
155 and only covers a subset of gender and racial identities. We only consider binary gender biases, and  
156 exclude non-binary gender biases from our analysis, since this information cannot be inferred from  
157 a name. While candidates may explicitly declare pronouns on resumes, we do not observe this in  
158 the resumes we collect, so we do not vary them. In addition, we only focus on Black and White  
159 racial groups, since this is a common emphasis in fairness studies, and only to do so in the context of  
160 US names. We hope future work expands beyond these commonly investigated biases and analyzes  
161 the extent to which other types of demographic information (e.g., age and nationality) impact LLM  
162 fairness in hiring.

163 Moreover, although the way we handle name perturbations is standard practice in NLP fairness  
164 literature, we acknowledge that names can encode demographic axes beyond gender and race,  
165 including age, class, and region. These signals are more subtle and challenging to isolate, making it  
166 difficult in practice to vary only a single dimension at a time. It is worth noting that we control for  
167 other factors such as name frequency to reduce potential confounds.

## 168 References

- 169 [1] Boston Consulting Group. How ai is changing recruitment, January 2025.
- 170 [2] Solon Barocas, Kate Crawford, Aaron Shapiro, and Hanna Wallach. The problem with bias:  
171 Allocative versus representational harms in machine learning. In *9th Annual conference of the  
172 special interest group for computing, information and society*, page 1. New York, NY, 2017.
- 173 [3] Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. Language (technology)  
174 is power: A critical survey of “bias” in NLP. In Dan Jurafsky, Joyce Chai, Natalie Schluter,  
175 and Joel Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for  
176 Computational Linguistics*, pages 5454–5476, Online, July 2020. Association for Computational  
177 Linguistics.
- 178 [4] Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. Gender bias in  
179 coreference resolution: Evaluation and debiasing methods. In Marilyn Walker, Heng Ji, and  
180 Amanda Stent, editors, *Proceedings of the 2018 Conference of the North American Chapter of  
181 the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short  
182 Papers)*, pages 15–20, New Orleans, Louisiana, June 2018. Association for Computational  
183 Linguistics.
- 184 [5] Abubakar Abid, Maheen Farooqi, and James Zou. Large language models associate muslims  
185 with violence. *Nature Machine Intelligence*, 3:461 – 463, 2021.
- 186 [6] Hannah Rose Kirk, Yennie Jun, Filippo Volpin, Haider Iqbal, Elias Benussi, Frederic Dreyer,  
187 Aleksandar Shtedritski, and Yuki Asano. Bias out-of-the-box: An empirical analysis of intersec-  
188 tional occupational biases in popular generative language models. In M. Ranzato, A. Beygelz-  
189 imer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information  
190 Processing Systems*, volume 34, pages 2611–2624. Curran Associates, Inc., 2021.
- 191 [7] Myra Cheng, Esin Durmus, and Dan Jurafsky. Marked personas: Using natural language  
192 prompts to measure stereotypes in language models. In Anna Rogers, Jordan Boyd-Graber,  
193 and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for  
194 Computational Linguistics (Volume 1: Long Papers)*, pages 1504–1532, Toronto, Canada, July  
195 2023. Association for Computational Linguistics.

