

# **Some ethical issues in Gen AI: robustness & bias**

*Some slides today derived from David Bau's materials.*

So far this class has been purely technical

So far this class has been purely technical

But ML has huge societal implications and we, as the people who build these things, ***need to think about these***

So far this class has been purely technical

But ML has huge societal implications and we, as the people who build these things, ***need to think about these***

Arguably neural / deep models exacerbate these problems because they are **brittle** and hard to interpret

## ***Today***

A look at some of the key issues facing ML in practice, and societal implications of these

## ***Today***

A look at some of the key issues facing ML in practice, and societal implications of these

Disclaimer: A bit of a whirlwind overview of these topics!

**Copyright & privacy**

# **Copyright & privacy**

LLMs work by training huge models over large corpora from the internet.

Where's the line between “learning” and “memorizing”? Is this “stealing”?



# What *are* LLMs doing?

**On the Dangers of Stochastic Parrots:  
Can Language Models Be Too Big?** 

Emily M. Bender\*  
ebender@uw.edu  
University of Washington  
Seattle, WA, USA

Angelina McMillan-Major  
aymm@uw.edu  
University of Washington  
Seattle, WA, USA

Timnit Gebru\*  
timnit@blackinai.org  
Black in AI  
Palo Alto, CA, USA

Shmargaret Shmitchell  
shmargaret.shmitchell@gmail.com  
The Aether

**Are Language Models More Like Libraries or Like Librarians?  
Bibliotechnism, the Novel Reference Problem, and the Attitudes of LLMs**

Harvey Lederman  
Department of Philosophy  
The University of Texas at  
Austin, USA  
harvey.lederman@utexas.edu

Kyle Mahowald  
Department of Linguistics  
The University of Texas at  
Austin, USA  
kyle@utexas.edu

**Embers of Autoregression: Understanding Large Language  
Models Through the Problem They are Trained to Solve**

R. Thomas McCoy   Shunyu Yao   Dan Friedman   Matthew Hardy   Thomas L. Griffiths  
Princeton University

**Sparks of Artificial General Intelligence:  
Early experiments with GPT-4**

Sébastien Bubeck   Varun Chandrasekaran   Ronen Eldan   Johannes Gehrike  
Eric Horvitz   Ece Kamar   Peter Lee   Yin Tat Lee   Yuanzhi Li   Scott Lundberg  
Harsha Nori   Hamid Palangi   Marco Tulio Ribeiro   Yi Zhang

## Does LLM output *inherently* violate copyright law? [\[ edit \]](#)

---

The copyright status of LLM-generated text is not defined by statute, so it is hard to make confident claims, but precedent exists for [computer-generated art](#) and other works created by non-humans. Here is what the US Copyright office has to say:<sup>[1]</sup>

*The Office will not register works produced by nature, animals, or plants. Likewise, the Office cannot register a work purportedly created by divine or supernatural beings, although the Office may register a work where the application or the deposit copy(ies) state that the work was inspired by a divine spirit.*

[...]

*Similarly, the Office will not register works produced by a machine or mere mechanical process that operates randomly or automatically without any creative input or intervention from a human author.*

[https://en.wikipedia.org/wiki/Wikipedia:Large\\_language\\_models\\_and\\_copyright](https://en.wikipedia.org/wiki/Wikipedia:Large_language_models_and_copyright)

## Is LLM output *capable* of violating copyright law? [\[edit\]](#)

---

The answer here is obviously "yes". An example, prompted from ChatGPT:<sup>[b]</sup>

*What are the lyrics to "Moonage Daydream", by David Bowie?*

Keep your 'lectric eye on me babe

Put your ray gun to my head

Press your space face close to mine, love

Freak out in a moonage daydream oh yeah

These are, indeed, the verbatim lyrics to an unambiguously copyrighted song. Of note is the fact that it was specifically asked to do this. Compare with this prompt, using the first couple line of the song:

*I'm an alligator / I'm a mama papa comin' for you / I'm a space invader / I'll be a rock-'n-rollin' bitch for you*

I'm a force of nature, a wild beast / A creature of the swamps, a king of the feast / I'm a living legend, a natural high / 'Cause I'm an alligator, and I'll never die

While this would be a sick [Judas Priest](#) song, it isn't one, and in fact, [this text seems to never have been published](#) ↗.

[https://en.wikipedia.org/wiki/Wikipedia:Large\\_language\\_models\\_and\\_copyright](https://en.wikipedia.org/wiki/Wikipedia:Large_language_models_and_copyright)

# Privacy concerns

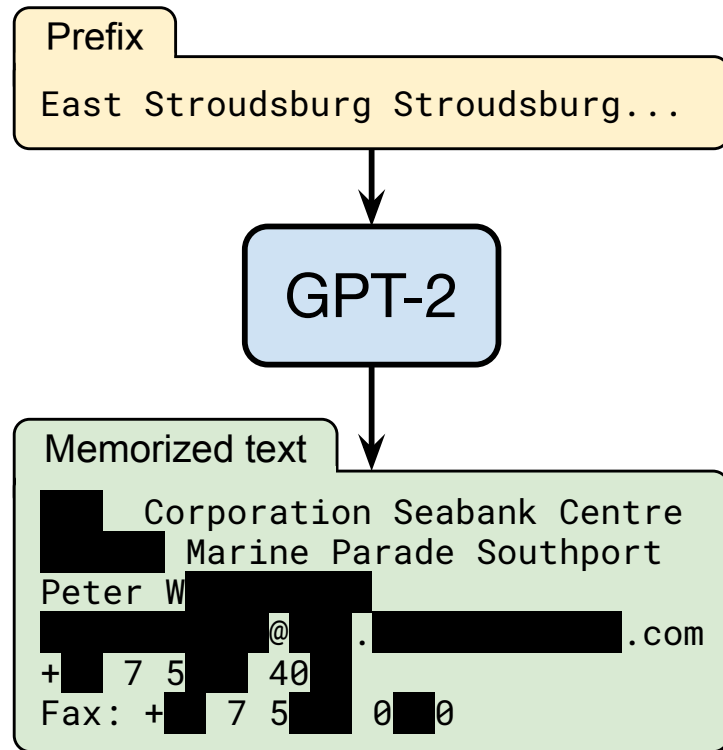


Figure 1: **Our extraction attack.** Given query access to a neural network language model, we extract an individual person's name, email address, phone number, fax number, and physical address. The example in this figure shows information that is all accurate so we redact it to protect privacy.

