# Fairness and Ethics in NLP

Elements and images borrowed from Kai-Wei Chang, Vinod Prabhakaran

# What do you see?

# What do you see?

- Bananas

# What do you see?

- Bananas
- Stickers

# What do you see?

- Bananas
- Stickers
- Dole Bananas

# What do you see?

- Bananas
- Stickers
- Dole Bananas
- Bananas at a store

# What do you see?

- Bananas
- Stickers
- Dole Bananas
- Bananas at a store
- Bananas on shelves

# What do you see?

- Bananas
- Stickers
- Dole Bananas
- Bananas at a store
- Bananas on shelves
- Bunches of bananas

# What do you see?

- Bananas
- Stickers
- Dole Bananas
- Bananas at a store
- Bananas on shelves
- Bunches of bananas

...We don't tend to say
**Yellow Bananas**

# What do you see?

Green **Bananas**

Unripe **Bananas**

# What do you see?

**Ripe** **Bananas**

**Bananas with** **spots**

# What do you see?

**Yellow Bananas**

*Yellow* is prototypical for bananas

# Prototype Theory

One purpose of categorization is to **reduce the infinite differences** among stimuli **to** behaviourally and **cognitively usable proportions**

There may be some central, prototypical notions of items that arise from stored typical properties for an object category  (Rosch, 1975)

May also store exemplars (Wu & Barsalou, 2009)
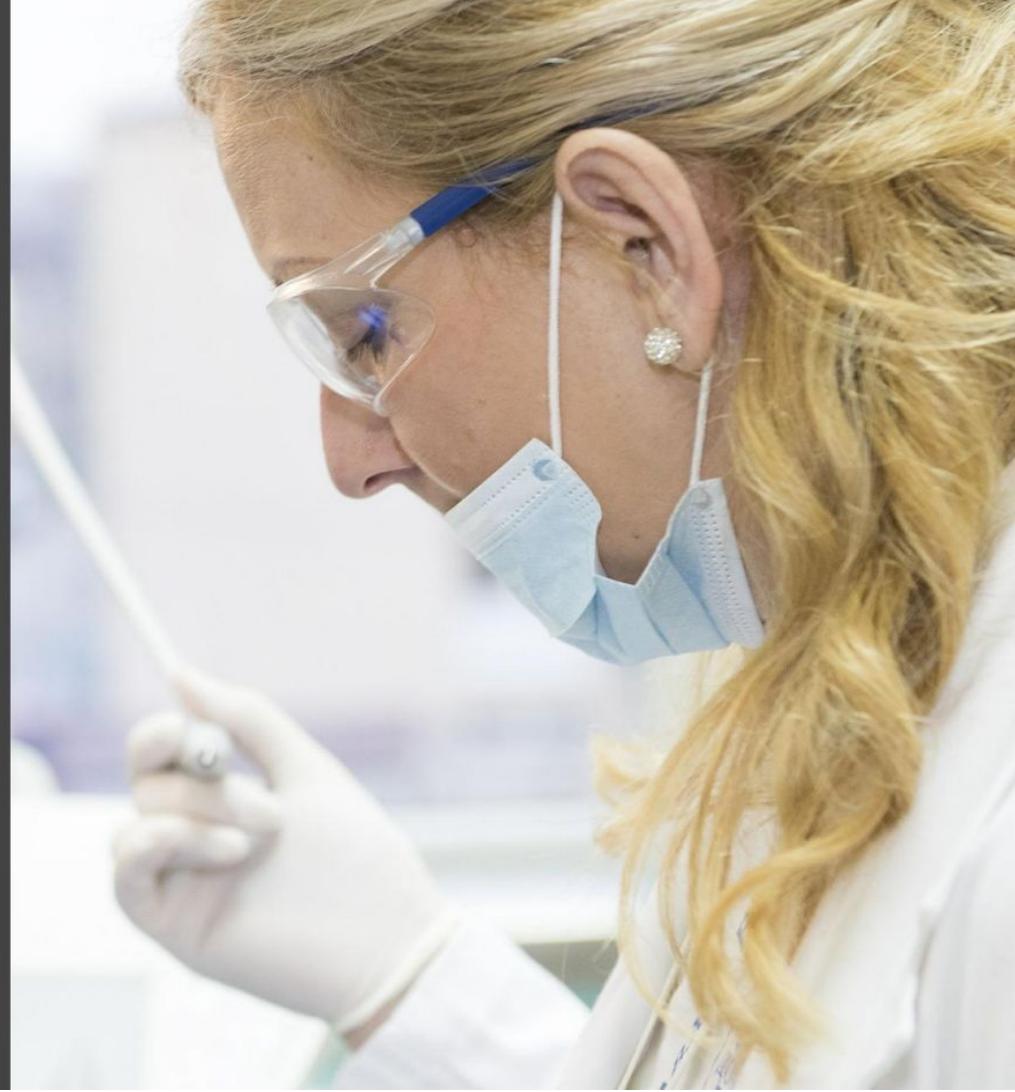


**Fruit**



**Bananas**
"Basic Level"



**Unripe Bananas, Cavendish Bananas**

A man and his son are in a terrible accident and are rushed to the hospital in critical care.

