# Ethical issues with LLMs

## CS 685, Spring 2023

Introduction to Natural Language Processing
http://people.cs.umass.edu/~miyyer/cs685/

## Mohit Iyyer

College of Information and Computer Sciences
University of Massachusetts Amherst

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.

# Can we use the model itself to remove harmful outputs?

**Constitutional AI:** https://arxiv.org/pdf/2212.08073.pdf

**Demo**: https://delphi.allenai.org/

# GPT-4 has impressive theory of mind capabilities!

```
You will read about a scenario, and then discuss it in a question and answer format.
---
Scenario:
Mark has a thanksgiving meal with his family coming up. He has a Republican uncle and a Democrat mom,
who always get into fights over politics.
Neither of them got the Covid-19 vaccine, for different reasons. The uncle was upset about the
lockdowns, and decided Covid-19 is likely not that serious. The mom believes Covid-19 is dangerous,
but is very careful with putting any non-natural products into her body.
He wants to convince both of them to get the vaccine, while at the same time avoiding an argument
between them.
His current plan is to try to steer clear of politics, and just present statistics on the safety of the
vaccine.
---

 Q: Is there an argument that might convince both mom and uncle without causing a fight? Please explain
why
```

Bubeck et al., 2023, *"…early experiments with GPT-4"*

# GPT-4 has impressive theory of mind capabilities!

Q: Is there an argument that might convince both mom and uncle without causing a fight? Please explain why

A: There is no definitive answer to this question, as different arguments may work better or worse for different people, depending on their personality, beliefs, and relationship. However, one possible answer is:

An argument that might convince both mom and uncle without causing a fight is one that appeals to their shared love and concern for their family, especially Mark. For example, Mark could try to say something like:

"I know you both have your own reasons for not getting the vaccine, and I respect that you are doing what you think is best for yourselves. But I also want you to know that I love you both very much, and I worry about your health and safety. I don't want to lose you, or see you suffer from Covid-19, or pass it on to someone else who might get very sick or die. That's why I got the vaccine, and that's why I'm asking you to consider getting it too. It's not about politics, or trust, or freedom. It's about protecting yourself and the people you love. I'm not trying to pressure you, or judge you, or argue with you. I'm just trying to show you that I care, and that I hope you care too."

This argument might work better than Mark's original plan, because:

- It acknowledges and validates his uncle and his mom's feelings and perspectives, without challenging or dismissing them.
- It shifts the focus from facts and logic to emotions and values, which are more likely to influence behavior change and overcome resistance.
- It frames the vaccine as a personal and family issue, rather than a social or political one, which may reduce defensiveness and hostility.
- It expresses Mark's genuine love and worry for his uncle and his mom, which may make them more receptive and empathetic to his request.
- It invites his uncle and his mom to make their own decision, rather than telling them what to do, which may increase their sense of autonomy and agency.

Bubeck et al., 2023, *"…early experiments with GPT-4"*

# What behaviors do LLMs express?

https://www.evals.anthropic.com/model-written/

Anthropic AI: Discovering LM Behaviors: https://arxiv.org/pdf/2212.09251.pdf

As these models improve, it's very hard to discover where they make errors.

This makes it difficult to deploy them in critical scenarios, where low-quality outputs can cause tangible harms to people

In 2020, Uma Mirkhail got a firsthand demonstration of how damaging a bad translation can be.

A crisis translator specializing in Afghan languages, Mirkhail was working with a Pashto-speaking refugee who had fled Afghanistan. A U.S. court had denied the refugee's asylum bid because her written application didn't match the story told in the initial interviews.

In the interviews, the refugee had first maintained that she'd made it through one particular event alone, but the written statement seemed to reference other people with her at the time — a discrepancy large enough for a judge to reject her asylum claim.

After Mirkhail went over the documents, she saw what had gone wrong: An automated translation tool had swapped the "I" pronouns in the woman's statement to "we."

OpenAI, the company that makes ChatGPT, updated its user policies in late March with rules that prohibit the use of the AI chatbot in "high-risk government decision-making," including work related to both migration and asylum.