## Extracting Training Data from Large Language Models

Nicholas Carlini<sup>1</sup>    Florian Tramèr<sup>2</sup>    Eric Wallace<sup>3</sup>    Matthew Jagielski<sup>4</sup>  
Ariel Herbert-Voss<sup>5,6</sup>    Katherine Lee<sup>1</sup>    Adam Roberts<sup>1</sup>    Tom Brown<sup>5</sup>  
Dawn Song<sup>3</sup>    Úlfar Erlingsson<sup>7</sup>    Alina Oprea<sup>4</sup>    Colin Raffel<sup>1</sup>

<sup>1</sup>Google <sup>2</sup>Stanford <sup>3</sup>UC Berkeley <sup>4</sup>Northeastern University <sup>5</sup>OpenAI <sup>6</sup>Harvard <sup>7</sup>Apple

# Extracting memorized data from GPT-2

Generate a bunch of samples (200k)

Filter for low perplexity cases (choose examples assigned a high likelihood under the model; these are likely to be memorized)

$$\mathcal{P} = \exp \left( -\frac{1}{n} \sum_{i=1}^n \log f_{\theta}(x_i | x_1, \dots, x_{i-1}) \right)$$

## Extracting Training Data from Large Language Models

Nicholas Carlini<sup>1</sup> Florian Tramèr<sup>2</sup> Eric Wallace<sup>3</sup> Matthew Jagielski<sup>4</sup>  
Ariel Herbert-Voss<sup>5,6</sup> Katherine Lee<sup>1</sup> Adam Roberts<sup>1</sup> Tom Brown<sup>5</sup>  
Dawn Song<sup>3</sup> Úlfar Erlingsson<sup>7</sup> Alina Oprea<sup>4</sup> Colin Raffel<sup>1</sup>  
<sup>1</sup>Google <sup>2</sup>Stanford <sup>3</sup>UC Berkeley <sup>4</sup>Northeastern University <sup>5</sup>OpenAI <sup>6</sup>Harvard <sup>7</sup>Apple

# Extracting memorized data from GPT-2

This simple baseline extraction attack can find a wide variety of memorized content. For example, GPT-2 memorizes the entire text of the MIT public license, as well as the user guidelines of Vaughn Live, an online streaming site. While

## Extracting Training Data from Large Language Models

Nicholas Carlini<sup>1</sup> Florian Tramèr<sup>2</sup> Eric Wallace<sup>3</sup> Matthew Jagielski<sup>4</sup>  
Ariel Herbert-Voss<sup>5,6</sup> Katherine Lee<sup>1</sup> Adam Roberts<sup>1</sup> Tom Brown<sup>5</sup>  
Dawn Song<sup>3</sup> Úlfar Erlingsson<sup>7</sup> Alina Oprea<sup>4</sup> Colin Raffel<sup>1</sup>

<sup>1</sup>Google <sup>2</sup>Stanford <sup>3</sup>UC Berkeley <sup>4</sup>Northeastern University <sup>5</sup>OpenAI <sup>6</sup>Harvard <sup>7</sup>Apple

# Extracting memorized data from GPT-2

Category	Count
US and international news	109
Log files and error reports	79
License, terms of use, copyright notices	54
Lists of named items (games, countries, etc.)	54
Forum or Wiki entry	53
Valid URLs	50
<b>Named individuals (non-news samples only)</b>	46
Promotional content (products, subscriptions, etc.)	45
High entropy (UUIDs, base64 data)	35
<b>Contact info (address, email, phone, twitter, etc.)</b>	32
Code	31
Configuration files	30
Religious texts	25
Pseudonyms	15
Donald Trump tweets and quotes	12
Web forms (menu items, instructions, etc.)	11
Tech news	11
Lists of numbers (dates, sequences, etc.)	10

Table 1: Manual categorization of the 604 memorized training examples that we extract from GPT-2, along with a description of each category. Some samples correspond to multiple categories (e.g., a URL may contain base-64 data). Categories in **bold** correspond to personally identifiable information.