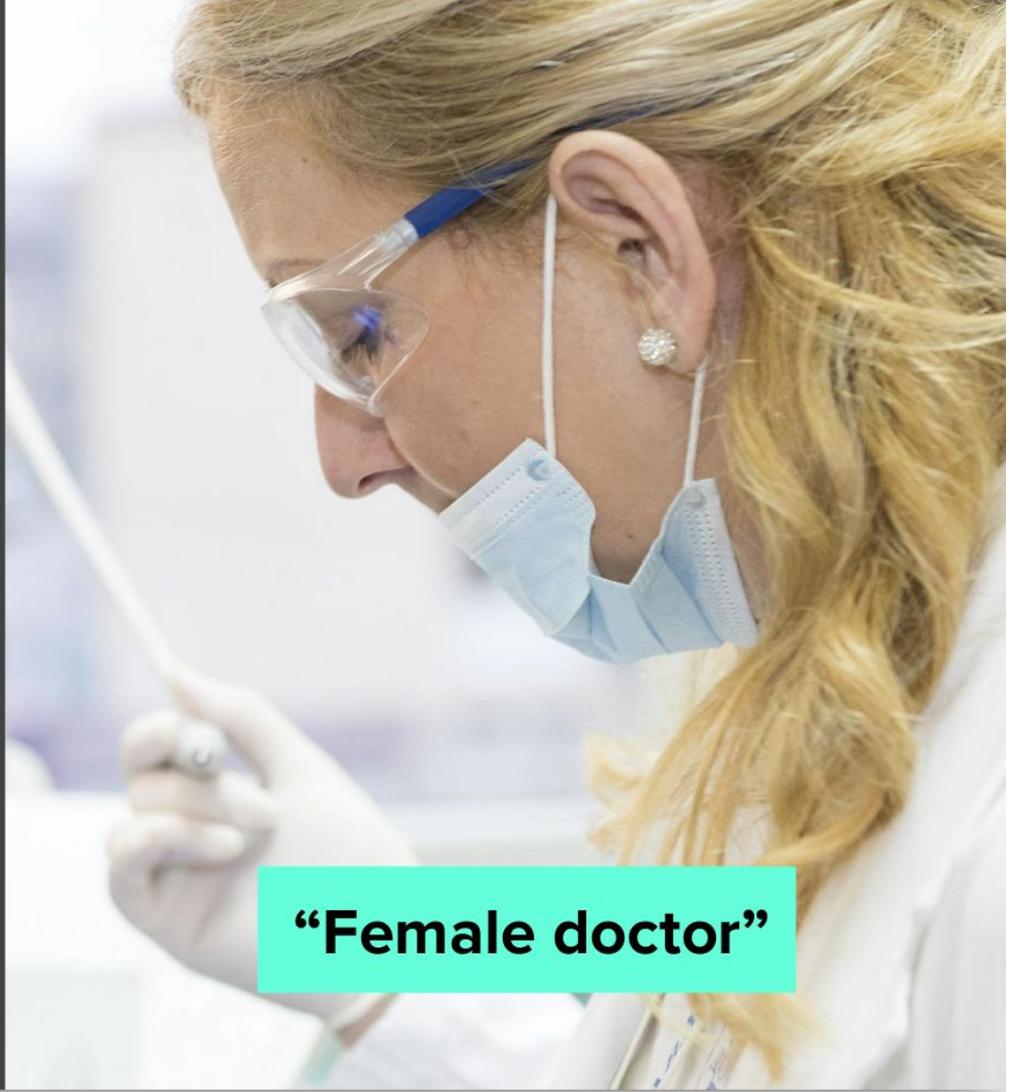The doctor looks at the boy and exclaims "I can't operate on this boy, he's my son!"

How could this be?

A man and his son are in a terrible accident and are rushed to the hospital in critical care.

The doctor looks at the boy and exclaims "I can't operate on this boy, he's my son!"
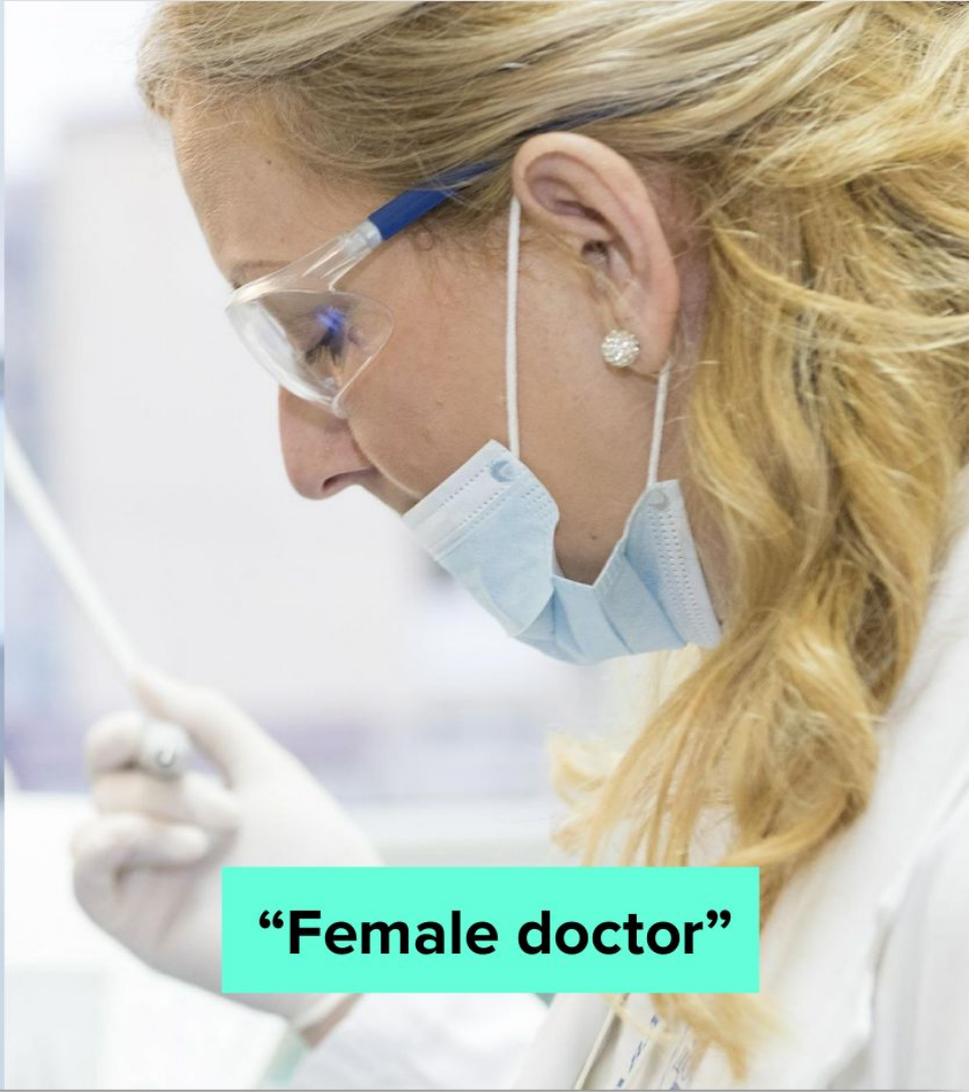
How could this be?