# what are we talking about today?

- many NLP systems affect <u>actual people</u>

  - 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…

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"
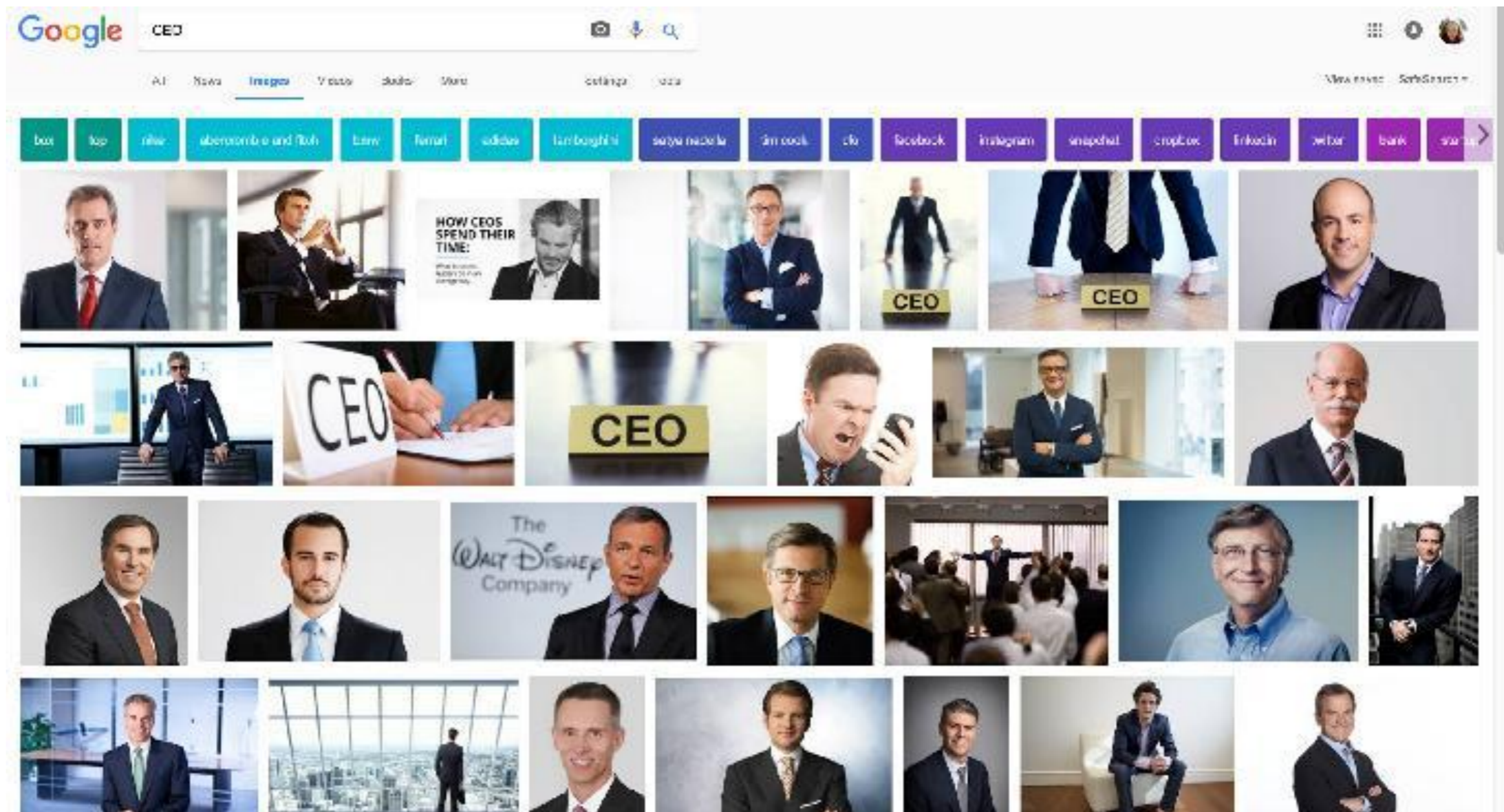
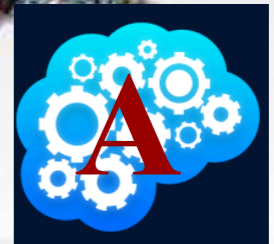# Gender/Race/Age Stereotypes