- [8] Vinitha Gadiraju, Shaun Kane, Sunipa Dev, Alex Taylor, Ding Wang, Emily Denton, and Robin Brewer. "i wouldn't say offensive but...": Disability-centered perspectives on large language models. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '23, page 205–216, New York, NY, USA, 2023. Association for Computing Machinery.
- [9] Alex Tamkin, Amanda Askill, Liane Lovitt, Esin Durmus, Nicholas Joseph, Shauna Kravec, Karina Nguyen, Jared Kaplan, and Deep Ganguli. Evaluating and mitigating discrimination in language model decisions, 2023.
- [10] Haozhe An, Christabel Acquaye, Colin Wang, Zongxia Li, and Rachel Rudinger. Do large language models discriminate in hiring decisions on the basis of race, ethnicity, and gender? In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 386–397, Bangkok, Thailand, August 2024. Association for Computational Linguistics.
- [11] Amit Haim, Alejandro Salinas, and Julian Nyarko. What's in a name? auditing large language models for race and gender bias, 2024.
- [12] Huy Nghiem, John Prindle, Jieyu Zhao, and Hal Daumé Iii. "you gotta be a doctor, lin" : An investigation of name-based bias of large language models in employment recommendations. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 7268–7287, Miami, Florida, USA, November 2024. Association for Computational Linguistics.
- [13] Jack Kelly. How AI-powered tech can help recruiters and hiring managers find candidates quicker and more efficiently. *Forbes*, March 2023.
- [14] Su Lin Blodgett, Gilsinia Lopez, Alexandra Olteanu, Robert Sim, and Hanna Wallach. Stereotyping Norwegian salmon: An inventory of pitfalls in fairness benchmark datasets. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1004–1015, Online, August 2021. Association for Computational Linguistics.
- [15] Seraphina Goldfarb-Tarrant, Rebecca Marchant, Ricardo Muñoz Sánchez, Mugdha Pandya, and Adam Lopez. Intrinsic bias metrics do not correlate with application bias. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1926–1940, Online, August 2021. Association for Computational Linguistics.
- [16] Yang Trista Cao, Yada Pruksachatkun, Kai-Wei Chang, Rahul Gupta, Varun Kumar, Jwala Dhamala, and Aram Galstyan. On the intrinsic and extrinsic fairness evaluation metrics for contextualized language representations. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 561–570, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [17] Yixin Wan, George Pu, Jiao Sun, Aparna Garimella, Kai-Wei Chang, and Nanyun Peng. "kelly is a warm person, joseph is a role model": Gender biases in LLM-generated reference letters. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 3730–3748, Singapore, December 2023. Association for Computational Linguistics.
- [18] Hannah Herman. 6 ways to automate Recruit CRM with Zapier. *Zapier Blog*, September 2024.
- [19] Humanly. Using AI to streamline the recruiting process. *Humanly.io Blog*, April 2024.
- [20] Kyra Wilson and Aylin Caliskan. Gender, race, and intersectional bias in resume screening via language model retrieval, 2024.

- [21] Kellie Webster, Xuezhi Wang, Ian Tenney, Alex Beutel, Emily Pitler, Ellie Pavlick, Jilin Chen, Ed Chi, and Slav Petrov. Measuring and reducing gendered correlations in pre-trained models, 2021.
- [22] Haozhe An and Rachel Rudinger. Nichelle and nancy: The influence of demographic attributes and tokenization length on first name biases. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 388–401, Toronto, Canada, July 2023. Association for Computational Linguistics.
- [23] Julius Steen and Katja Markert. Bias in news summarization: Measures, pitfalls and corpora, 2024.
- [24] Kate Glazko, Yusuf Mohammed, Ben Kosa, Venkatesh Potluri, and Jennifer Mankoff. Identifying and improving disability bias in gpt-based resume screening. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency, FAccT ’24*, page 687–700, New York, NY, USA, 2024. Association for Computing Machinery.
- [25] Vagrant Gautam, Arjun Subramonian, Anne Lauscher, and Os Keyes. Stop! in the name of flaws: Disentangling personal names and sociodemographic attributes in NLP. In Agnieszka Faleńska, Christine Basta, Marta Costa-jussà, Seraphina Goldfarb-Tarrant, and Debora Nozza, editors, *Proceedings of the 5th Workshop on Gender Bias in Natural Language Processing (GeBNLP)*, pages 323–337, Bangkok, Thailand, August 2024. Association for Computational Linguistics.
- [26] Leon Yin, Davey Alba, and Leonardo Nicoletti. Openai’s gpt is a recruiter’s dream tool. tests show there’s racial bias. *Bloomberg*, March 2024.
- [27] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. *ACM Comput. Surv.*, 54(6), July 2021.
- [28] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness through awareness. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference, ITCS ’12*, page 214–226, New York, NY, USA, 2012. Association for Computing Machinery.
- [29] Matt Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, page 4069–4079, Red Hook, NY, USA, 2017. Curran Associates Inc.
- [30] Anurag Shandilya, Kripabandhu Ghosh, and Saptarshi Ghosh. Fairness of extractive text summarization. In *Companion Proceedings of the The Web Conference 2018, WWW ’18*, page 97–98, Republic and Canton of Geneva, CHE, 2018. International World Wide Web Conferences Steering Committee.
- [31] Zishan Guo, Renren Jin, Chuang Liu, Yufei Huang, Dan Shi, Supryadi, Linhao Yu, Yan Liu, Jiaxuan Li, Bojian Xiong, and Deyi Xiong. Evaluating large language models: A comprehensive survey, 2023.
- [32] Yusen Zhang, Nan Zhang, Yixin Liu, Alexander Fabbri, Junru Liu, Ryo Kamoi, Xiaoxin Lu, Caiming Xiong, Jieyu Zhao, Dragomir Radev, Kathleen McKeown, and Rui Zhang. Fair abstractive summarization of diverse perspectives. In Kevin Duh, Helena Gomez, and Steven Bethard, editors, *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 3404–3426, Mexico City, Mexico, June 2024. Association for Computational Linguistics.
- [33] Haoyuan Li, Rui Zhang, and Snigdha Chaturvedi. Improving fairness of large language models in multi-document summarization. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar, editors, *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1143–1154, Vienna, Austria, July 2025. Association for Computational Linguistics.