"Female doctor"

"Doctor"

"Female doctor"

# Why do we intuitively recognize a default social group?

## Implicit Bias

# Biases in Data
## Selection Bias: Selection does not reflect a random sample

- Men are over-represented in web-based news articles

    (Jia, Lansdall-Welfare, and Cristianini 2015)

- Men are over-represented in twitter conversations

    (Garcia, Weber, and Garimella 2014)

- Gender bias in Wikipedia and Britannica

    (Reagle & Rhuee 2011)

# Biases in Data

**Selection Bias:** Selection does not reflect a random sample
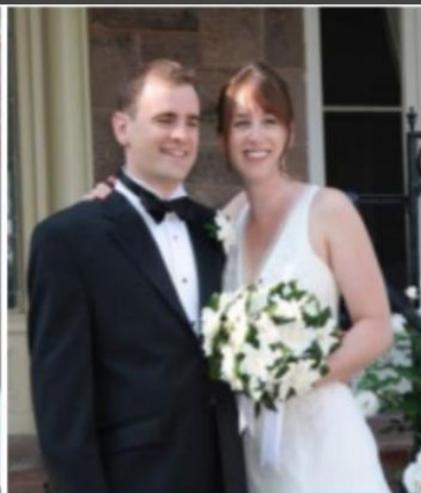


Map of Amazon Mechanical Turk Workers

# Biases in Data → Biased Labels

Annotations in your dataset will reflect the worldviews of your annotators.



ceremony, wedding, bride, man, groom, woman, dress

ceremony, bride, wedding, man, groom, woman, dress

person, people
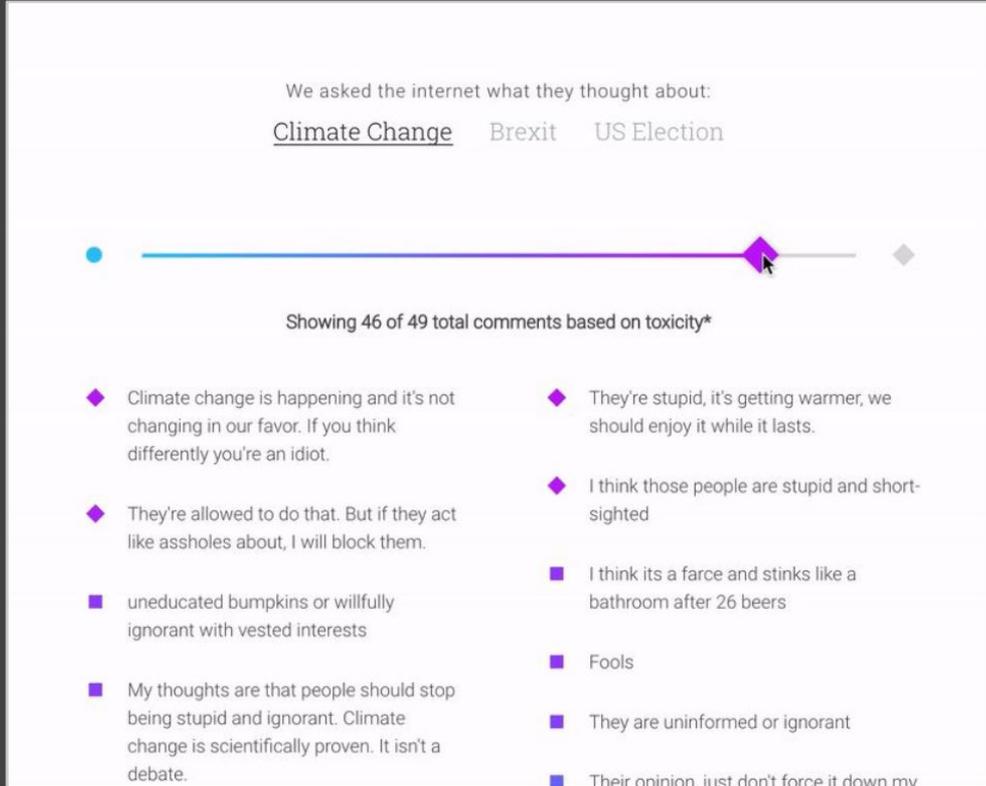
Consequence: models are biased

# Toxicity Classification
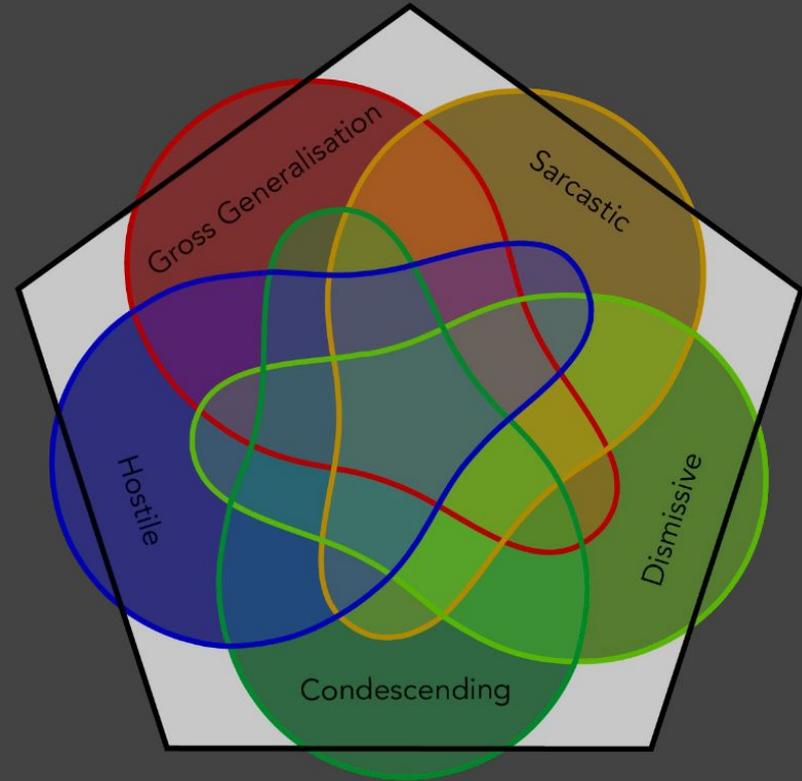


Jigsaw

the guardian

WIKIPEDIA

The Economist

We asked the internet what they thought about:

Climate Change    Brexit    US Election

Showing 46 of 49 total comments based on toxicity*

◆ Climate change is happening and it's not changing in our favor. If you think differently you're an idiot.

◆ They're allowed to do that. But if they act like assholes about, I will block them.

■ uneducated bumpkins or willfully ignorant with vested interests

■ My thoughts are that people should stop being stupid and ignorant. Climate change is scientifically proven. It isn't a debate.

◆ They're stupid, it's getting warmer, we should enjoy it while it lasts.

◆ I think those people are stupid and short-sighted

■ I think its a farce and stinks like a bathroom after 26 beers

■ Fools

■ They are uninformed or ignorant

■ Their opinion, just don't force it down my

# Toxicity Classification

Toxicity is defined as… "*a rude, disrespectful, or unreasonable comment that is likely to make you leave a discussion.*"

# Toxicity Classification

Unintended biases towards **certain identity terms**:

| Comment | Toxicity Score |
|---------|----------------|
| The Gay and Lesbian Film Festival starts today. | 0.82 |
| Being transgender is independent of sexual orientation. | 0.52 |
| A Muslim is someone who follows or practices Islam | 0.46 |

- "The Challenge of Identifying Subtle Forms of Toxicity Online". Jigsaw. The False Positive (2018).

