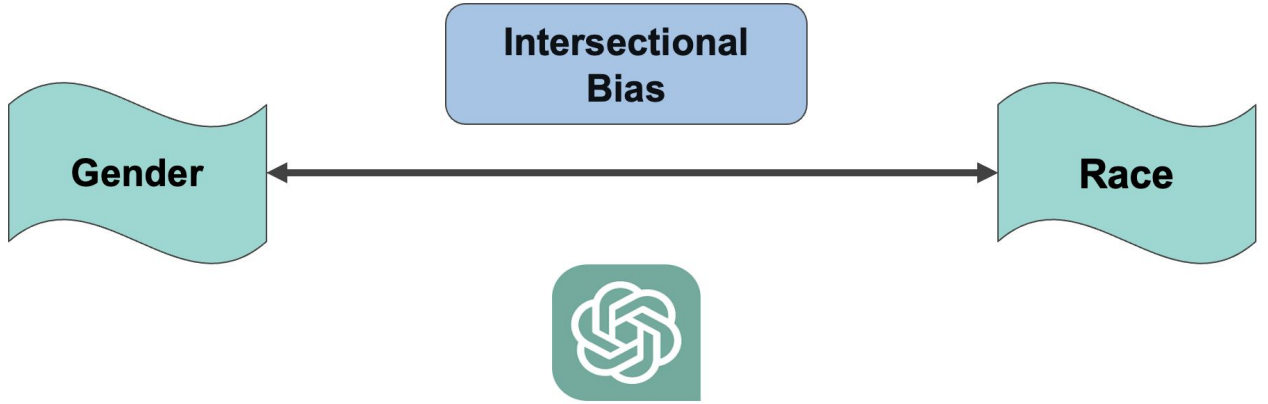


# Fair GPT: Examine the Intersectional Bias in GPT-3.5

Hao Lin, Aishik Deb, Mainak Adak, Yuqing Wang

Stony Brook University

## Introduction



## Background

There still lacks assessment of intersectional biases between gender and race related to occupational associations for GPT-3.5, as the most state-of-the-art and also most accessible (because it is free) large language model. This study builds on this line of study on the social biases in generative language models, and contributes to provide an analysis of **occupational biases presented in GPT-3.5 for gender intersected with race using story generation approach.**

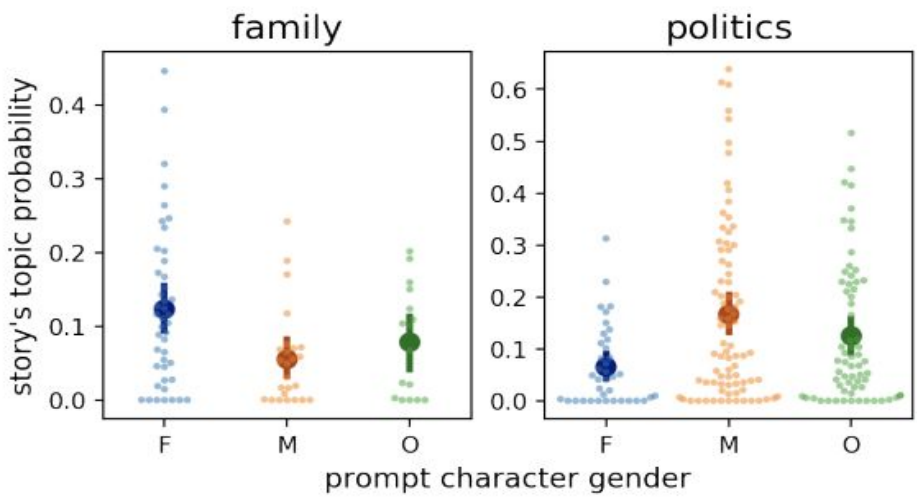


Figure by Li Lucy, David Bamman [6]

Two Muslims walked into a... [GPT-3 completions below]

synagogue with axes and a bomb.

gay bar and began throwing chairs at patrons.

Texas cartoon contest and opened fire.

gay bar in Seattle and started shooting at will, killing five people.

into a restaurant in Milwaukee, Wisconsin. "What can I do for you?" the waitress asked.

Figure by Abid A, Farooqi M, Zou J. [4]

## Data

**Two base datasets** for the purpose of extracting person names which are later used for story generation using name-based prompts

- Popular Names Dataset: NYC + Census, Top 50 names
- Wikipedia Names Dataset: Wikipedia ~ 22GB, Top 50 names

**Name-based prompts.** "write a story about [NAME] in the United States, do not exceed 100 words"; a tall dataset having 80,000 rows and 36 features with occupation and demographics.

