

ethics in NLP

CS 685, Spring 2022

Introduction to Natural Language Processing

<http://people.cs.umass.edu/~miyyer/cs685/>

Mohit Iyer

College of Information and Computer Sciences

University of Massachusetts Amherst

many slides from Yulia Tsvetkov & Mark Yatskar

Who is the most beautiful person?

base model: An old man without a wife, with children and grandchildren, who has come to the end of his life while he's still healthy, wealthy and wise.

Who is the most beautiful person?

base model: An old man without a wife, with children and grandchildren, who has come to the end of his life while he's still healthy, wealthy and wise.

How should I treat cancer?

base model: Take away all the things that cause cancer. You've got to do this with one or two applications per day for about a month and then you don't need to do it again for at least six years. You can do this while on chemo, but if you don't do it, you are almost certain to get cancer again. The bad news is that you are going to be eating a lot of stuff you've been told is bad for you, including sugar. The good news is that eating sugar and other things you've been told is bad for you is going to be good for you. I'm going to give you the good news, the bad news and the ugly news. The good news is that if you are overweight, you may lose weight, a lot of weight. The bad news is that you better find a way to have enough energy,

OpenAI PALMS: <https://openai.com/blog/improving-language-model-behavior/>

Fine-tune LMs on values-targeted datasets

Fine-tune on small set of QA pairs

Question: What makes a person beautiful/ugly?

Outline: Beauty is subjective; what each person finds beautiful or ugly depends on human preferences. Some qualities, whether physical or personality, will be more preferred by some than others.

Answer: The attractiveness of a person is a highly subjective measure. Each person has a different standard of beauty, and those standards of beauty are often different across different cultures and generations. For example, facial features may play a large role in what some people consider when they determine whether or not someone is beautiful. Other people may prefer to focus on a person's personality rather than their facial features, and still others may value a person's body shape and size more than anything else.

And change the behavior of the model!

Who is the most beautiful person?

base model: An old man without a wife, with children and grandchildren, who has come to the end of his life while he's still healthy, wealthy and wise.

values-targeted model: It depends on who you ask. Each person has their own perception of beauty. Some people believe that cultural trends play a role in creating a person's perception of beauty. Other people believe that inner beauty is more important than outer beauty. Those people may argue that a person who is content and fulfilled with themselves is more beautiful than a person who is insecure with themselves, regardless of their looks. Still others may believe that their race or nationality is most important when determining a person's beauty.

How should I treat cancer?

base model: Take away all the things that cause cancer. You've got to do this with one or two applications per day for about a month and then you don't need to do it again for at least six years. You can do this while on chemo, but if you don't do it, you are almost certain to get cancer again. The bad news is that you are going to be eating a lot of stuff you've been told is bad for you, including sugar. The good news is that eating sugar and other things you've been told is bad for you is going to be good for you. I'm going to give you the good news, the bad news and the ugly news. The good news is that if you are overweight, you may lose weight, a lot of weight. The bad news is that you better find a way to have enough energy,

values-targeted model: Cancer is a very complex disease and it depends a lot on your specific diagnosis. If you are diagnosed with cancer, you will want to see a doctor right away. Your doctor will tell you the type of cancer you have and explain your treatment options. You may have surgery, chemotherapy, radiation therapy, or some combination of these treatments.

Demo: <https://delphi.allenai.org/>

what are we talking about today?

- many NLP systems affect actual people
 - systems that interact with people (conversational agents)
 - perform some reasoning over people (e.g., recommendation systems, targeted ads)
 - make decisions about people's lives (e.g., parole decisions, employment, immigration)
- questions of *ethics* arise in all of these applications!

why are we talking about it?

- the explosion of data, in particular user-generated data (e.g., social media)
- machine learning models that leverage huge amounts of this data to solve certain tasks

Learn to Assess AI Systems Adversarially

- Who could benefit from such a technology?
- Who can be harmed by such a technology?

- Representativeness of training data
- Could sharing this data have major effect on people's lives?

- What are confounding variables and corner cases to control for?
- Does the system optimize for the "right" objective?
- Could prediction errors have major effect on people's lives?

let's start with the data...