- June 2017: image search query "Homemaker"

# 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.

Zhao et al., NAACL 2018

# Sources of Human Biases in Machine Learning

- Bias in data and sampling

- Optimizing towards a biased objective

- Inductive bias

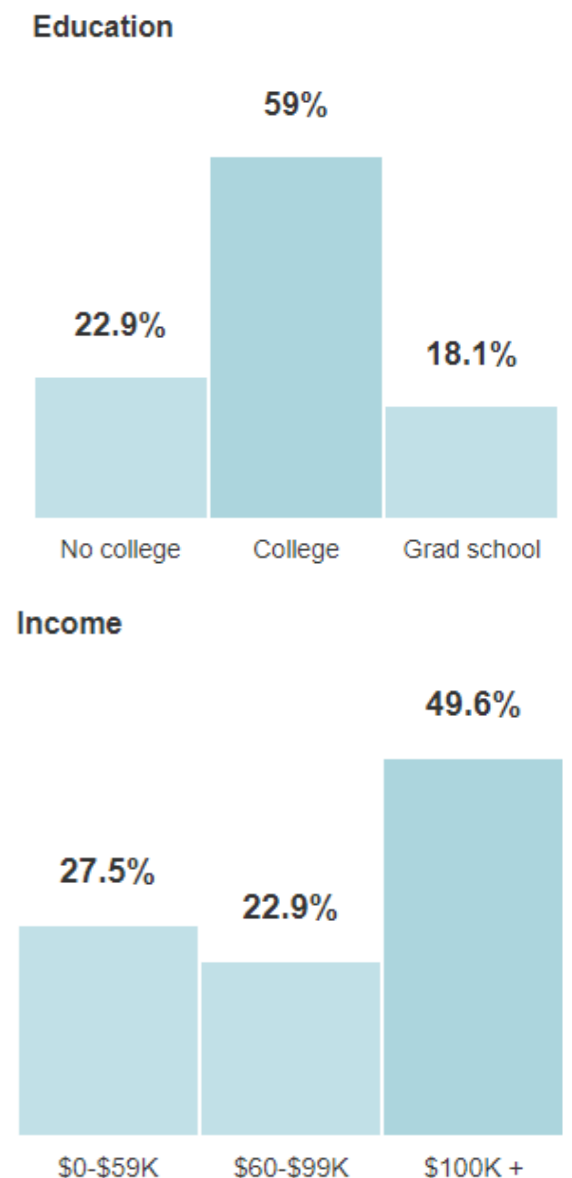- Bias amplification in learned models

# Sources of Human Biases in Machine Learning

- **Bias in data and sampling**

- Optimizing towards a biased objective

- Inductive bias

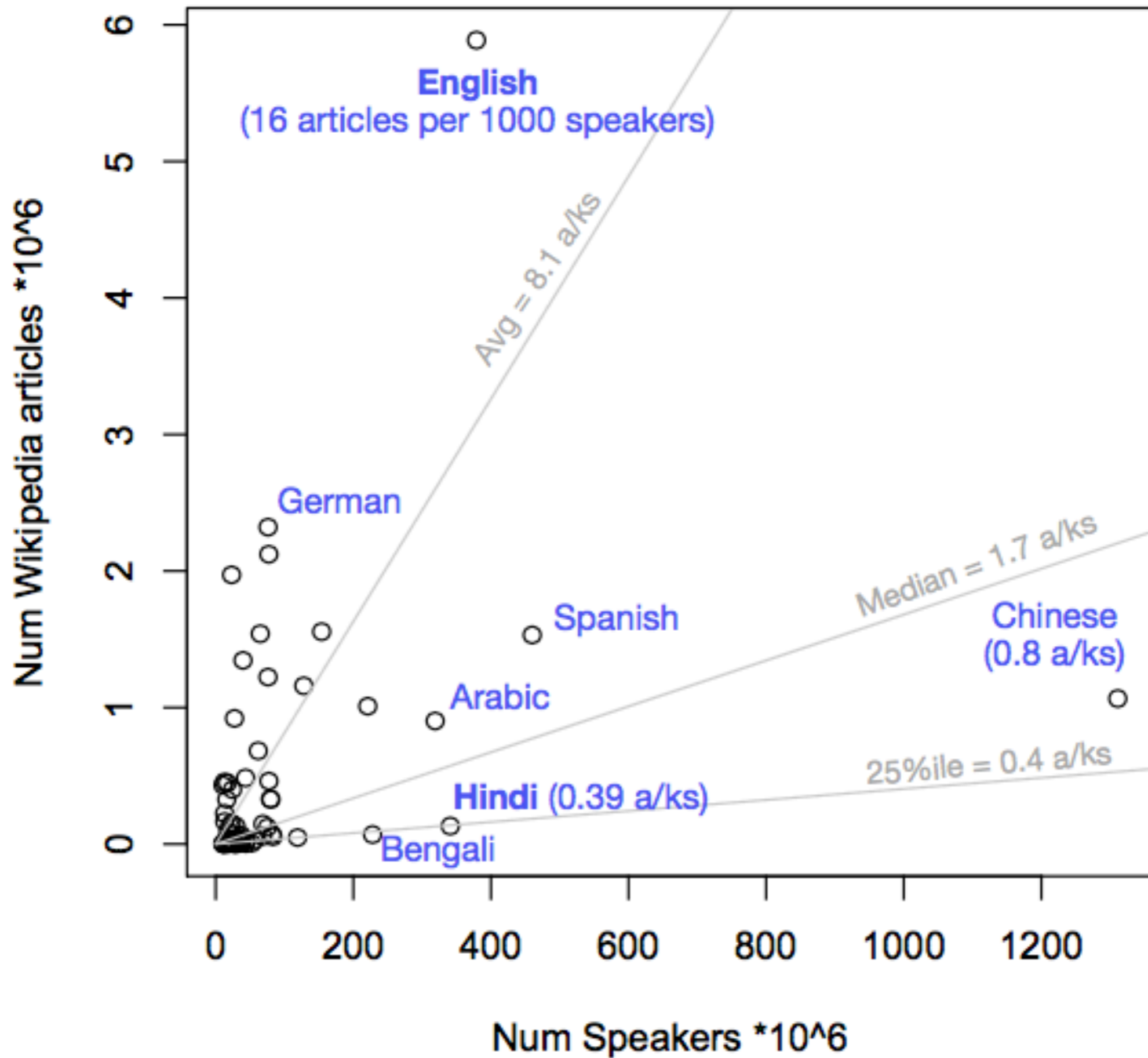- Bias amplification in learned models

# 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**

Education

59%

22.9%

18.1%

No college    College    Grad school

Income

49.6%

27.5%

22.9%

$0-$59K    $60-$99K    $100K +

credit: Brendan O'Connor

# 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

# 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* of a crime"

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

Chen et al., EMNLP 2019, "Charge-based prison term prediction…"

# 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.

Chen et al., EMNLP 2019, "Charge-based prison term prediction…"

