# Word Embedding Bias in Large Language Models

## Poomrapee Chuthamsatid - Supervised by Dr. Alex Thomo

University of Victoria, Department of Software Engineering, March 2025

## BACKGROUND

The rapid development of large language models (LLMs) has expanded natural language processing (NLP) applications, from text generation to chatbots.

- Word embeddings are the core of these systems, converting words into numeric vectors based on their statistics usage patterns in text corpora..
- However, embeddings often reflect societal biases, reinforcing stereotypes [1].
- For instance, they may link professions like nurses to women and engineers to men.

## OBJECTIVES

- Analyze gender and race biases in modern LLMs.
  - OpenAI, Cohere, Google, Microsoft, and BGE
- Examine bias in word embeddings and their impact on real-world applications.
  - Tech Industry and Higher Education
- Address biases to ensure fairer and more ethical AI systems.

## DATA SET

| | |
|---|---|
| **Test Word Sets** | The most frequent 100,000 words from the GloVe embedding dataset |
| **Word Embedding Models** | OpenAI, Cohere, Google, Microsoft, and BGE embedding models. |
| **Stimuli Words (Attribute Sets)** | <table> |
| **Big Tech Words** | Big Tech companies based on [3]. Such as Google, Amazon, and Facebook |
| **Top University Words** | The top 50 universities from the 2024 Times Higher Education rankings |

Stimuli Words (Attribute Sets) table:

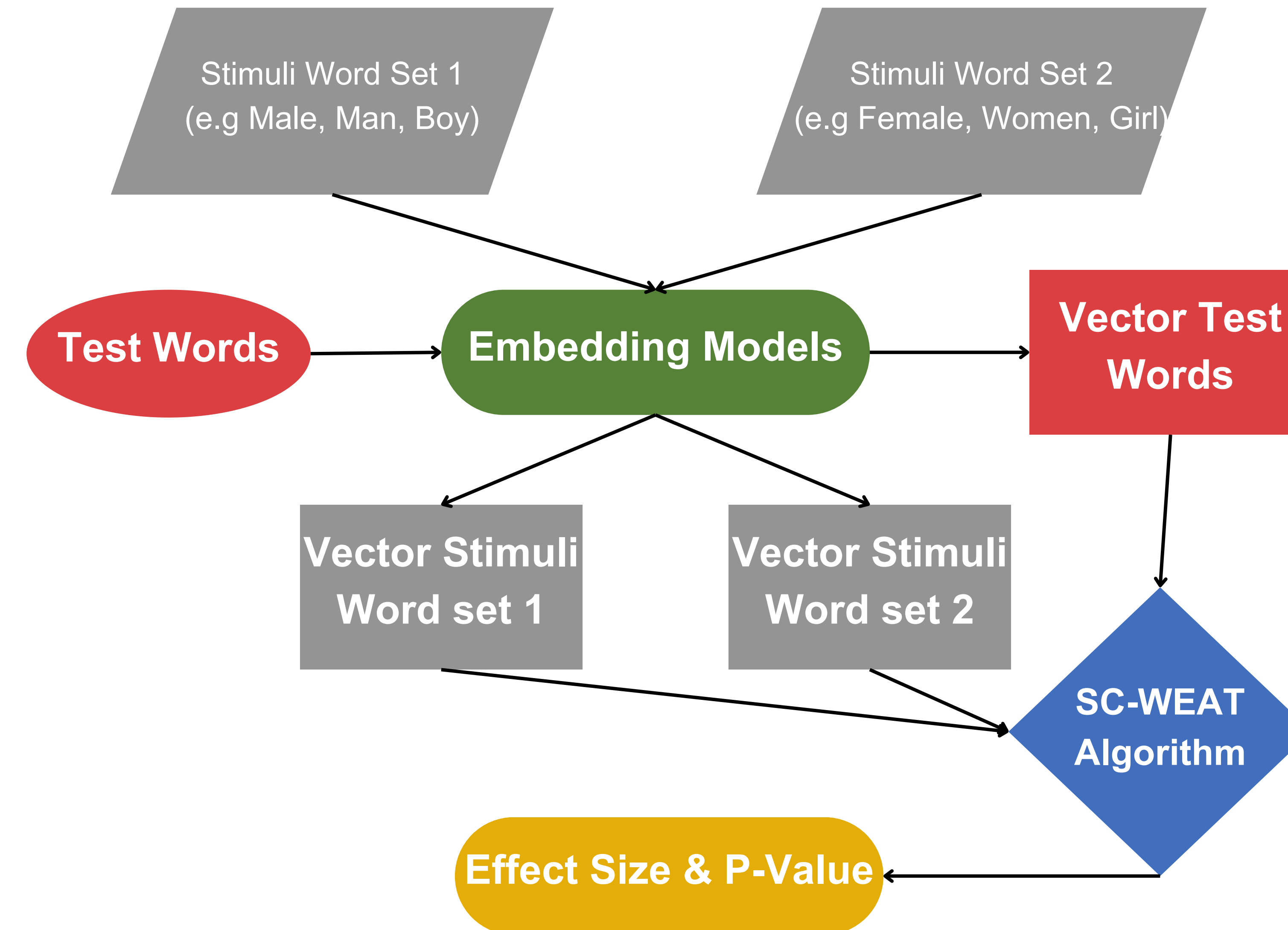| Category | Stimuli Group | Stimuli Words |
|---|---|---|
| Gender | Female | Female, Woman, Girl, Hers, Sister, She, Her, Daughter |
| | Male | Male, Man, Boy, Brother, He, Him, His, Son |
| Race | White | American, Australian, British, Canadian, White, Caucasian, European, French, German, Italian |
| | Asian | Asian, Chinese, Japanese, Indonesian, Indian, Korean, Pakistani, Thai, Filipino, Brown |
| | Black | African, African-American, Black, Congolese, Egyptian, Ethiopian, Haitian, Jamaican, Kenyan, Nigerian |