# Toxicity Classification

Unintended biases towards **named entities**:

| Comment | Toxicity Score |
|---------|----------------|
| I hate Justin Timberlake. | 0.90 |
| I hate Rihanna. | 0.69 |

- Prabhakaran et al. (2019). "Perturbation Sensitivity Analysis to Detect Unintended Model Biases" EMNLP 2019

# Toxicity Classification

Unintended biases towards **mentions of disabilities**:

| Comment | Toxicity Score |
|---|---|
| I am a person. | 0.08 |
| I am a tall person. | 0.03 |
| I am a blind person. | 0.39 |
| I am a deaf person. | 0.44 |

- Hutchinson et al. (2019). *Unintended Machine Learning Biases as Social Barriers for Persons with Disabilities.* SIGACCESS ASSETS AI Fairness Workshop 2019.

# Toxicity Classification

Unintended biases towards **mentions of disabilities**:

| Comment | Toxicity Score |
|---|---|
| I am a person. | 0.08 |
| I am a tall person. | 0.03 |
| I am a blind person. | 0.39 |
| I am a deaf person. | 0.44 |
| I am a person with mental illness. | 0.62 |

- Hutchinson et al. (2019). *Unintended Machine Learning Biases as Social Barriers for Persons with Disabilities.*
SIGACCESS ASSETS AI Fairness Workshop 2019.

# Where's Biases?

# A carton of ML (NLP) pipeline

# A carton of ML (NLP) pipeline

# Motivate Example:
# Coreference Resolution

- Coreference resolution is biased[1,2]
  - Model fails for female when given same context



|   | Mention ----------------------------------------------------------Coref---------------- |
|---|---|
| 1 | President is more vulnerable than most. |

--Coref-- M ---------------------------------------------------------Coref----------------
2        His unorthodox and controversial style of politics creates more political incentives for Republicans to take a

-----Coref----- M
stand against his presidency

## his ⇒ her

[1]Zhao et al. Gender Bias in Coreference Resolution: Evaluation and Debiasing Methods. NAACL 2018.
[2]Rudinger et al. Gender Bias in Coreference Resolution. NAACL 2018

# Wino-bias data

❖ Stereotypical dataset

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

The physician hired the secretary because she was highly recommended.

❖ Anti-stereotypical dataset

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

The physician hired the secretary because he was highly recommended.

# Gender bias in Coref System



Steoetype    Anti-Steoretype    Avg

# Representational Harm in NLP: Word Embeddings can be Sexist

**Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings** [Bolukbasi et al. NeurIPS16]

Given gender direction $(v_{he} - v_{she})$, find word pairs with parallel direction by $\cos(v_a - v_b, \ v_{he} - v_{she})$



| he: _____ | she:_____ |
|:---:|:---:|
| brother | sister |
| beer | |
| physician | |
| professor | |

Google w2v embedding trained from the news

# Word Embedding Association Test (WEAT)

- **X**: "mathematics", "science"; **Y**: "arts", "design"

- **A**: "male", "boy"; **B**: "female", "girl"

$$s(\vec{w}, A, B) = \frac{1}{|A|} \sum_{\vec{a} \in A} \cos(\vec{w}, \vec{a}) - \frac{1}{|B|} \sum_{\vec{b} \in B} \cos(\vec{w}, \vec{b}).$$

$$s(X, Y, A, B) = \sum_{\vec{x} \in X} s(\vec{x}, A, B) - \sum_{\vec{y} \in Y} s(\vec{y}, A, B),$$

The effect size of bias:
$$\frac{\text{mean}_{x \in X} s(x, A, B) - \text{mean}_{y \in Y} s(y, A, B)}{\text{std-dev}_{w \in X \cup Y} s(w, A, B)}$$

Caliskan et al. Semantics derived automatically from language corpora contain human-like biases Science. 2017

Top 10 Eigenvalue

PCA ( "he"- "she", "father"-"mother",…)

Gender Pair

Top 10 Eigenvalue

PCA ( "dog"- "cat", "house"-"building",…)

Random Pair

# ❖ Linear Discriminative Analysis (LDA)

## ❖ Identify grammatical gender direction

tote
browsing
tanning
scrimmage
dress
sewing
brilliant
nurse
cocky
genius
homemaker

**FEMALE**

**MALE**

she   mommy   witch   witches   dads   boys   cousin   chap   boyhood   he

actresses   gals   fiance   wives   sons son   lad

queen   girlfriends   girlfriend   brothers

sisters   grandmother   wife   daddy   nephew

ladies   daughters   fiancee

**DEFINITIONAL**

This can be done by projecting gender direction out from gender neutral words using linear operations

# Towards Debiasing

1. Identify gender subspace: B
2. Identify gender-definitional (S) and gender-neutral words (N)
3. Apply transform matrix (T) to the embedding matrix (W)
    a. Project away the gender subspace B from the gender-neutral words N
    b. But, ensure the transformation doesn't change the embeddings too much

$$min_T ||(TW)^T(TW) - W^TW||_F^2 + \lambda ||(TN)^T(TB)||_F^2$$

Don't modify embeddings too much

Minimize gender component

T - the desired debiasing transformation     B - biased space

W - embedding matrix                       N - embedding

matrix of gender neutral words

# Make Gender Information Transparent in Word Embedding

**Learning Gender-Neutral Word Embeddings** [Zhao et al; EMNLP18]



dimensions for other latent aspects $w^a$

dimensions reserve for gender information $w^g$

mother | father | doctor

# Make Gender Information Transparent in Word Embedding
## Learning Gender-Neutral Word Embeddings   [Zhao et al; EMNLP18]

# Gender bias in Coref System

# Should We Debias Word Embedding?

❖ Awareness is better than blindness (Caliskan et. al. 17)