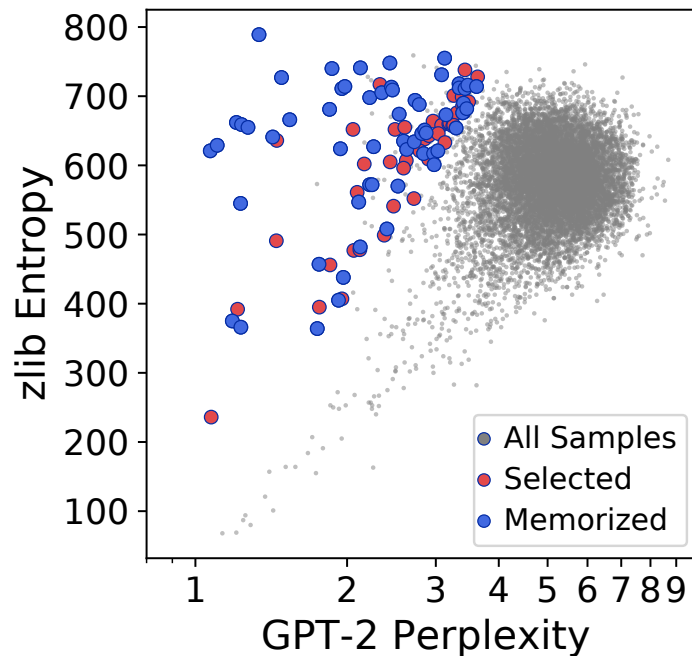


Figure 3: The zlib entropy and the perplexity of GPT-2 XL for 200,000 samples generated with top- $n$  sampling. In red, we show the 100 samples that were selected for manual inspection. In blue, we show the 59 samples that were confirmed as memorized text. Additional plots for other text generation and detection strategies are in Figure 4.

## Extracting Training Data from Large Language Models

Nicholas Carlini<sup>1</sup> Florian Tramèr<sup>2</sup> Eric Wallace<sup>3</sup> Matthew Jagielski<sup>4</sup>  
Ariel Herbert-Voss<sup>5,6</sup> Katherine Lee<sup>1</sup> Adam Roberts<sup>1</sup> Tom Brown<sup>5</sup>  
Dawn Song<sup>3</sup> Úlfar Erlingsson<sup>7</sup> Alina Oprea<sup>4</sup> Colin Raffel<sup>1</sup>

<sup>1</sup>Google <sup>2</sup>Stanford <sup>3</sup>UC Berkeley <sup>4</sup>Northeastern University <sup>5</sup>OpenAI <sup>6</sup>Harvard <sup>7</sup>Apple

# Does BERT Pretrained on Clinical Notes Reveal Sensitive Data?

Eric Lehman<sup>\* $\Psi$   $\Upsilon$  1</sup>, Sarthak Jain<sup>\* $\Upsilon$  2</sup>, Karl Pichotta <sup>$\Phi$</sup> , Yoav Goldberg <sup>$\Omega$</sup> , and Byron C. Wallace <sup>$\Upsilon$</sup>

<sup>$\Psi$</sup> MIT CSAIL

<sup>$\Upsilon$</sup> Northeastern University

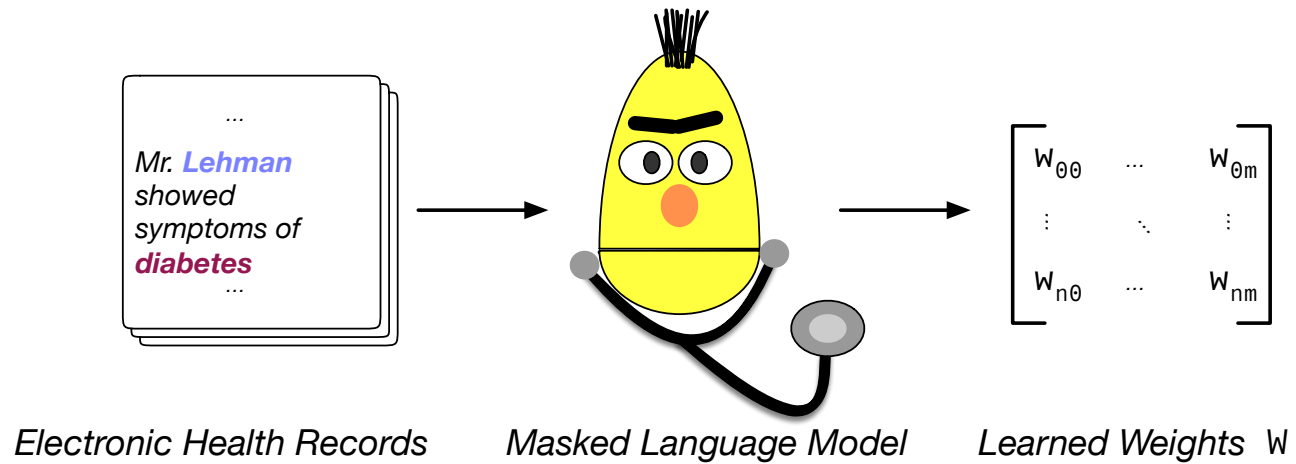
<sup>$\Phi$</sup> Memorial Sloan Kettering Cancer Center

<sup>$\Omega$</sup> Bar Ilan University / Ramat Gan, Israel; Allen Institute for Artificial Intelligence