# Sources of Human Biases in Machine Learning

- Bias in data and sampling

- Optimizing towards a biased objective

- Inductive bias
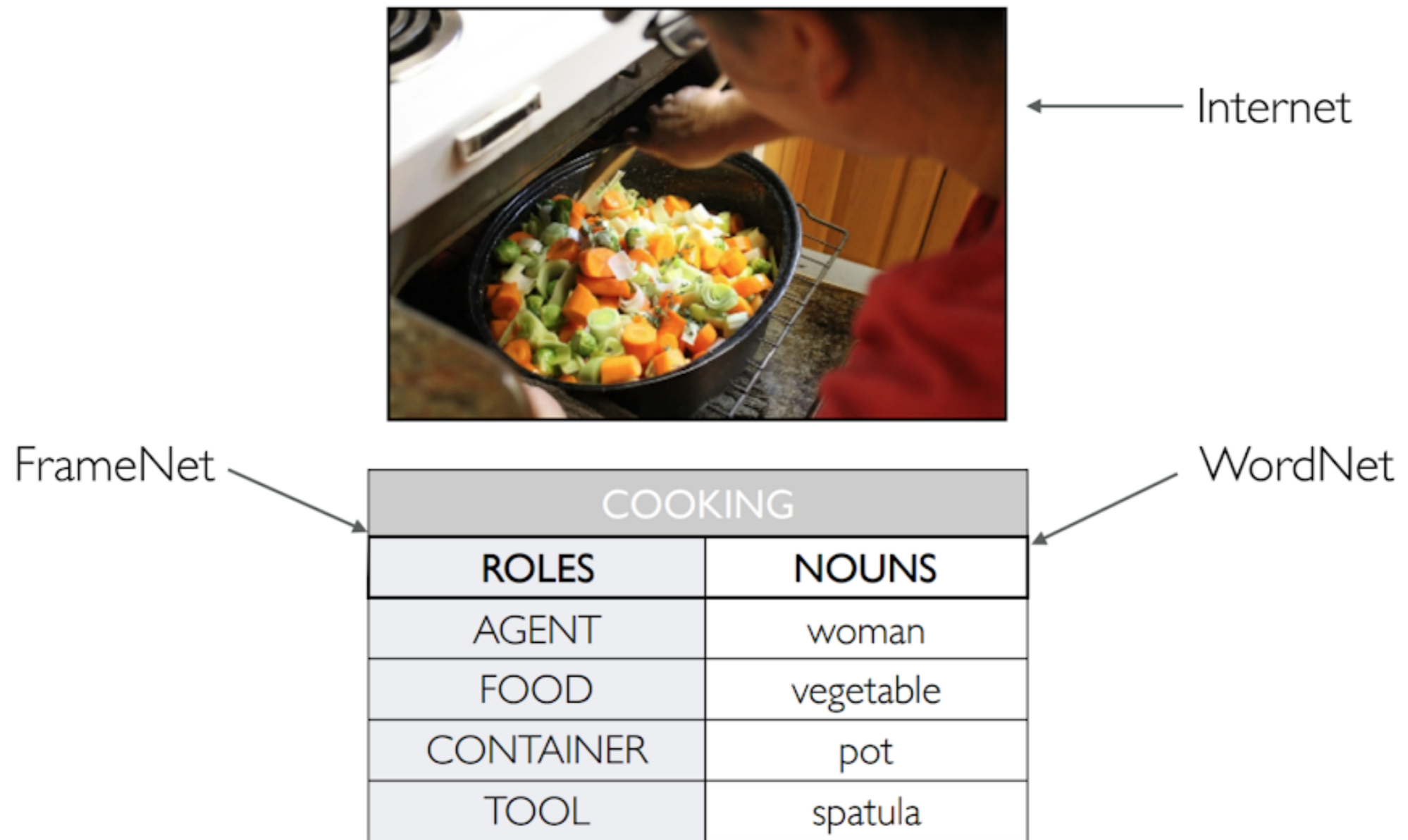
- **Bias amplification in learned models**