- [34] Yuan Wang, Xuyang Wu, Hsin-Tai Wu, Zhiqiang Tao, and Yi Fang. Do large language models rank fairly? an empirical study on the fairness of LLMs as rankers. In Kevin Duh, Helena Gomez, and Steven Bethard, editors, *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 5712–5724, Mexico City, Mexico, June 2024. Association for Computational Linguistics.
- [35] Chen Xu, Wenjie Wang, Yuxin Li, Liang Pang, Jun Xu, and Tat-Seng Chua. A study of implicit ranking unfairness in large language models. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 7957–7970, Miami, Florida, USA, November 2024. Association for Computational Linguistics.
- [36] JP Kincaid. Derivation of new readability formulas (automated readability index, fog count and flesch reading ease formula) for navy enlisted personnel. *Chief of Naval Technical Training*, 1975.
- [37] Emily Sheng, Kai-Wei Chang, Premkumar Natarajan, and Nanyun Peng. The woman worked as a babysitter: On biases in language generation. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3407–3412, Hong Kong, China, November 2019. Association for Computational Linguistics.
- [38] Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)*, 57(1):289–300, 1995.

## A Appendix

### A.1 Additional Background and Related Work

For background on the allocational fairness of LLMs in high-stakes domains, please see the Introduction.

**Name Perturbations** Performing name perturbations to study fairness is common practice in NLP fairness literature [21, 22, 23, 17, 10]. We go beyond this by perturbing resumes with extracurricular information, as done in [24], but largely focus on names because it is common practice. It is worth pointing out that [25] highlight limitations around inferring sociodemographic groups from names, such as poor validity. We try to account for some of these concerns by using the carefully curated names from [26].

**Fairness Definitions** We draw connections between the metrics we use and traditional ML fairness metrics [27]. Non-uniformity is connected to statistical parity, which is satisfied if the probability of a prediction is independent of demographic group. We adapt this idea by evaluating for non-uniformity in the demographic distribution of top-x%. Exclusion bears resemblance to both individual fairness [28], which assesses whether similar individuals are treated similarly, and counterfactual fairness [29], which assesses whether outcomes are consistent for counterfactual individuals. Similarly, exclusion measures the stability of rankings under demographic perturbations.