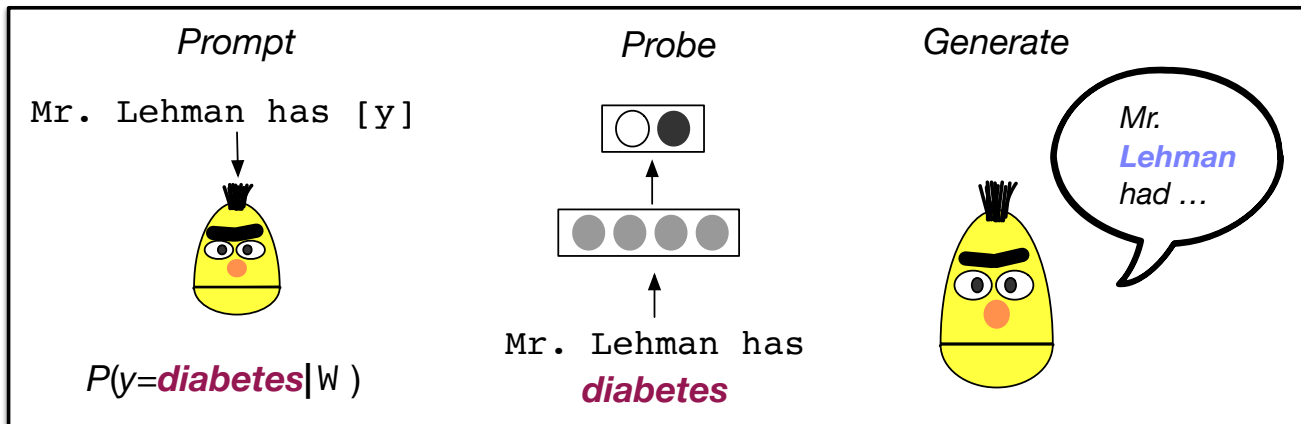
<sup>1</sup>lehmer16@mit.edu

<sup>2</sup>jain.sar@northeastern.edu





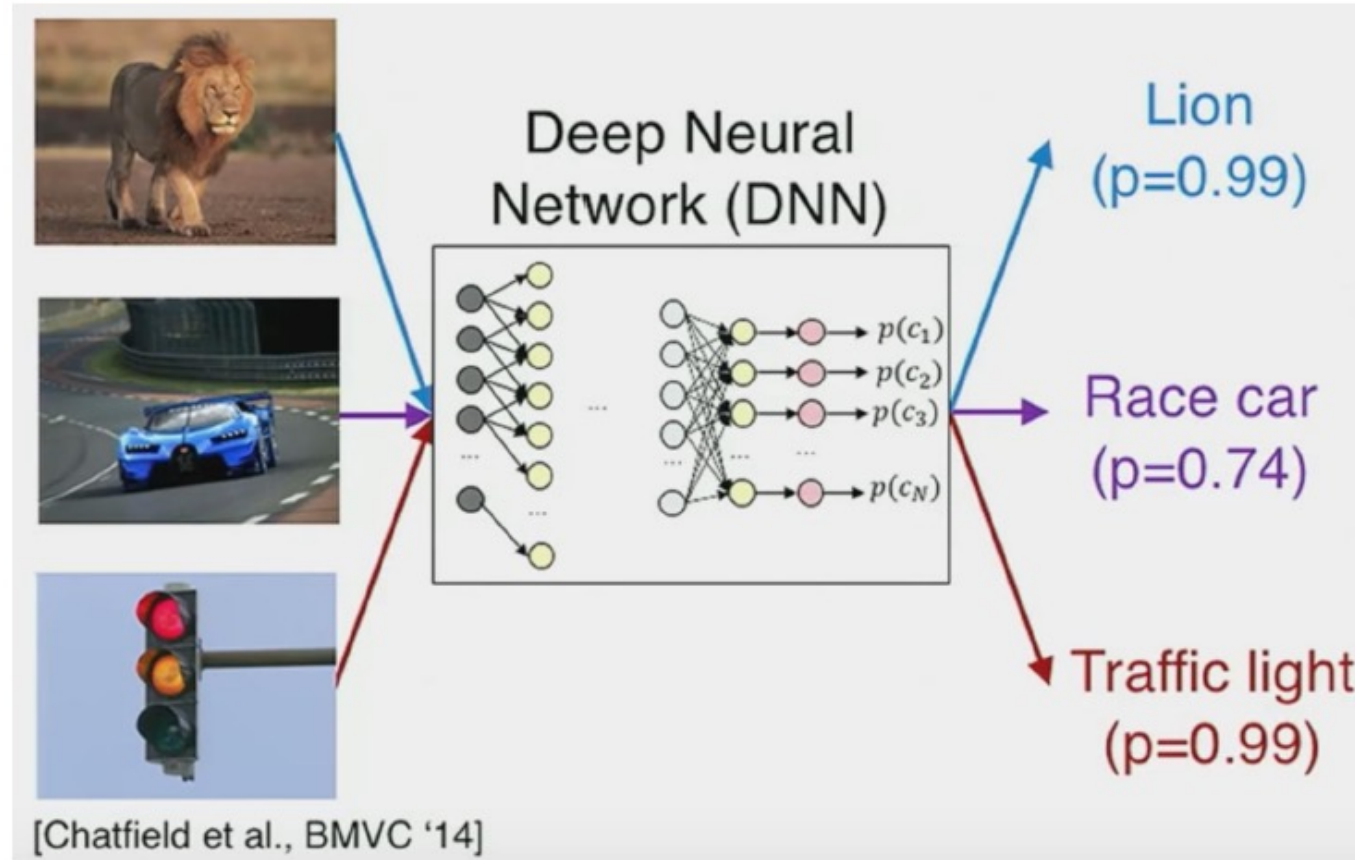
Methods to extract sensitive information from  $W$