## WORK FLOW



### SC-WEAT [2]

- Measures bias using cosine similarity between word vectors.

$$ES(\vec{w}, A, B) = \frac{\text{mean}_{a \in A} \cos(\vec{w}, \vec{a}) - \text{mean}_{b \in B} \cos(\vec{w}, \vec{b})}{\text{std\_dev}_{x \in A \cup B} \cos(\vec{w}, \vec{x})}$$
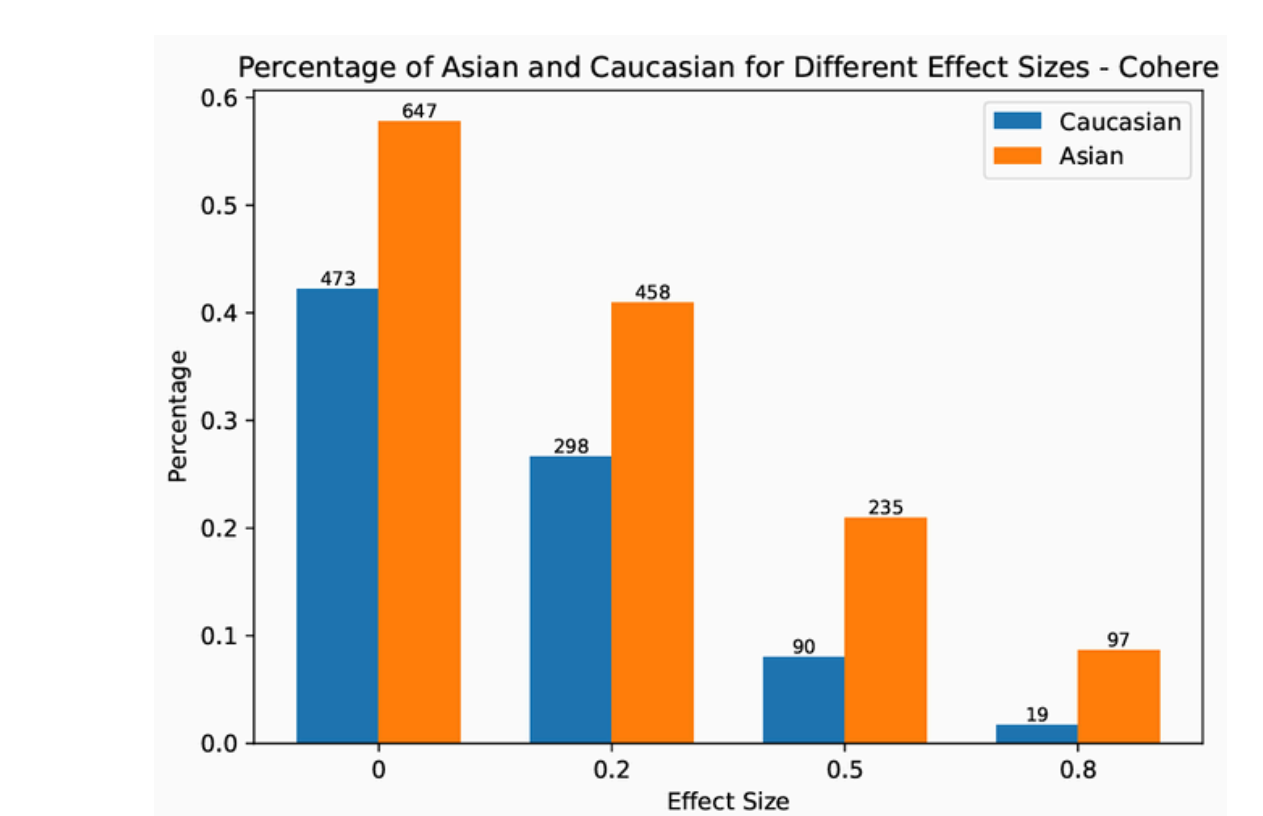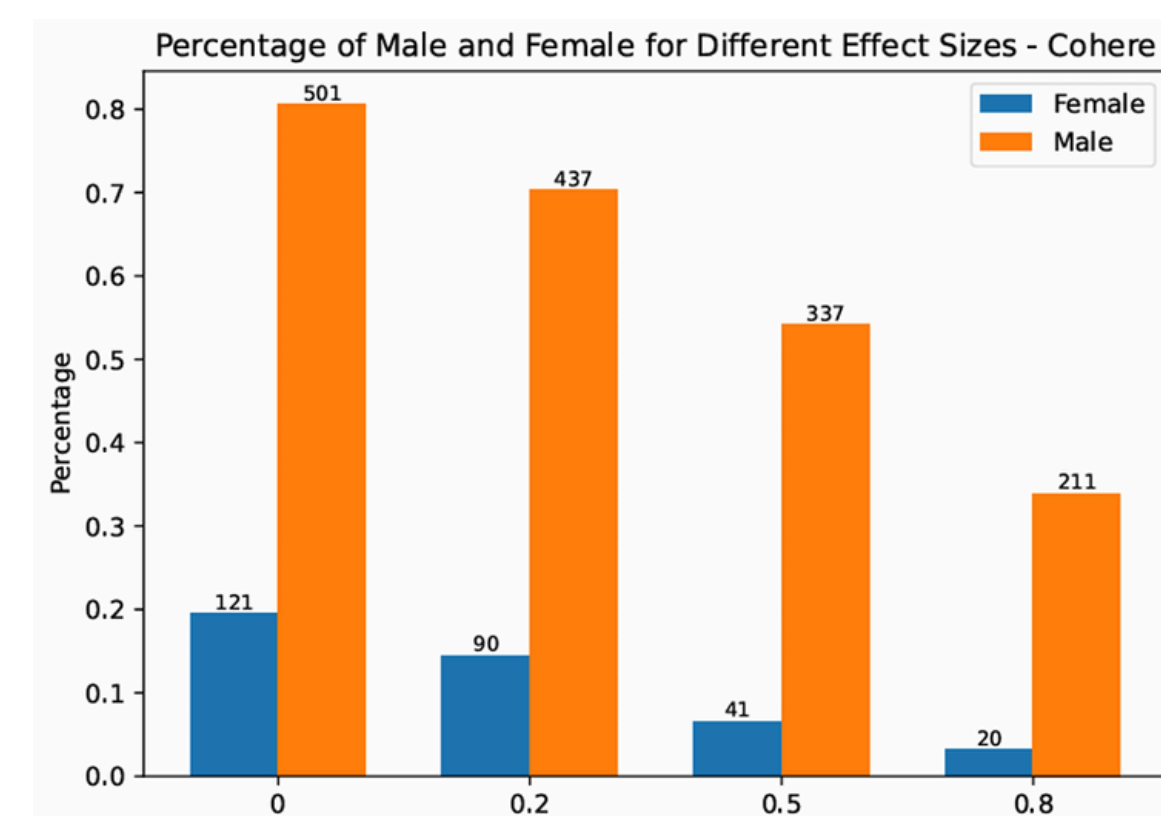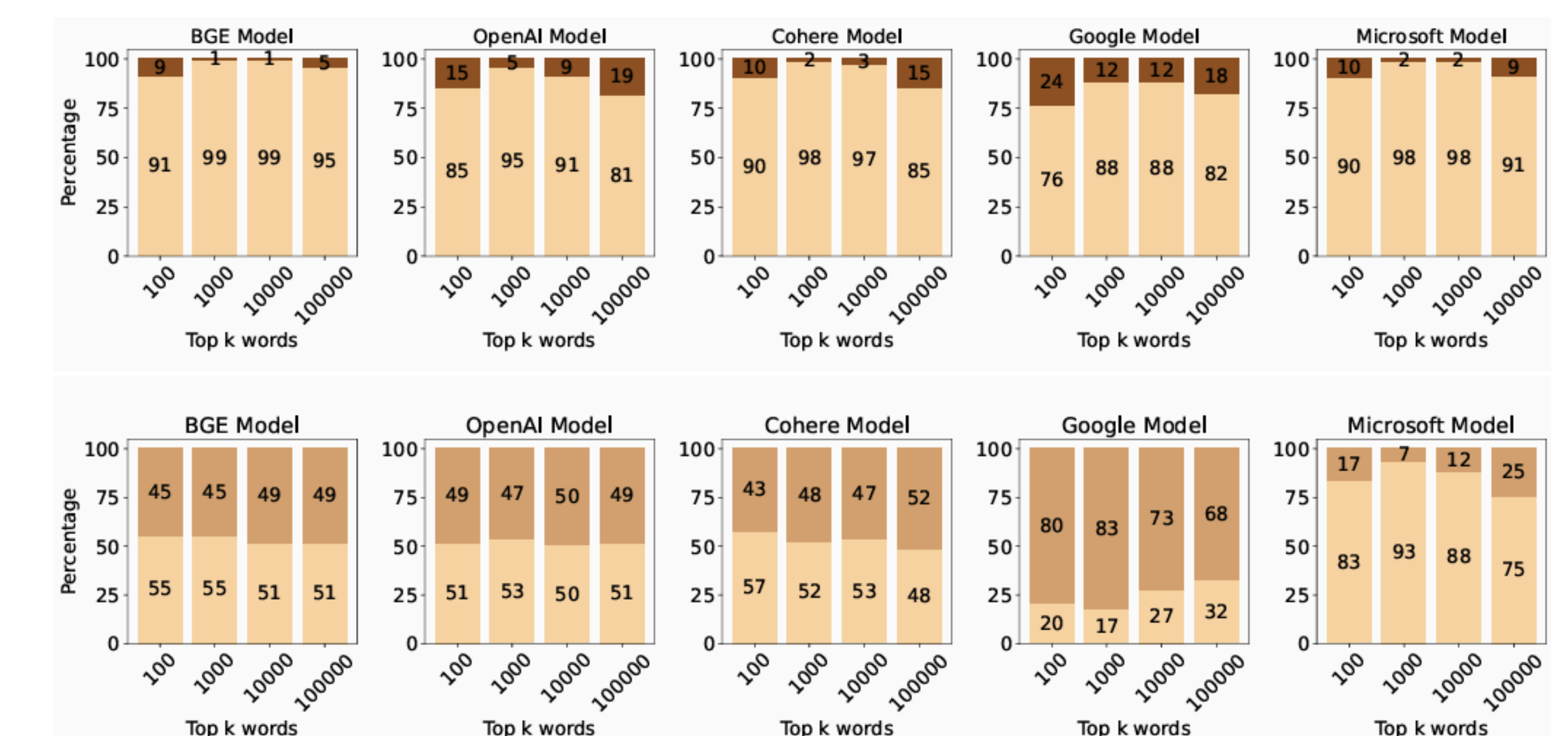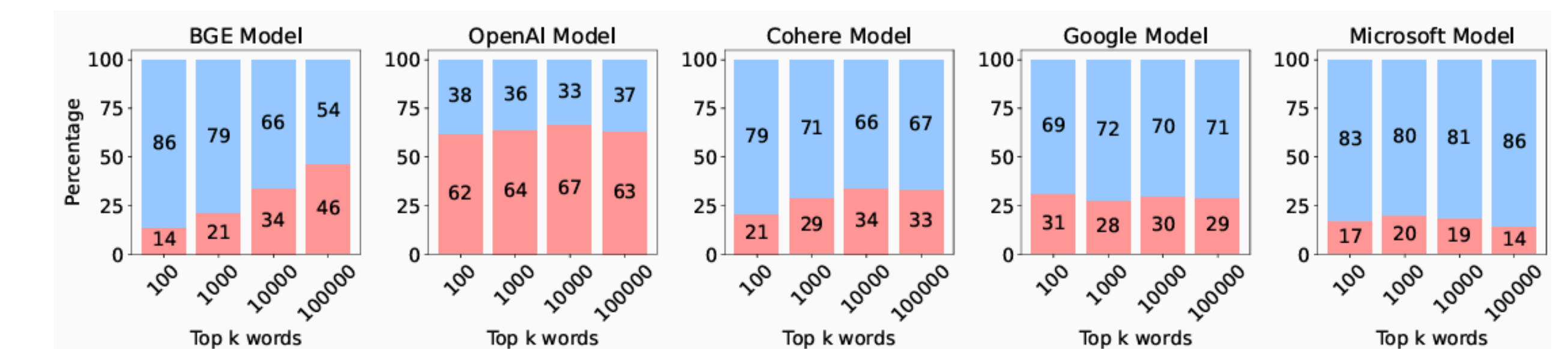


### Analysis

- Bias Analysis by Frequency Range and Effect Size
- Semantic Categories of Gender- and Race-Associated Words
- Bias in Big Tech and Higher Education Contexts

## RESULTS



Gender Association of Top Words: Male is light blue, female is pink



Race Association of Top Words. 1st row: White (lighter color) vs Black (darker color); 2nd row: White (lighter color) vs Asian (darker color)



Big Tech Association by Gender          Top University Association by Race

## CONCLUSION

- **Male** group association dominates in most models
- **Black** group consistently underrepresented
- **Male / Asian** groups dominate in Big Tech
- **Male / Caucasian** groups dominate in Higher Education

## REFERENCES

[1] Eric Michael Smith et al. ""I'm sorry to hear that": Finding new biases in language models with a holistic descriptor dataset". In: Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics, 2022, pp. 9180–9211.

[2] Aylin Caliskan et al. "Gender bias in word embeddings: A comprehensive anal ysis of frequency, syntax, and semantics". In: Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society. 2022, pp. 156–170.

[3] Mohamed Abdalla and Moustafa Abdalla. "The Grey Hoodie Project: Big to bacco, big tech, and the threat on academic integrity". In: Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society. 2021, pp. 287–297.

University of Victoria