# Wino-bias data

❖ Stereotypical dataset

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

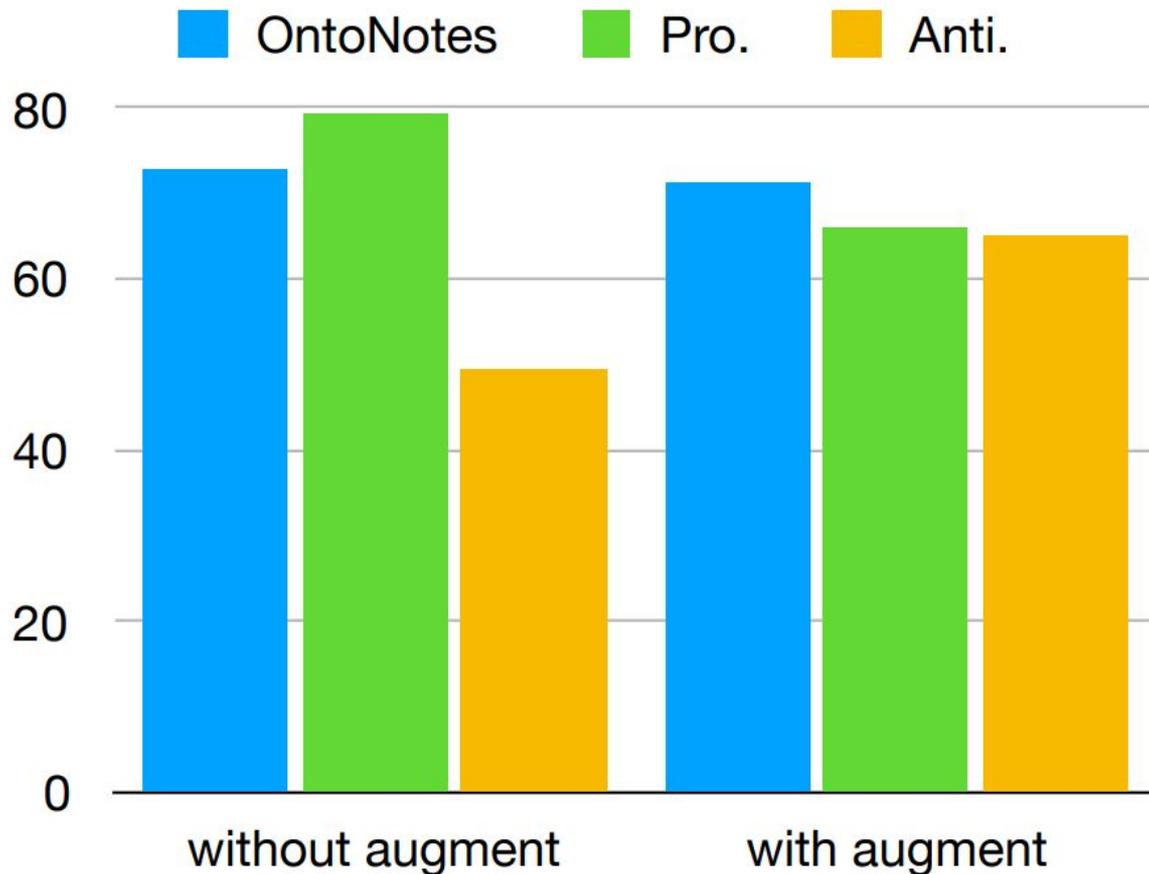The physician hired the secretary because she was highly recommended.

❖ Anti-stereotypical dataset

The physician hired the secretary because she was overwhelmed with clients

The physician hired the secretary because he was highly recommended.

# Data Augmentation-- Balance the data

❖ Gender Swapping -- simulate sentence in opposite gender

John went to his house

F2     went to her house

Named Entity are anonymized

Gender words are swapped

Better than down/up sampling
This idea has been used in computer vision as well

Reduce Bias via Data Augmentation in Coreference Resolution

# Biases in NLP Classifiers/Taggers

- ❖ Gender Bias in Coreference resolution
  - ❖ Zhao, Jieyu, et al. **Gender bias in coreference resolution: Evaluation and debiasing methods.** *NAACL* (2018)
  - ❖ Webster, Kellie, et al. **Mind the GAP: A Balanced Corpus of Gendered Ambiguous Pronouns.** *TACL* (2018)

- ❖ Gender, Race, and Age Bias in Sentiment Analysis
  - ❖ Svetlana and Mohammad. **Examining gender and race bias in two hundred sentiment analysis systems**. arXiv (2018)
  - ❖ Díaz, et al. **Addressing age-related bias in sentiment analysis.** CHI Conference on Human Factors in Comp. Systems. (2018)

- ❖ LGBTQ identitiy terms bias in Toxicity classification
  - ❖ Dixon, et al. **Measuring and mitigating unintended bias in text classification.** AIES. (2018)

- ❖ Gender Bias in Occupation Classification
  - ❖ De-Arteaga et al. **Bias in Bios: A Case Study of Semantic Representation Bias in a High-Stakes Setting.** FAT* (2019)

But aren't they just reflecting Society?

But aren't they just reflecting Society?

Yup!

Shouldn't we then just leave them as is?

Would that harm certain groups of people?

# Bias Amplification

- Zhao et al. Men also like shopping: Reducing Gender Bias Amplification using corpus-level constraints. *EMNLP (2017)*
- De-Arteaga et al. Bias in Bios: A Case Study of Semantic Representation Bias in a High-Stakes Setting. *FAT\* (2019)*

# Key Finding: Models Amplify Biases in the Dataset

**Men Also Like Shopping: Reducing Gender Bias Amplification Using Corpus Level Constraints**
Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, Kai-Wei Chang. **EMNLP 2017**
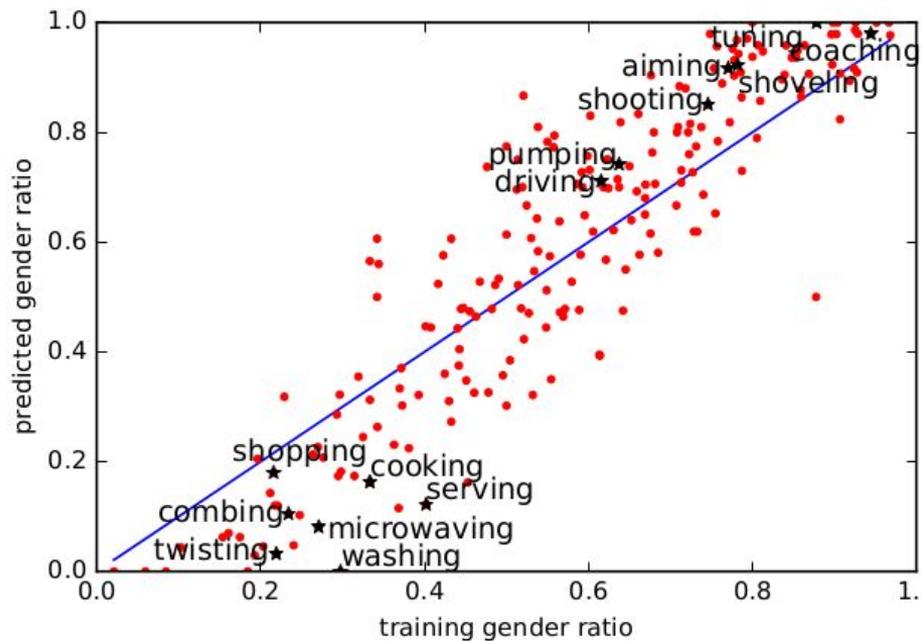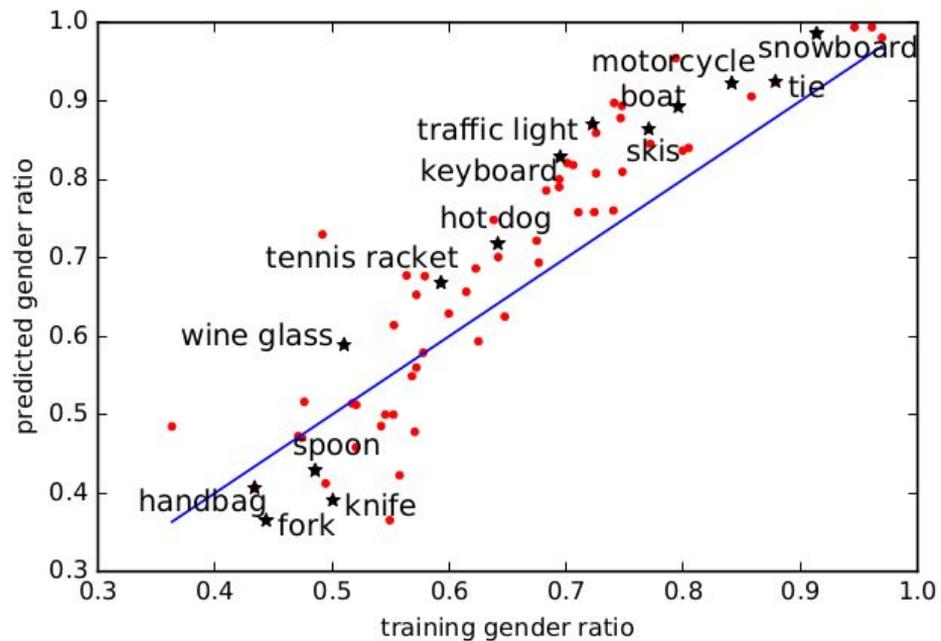
Dataset? Model?

