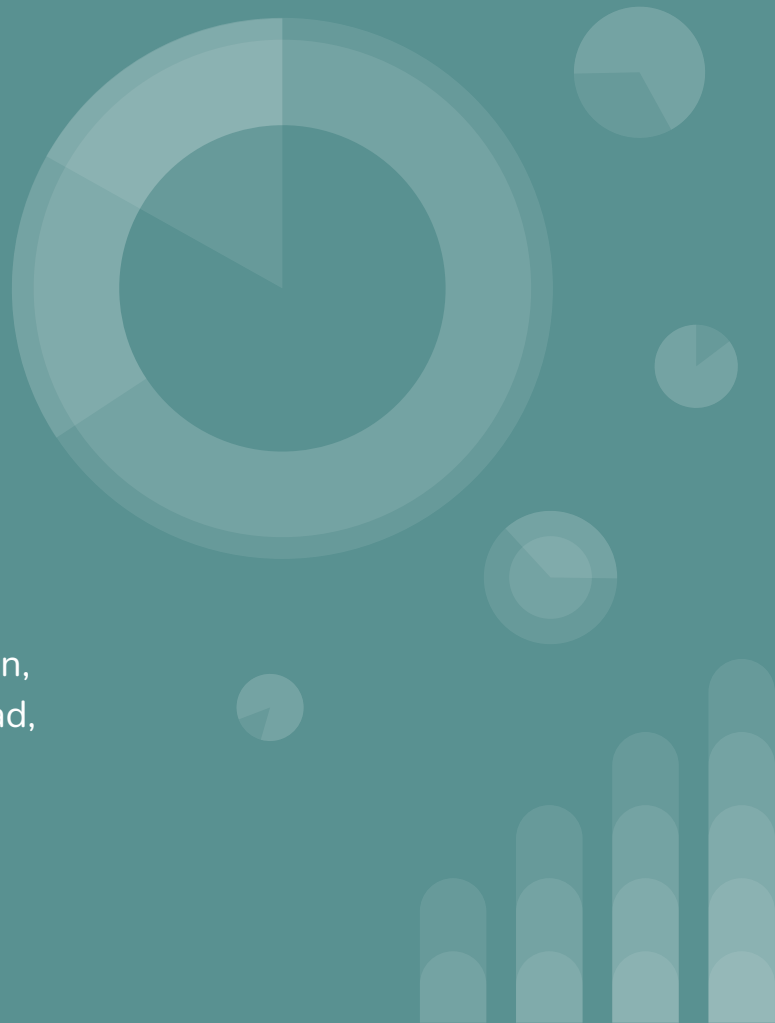# ChatGPT Bias

Hiresh Poosarla, Rishi Sinha, Hayden Fu, Alexander Suen, Aamod Gandhi, Vincent Lo, Harshaan Chugh, Lavi Avigad, Harish Senthilkumar, Karthik Subramanian

Mentor: Phil Mui

# Introduction & Background

- What is ChatGPT and how does it work?
- What are some issues and problems with ChatGPT?
- How do we aim to analyze and quantify ChatGPT's (political) bias?

# Methodology



**THE BIPARTISAN PRESS**

- Human-Generated Questions
- GPT Questions
- Bipartisan Press API
- Statistical Z-Test
  - Z-score

# Results

| Topic | Standard Deviation | Z- Test Statistically Significant (Y/N) | Z-Score | Z P-Val | Original Mean Score | Scaled Mean Score |
|-------|-------------------|----------------------------------------|---------|---------|---------------------|-------------------|
| Abortion | 7.357498498 | N | 1.371464039 | 0.1702303462 | 1.842272967 | 0.03838068681 |
| Gun Control | 5.284963183 | N | 1.988457776 | 4.68E-02 | -6.738644367 | -0.1403884243 |
| Climate Change | 2.658053943 | Y | 13.88575102 | 7.73E-44 | -6.128644367 | -0.127680091 |
| Animal Testing | 1.89695981 | Y | -6.44295277 | 1.17E-10 | -2.2314258 | -0.0464880375 |
| Healthcare | 4.566345259 | N | -2.83308651 | 0.004610090769 | -2.361935067 | -0.04920698056 |
| Freedom of Religion | 4.542296182 | N | -2.636871064 | 0.008367462923 | -2.186773067 | -0.04555777222 |
| Death Penalty | 3.89821359 | N | -2.193564431 | 0.02826673786 | -1.5611887 | -0.03252476458 |
| Gender | 4.144928963 | Y | -10.89901562 | 1.17E-27 | -8.2479067 | -0.1718313896 |
| Racism + Police | 7.350375827 | N | -2.701506409 | 0.006902614992 | -3.625391567 | -0.07552899097 |
| Marijuana | 3.572207308 | Y | -4.035141861 | 5.46E-05 | -2.696034267 | -0.05616738056 |
| Marriage Equality | 5.948609917 | Y | -3.032669015 | 0.00242401311 | -3.293668433 | -0.06861809236 |

# Conclusion

We used the Bipartisan Press API to score the bias of responses from ChatGPT to a mix of human and GPT-generated questions.

ChatGPT is (relatively) unbiased!

- Most scores were negative (left-leaning)
- Overall, ChatGPT is not polarized in either direction

# Future Works

- Prompt generation to reduce ChatGPT's bias
- Linguistic correlation for prompt <-> bias
  - Prompt generation models
- Eventually generalize this to other GPT models
  - And eventually all other future large-language models

# Thank you!