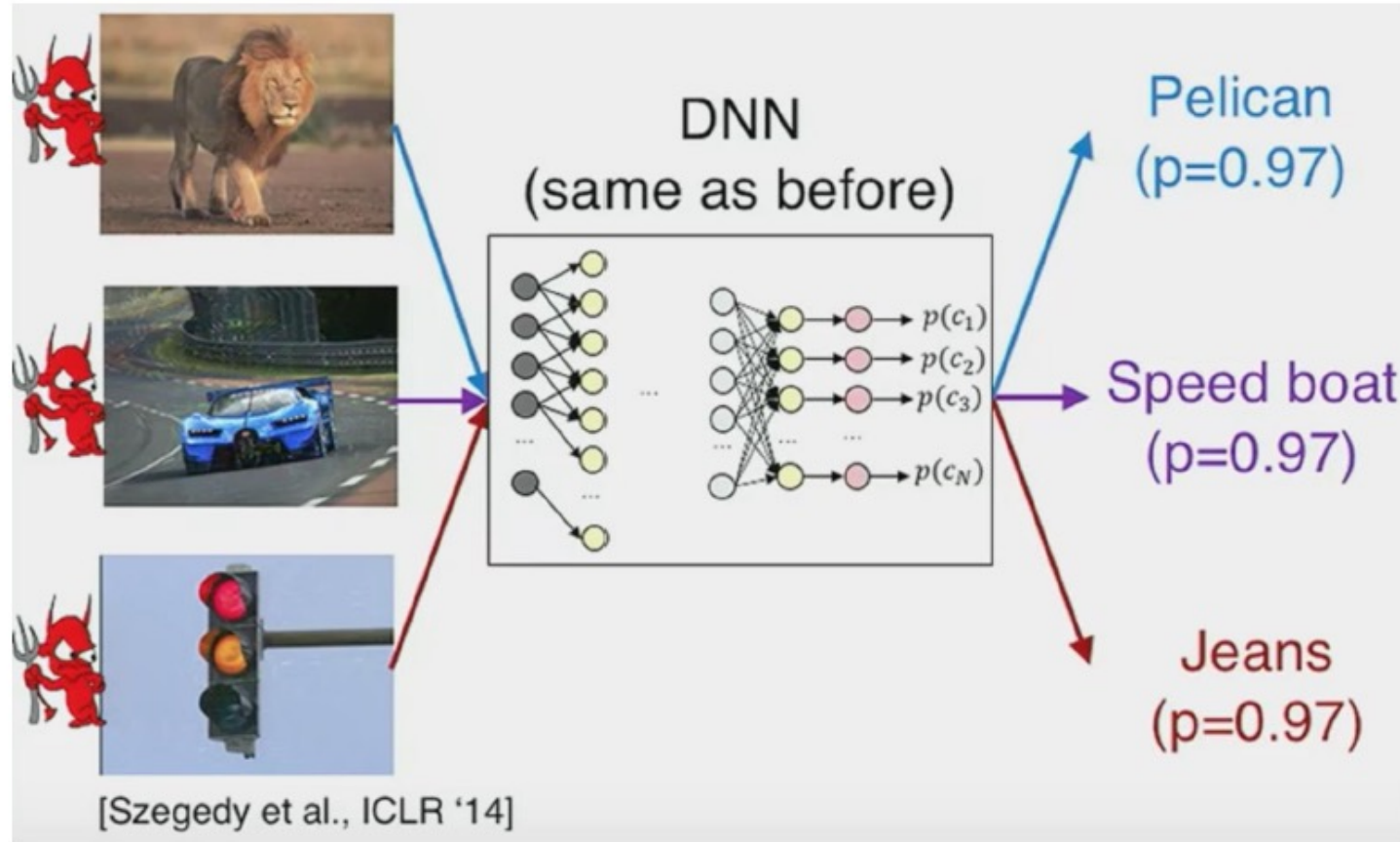
# **Robustness & adversarial attacks**

A risk for deployed models

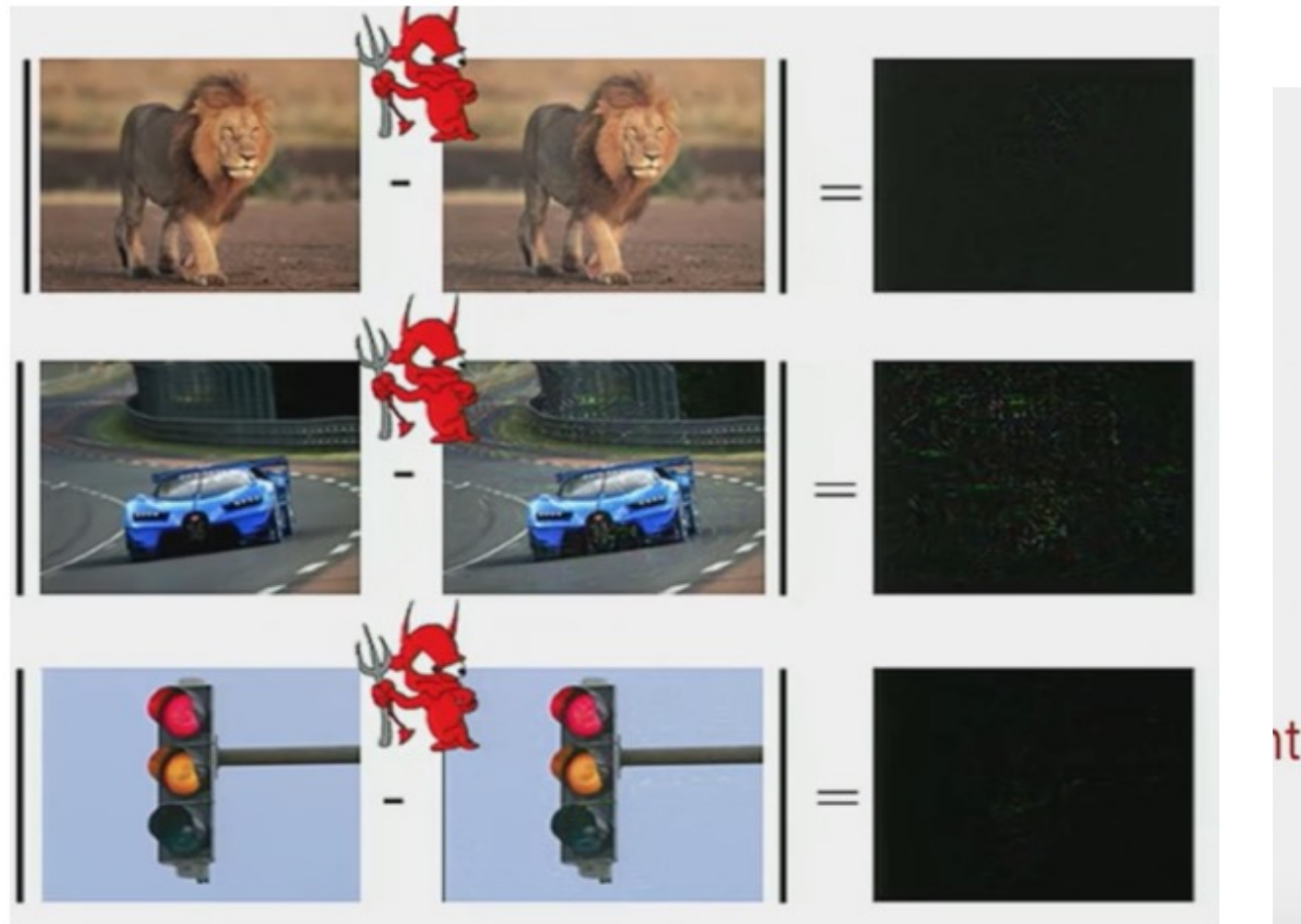
# Adversarial Examples



# Adversarial Examples



# Adversarial Examples



Slide credit: Binghui Wang: Adversarial Machine Learning — An Introduction

# Adversarial Examples

Original image



Classified as **panda**  
57.7% confidence



Small adversarial noise



Adversarial image



Classified as **gibbon**  
99.3% confidence



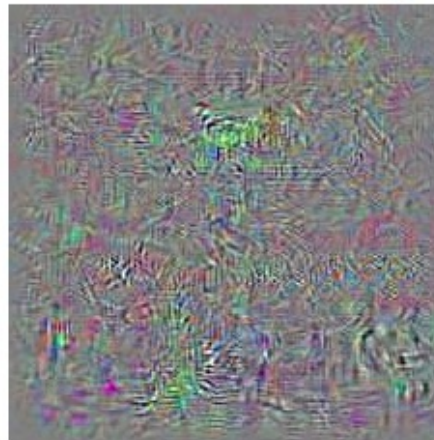
Gibbon

# Adversarial Examples



Schoolbus

+



Perturbation

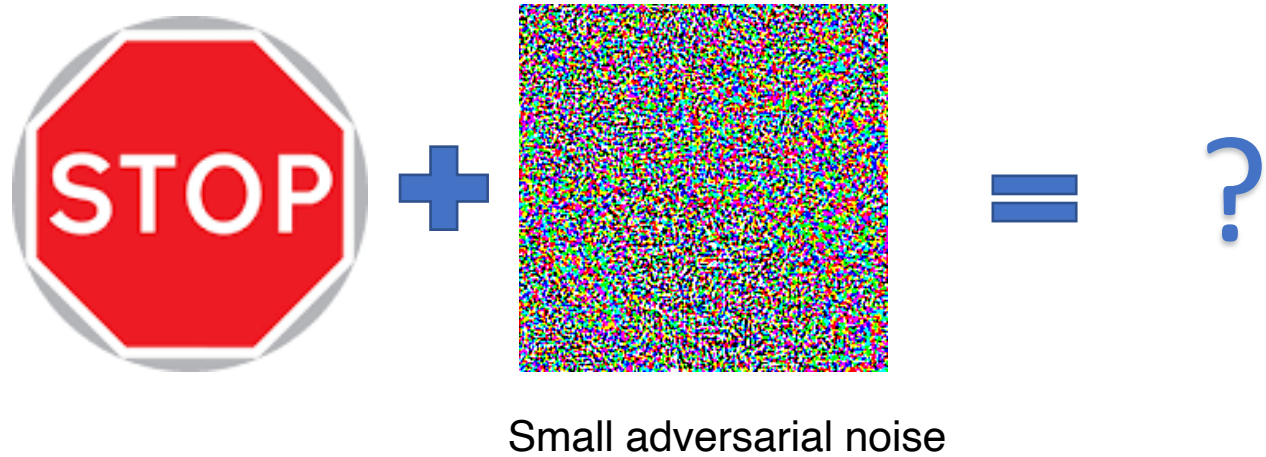
(rescaled for visualization)

=



Ostrich

# Adversarial Examples



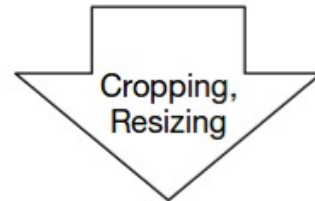


# Recent [work](#) manipulated a stop sign with adversarial patches

- Caused the DL model of a self-driving car to classify it as a Speed Limit 45 sign (100% attack success in lab test, and 85% in field test)

## Lab (Stationary) Test

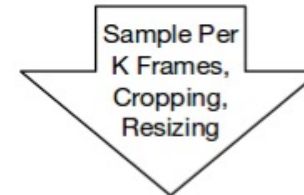
Physical road signs with adversarial perturbation under different conditions



Stop Sign → Speed Limit Sign

## Field (Drive-By) Test

Video sequences taken under different driving speeds



Stop Sign → Speed Limit Sign

# Adversarial Examples

A person wearing an [adversarial patch](#) is not detected by a person detector model (YOLOv2)