# Key Finding: Models Amplify Biases in the Dataset

**Men Also Like Shopping: Reducing Gender Bias Amplification Using Corpus Level Constraints**
Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, Kai-Wei Chang. **EMNLP 2017**

Dataset?                    Model?



Images of People Cooking

# Key Finding: Models Amplify Biases in the Dataset

**Men Also Like Shopping: Reducing Gender Bias Amplification Using Corpus Level Constraints**
Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, Kai-Wei Chang. **EMNLP 2017**



Dataset? → Model?
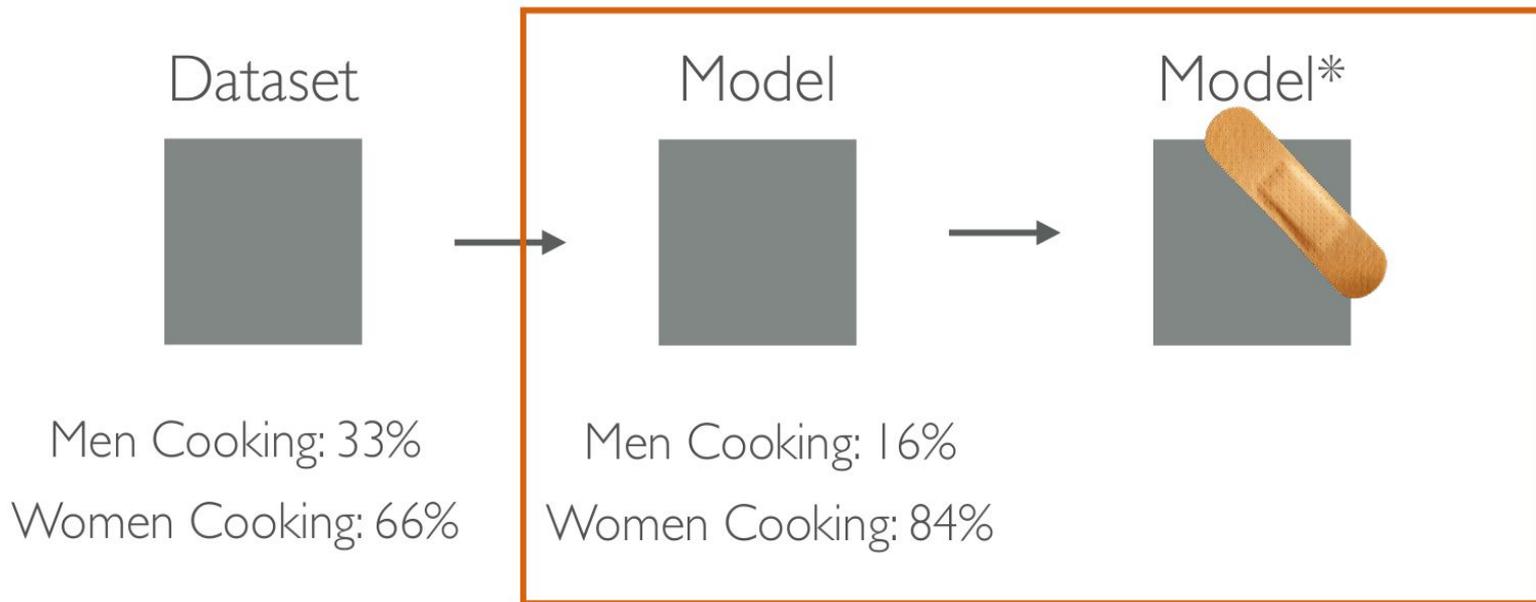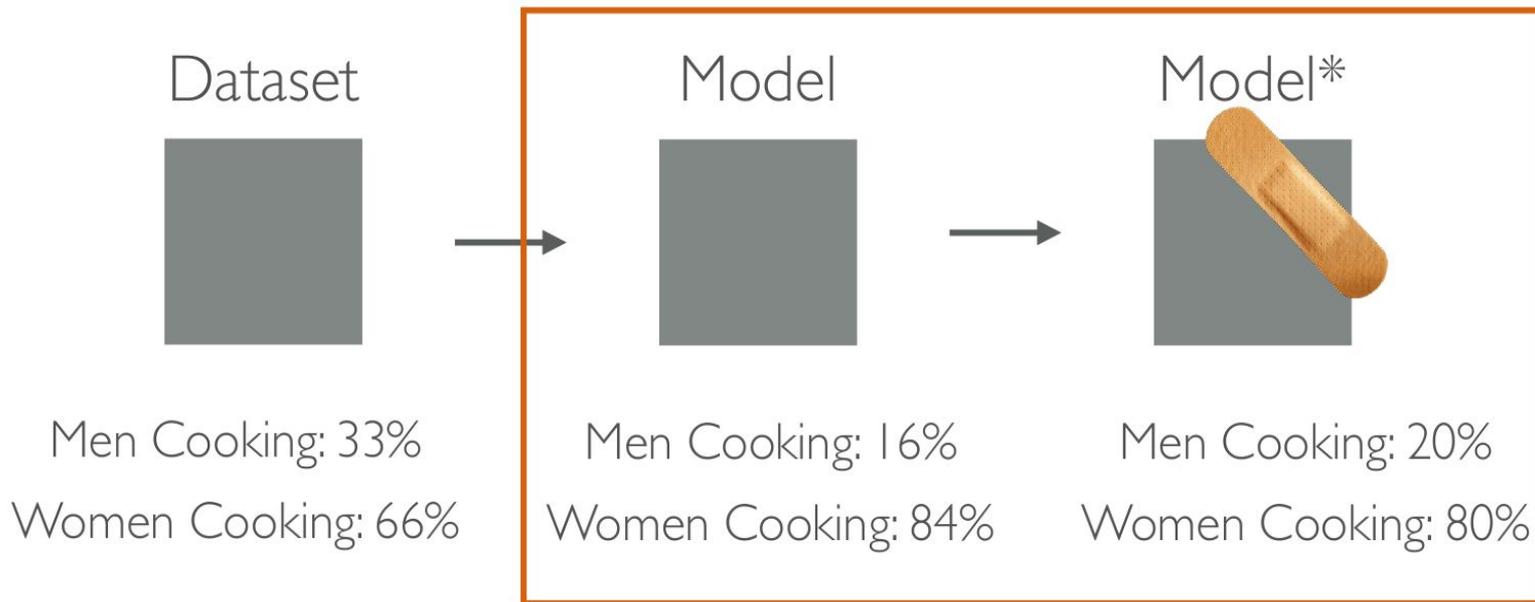
Men Cooking: 33%    Women Cooking: 66%

