# Closing Gender Gap: Analyzing the Pipeline Model Output and Building Extension That Detects Keywords

Ruiqi Wang    Advisor: Haiyi Zhu

**Carnegie Mellon University**

## Introduction

### Background
Previously, Ruotong Wang has developed a pipeline model that detects potential gender stereotypes and biases. A list of keywords was generated from the pipeline mode. My work focuses on the use the keyword list.

### Abstract
In the first past, I developed a chrome extension that highlights the keywords on website as a visualizing tool. In the second part, I analyze how the keyword list overlaps with gender biases from previous literature.

## Research

### Chrome Extension

I developed a chrome extension that highlights all the keywords generated from the pipeline model. Red represents keywords that have negative coefficients, which are female-indicative, while blue is for male.



### Analysis of the PipelineModel Output

To quantitatively evaluate the result from the pipeline model, we perform the following evaluation: First, we identified a list of gender stereotype indicators from previous literature. Then, we used an automated method to determine how our keywords overlap with an existing stereotype.

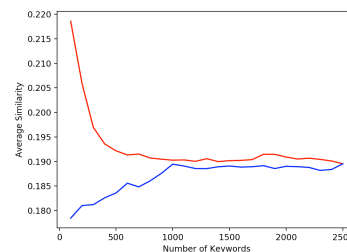**Identify gender stereotype categorization lexicon in previous literature**
We manually collect 39 literatures on Google scholar using the keywords "Wikipedia gender stereotype", "gender bias", etc. Then, we list all possible types of bias and categorize them into eight dimensions ("powerful", "wealthy"," sporty", etc.).

**Calculate the similarity between the SVM output and seed words**
We quantify the similarity between our pipeline model output and seed words using the token similarity given by the word embedding model in SpaCy. We repeatedly carry out this calculation for the set of words with high absolute coefficients.

**Result**
The plot on the left captures the similarity between top keywords and the seed of appearance bias. The red line represents female, while the blue line represents male. The fact that the red line is above the blue line indicates that our female-indicative keywords are more related to appearance, which corresponds to the previous literature that people tend to focus on appearance when describing female.



## Conclusion

### Discussion
The chrome extension serves as a visualizing tool for the pipeline model, which helps people quickly capture the keywords that might indicate gender bias. The analysis of SVM output, which is based on SpaCy similarity, shows that the SVM model output fits with the gender bias from previous literature.

### Future Work
Many interesting experiments have been left due to lack of time. Below are some ideas: (1). The chrome extension could be further developed by adding analysis to each highlighted word. For example, the bias categorization and SVM coefficient can be added; (2). It could be interesting to carry out interviews with Wikipedia editors and collect their feedback on this extension.

### Acknowledgement