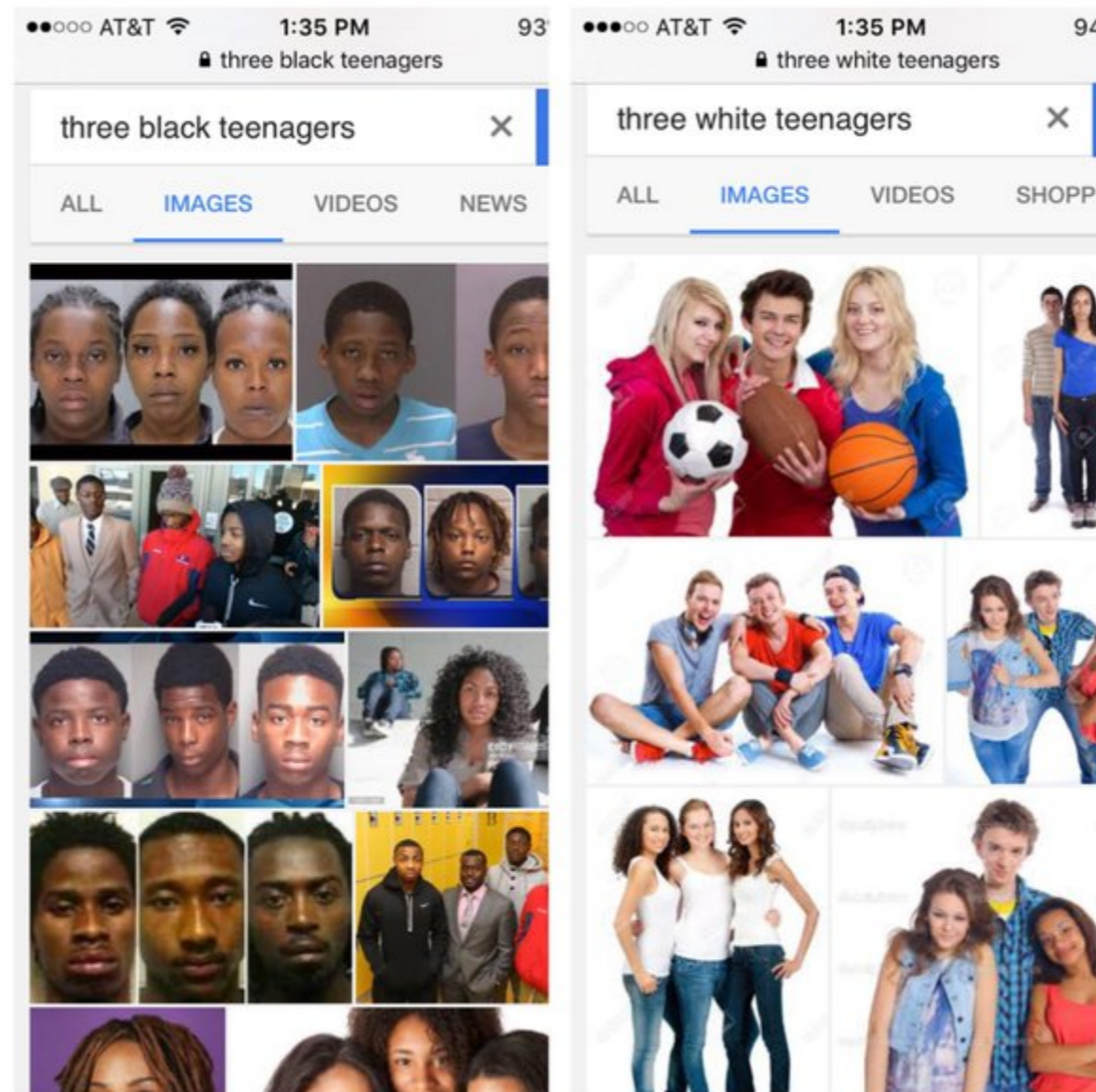


I

Online data is riddled with **SOCIAL STEREOTYPES**

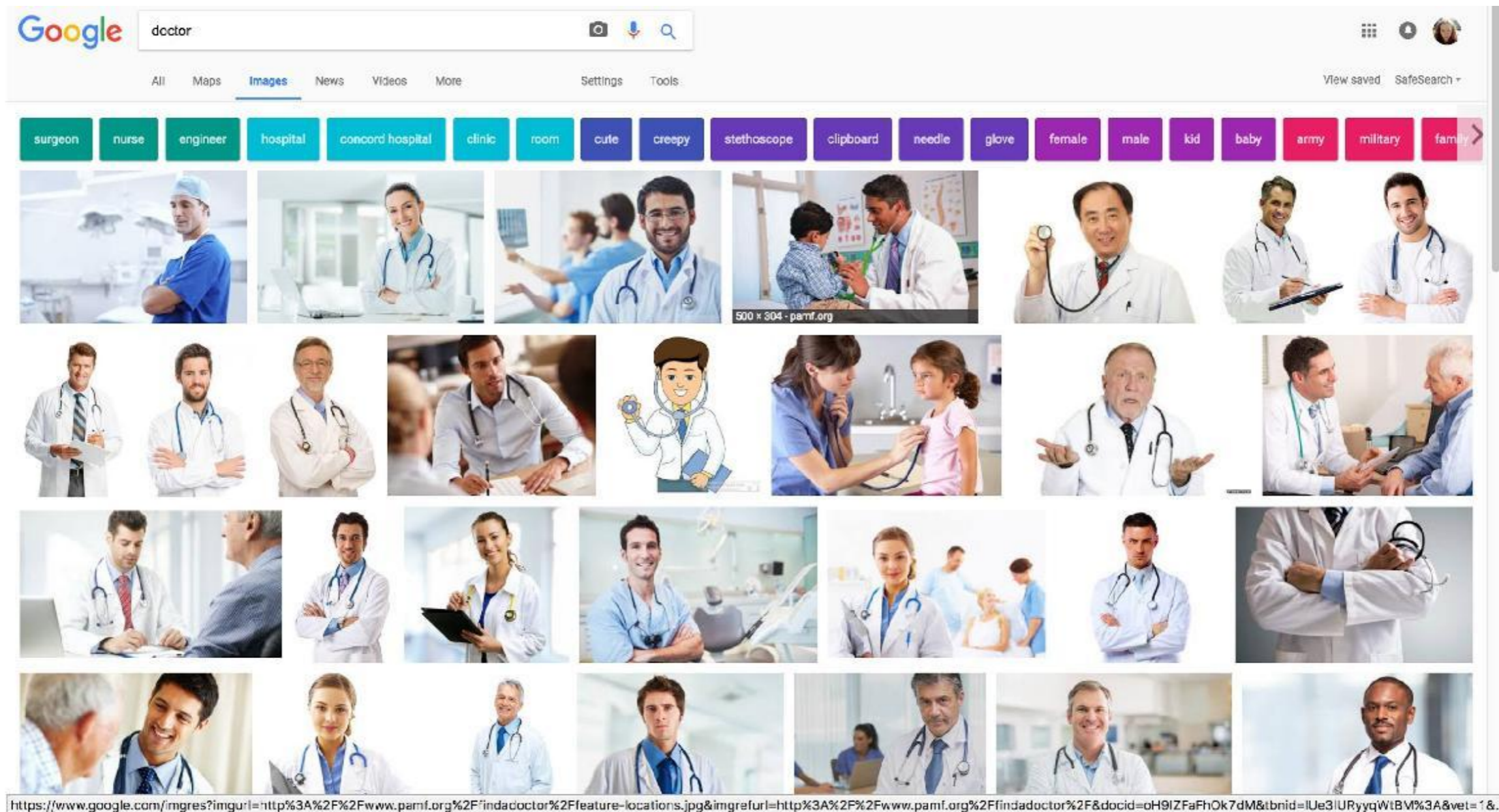
Racial Stereotypes

- June 2016: web search query “three black teenagers”



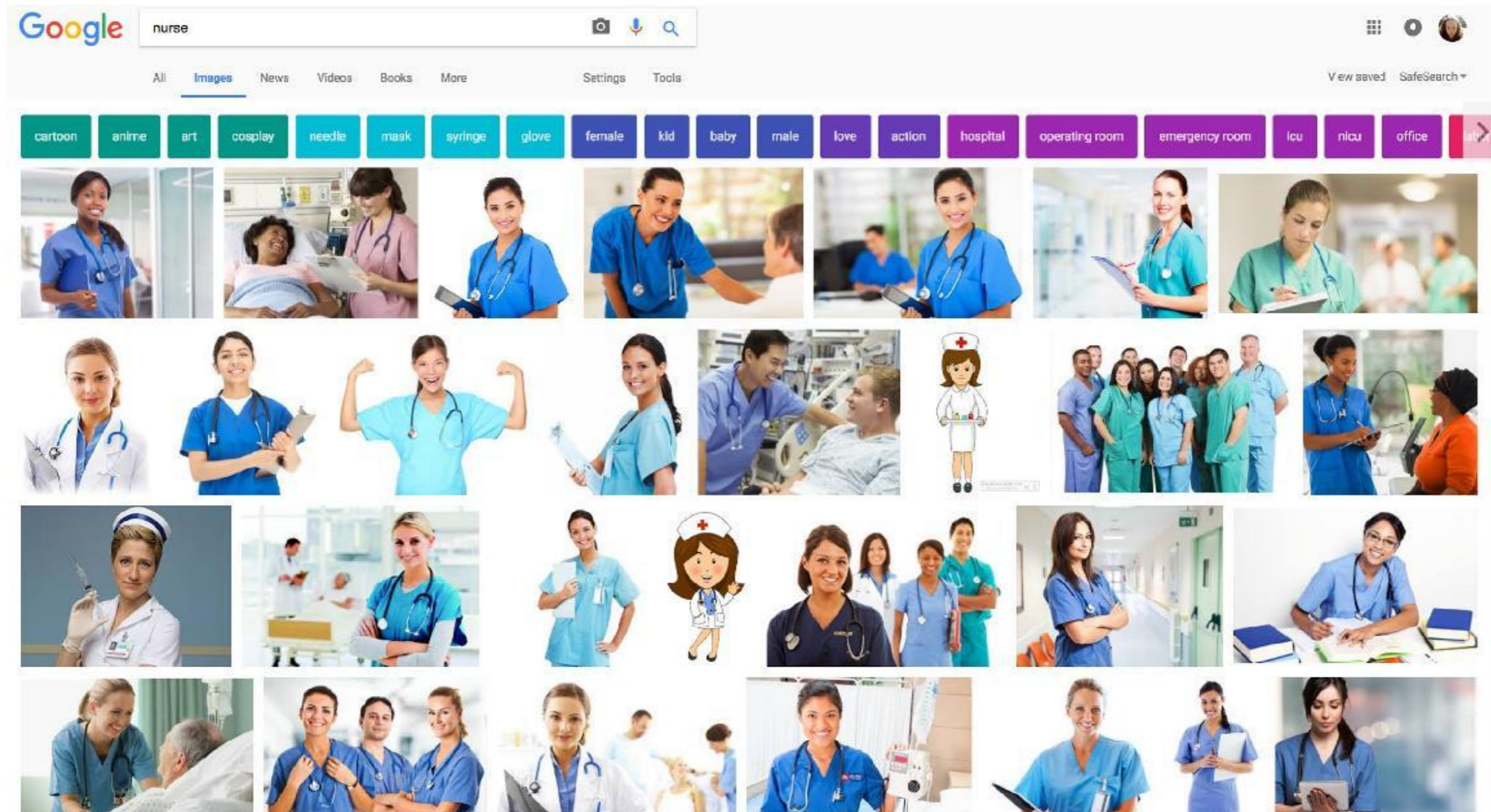
Gender/Race/Age Stereotypes

- June 2017: image search query “Doctor”