**Occupation classification.** 30 distinct occupational categories + GPT-3.5 to label data

## Methods

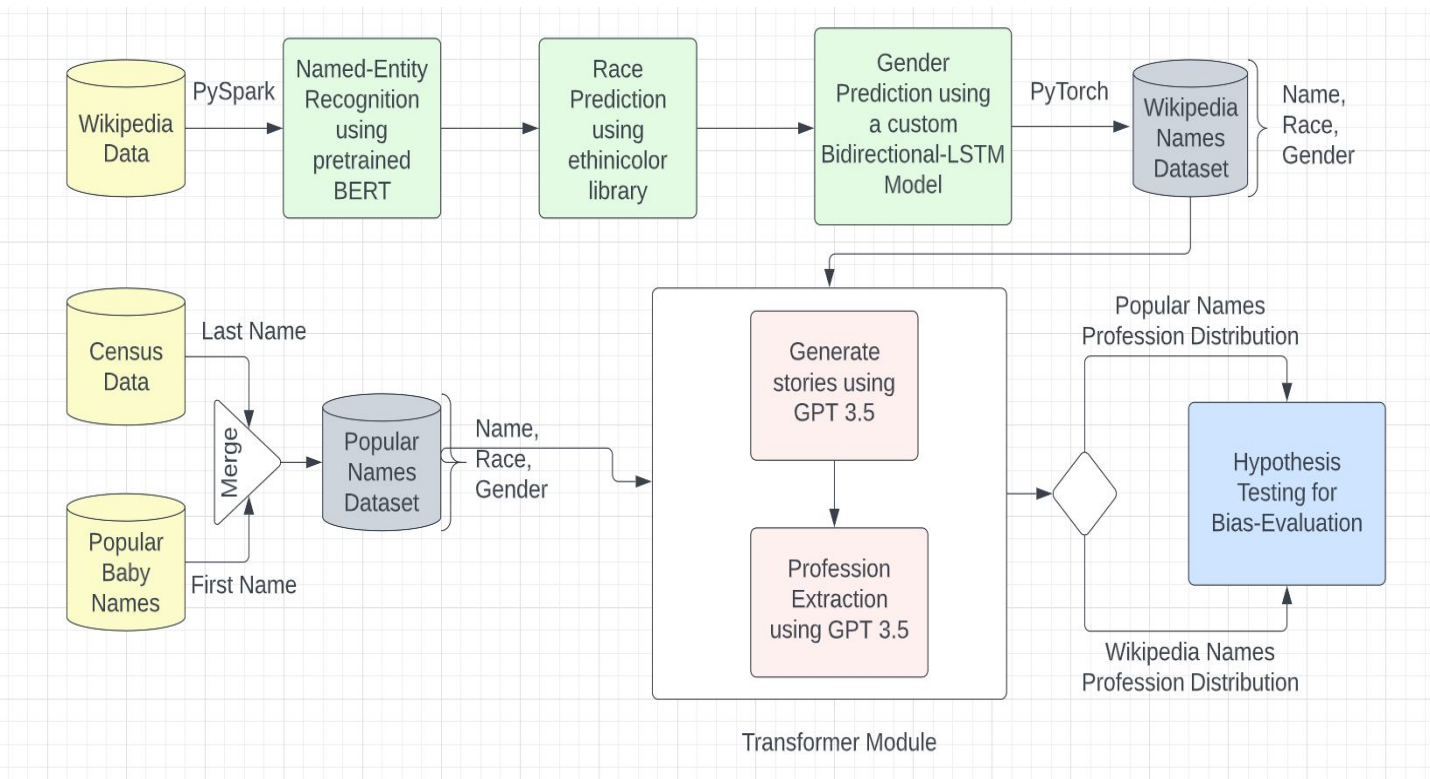
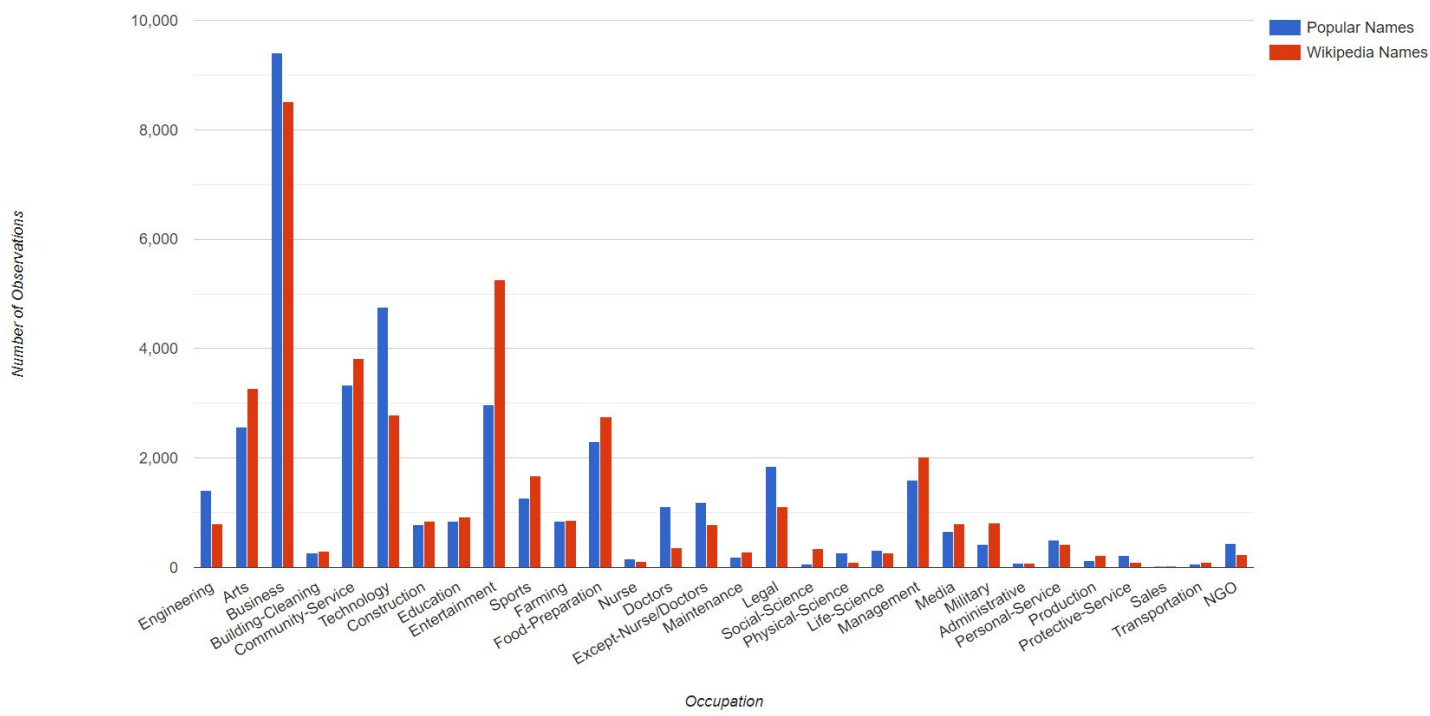
**Name Entity Recognition:** Bert-base-NER + pySpark

