# MITIGATING GENDER BIAS IN NLP: LITERATURE REVIEW

ACL 2019
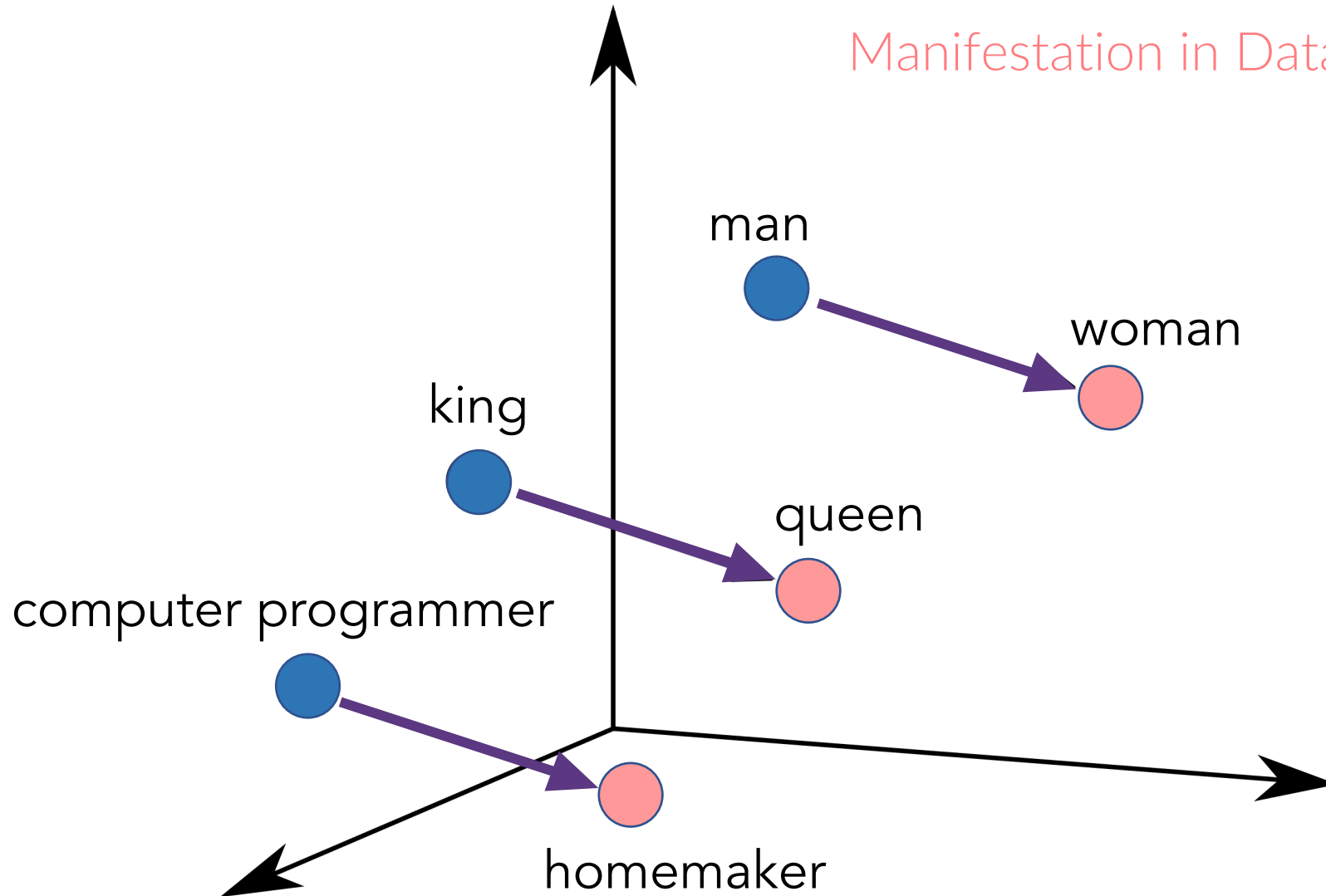
**UC Santa Barbara, UCLA**

Tony Sun · Andrew Gaut · Shirlyn Tang · Yuxin Huang

Mai ElSherief · Jieyu Zhao · Diba Mirza · Elizabeth Belding · Kai-Wei Chang · William Yang Wang

# Gender Bias Origins

# What is Gender Bias?

"Gender bias is the preference or prejudice toward one gender over another"

Allocation vs Representation Bias
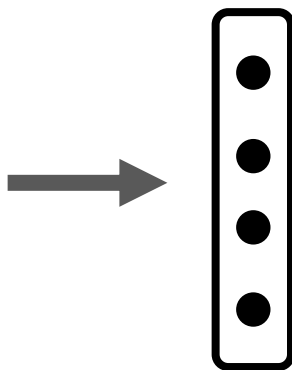
# Gender Bias in Machine Translation