# 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*
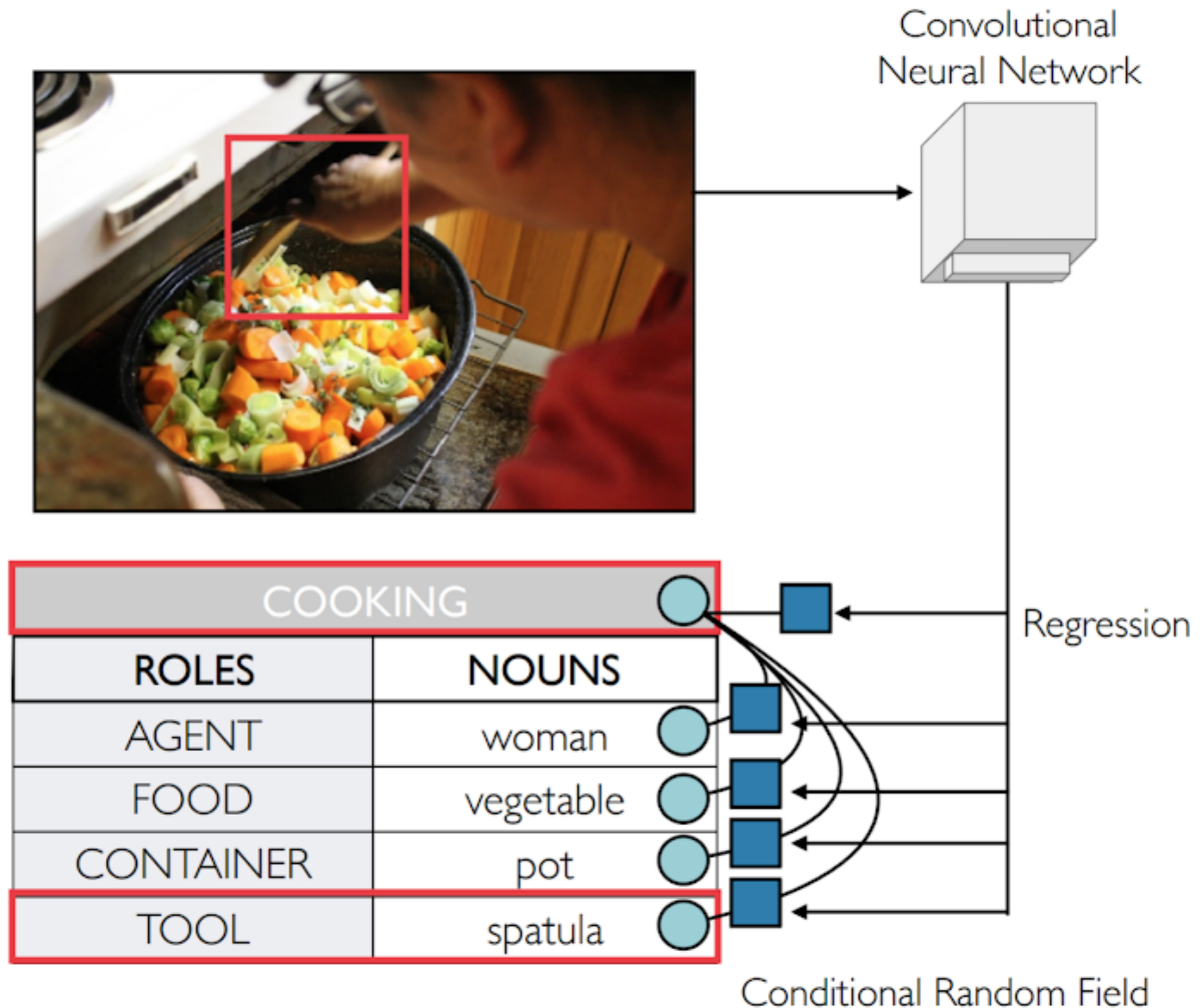