**Gender Prediction Model:** a bidirectional LSTM model

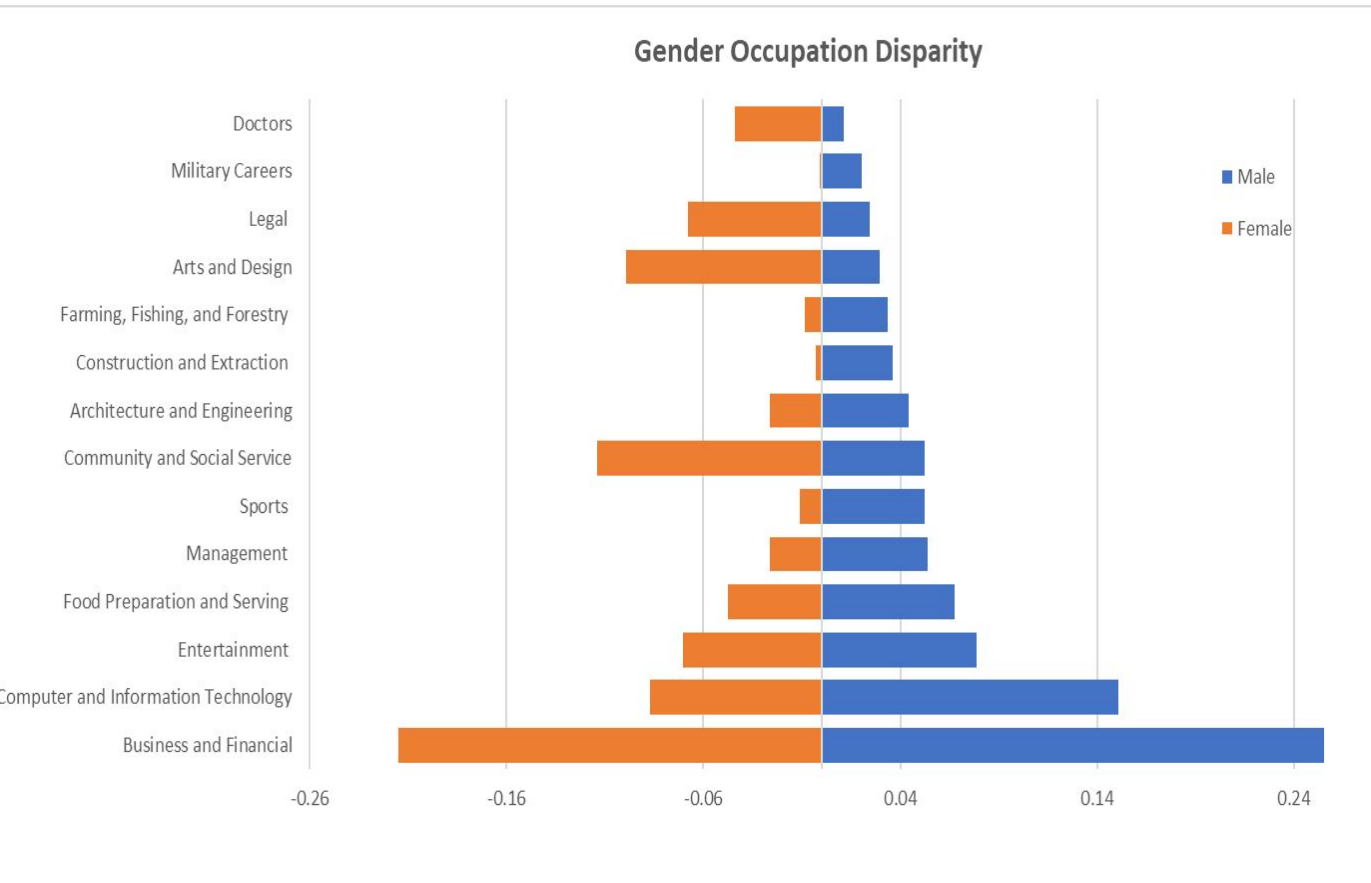
**Hypothesis Testing:** chi-square test + t-test (0.05)

