

# Quantifying Gender Bias in Different Corpora

Marzieh Babaeianjelodar

Josh Gordon

Jeanna Matthews

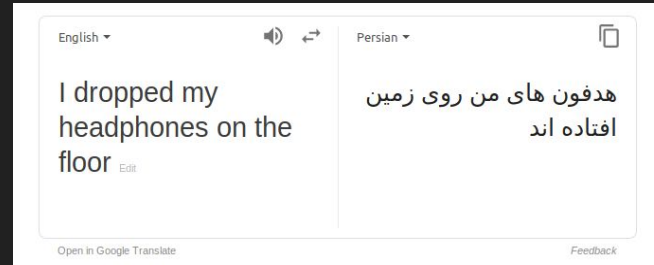
# Introduction

- A part of my thesis research
- A group of researchers
- Accepted in FATES on the Web Conference



# Machine Learning Systems

- Machine Learning ?
  - Train examples
  - Learns like a human (e.g. Amazon Alexa)
  - Can predict the future
  - Examples: Machine Translation, Sentiment Analysis
  - Computer Vision (train images)



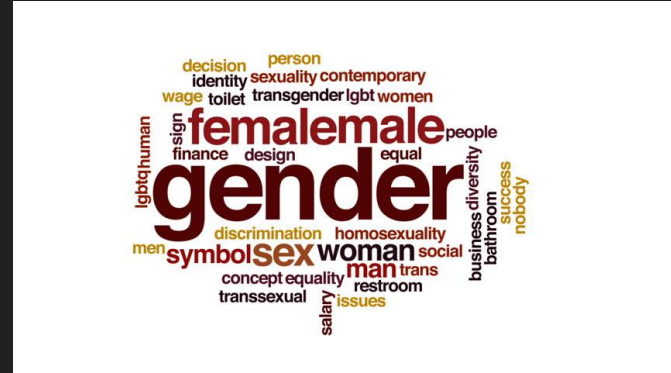
# Effect of NLP in our lives

- Natural Language Processing: NLP
- NLP: Making a computer understand human language
- Used in different applications:
  - Resume parsing
  - Predicting the oscars (based on movie ranks)
  - Computer Science students (applications for universities)
- NLP systems can carry Bias:
  - Through the algorithm
  - Through the data
  - Hiring more male students rather than female students



# Data

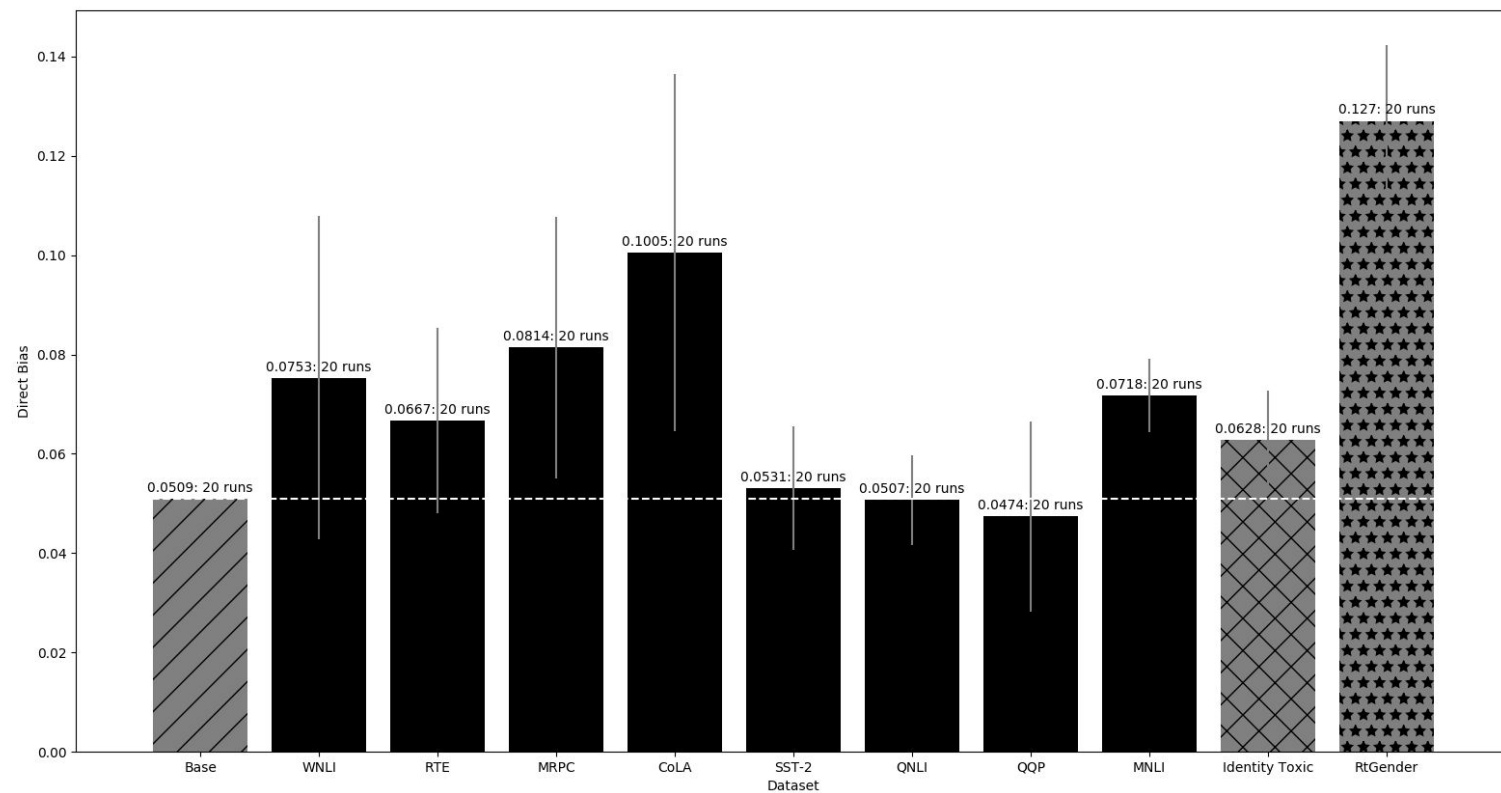
- Body of text to feed to the NLP systems:
  - GLUE: General Language Understanding Evaluation
  - RTGender: labeled based on (source of response to social media comments) f and m
  - Toxic identity: labeled based on attacking race, gender, religion, homosexuality, etc
- Importance: How the system carries these biases



# Word Embeddings and Gender Bias

- Words are represented as positions in space
- The distribution is learned from the data
  - Position carries semantic meaning
  - Position can encode gender
- Can compute “gender direction”
- Bias is represented by average similarity to this direction

# Our Results



# Importance of Thorough Testing

- Biases show up in seemingly innocuous systems
- Can be unclear who is accountable
  - Wider issues with giving computers agency
- Debiasing algorithms



Thank You for your Attention