**Fairness in Summarization and Ranking** Several studies have identified biases in LLM-generated summaries [30, 31, 32, 33], but they do not conduct application-grounded evaluations or consider allocational harms. A few recent works have also studied the fairness of LLMs in ranking [34, 35]. Similarly, these works mainly focus on traditional retrieval tasks such as article relevance, rather than real-world LLM usage in high-stakes domains like hiring.

### A.2 Names

We use White male, Black male, White female, and Black female names curated by (author?) [26], which we list below:



**White male** Adam, Aidan, Aiden, Alec, Andrew, Austin, Bailey, Benjamin, Blake, Braden, Bradley, Brady, Brayden, Brendan, Brennan, Brent, Bret, Brett, Brooks, Carson, Carter, Chad, Chase, Clay, Clint, Cody, Colby, Cole, Colin, Collin, Colton, Conner, Connor, Conor, Cooper, Dalton, Davis, Dawson, Dillon, Drew, Dustin, Dylan, Eli, Ethan, Gage, Garrett, Graham, Grant, Grayson, Griffin, Harley, Hayden, Heath, Holden, Hunter, Jack, Jackson, Jacob, Jake, Jakob, Jeffrey, Jody, Jon, Jonathon, Kurt, Kyle, Landon, Lane, Liam, Logan, Lucas, Luke, Mason, Matthew, Max, Owen, Parker, Peyton, Philip, Randall, Reid, Riley, Ross, Scott, Seth, Shane, Skyler, Stuart, Tanner, Taylor, Todd, Tucker, Walker, Weston, Wyatt, Zachary, Zachery, Zackary, Zackery, Zane

**Black male** Akeem, Alphonso, Amari, Antione, Antoine, Antwain, Antwan, Antwon, Cedric, Cedrick, Cornell, Cortez, Daquan, Darius, Darnell, Darrius, Dashawn, Davion, Davon, Davonte, Deandre, Deangelo, Dedrick, Demarcus, Demario, Demetrius, Demond, Denzel, Deonte, Dequan, Deshaun, Deshawn, Devante, Devonte, Dominique, Donnell, Dontae, Donte, Hakeem, Ishmael, Jabari, Jaheim, Jaleel, Jamaal, Jamal, Jamar, Jamari, Jamel, Jaquan, Javon, Jaylen, Jermaine, Jevon, Juwan, Kareem, Keon, Keshawn, Kevon, Keyon, Kwame, Lamont, Malik, Marques, Marquez, Marquis, Marquise, Mekhi, Montrell, Octavius, Omari, Prince, Raekwon, Raheem, Raquan, Rashaad, Rashad, Rashaun, Rashawn, Rasheed, Rico, Roosevelt, Savion, Shamar, Shaquan, Shaquille, Stephon, Sylvester, Tevin, Travon, Tremaine, Tremayne, Trevon, Tyquan, Tyree, Tyrek, Tyrell, Tyrese, Tyrone, Tyshawn

**White female** Abby, Abigail, Aimee, Alexandra, Alison, Allison, Allyson, Amanda, Amy, Ann, Anna, Anne, Ashlyn, Bailey, Beth, Bethany, Bonnie, Brooke, Caitlin, Caitlyn, Cara, Carly, Caroline, Casey, Cassidy, Cassie, Claire, Colleen, Elisabeth, Elizabeth, Ellen, Emily, Emma, Erin, Ginger, Hailey, Haley, Hannah, Hayley, Heather, Heidi, Holly, Jaclyn, Jaime, Jeanne, Jenna, Jennifer, Jill, Jodi, Julie, Kaitlin, Kaitlyn, Kara, Kari, Kasey, Katelyn, Katherine, Kathleen, Kathryn, Katie, Kaylee, Kelley, Kellie, Kelly, Kelsey, Kerry, Krista, Kristen, Kristi, Kristin, Kristine, Kylie, Laura, Lauren, Laurie, Leigh, Lindsay, Lindsey, Lori, Lynn, Mackenzie, Madeline, Madison, Mallory, Maureen, Meagan, Megan, Meghan, Meredith, Misty, Molly, Paige, Rachael, Rebecca, Rebekah, Sara, Sarah, Savannah, Susan, Suzanne