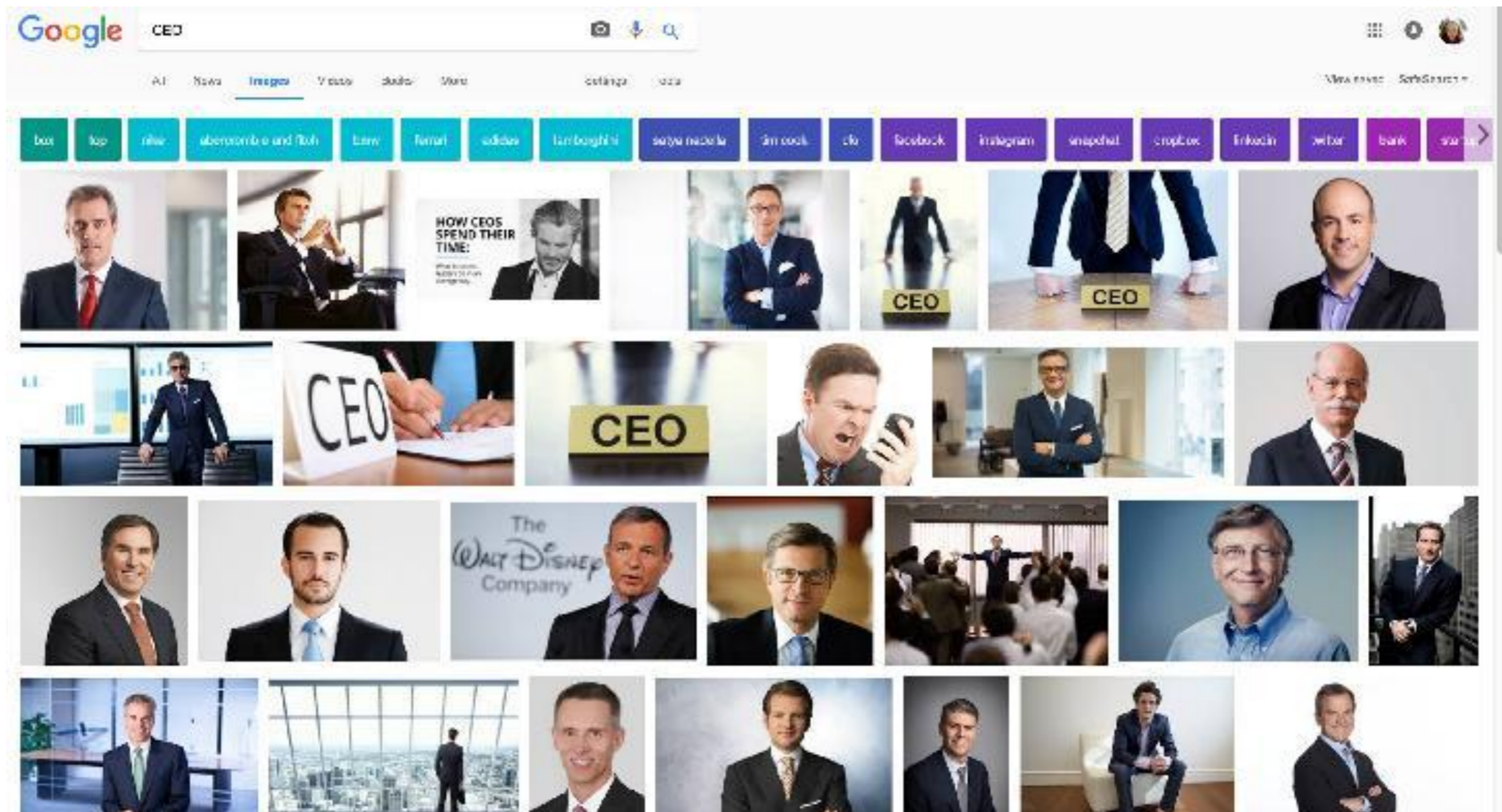
Gender/Race/Age Stereotypes

- June 2017: image search query “Nurse”



Gender/Race/Age Stereotypes

- June 2017: image search query “CEO”





Consequence: models are biased


Gender Biases on the Web

- The dominant class is often portrayed and perceived as relatively more professional ([Kay, Matuszek, and Munson 2015](#))
- Males are over-represented in the reporting of web-based news articles ([Jia, Lansdall-Welfare, and Cristianini 2015](#))
- Males are over-represented in twitter conversations ([Garcia, Weber, and Garimella 2014](#))
- Biographical articles about women on Wikipedia disproportionately discuss romantic relationships or family-related issues ([Wagner et al. 2015](#))
- IMDB reviews written by women are perceived as less useful ([Otterbacher 2013](#))


Biased NLP Technologies

- Bias in word embeddings ([Bolukbasi et al. 2017](#); [Caliskan et al. 2017](#); [Garg et al. 2018](#))
- Bias in Language ID ([Blodgett & O'Connor. 2017](#); [Jurgens et al. 2017](#))
- Bias in Visual Semantic Role Labeling ([Zhao et al. 2017](#))
- Bias in Natural Language Inference ([Rudinger et al. 2017](#))
- Bias in Coreference Resolution ([At NAACL: Rudinger et al. 2018](#); [Zhao et al. 2018](#))
- Bias in Automated Essay Scoring ([At NAACL: Amorim et al. 2018](#))

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




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



The physician hired the secretary because she was highly recommended.



The physician hired the secretary because he was highly recommended.



Sources of Human Biases in Machine Learning

- Bias in data and sampling
- Optimizing towards a biased objective
- Inductive bias
- Bias amplification in learned models

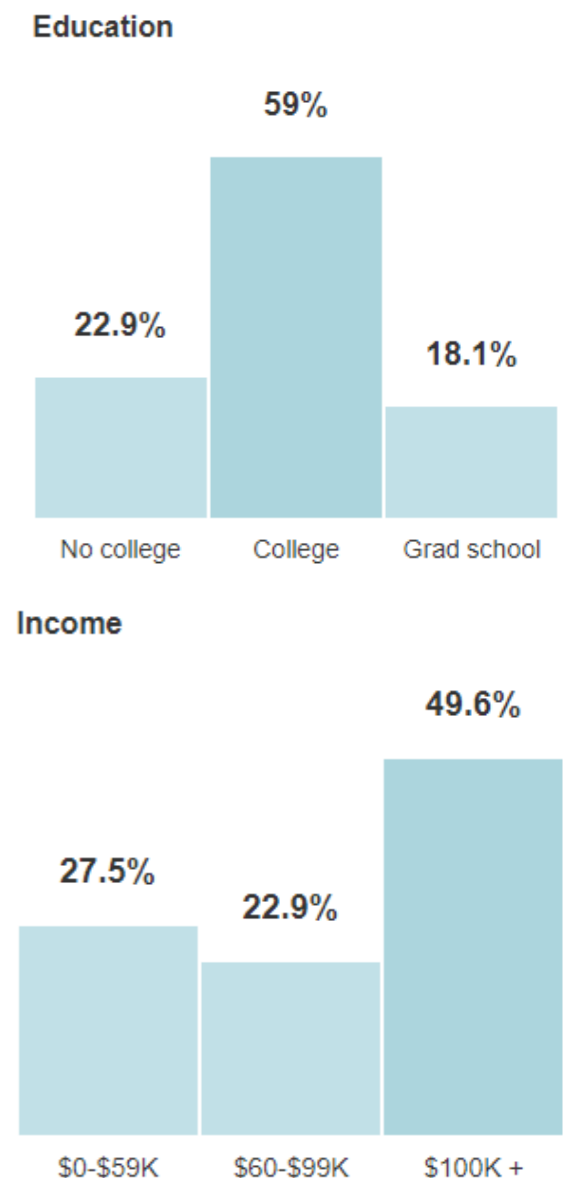
Sources of Human Biases in Machine Learning

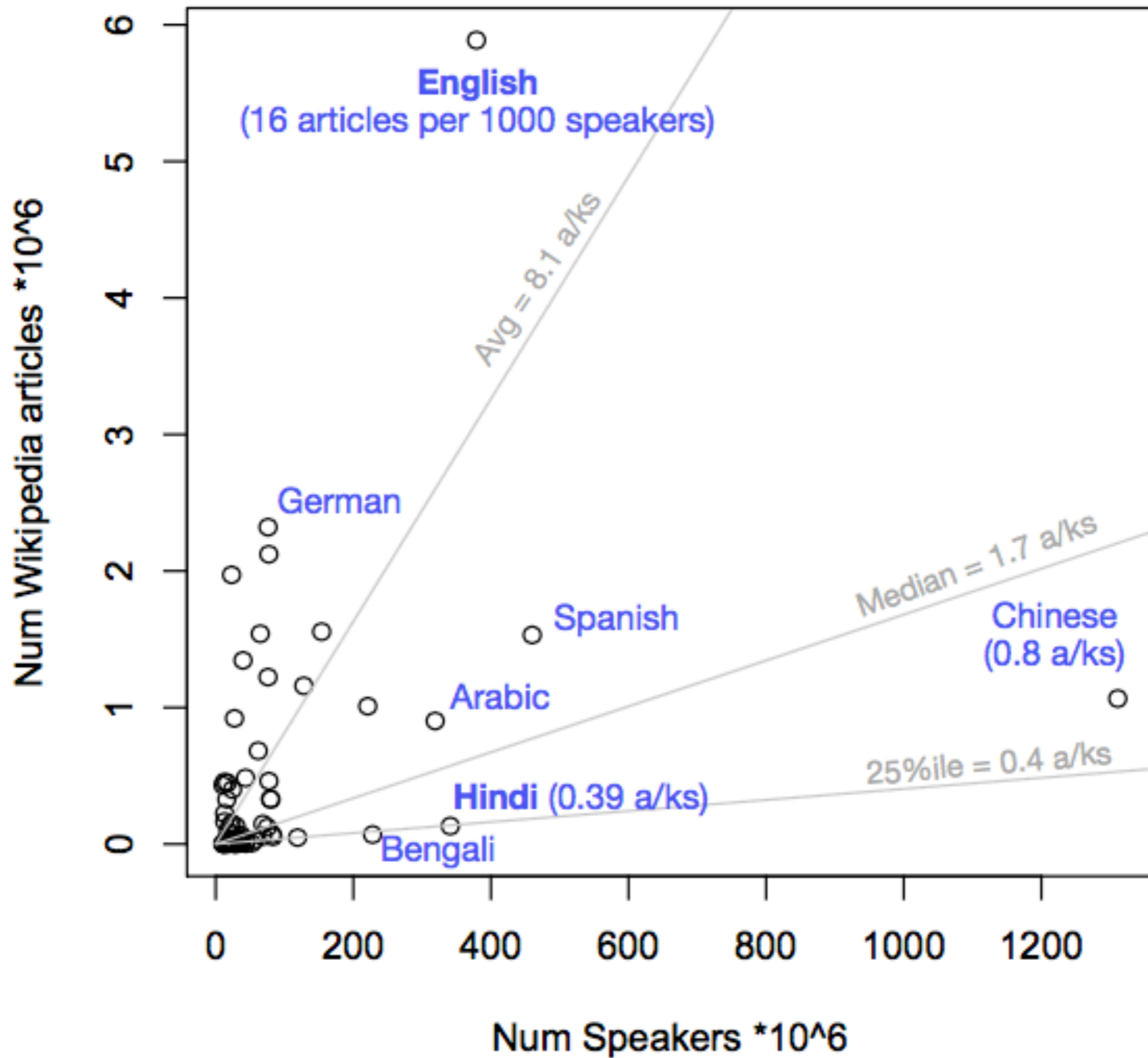
- **Bias in data and sampling**
- Optimizing towards a biased objective
- Inductive bias
- Bias amplification in learned models