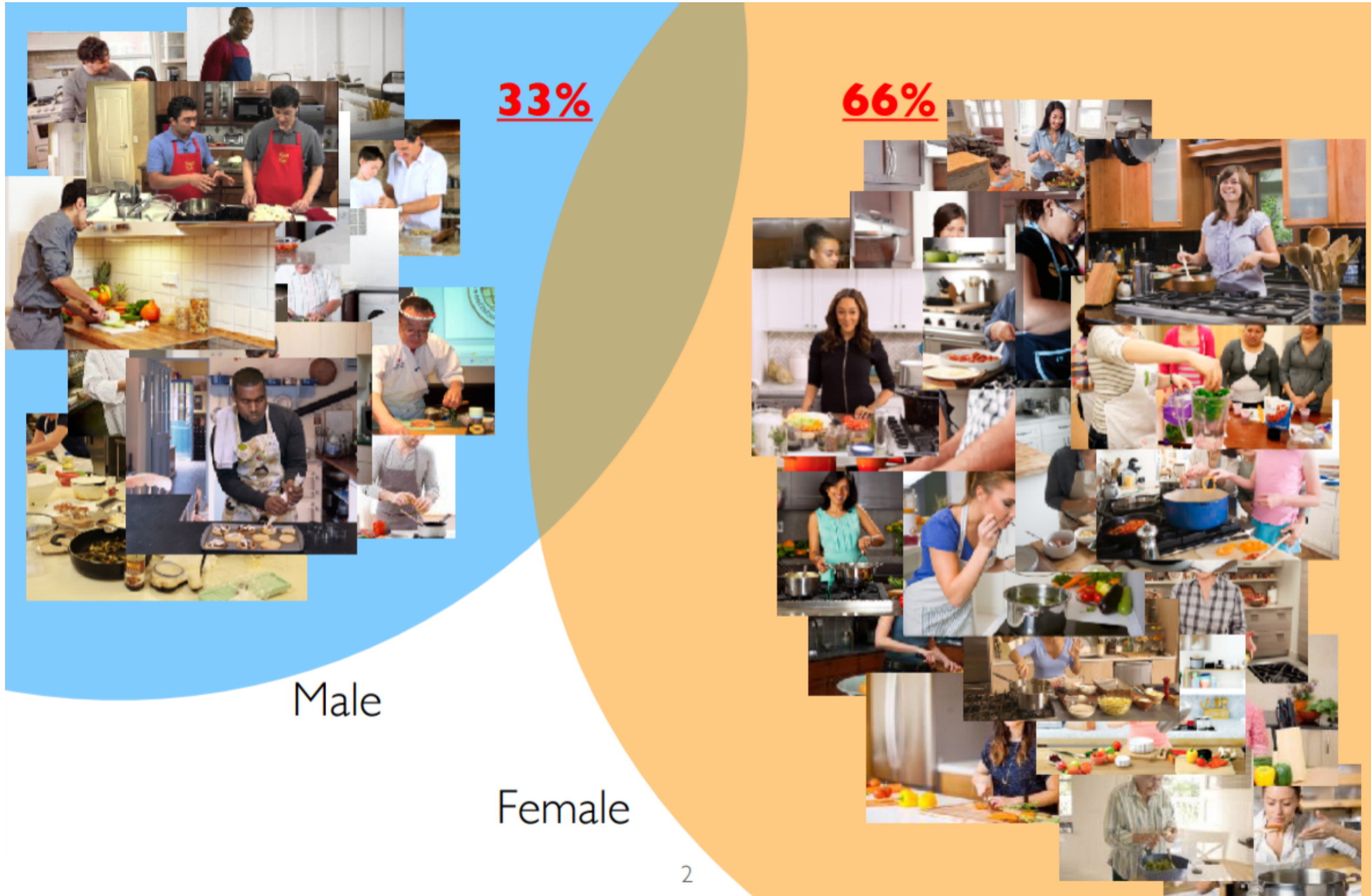
# imSitu Visual Semantic Role Labeling (vSRL)