User Input    "*He* is a nurse. *She* is a doctor."
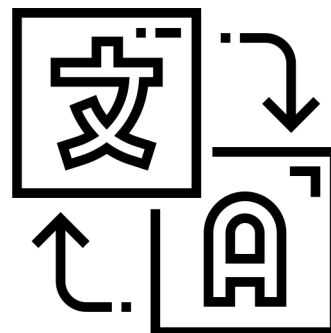
Machine Translation Model

Machine Translation Output

Ő ápolónő. Ő egy orvos

"*She* is a nurse. *He* is a doctor."

Text-based Dataset

Pre-processed data

Propagation of gender bias

# Gender Bias in NLP

Word Embeddings

Language Modeling

Coreference Resolution

Machine Translation

Speech Recognition

Sentiment Analysis

Caption Generation

# Gender Debiasing Pipeline

Task Specific
Training Set

NLP Algorithm

Gender Bias
Evaluation Test Set

Debiasing
Gender

Observing Gender Bias in NLP
Algorithm's Predictions
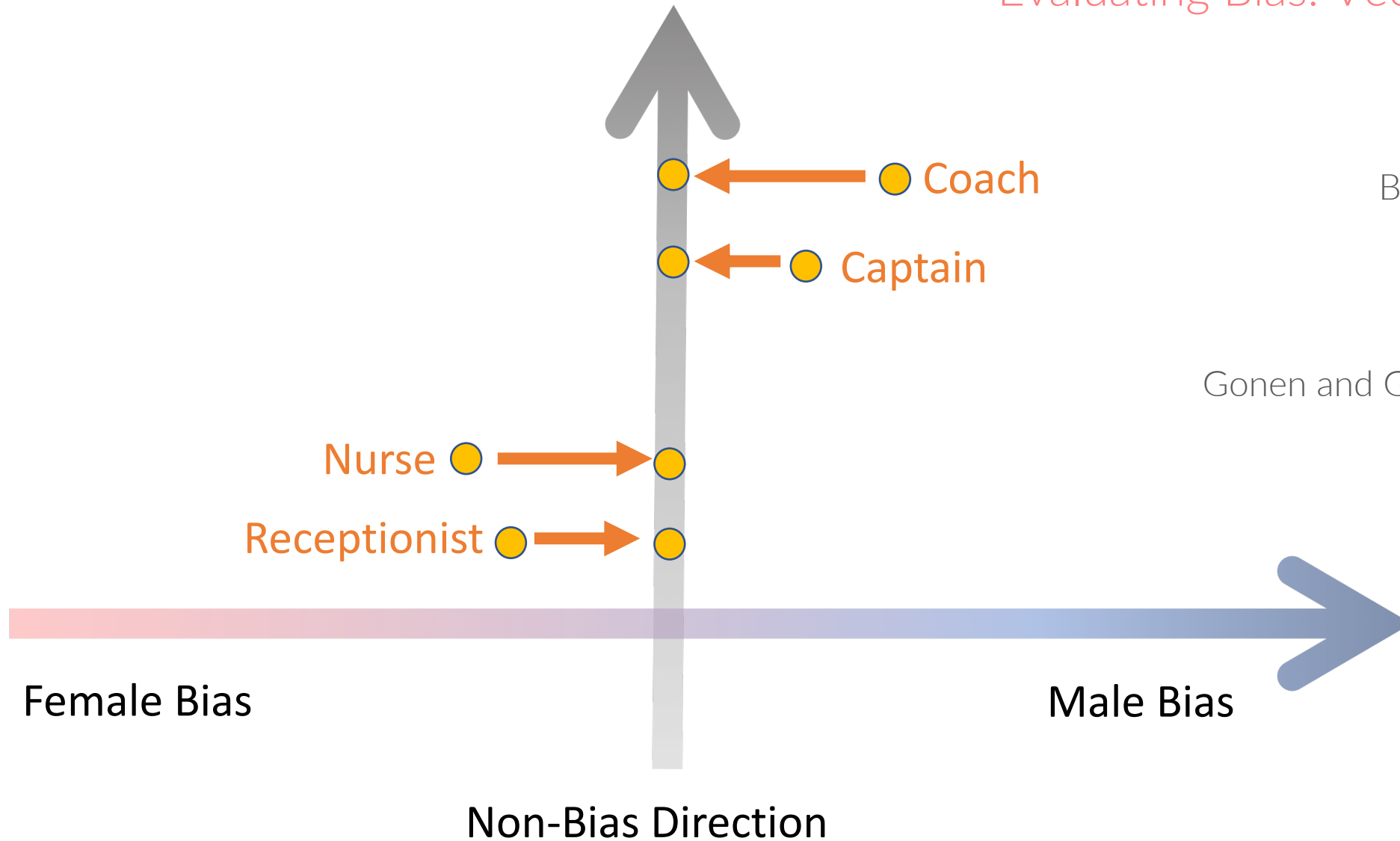
Gender Bias
Observation