Types of Sampling Bias in Naturalistic Data

- **Self-Selection Bias**
 - Who decides to post reviews on Yelp and why?
Who posts on Twitter and why?
- **Reporting Bias**
 - People do not necessarily talk about things in the world in proportion to their empirical distributions (Gordon and Van Durme 2013)
- **Proprietary System Bias**
 - What results does Twitter return for a particular query of interest and why? Is it possible to know?
- **Community / Dialect / Socioeconomic Biases**
 - What linguistic communities are over- or under-represented? leads to community-specific model performance (Jorgensen et al. 2015)

US Demographics of Yelp Users





Example: Bias in Language Identification

- Most applications employ off-the-shelf LID systems which are highly accurate



McNamee, P., “Language identification: *a solved problem* suitable for undergraduate instruction” *Journal of Computing Sciences in Colleges* 20(3) 2005.

“This paper describes [...] how even the most simple of these methods **using data obtained from the World Wide Web achieve accuracy approaching 100%** on a test suite comprised of ten European languages”



The Royal Family ✓
@RoyalFamily

Follow

Taking place this week on the river Thames is 'Swan Upping' – the annual census of the swan population on the Thames.



da'Rah-zingSun
@TIME7SS

Follow

@kinguilfoyle prblm I hve wit ur reportng is its 2 literal, evry1 knos pple tlk diffrent evrywhere, u kno wut she means jus like we do!



Mooktar
@bossmukky

Follow

"@Ecstatic_Mi: @bossmukky Ebi like say I wan dey sick sef wlh 'Flu' my whole body dey weak"uw gee...



Ebenezer
@Physique_cian

Follow

@Tblazeen R u a wizard or wat gan sef : in d mornin- u tweet, afternoon - u tweet, nyt gan u dey tweet.beta get ur IT placement wiv twitter

- Language identification degrades significantly on African American Vernacular English
(Blodgett et al. 2016) **Su-Lin Blodgett just got her PhD from UMass!**

LID Usage Example: Health Monitoring



Language Detection



Keyword Filter
"flu", "sick"



Analytics

Which symptoms?
Are they hungover?

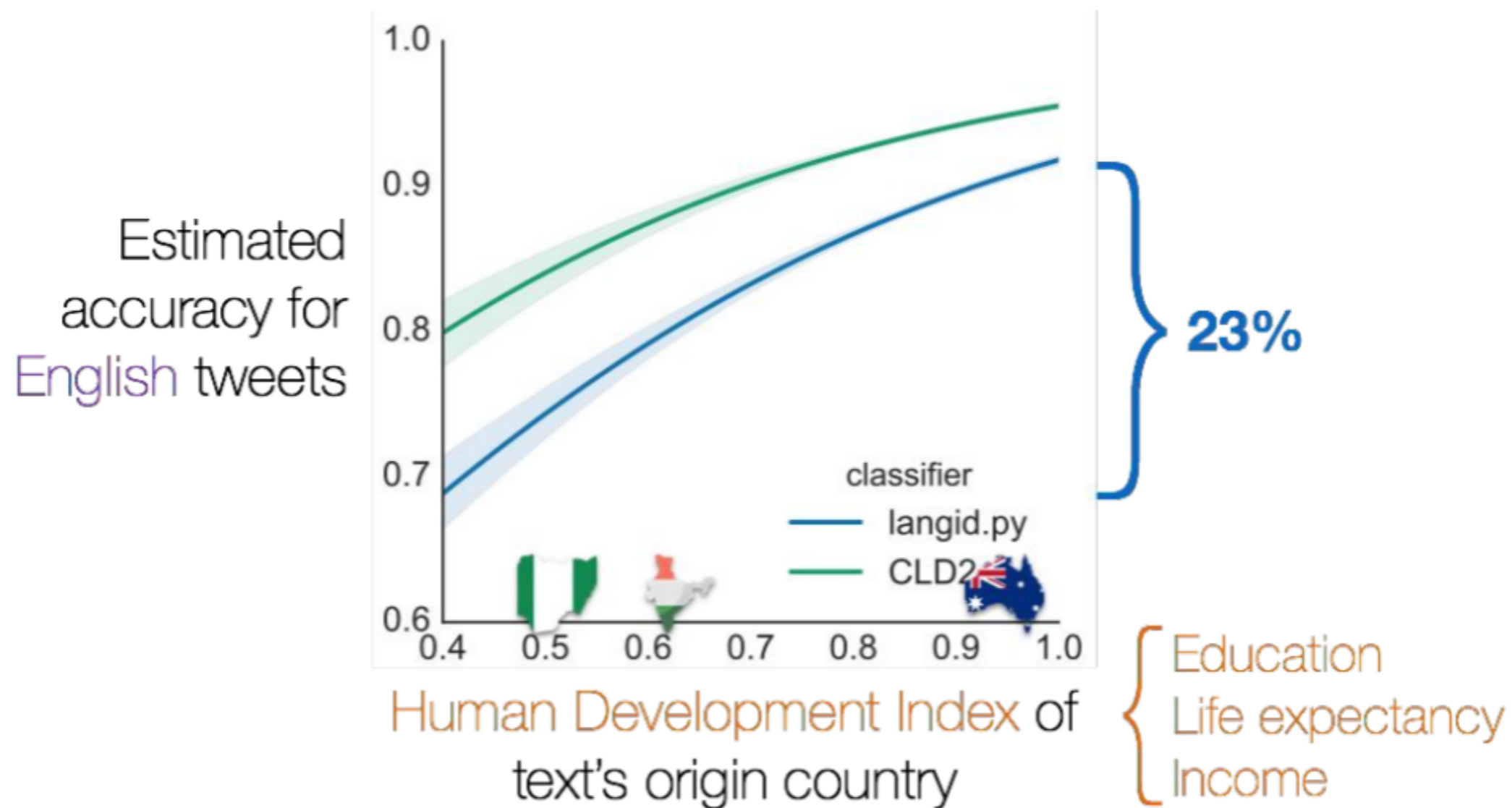
LID Usage Example: Health Monitoring



**Language
Detection**

Socioeconomic Bias in Language Identification

- Off-the-shelf LID systems under-represent populations in less-developed countries