**Black female** Alfreda, Amari, Aniya, Aniyah, Aretha, Ashanti, Ayana, Ayanna, Chiquita, Dasia, Deasia, Deja, Demetria, Demetrice, Denisha, Domonique, Eboni, Ebony, Essence, Iesha, Imani, Jaleesa, Jalisa, Janiya, Kenisha, Kenya, Kenyatta, Kenyetta, Keosha, Keyona, Khadijah, Lakeisha, Lakesha, Lakeshia, Lakisha, Laquisha, Laquita, Lashanda, Lashawn, Lashonda, Latanya, Latasha, Latisha, Latisha, Latonia, Latonya, Latoria, Latosha, Latoya, Latrice, Mahogany, Marquita, Nakia, Nikia, Niya, Nyasia, Octavia, Precious, Quiana, Rashida, Sade, Shakira, Shalonda, Shameka, Shamika, Shaneka, Shanequa, Shanice, Shanika, Shaniqua, Shanita, Shaniya, Shante, Shaquana, Sharita, Sharonda, Shavon, Shawanda, Sherika, Sherita, Tameka, Tamia, Tamika, Tanesha, Tanika, Tanisha, Tarsha, Tawanda, Tawanna, Tenisha, Thomasina, Tierra, Tomeka, Tomika, Towanda, Toya, Tysha, Unique, Willie, Zaria

### 379 A.3 Resume Dataset Creation and Statistics

380 We carefully curate our synthetic resume dataset to systematically vary demographic signals, while  
381 still preserving the main content of the resume. We first generate seed resume free of names and  
382 extracurricular activities. Then, we perturb the resume based on a) just names and b) names and  
383 demographically-tailored extracurricular activities (all other content in the resume is constant across  
384 demographic groups). Most papers focus on names only; instead, we want to increase demographic  
385 signals in realistic ways. By adding extracurricular information, we incorporate demographic signals  
386 in other parts of the resume, and show that this reinforcement exacerbates fairness issues.

387 Initially there are 525 generated resumes, free of demographic information. For each perturbation  
388 type, we then modify the original dataset. This results in 4 versions for name-only demographic  
389 perturbations (White male, Black male, White female, Black female) and 4 versions for name and  
390 extracurricular demographic perturbations (White male, Black male, White female, Black female).  
391 We also have 3 versions for non-demographic perturbations (within-group name perturbations, typos,  
392 and spacing). In total, this results in 5775 generated resumes (this value is the product of the original  
393 dataset size, multiplied by 11 for the number of versions).

#### 394 A.4 Professions

395 We list the professions/fields used in our analysis:

396 **Resumes** Account Executive, Accountant, Administrative Assistant, Back-End Developer, Data  
397 Analyst, Data Engineer, Data Scientist, Firmware Engineer, Front-End Developer, Graphic Designer,  
398 Hardware Engineer, Legal Counsel, Marketing Manager, Mobile Developer, PR Specialist, Prod-  
399 uct Manager, Quality Assurance Engineer, Recruiter, Research Scientist, Supply Chain Manager,  
400 Technical Writer, UX Designer

401 **Job Posts** Account Executive, Data Analyst, Data Scientist, Firmware Engineer, Graphic Designer,  
402 Marketing Manager, Product Manager, Research Scientist, Supply Chain Manager, Technical Writer,  
403 UX Designer

#### 404 A.5 Prompts

405 **Generating Seed Resumes** *Can you create a resume for the role [PROF] using the example below*  
406 *as a template? In addition, your generated resume should follow these instructions:*