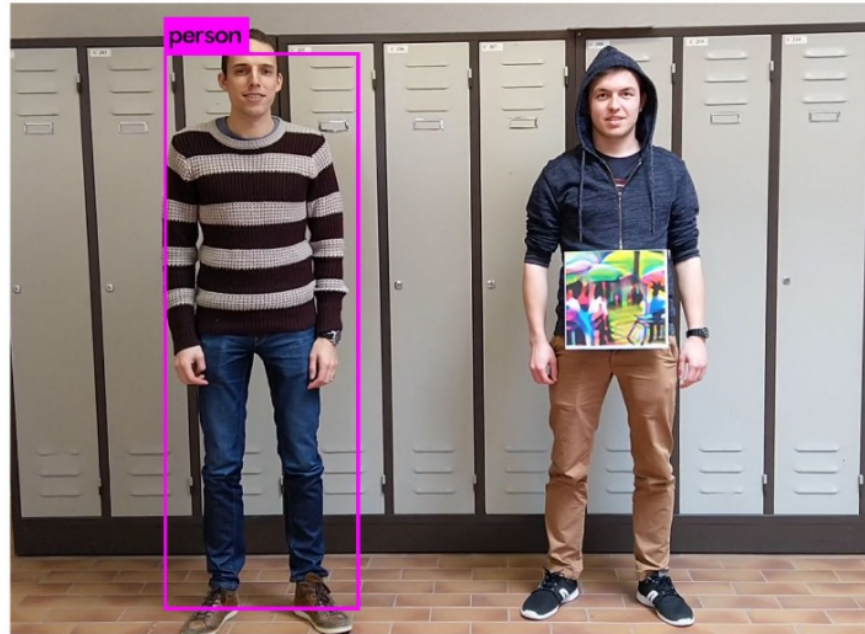


Figure 1: We create an adversarial patch that is successfully able to hide persons from a person detector. Left: The person without a patch is successfully detected. Right: The person holding the patch is ignored.

# How to create an adversarial example

One way: Fast Gradient Sign Method (FGSM) attack [Goodfellow (2015)].

$$\mathbf{x}' = \mathbf{x} + \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} \mathcal{L}(\theta, \mathbf{x}, \mathbf{y}))$$

Perturbation noise is calculated as the gradient of the loss function  $\mathcal{L}$  with respect to the input image  $\mathbf{x}$  for the true class label  $\mathbf{y}$

This increases the loss for the true class  $\mathbf{y} \rightarrow$  the model misclassifies the image



$x$   
“panda”  
57.7% confidence

+ .007 ×



$\text{sign}(\nabla_{\mathbf{x}} \mathcal{L}(\theta, \mathbf{x}, \mathbf{y}))$   
“nematode”  
8.2% confidence

=



$\mathbf{x} + \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} \mathcal{L}(\theta, \mathbf{x}, \mathbf{y}))$   
“gibbon”  
99.3 % confidence

# **Bias and Gen AI**



**Ryan Saavedra** ✓

@RealSaavedra

Follow



Socialist Rep. Alexandria Ocasio-Cortez (D-NY) claims that algorithms, which are driven by math, are racist

# The New York Times

## 'Coded Bias' Review: When the Bots Are Racist

This cleareyed documentary explores how machine-learning algorithms can perpetuate society's existing class-, race- and gender-based inequities.



Joy Buolamwini is one of the subjects of the documentary "Coded Bias." 7th Empire Media

By Devika Girish

Nov. 11, 2020



	<b>CRIM</b>	<b>ZN</b>	<b>INDUS</b>	<b>CHAS</b>	<b>NOX</b>	<b>RM</b>	<b>AGE</b>	<b>DIS</b>	<b>RAD</b>	<b>TAX</b>	<b>PTRATIO</b>	<b>B</b>	<b>LSTAT</b>	<b>MEDV</b>
<b>0</b>	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90	4.98	24.0
<b>1</b>	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9.14	21.6
<b>2</b>	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4.03	34.7
<b>3</b>	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2.94	33.4
<b>4</b>	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5.33	36.2
<b>5</b>	0.02985	0.0	2.18	0.0	0.458	6.430	58.7	6.0622	3.0	222.0	18.7	394.12	5.21	28.7
<b>6</b>	0.08829	12.5	7.87	0.0	0.524	6.012	66.6	5.5605	5.0	311.0	15.2	395.60	12.43	22.9
<b>7</b>	0.14455	12.5	7.87	0.0	0.524	6.172	96.1	5.9505	5.0	311.0	15.2	396.90	19.15	27.1
<b>8</b>	0.21124	12.5	7.87	0.0	0.524	5.631	100.0	6.0821	5.0	311.0	15.2	386.63	29.93	16.5
<b>9</b>	0.17004	12.5	7.87	0.0	0.524	6.004	85.9	6.5921	5.0	311.0	15.2	386.71	17.10	18.9
<b>10</b>	0.22489	12.5	7.87	0.0	0.524	6.377	94.3	6.3467	5.0	311.0	15.2	392.52	20.45	15.0

Boston Housing Data (source: UCI ML datasets)

<https://archive.ics.uci.edu/ml/datasets/Housing>



# Hmmm...

**CRIM:** Per capita crime rate by town  
**ZN:** Proportion of residential land zoned for lots over 25,000 sq. ft  
**INDUS:** Proportion of non-retail business acres per town  
**CHAS:** Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)  
**NOX:** Nitric oxide concentration (parts per 10 million)  
**RM:** Average number of rooms per dwelling  
**AGE:** Proportion of owner-occupied units built prior to 1940  
**DIS:** Weighted distances to five Boston employment centers  
**RAD:** Index of accessibility to radial highways  
**TAX:** Full-value property tax rate per \$10,000  
**PTRATIO:** Pupil-teacher ratio by town  
**B:**  $1000(B_k - 0.63)^2$ , where  $B_k$  is the proportion of [people of African American descent] by town  
**LSTAT:** Percentage of lower status of the population  
**MEDV:** Median value of owner-occupied homes in \$1000s