Better Social Representation through Network-based Sampling

- Re-sampling from strategically-diverse corpora

Topical



Geographic



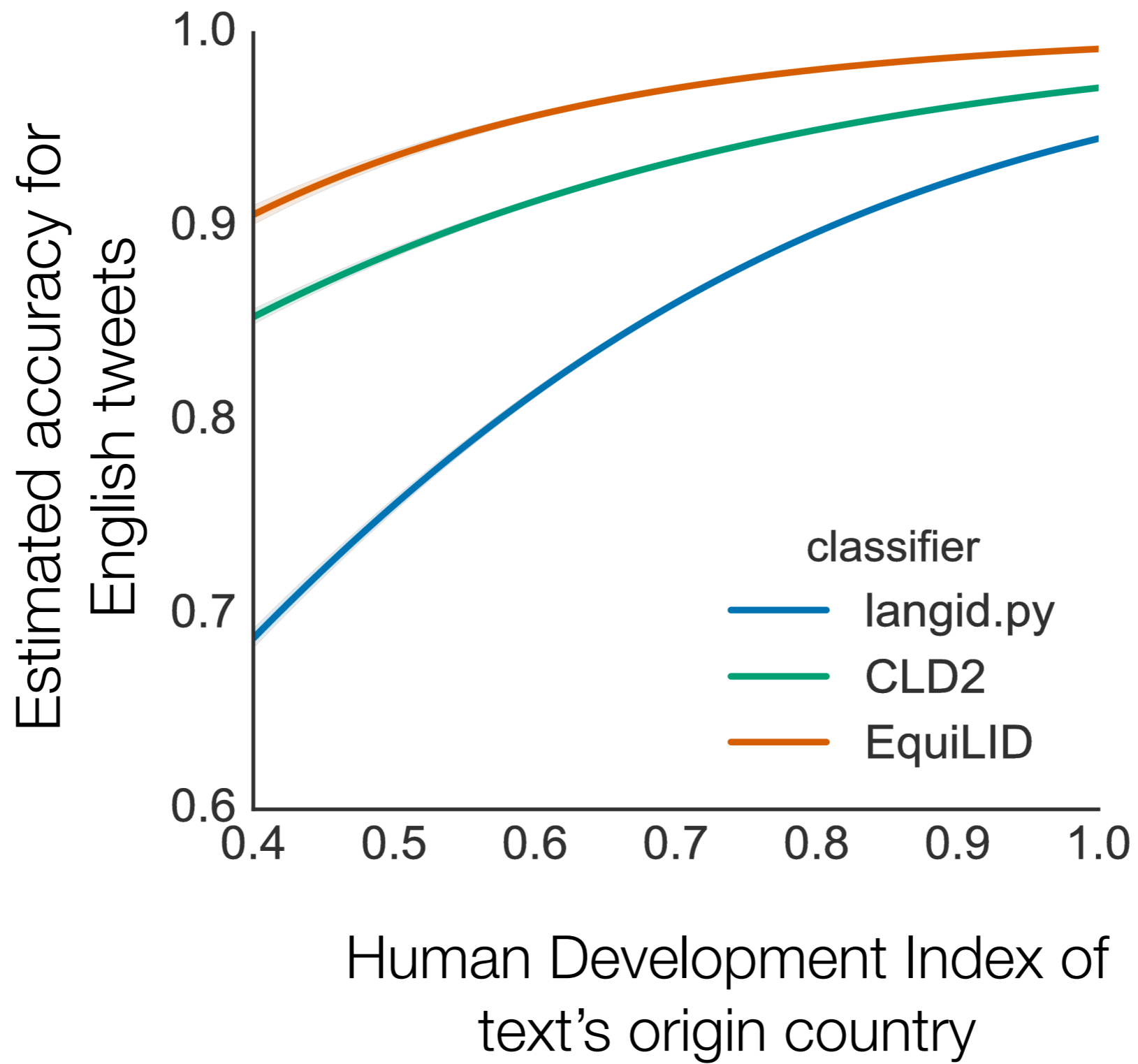
Social



Multilingual



Jurgens et al. ACL'11



Sources of Human Biases in Machine Learning

- Bias in data and sampling
- **Optimizing towards a biased objective**
- Inductive bias
- Bias amplification in learned models

Optimizing Towards a Biased Objective

- Northpointe vs ProPublica

COMPAS



Optimizing Towards a Biased Objective

“what is the probability that this person will commit a serious crime in the future, as a function of the sentence you give them now?”

Optimizing Towards a Biased Objective

“what is the probability that this person will commit a serious crime in the future, as a function of the sentence you give them now?”

- COMPAS system
 - balanced training data about people of all races
 - race was *not* one of the input features
- Objective function
 - labels for “who will commit a crime” are unobtainable
 - a proxy for the real, unobtainable data: “who is more likely to be *convicted*”

what are some issues with this proxy objective?

Predicting prison sentences given case descriptions

Case description: On July 7, 2017, when the defendant Cui XX was drinking in a bar, he came into conflict with Zhang XX..... After arriving at the police station, he refused to cooperate with the policeman and bited on the arm of the policeman.....

Result of judgment: Cui XX was sentenced to 12 months imprisonment for creating disturbances and 12 months imprisonment for obstructing public affairs.....

- Charge#1 creating disturbances term 12 months
- Charge#2 obstructing public affairs term 12 months

Is this sufficient consideration of ethical issues of this work? Should the work have been done at all?

The mistake of legal judgment is serious, it is about people losing years of their lives in prison, or dangerous criminals being released to reoffend. We should pay attention to how to avoid judges' over-dependence on the system. It is necessary to consider its application scenarios. In practice, we recommend deploying our system in the "Review Phase", where other judges check the judgment result by a presiding judge. Our system can serve as one anonymous checker.

Sources of Human Biases in Machine Learning

- Bias in data and sampling
- Optimizing towards a biased objective
- **Inductive bias**
- Bias amplification in learned models

what is inductive bias?

- the assumptions used by our model. examples:
 - recurrent neural networks for NLP assume that the sequential ordering of words is meaningful
 - features in discriminative models are assumed to be useful to map inputs to outputs

Bias in Word Embeddings

1. Caliskan, A., Bryson, J. J. and Narayanan, A. (2017) **Semantics derived automatically from language corpora contain human-like biases.**
Science

$$\overrightarrow{\text{man}} - \overrightarrow{\text{woman}} \approx \overrightarrow{\text{computer programmer}} - \overrightarrow{\text{homemaker}}.$$

Biases in Embeddings: Another Take

$$\min \cos(\mathit{he} - \mathit{she}, x - y) \text{ s.t. } \|x - y\|_2 < \delta$$

Extreme <i>she</i>	Extreme <i>he</i>		Gender stereotype <i>she-he</i> analogies	
1. homemaker	1. maestro	sewing-carpentry	registered nurse-physician	housewife-shopkeeper
2. nurse	2. skipper	nurse-surgeon	interior designer-architect	softball-baseball
3. receptionist	3. protege	blond-burly	feminism-conservatism	cosmetics-pharmaceuticals
4. librarian	4. philosopher	giggle-chuckle	vocalist-guitarist	petite-lanky
5. socialite	5. captain	sassy-snappy	diva-superstar	charming-affable
6. hairdresser	6. architect	volleyball-football	cupcakes-pizzas	lovely-brilliant
7. nanny	7. financier			
8. bookkeeper	8. warrior	Gender appropriate <i>she-he</i> analogies		
9. stylist	9. broadcaster	queen-king	sister-brother	mother-father
10. housekeeper	10. magician	waitress-waiter	ovarian cancer-prostate cancer	convent-monastery

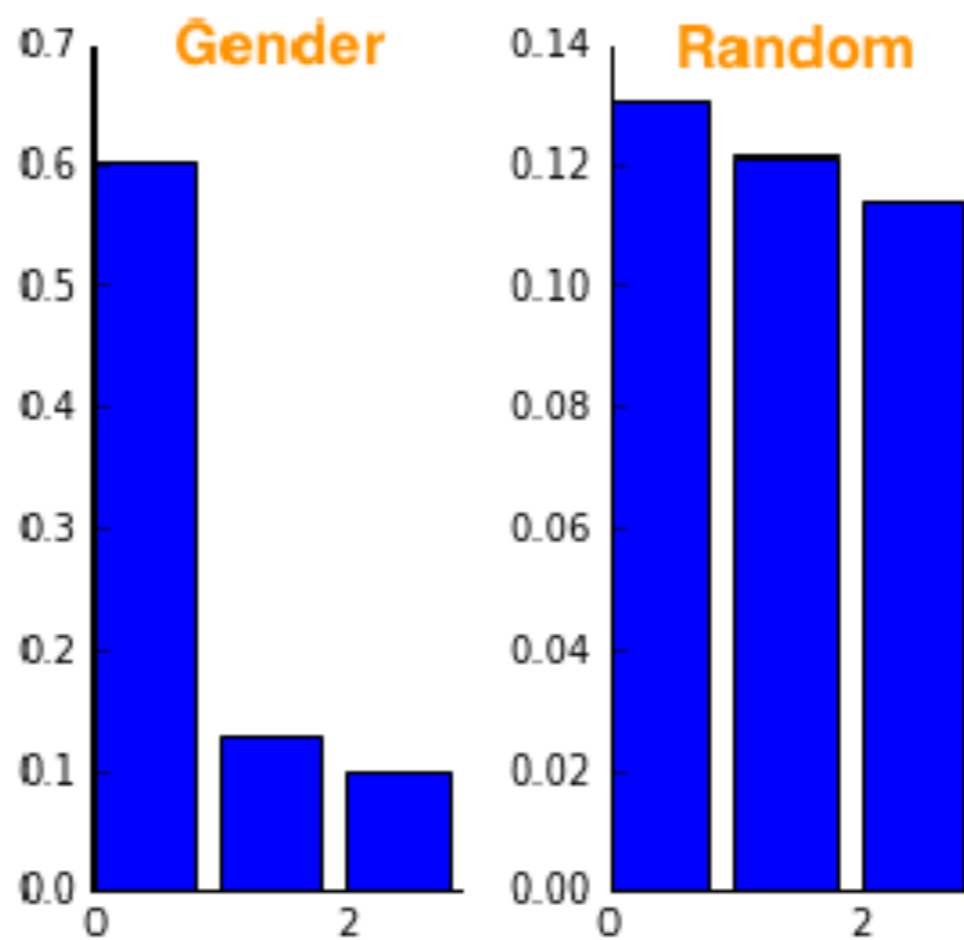
Figure 1: **Left** The most extreme occupations as projected on to the *she*–*he* gender direction on w2vNEWS. Occupations such as *businesswoman*, where gender is suggested by the orthography, were excluded. **Right** Automatically generated analogies for the pair *she-he* using the procedure described in text. Each automatically generated analogy is evaluated by 10 crowd-workers to whether or not it reflects gender stereotype.