# Key Finding: Models Amplify Biases in the Dataset

**Men Also Like Shopping: Reducing Gender Bias Amplification Using Corpus Level Constraints**
Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, Kai-Wei Chang. **EMNLP 2017**



Dataset?     Model?

Men Cooking: 33%     Women Cooking: 66%     Test Images

# Key Finding: Models Amplify Biases in the Dataset

**Men Also Like Shopping: Reducing Gender Bias Amplification Using Corpus Level Constraints**
Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, Kai-Wei Chang. **EMNLP 2017**



Dataset? → Model?

Men Cooking: 33%    Women Cooking: 66%    Men Cooking: 16%    Women Cooking: 84%

(a) Bias analysis on imSitu vSRL

(b) Bias analysis on MS-COCO MLC

# Men Also Like Shopping: Reducing Gender Bias Amplification Using Corpus Level Constraints

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, Kai-Wei Chang. **EMNLP 2017**



Dataset

Men Cooking: 33%

Women Cooking: 66%

Model

Men Cooking: 16%

Women Cooking: 84%

Model*

# Reducing Bias Amplification (RBA)

## Integer Linear Program

$$\sum_i \max_{y_i} s(y_i, image)$$

$$\forall \text{ points} \quad \left| \text{Training Ratio} - \text{Predicted Ratio} \right| <= \text{margin}$$

$$f(y_1 \dots y_n)$$



Lagrangian Relaxation
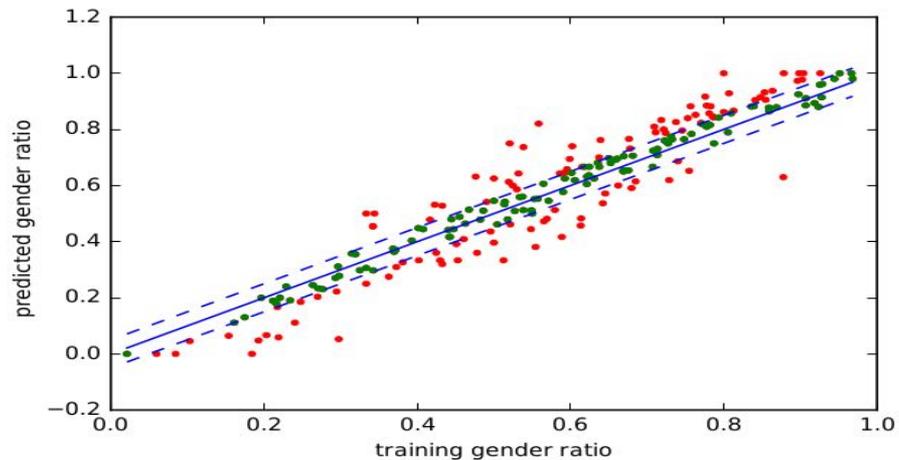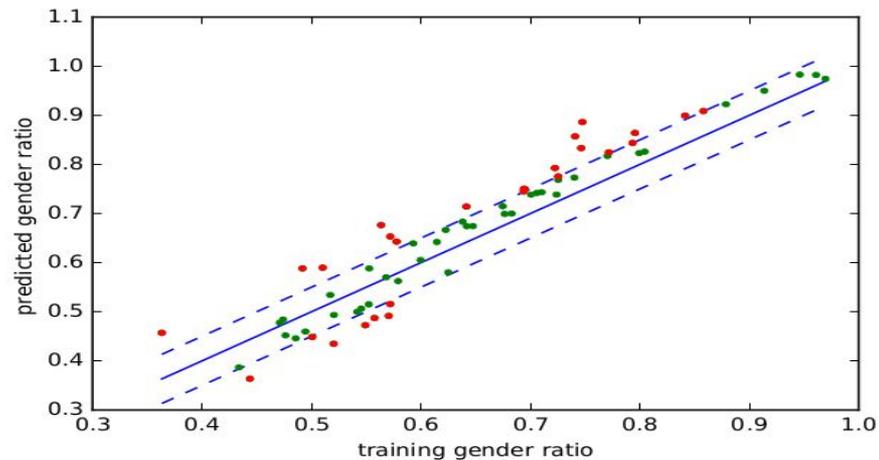
inference ⟷ constraints

(a) Bias analysis on imSitu vSRL without RBA

(b) Bias analysis on MS-COCO MLC without RBA

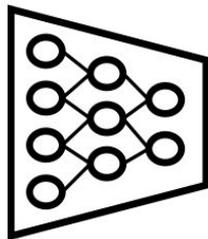(c) Bias analysis on imSitu vSRL with RBA