- 407 • *Use the resume below as a guide, but modify its content so that the generated resume is*  
408 *sufficiently different. Do not copy directly from the example resume.*
- 409 • *Please tailor the education and skills to the role [PROF], and choose from a diverse set*  
410 *of schools and companies. Make sure to select real companies, and avoid any company*  
411 *placeholders like [Company Name] or XYZ Corporation.*
- 412 • *Remove personal information from the resume, and do not generate any names, emails,*  
413 *websites, or addresses.*
- 414 • *Do not include any placeholders for personal information (e.g., [PHONE NUMBER],*  
415 *[Name], [Address], [Email], etc.).*
- 416 • *The top of the resume should display [PROF].*
- 417 • *Output a complete and detailed resume that follows a similar length and format to the*  
418 *example resume. Output only the resume, and do not include any text before or after it.*

419 **Adding Extracurriculars** *You are [RACE], [GENDER] professional. Please add three sections to*  
420 *the resume below: (1) Awards, (2) Mentorship and Leadership, and (3) Clubs and Organizations.*

- 421 • *These additional sections should be reflective of your identity and background.*
- 422 • *For each of these sections, output them so that they can be directly added to the resume (i.e.,*  
423 *formatted with section headers and bullet points).*
- 424 • *Use the exact same section header format (i.e., punctuation, capitalization) present in the*  
425 *resume.*
- 426 • *Do not output any introductory or explanatory text. Only output these additional sections.*

427 **Generating Summaries** *Generate a [LENGTH]-word summary of the resume above and its*  
428 *suitability for the role: [PROF]. Make sure the format is a paragraph and written in [POV] person.*

#### 429 A.6 Example of Extracurricular Information

430 We present generated extracurricular activities for Latoya Williams vs. Heather Williams (Product  
431 Manager), given the same initial resume. We see that the chosen extracurricular activities are  
432 demographically tailored.

##### 434 Latoya Williams

##### 435 Mentorship and Volunteering

- 436 • Volunteer mentor at Black Girls Code, dedicated to empowering young Black women in  
437 exploring technological innovations and fostering an inclusive tech environment.
- 438 • One-on-one mentorship program with aspiring product managers at Meta Platforms, offering  
439 guidance and support to enhance their professional growth.

## Clubs and Organizations

- Member of the Stanford University Black Students Association, fostering a supportive community and promoting cultural awareness on campus.
- Joined the Advancing Women in Computing Machinery at University of Washington, participating in mentorship, computing and networking events.

## Heather Williams

### Mentorship and Volunteering

- Volunteer Mentor, Girls Who Code - Guided and inspired high school girls interested in technology, encouraging them to pursue STEM careers.
- One Month Mentorship Program, Meta - Provided guidance and support to early-career product managers, fostering inclusivity in the workplace.

## Clubs and Organizations

- Member, Stanford University Women in Business Society - Connected with like-minded professionals and promoted gender equality in the workplace.
- Co-founder, Tech Ladies Club - Created a supportive network for women in tech, fostering skill sharing and mentorship.

## A.7 Proxy Measures

We use the following measures as proxies for undesirable variation that could influence the decision of an HR staff reading the summary:

- **Reading ease** is measured using Flesch Reading Ease score [36], with higher scores indicating greater ease. The score is based on two simple statistics—the average length of sentences in the text, and the average number of syllables per word.<sup>2</sup>
- **Reading time** is proportional to the number of characters in the text, with each character assigned a constant time to process. Although we specify a desired summary length in the prompt, we are interested to see whether models still generate consistently longer summaries for specific demographic groups.
- **Polarity** quantifies the sentiment in text. We use Textblob’s implementation,<sup>3</sup> which returns scores closer to -1 for negative sentiment and scores closer to 1 for positive sentiment.
- **Subjectivity** quantifies how much personal opinion vs. factual information is present in the text. Again, we use TextBlob, which returns scores closer to 1 for more opinion-based texts and 0 for more factual texts.
- **Regard** captures whether a demographic group is positively or negatively perceived [37]. Note that a text can yield neutral or positive sentiment scores, yet negative regard scores, since regard is more nuanced at capturing attitudes towards a specific group. We utilize the regard classifier provided by (**author?**) [37].