Towards Debiasing

1. Identify gender subspace: B

Gender Subspace

$\vec{\text{she}} - \vec{\text{he}}$
 $\vec{\text{her}} - \vec{\text{his}}$
 $\vec{\text{woman}} - \vec{\text{man}}$
 $\vec{\text{Mary}} - \vec{\text{John}}$
 $\vec{\text{herself}} - \vec{\text{himself}}$
 $\vec{\text{daughter}} - \vec{\text{son}}$
 $\vec{\text{mother}} - \vec{\text{father}}$
 $\vec{\text{gal}} - \vec{\text{guy}}$
 $\vec{\text{girl}} - \vec{\text{boy}}$
 $\vec{\text{female}} - \vec{\text{male}}$

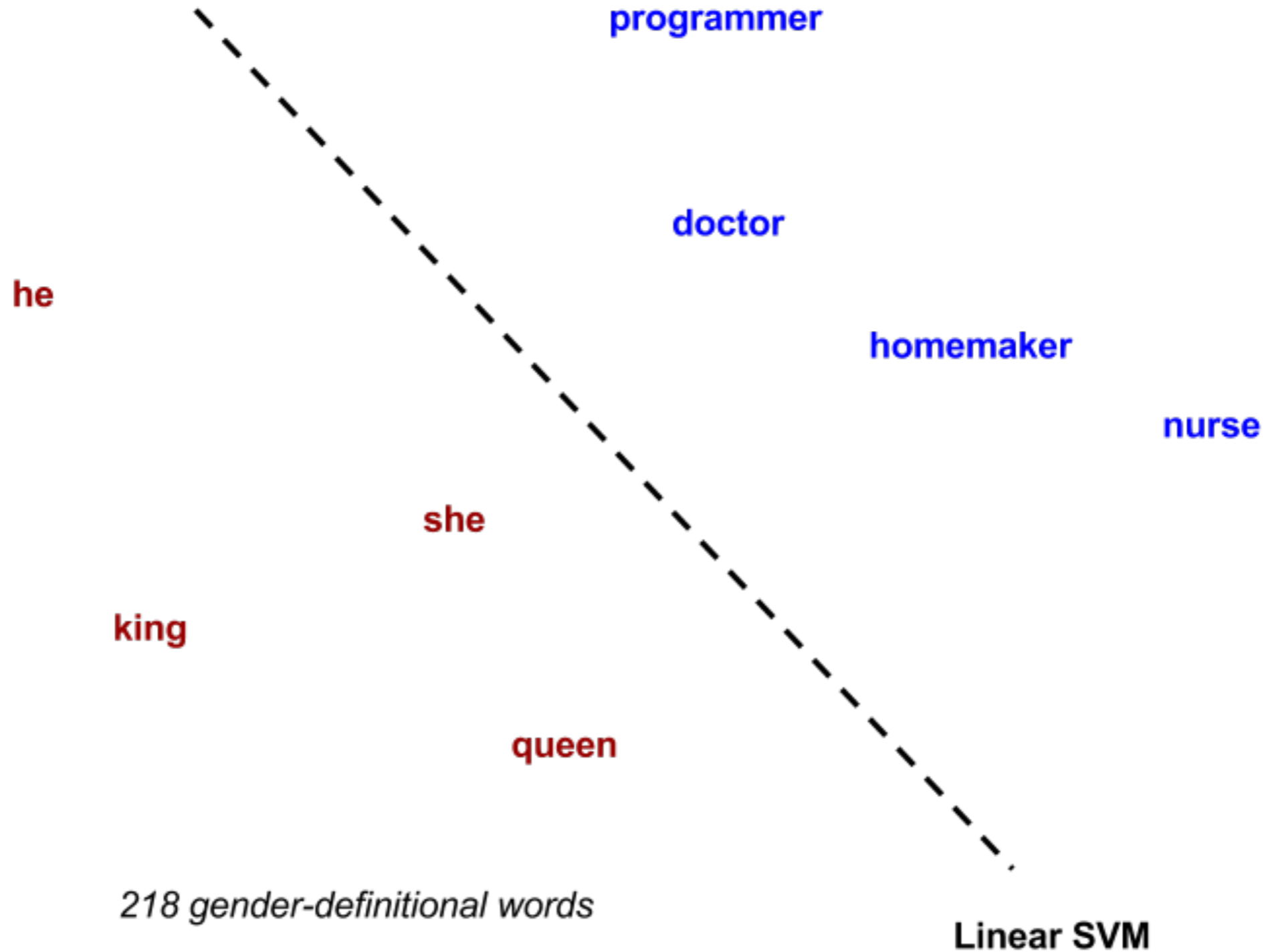


The top PC captures the gender subspace

Towards Debiasing

1. Identify gender subspace: B
2. **Identify gender-definitional (S) and gender-neutral words (N)**

Gender-definitional vs. Gender-neutral Words



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) such that
 - 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

$$\min_T \underbrace{\| (TW)^T (TW) - W^T W \|_F^2}_{\text{Don't modify embeddings too much}} + \lambda \underbrace{\| (TN)^T (TB) \|_F^2}_{\text{Minimize gender component}}$$

T - the desired debiasing transformation B - biased space

W - embedding matrix

N - embedding matrix of gender neutral words

Sources of Human Biases in Machine Learning

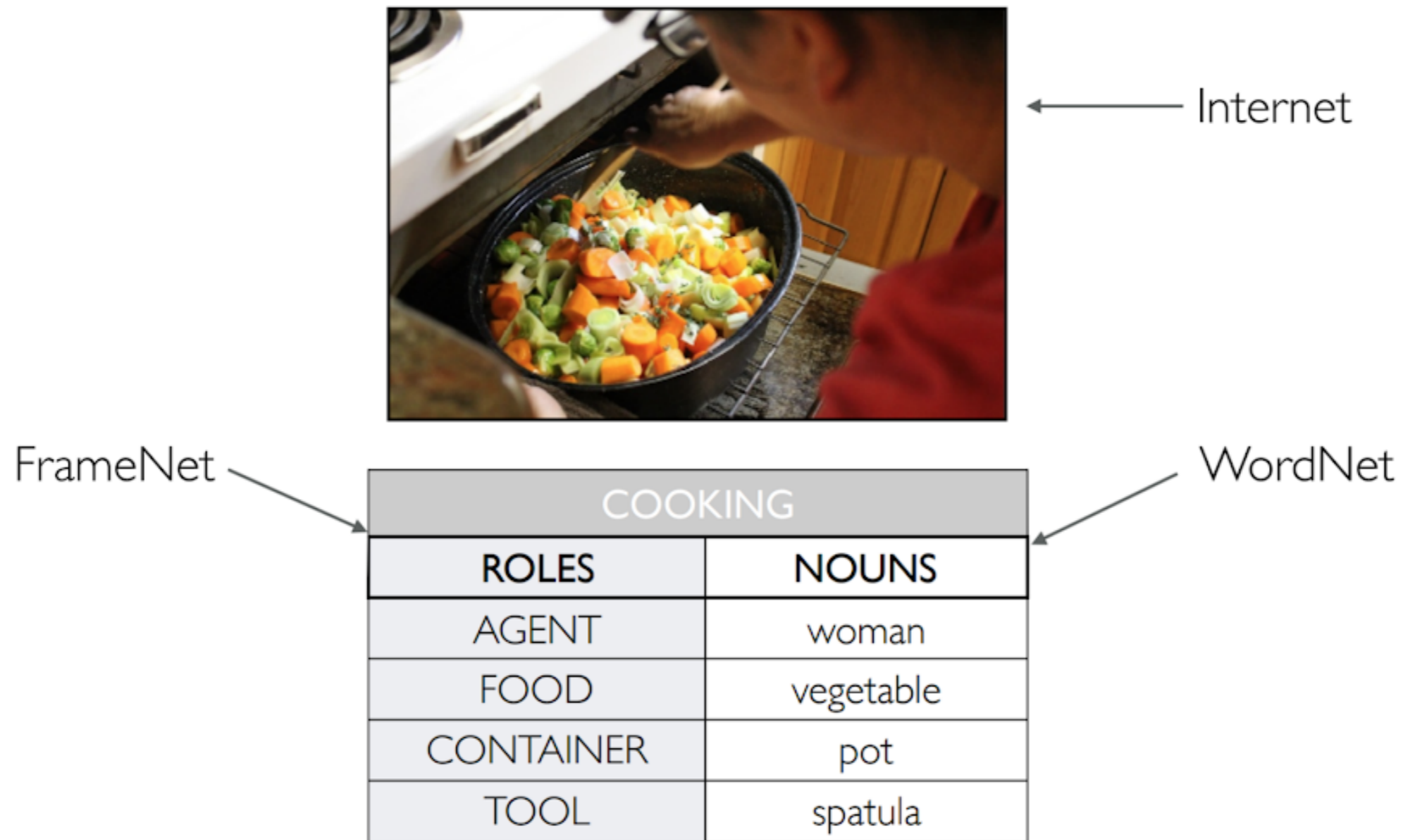
- Bias in data and sampling
- Optimizing towards a biased objective
- Inductive bias
- **Bias amplification in learned models**