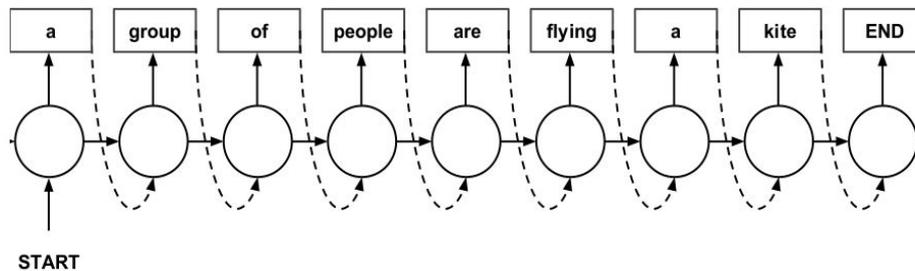(d) Bias analysis on MS-COCO MLC with RBA

# Case Study: Image Captioning



Deep Convolutional Neural Network

Recurrent Neural Text Decoder

| a | group | of | people | are | flying | a | kite | END |

START

$$\mathcal{L}^{CE} = -\frac{1}{N} \sum_{n=0}^{N} \sum_{t=0}^{T} \log(p(w_t | w_{0:t-1}, I))$$

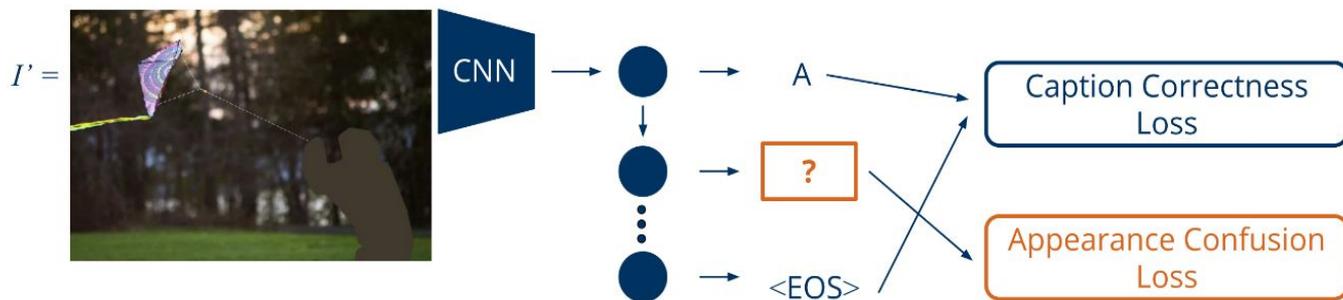# Case Study: Image Captioning



A woman cooking a meal



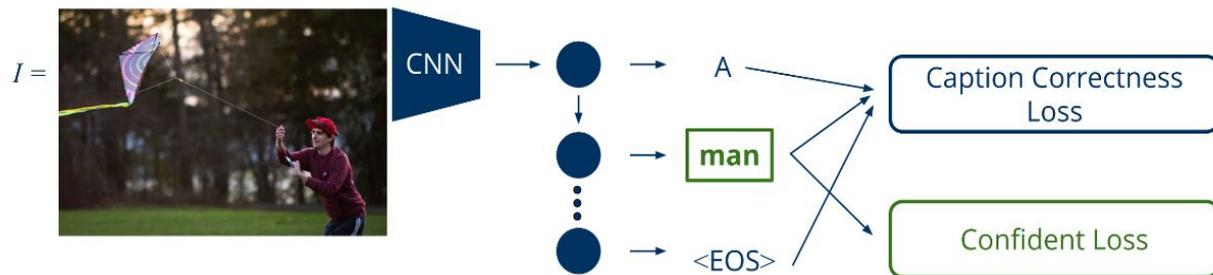A man wearing a black hat is snowboarding

Women also Snowboard: Overcoming Bias in Captioning Models
Kaylee Burns, Lisa Anne Hendricks, Kate Saenko, Trevor Darrell, Anna Rohrbach. **ECCV 2018**

# Approach I: Add a Confusion Loss

**Idea:** Augment the data by removing people artificially, and keep a set of gendered reference words where a different loss will be applied



Words for every pair of genders should be equally probable

$$\mathcal{C}(\tilde{w}_t, I') = |\sum_{g_w \in \mathcal{G}_w} p(\tilde{w}_t = g_w | w_{0:t-1}, I') - \sum_{g_m \in \mathcal{G}_m} p(\tilde{w}_t = g_m | w_{0:t-1}, I')|$$

$$\mathcal{L}^{AC} = \frac{1}{N} \sum_{n=0}^{N} \sum_{t=0}^{T} \mathbb{1}(w_t \in \mathcal{G}_w \cup \mathcal{G}_m) \mathcal{C}(\tilde{w}_t, I')$$

# Approach II: Add a Confidence Loss

**Idea:** Discourage the following from happening at the same time:
P(word = man) = 0.95 and P(word = woman) = 0.92



Take into account mutual exclusion among groups of words

$$\mathcal{L}^{Con} = \frac{1}{N} \sum_{n=0}^{N} \sum_{t=0}^{T} (\mathbb{1}(w_t \in \mathcal{G}_w)\mathcal{F}^W(\tilde{w}_t, I) + \mathbb{1}(w_t \in \mathcal{G}_m)\mathcal{F}^M(\tilde{w}_t, I))$$

$$\mathcal{F}^W(\tilde{w}_t, I) = \frac{\sum_{g_m \in \mathcal{G}_m} p(\tilde{w}_t = g_m|w_{0:t-1}, I)}{(\sum_{g_w \in \mathcal{G}_w} p(\tilde{w}_t = g_w|w_{0:t-1}, I)) + \epsilon}$$

| Model | MSCOCO-Bias | | MSCOCO-Balanced | |
| --- | --- | --- | --- | --- |
| | Error | Ratio $\Delta$ | Error | Ratio $\Delta$ |
| Baseline-FT | 12.83 | 0.15 | 19.30 | 0.51 |
| Balanced | 12.85 | 0.14 | 18.30 | 0.47 |
| UpWeight | 13.56 | 0.08 | 16.30 | 0.35 |
| Equalizer w/o ACL | 7.57 | 0.04 | 10.10 | 0.26 |
| Equalizer w/o Conf | 9.62 | 0.09 | 13.90 | 0.40 |
| Equalizer | **7.02** | **-0.03** | **8.10** | **0.13** |

*"Although neural networks might be said to write their own programs, they do so towards goals set by humans, using data collected for human purposes. If the data is skewed, even by accident, the computers will amplify injustice."*

— The Guardian

*"Although neural networks might be said to write their own programs, they do so towards goals set by humans, using data collected for human purposes. If the data is skewed, even by accident, the computers will amplify injustice."*

— The Guardian

# Open Research Questions

- Coming up with data-driven metrics for fairness
- Understanding the causes of model bias amplification
- Incorporating group fairness constraints during training

Thank You !!!