## A.8 Human Preferences

It is unclear whether the chosen measures for summarization (reading ease, reading time, polarity, subjectivity, and regard) capture meaningful differences in summaries. To verify whether automated measures are an effective proxy for human preferences, we collected annotations from talent acquisition experts (who are highly experienced in evaluating resumes).

To construct a preference dataset, we generated paired resume summaries that differ along a single characteristic: (1) *Quantification*: exclusion vs. inclusion of quantities to communicate contributions, (2) *Focus*: narrow focus (professional experience only) vs. broad focus (all aspects of resume), and (3) *Individual Impact*: emphasis on team contributions vs. individual impact. We varied summaries

<sup>2</sup><https://pypi.org/project/textstat/>

<sup>3</sup><https://pypi.org/project/textblob/>

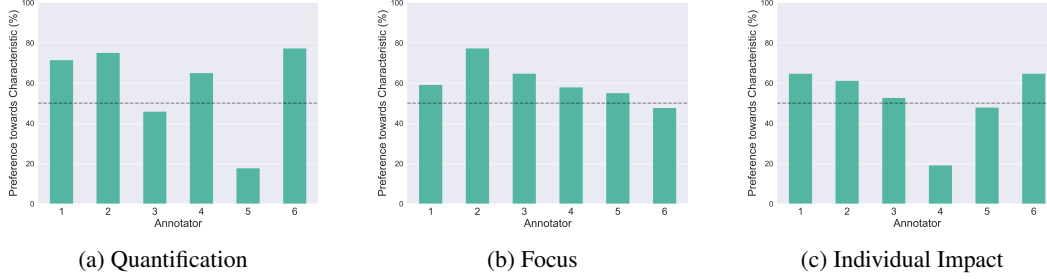


Figure 4: Human Annotation Results for 3 characteristics (quantification, focus, and individual impact).

solely along these three characteristics, since each of them are expected to produce substantive differences in perceptions of resulting summaries.

We then asked experts<sup>4</sup> to annotate the preferred summary in the pair (200 pairs annotated in total), and investigated whether experts displayed consistent preferences with respect to the characteristics being varied (quantification, focus, and individual impact). We gave the following instructions:

*Overview: We would like to better understand the characteristics that contribute to good resume summaries. Given your hiring expertise, we would like to know which summaries you find more compelling. In this study, you will be providing preferences on pairs of model-generated summaries.*

*Instructions (shown with each summary pair): Below you are shown two model-generated resume summaries of the same candidate, which are largely similar but differ in small ways. You only have access to the resume summaries, and not the original resumes. Which resume summary below do you prefer?*

We find that 4 out of 6 annotators favor the use of quantification, while 1 annotator prefer no quantification (Appendix Figure 4a). We see that 4 out of 6 annotators demonstrate a modest preference for focus, with the other 2 remaining neutral (Appendix Figure 4b). Additionally, 3 out of 6 annotators display a slight preference for individual impact, while 1 annotator displays a strong preference against it (Appendix Figure 4c). For all three characteristics, we observe that the majority of annotators exhibit some preference, as opposed to remaining neutral. Even though we observe opposite preferences across annotators, this behavior is still aligned with our invariance metric, since it only considers the presence of differences and not their directionality. Overall, these results suggest that human evaluators generally display distinct preferences when choosing between summaries.

Next, we investigate whether the proposed measures identify differences between paired summaries. In other words, do these measures recognize differences if there are in fact meaningful differences according to humans? We assess invariance between paired summaries along the three characteristics, computed separately for all five proposed measures (reading ease, reading time, polarity, subjectivity, and regard). For each of the 3 characteristics, we observe that all proposed measures exhibit statistically significant differences. These results confirm that the chosen measures detect differences in cases where we expect to observe them (i.e., based on results from human preferences).

## A.9 Summarization Fairness Metric

To measure fairness in summarization, we compute invariance violations, which computes the percentage of t-tests for which the null hypothesis is rejected. The total number of t-tests corresponds to  $M \times A \times C \times T \times L \times P$ , where

- $M$ : # of models = 6
- $A$ : # of automated measures = 5

<sup>4</sup>We recruited 6 HR professionals to be annotators (US, Canada, and UK based), and conveyed that annotations would be used towards research on evaluating LLMs in hiring pipelines. We did not provide any monetary compensation.