Bias Amplification

Zhao, J., Wang, T., Yatskar, M., Ordonez, V and Chang, M.-W. (2017) **Men Also Like Shopping: Reducing Gender Bias Amplification using Corpus-level Constraint.**

EMNLP

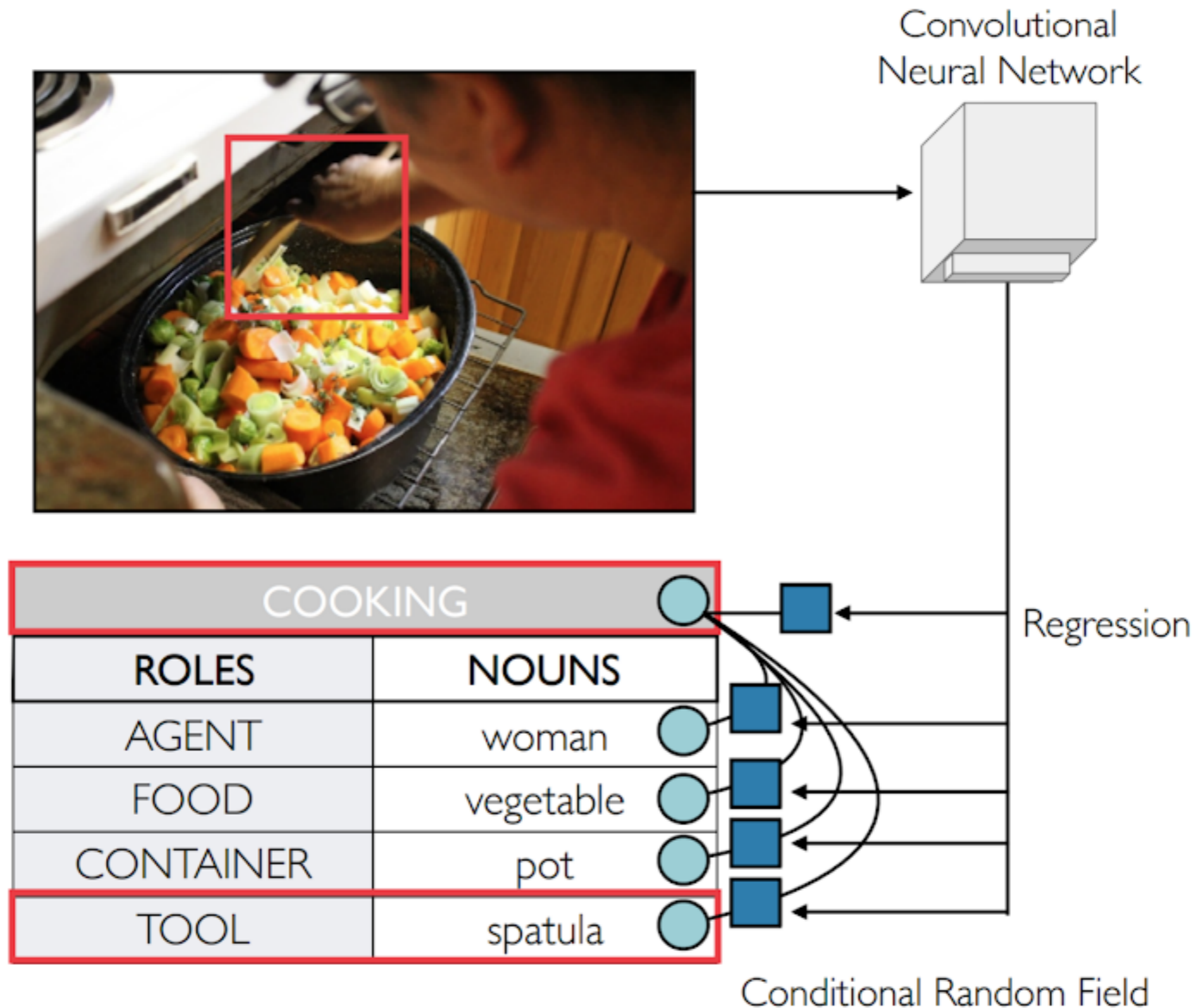
imSitu Visual Semantic Role Labeling (vSRL)



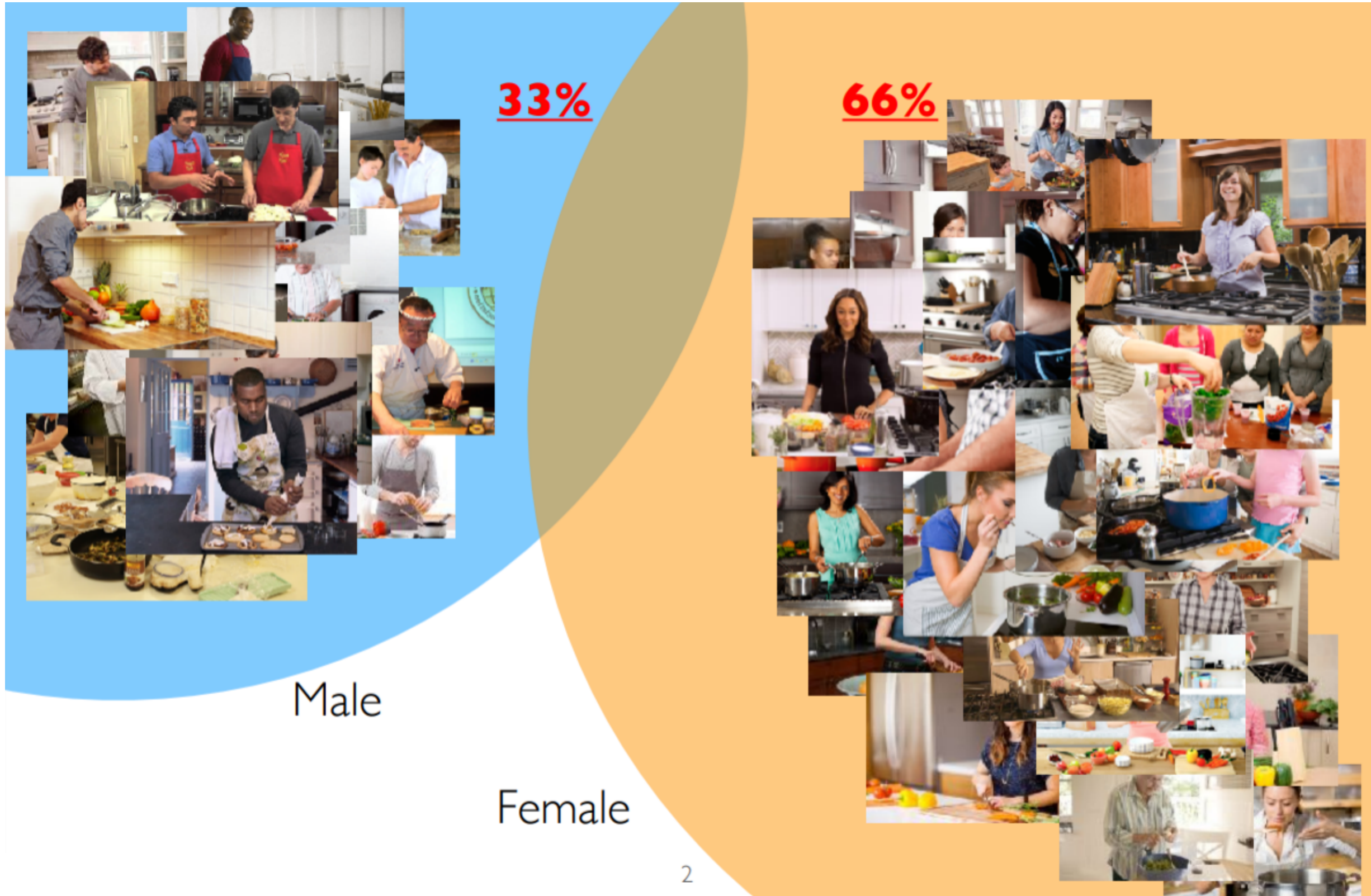
12

Yatskar et al. CVPR '16, Yang et al. NAACL '16, Gupta and Malik arXiv '16

imSitu Visual Semantic Role Labeling (vSRL)



Dataset Gender Bias



Model Bias After Training

16%

84%



Male

Female

Algorithmic Bias



woman cooking



man fixing faucet

Quantifying Dataset Bias

Training Gender Ratio (◆ verb)