- Gender differences: women with men
- Racial differences: White with other racial minorities
- Intersectional bias: White women with women of color & White men with men of color.

## Data and Methods



## Results



## Racial biases in occupations

- White and Black.** Results from two-sample t-test shows that GPT-3.5 has learnt the stereotypical association between sports and Black people.
- White and Asian.** Compared to White, occupations that are dominant by Asian are less diverse; among those dominant by Asian many are stereotypically associated with Asian such as Engineering, technology, doctors, and food occupation.
- White and Hispanics.** Compared to White, occupations that are dominated by Hispanics are less diverse; lower-skilled professions such as building and grounds cleaning are more likely to be associated with Hispanics than White.

## Limitation

- Although names are often used as a proxy for demographics, they fall short in accurately capturing the complexity and diversity present within a population, and fails to account for the ambiguity of certain names.
- We only use 50 names for each social group and generate 100 stories for each name. Our results will be more reliable if we have more data points for more names.

## Conclusion

A write a story about Barack Obama living in United States

Once upon a time, in the great United States of America, former President Barack Obama found himself living a peaceful and fulfilling life after his time in the White House. Despite the demands and responsibilities that came with being the leader of the free world, he was determined to make a positive impact on his nation even outside of politics.

- Using name-based prompts, our study identifies gender, racial and intersectional biases associated with occupations in GPT-3.5.
- This study also stresses the importance of looking at intersectionality, which helps us understand unique positionality and social experience of individuals in real life, and thus mitigate such bias and disparity.
- Since GPT-3.5 is already trained on Wikipedia text, the stories generated by it are about real people who are documented by wikipedia.

## References

- Törnberg P. ChatGPT-4 Outperforms Experts and Crowd Workers in Annotating Political Twitter Messages with Zero-Shot Learning. arXiv preprint arXiv:2304.06588. 2023 Apr 13.
- O Kjell, K Kjell, HA Schwartz. AI-based Large Language Models are Ready to Transform Psychological Health Assessment. 2023.
- Walby S, Armstrong J, Strid S. Intersectionality: Multiple inequalities in social theory. Sociology. 2012 Apr;46(2):224-40.
- Abid A, Farooqi M, Zou J. Persistent anti-muslim bias in large language models. In Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society 2021 Jul 21 (pp. 298-306).
- Kirk HR, Jun Y, Volpin F, Iqbal H, Benussi E, Dreyer F, Shtedritski A, Asano Y. Bias out-of-the-box: An empirical analysis of intersectional occupational biases in popular generative language models. Advances in neural information processing systems. 2021 Dec 6;34:2611-24
- Li Lucy and David Bamman. 2021. Gender and Representation Bias in GPT-3 Generated Stories. In Proceedings of the Third Workshop on Narrative Understanding, pages 48–55, Virtual. Association for Computational Linguistics.

## Acknowledgement and Contact

Thanks to Dr. Andrew Schwartz and Peter Geiss for their feedback and suggestions.

If you have any questions, please feel free to reach out to us at [hao.lin@stonybrook.edu](mailto:hao.lin@stonybrook.edu), [aideb@cs.stonybrook.edu](mailto:aideb@cs.stonybrook.edu), [madak@cs.stonybrook.edu](mailto:madak@cs.stonybrook.edu), [yuqing.wang.1@stonybrook.edu](mailto:yuqing.wang.1@stonybrook.edu)