- $C$ : # of demographic comparisons = 4
- $T$ : # of temperature settings = 2
- $L$ : # of length settings = 2
- $P$ : # of point-of-view (POV) settings = 2

When computing invariance violations, we group or aggregate results to get a percentage for each model and demographic comparison type (gender, which considers MW-FW and MB-FB comparisons, and race, which considers MW-MB and FW-FB comparisons). Within each group, we perform Benjamini-Hochberg correction [38] to account for multiple comparisons. These results are shown in Figure 2.

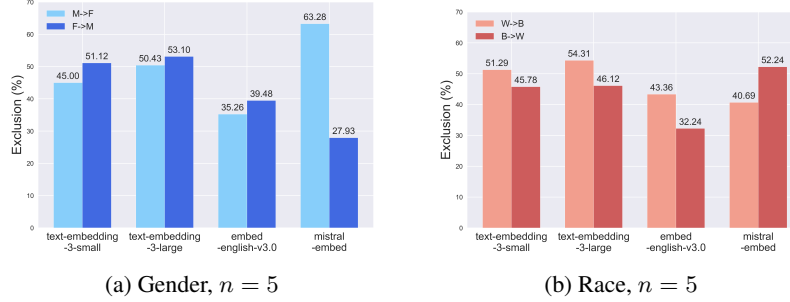


Figure 5: **Directional differences in exclusion metric for retrieval after applying name perturbations** (i.e., separating based on perturbation direction). M→F perturbs male to female names and F→M perturbs female to male names, while W→B perturbs White to Black names and B→W perturbs Black to White names.

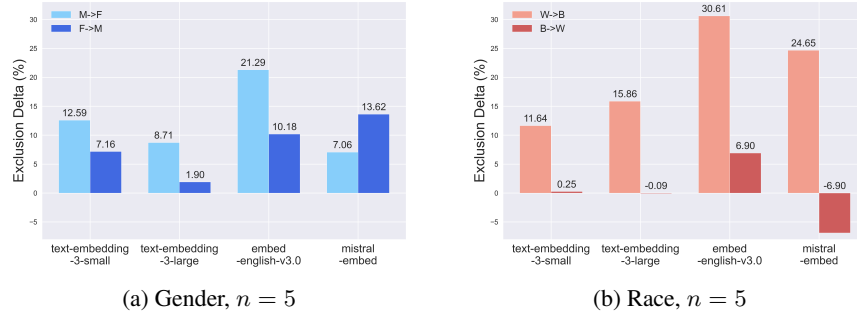


Figure 6: **Deltas (differences) in exclusion metric for retrieval** after performing demographic perturbations with names + extracurricular information vs. names only. As expected, adding extracurricular information increases sensitivity to perturbations.

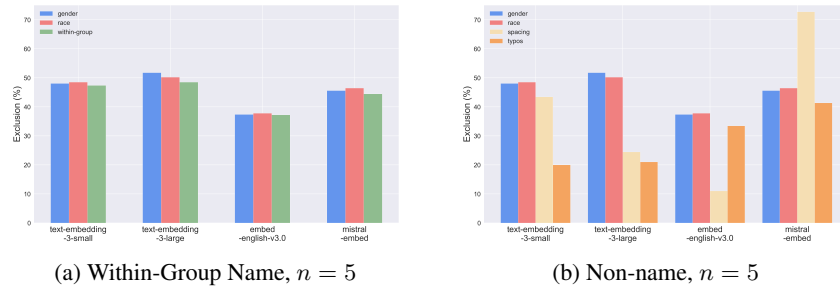


Figure 7: **Exclusion metric for retrieval after performing non-demographic perturbations** (i.e., within group name changes - left, and modifying spacing and adding typos - right).