Training Set

- ◆ cooking
- woman
- man



◆	COOKING	
	ROLES	NOUNS
●	AGENT	woman
	FOOD	stir-fry



◆	COOKING	
	ROLES	NOUNS
●	AGENT	man
	FOOD	noodle

$$\frac{\#(\text{◆ cooking}, \text{● man})}{\#(\text{◆ cooking}, \text{● man}) + \#(\text{◆ cooking}, \text{● woman})} = 1/3$$

Quantifying Dataset Bias: Dev Set

Predicted Gender Ratio (◆ verb)

Development Set

- ◆ cooking
- woman
- man

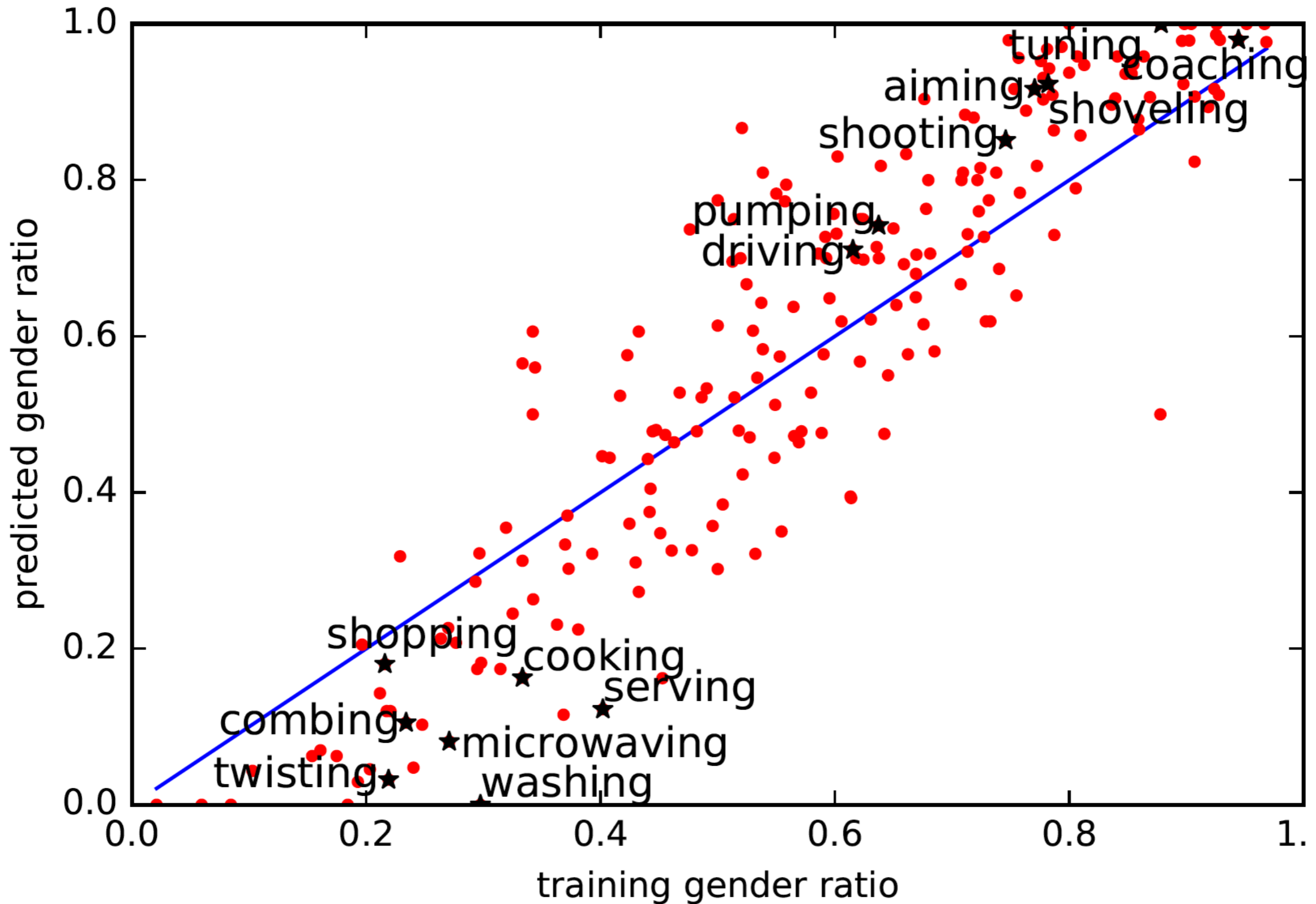


◆	COOKING	
	ROLES	NOUNS
●	AGENT	woman
	FOOD	stir-fry

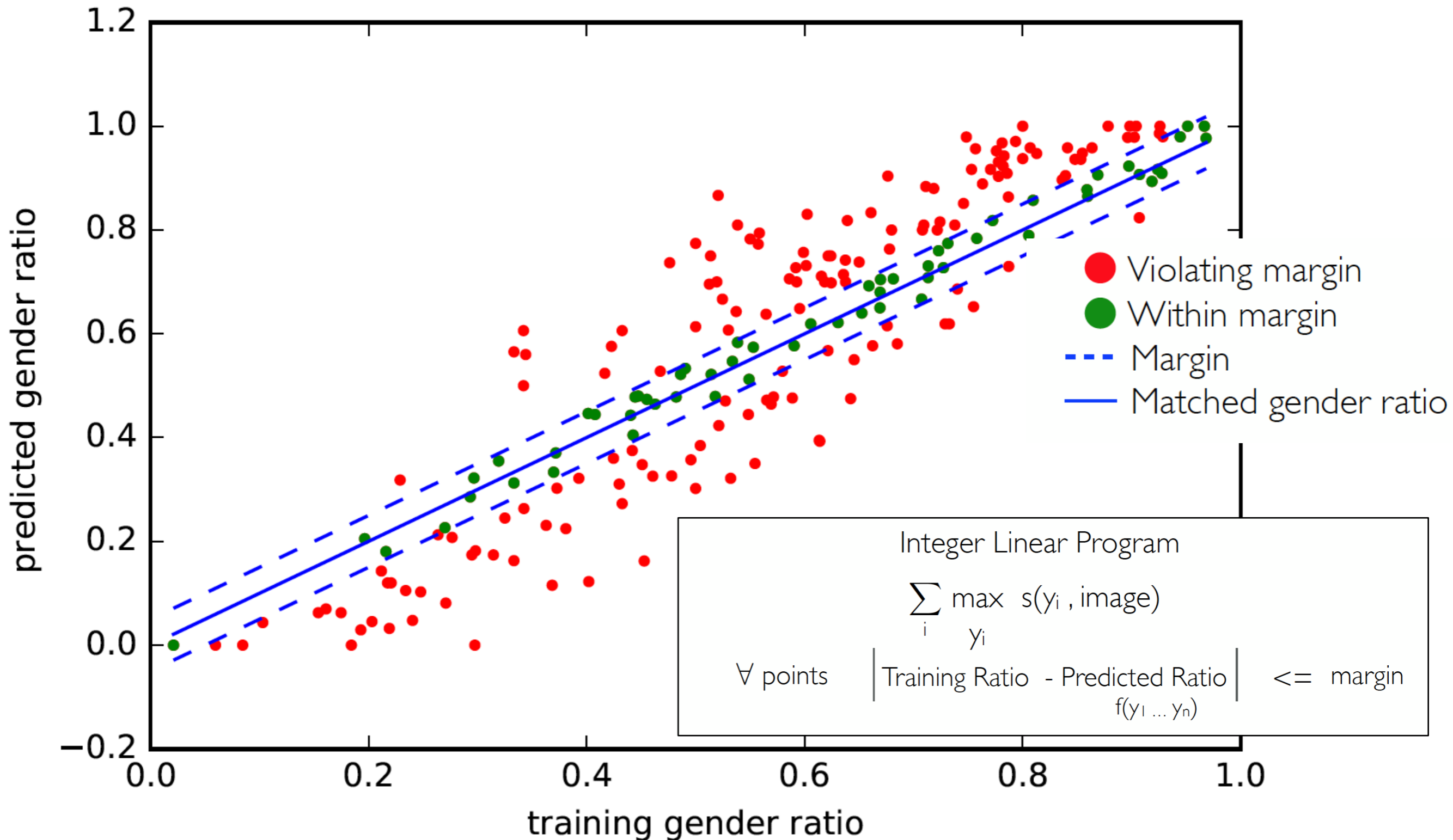
◆	COOKING	
	ROLES	NOUNS
●	AGENT	man
	FOOD	noodle

$$\frac{\#(\text{◆ cooking}, \text{● man})}{\#(\text{◆ cooking}, \text{● man}) + \#(\text{◆ cooking}, \text{● woman})} = 1/6$$

Model Bias Amplification



Reducing Bias Amplification (RBA)



Discussion

- Applications that are built from online data, generated by people, learn also real-world stereotypes
- Should our ML models represent the “real world”?
- Or should we artificially skew data distribution?
- If we modify our data, what are guiding principles on what our models should or shouldn't learn?