6

# Gender Bias Evaluation

Method Categorizations

1. Vector Spaces

2. Gender Bias Evaluation Test Set (GBETs)

# Gender Bias in Word Embeddings

Evaluating Bias: Vector Spaces



Bolukbasi et al., 2016

Caliskan et al, 2017

Manizini et al., 2017

Gonen and Goldberg et al., 2019

Coach

Captain

Nurse

Receptionist

Female Bias

Male Bias

Non-Bias Direction

- New data sets for evaluating gender bias
  - Traditional data sets lack gender-specific information
  - Eliminate confounding variables

Rudinger et al., 2018

Zhao et al., 2018

Webster et al., 2018

Kiritchenko and Mohammad, 2018

He called his mother ⟶ She called her father



The physician hired the secretary because he was overwhelmed with clients.

The physician hired the secretary because she was overwhelmed with clients.

Sentences from WinoBias. (Zhao et al., 2018)

9

# We know there's bias. Now what?

1. Inference vs Retraining

2. Training Data vs Algorithm

3. Key Debiasing Methods

# Categorizing Debiasing Methods

## Retraining

Fully retrain the model

## Inference

Adjust the model's predictions at test time

# Categorizing Debiasing Methods

## Training Data

Debias the training data.

Always Retraining

## Algorithm

Debias the algorithm.

Retraining or Inference
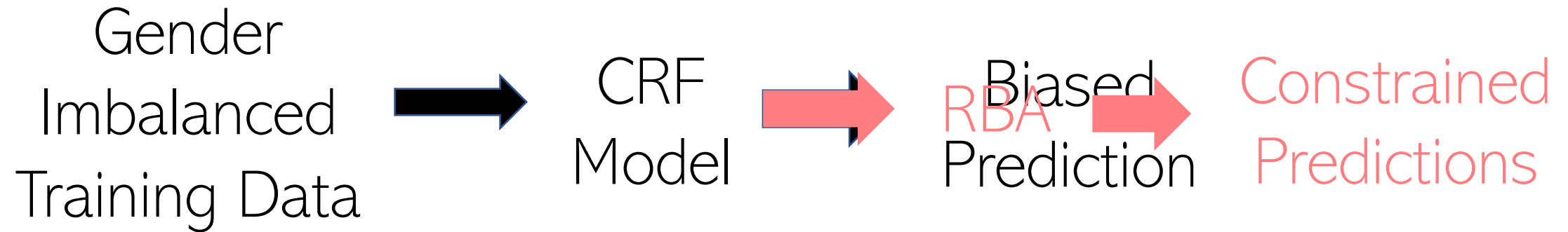
# Debiasing via Data Augmentation

Zhao et al., 2018

## Original Corpus
- He is a doctor.
- The doctor called his mom.

## Genderswapped Corpus
- She is a doctor.
- The doctor called her dad.

**Coreference Resolution Model**

# Constraining Algorithm Predictions

Gender Imbalanced Training Data → CRF Model → RBA Biased Prediction → Constrained Predictions

Five examples from the imSitu visual semantic role labeling dataset (Zhao et al., 2018)

14

# Our Contributions

- We provide a comprehensive literature review of gender bias in NLP
- Critically discuss issues with the purpose of identifying optimizations, knowledge gaps, and directions for future research
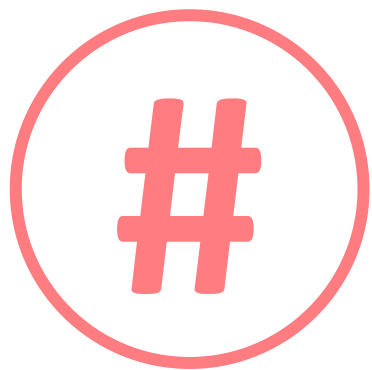
# Future Directions
Ideas Worth Exploring

1.  Non-English Languages

2.  Non-binary Bias

3.  Interdisciplinary Communication

# Thanks for Listening

ACL 2019

# Q & A

Mitigating Gender Bias in NLP: Literature Review

ajg@ucsb.edu          tonysun@ucsb.edu