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

# imSitu Visual Semantic Role Labeling (vSRL)

# Dataset Gender Bias



33%

66%

Male

Female

2

# 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 |

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

# Quantifying Dataset Bias: Dev Set



Predicted Gender Ratio (◆ verb)

Development Set
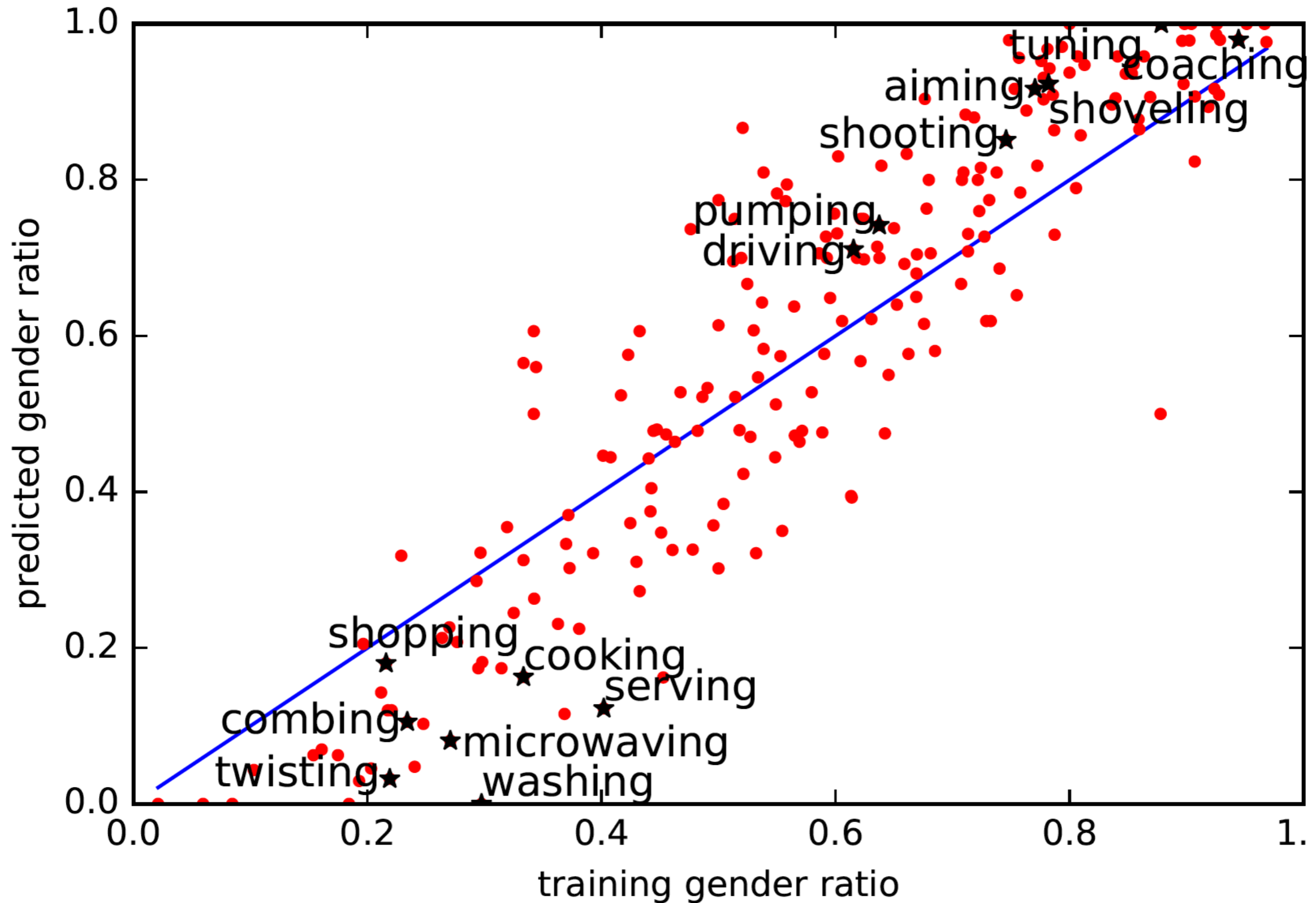
◆ cooking

● woman

● man

COOKING

| ROLES | NOUNS |
|---|---|
| ● AGENT | woman |
| FOOD | stir-fry |

COOKING

| ROLES | NOUNS |
|---|---|
| ● AGENT | man |
| FOOD | noodle |

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

# Model Bias Amplification

predicted gender ratio (y-axis), training gender ratio (x-axis)

tuning
aiming
shoveling
coaching
shooting
pumping
driving
shopping
cooking
serving
combing
microwaving
twisting
washing

# 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?