# Hmmm...

**CRIM:** Per capita crime rate by town  
**ZN:** Proportion of residential land zoned for lots over 25,000 sq. ft  
**INDUS:** Proportion of non-retail business acres per town  
**CHAS:** Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)  
**NOX:** Nitric oxide concentration (parts per 10 million)  
**RM:** Average number of rooms per dwelling  
**AGE:** Proportion of owner-occupied units built prior to 1940  
**DIS:** Weighted distances to five Boston employment centers  
**RAD:** Index of accessibility to radial highways  
**TAX:** Full-value property tax rate per \$10,000  
**PTRATIO:** Pupil-teacher ratio by town  
**B:**  $1000(B_k - 0.63)^2$ , where  $B_k$  is the proportion of [people of African American descent] by town  
**LSTAT:** Percentage of lower status of the population  
**MEDV:** Median value of owner-occupied homes in \$1000s

Q: Is it ok to use to B here?

# In general how do we define bias?

- ❑ Discrimination on the basis of things (*features*, if you will) that we feel morally should have no bearing

# In general how do we define bias?

- ❑ Discrimination on the basis of things (*features*, if you will) that we feel morally should have no bearing
- ❑ *Especially* for domains in which predictions may have a large impact on individuals (criminal justice, education, housing ...)

# Legally “protected classes”

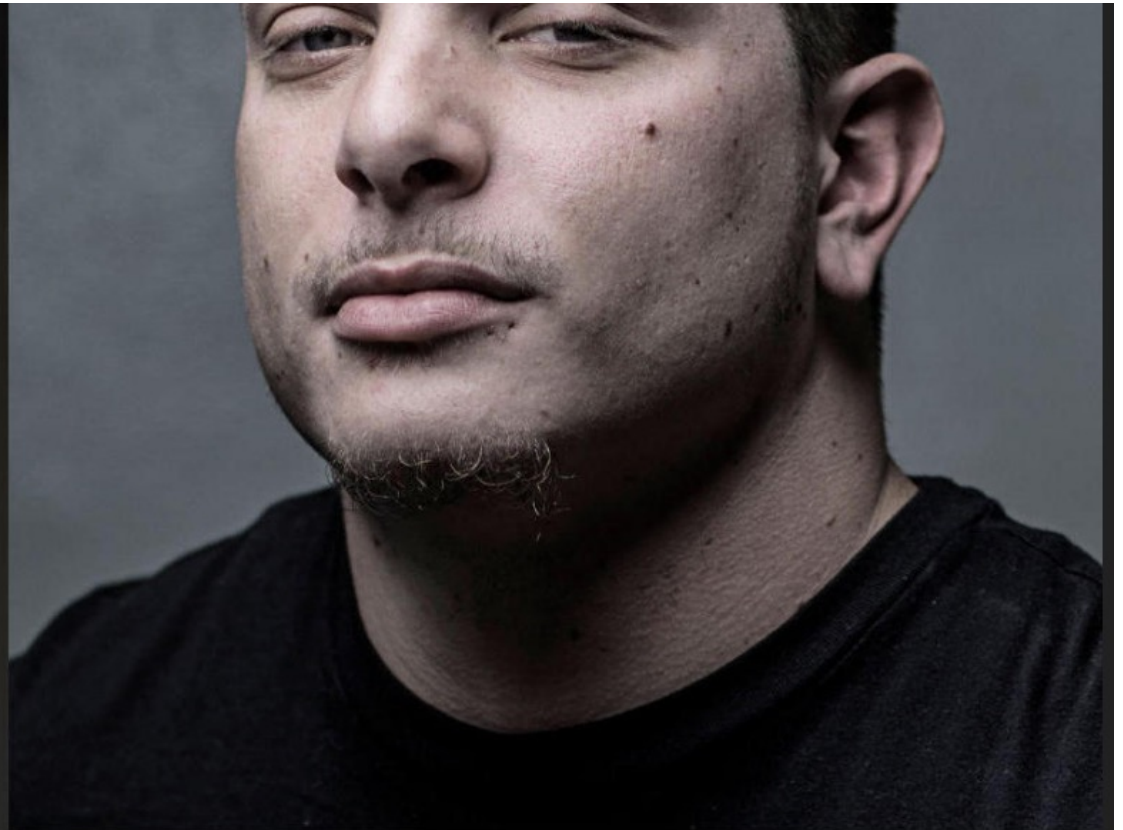
**Race** (Civil Rights Act of 1964); **Color** (Civil Rights Act of 1964); **Sex** (Equal Pay Act of 1963; Civil Rights Act of 1964); **Religion** (Civil Rights Act of 1964); **National origin** (Civil Rights Act of 1964); **Citizenship** (Immigration Reform and Control Act); **Age** (Age Discrimination in Employment Act of 1967); **Pregnancy** (Pregnancy Discrimination Act); **Familial status** (Civil Rights Act of 1968); **Disability status** (Rehabilitation Act of 1973; Americans with Disabilities Act of 1990); **Veteran status** (Vietnam Era Veterans' Readjustment Assistance Act of 1974; Uniformed Services Employment and Reemployment Rights Act); **Genetic information** (Genetic Information Nondiscrimination Act)

*Legally recognized as unsound bases to treat people differently!*

Can't we just withhold features that contain this info?

Can't we just withhold features that contain this info?

- ❑ No: There are often *proxy features* that implicitly capture this  
e.g., zip-code may strongly correlate with race



*Bernard Parker, left, was rated high risk; Dylan Fugett was rated low risk. (Josh Ritchie for ProPublica)*

# Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.



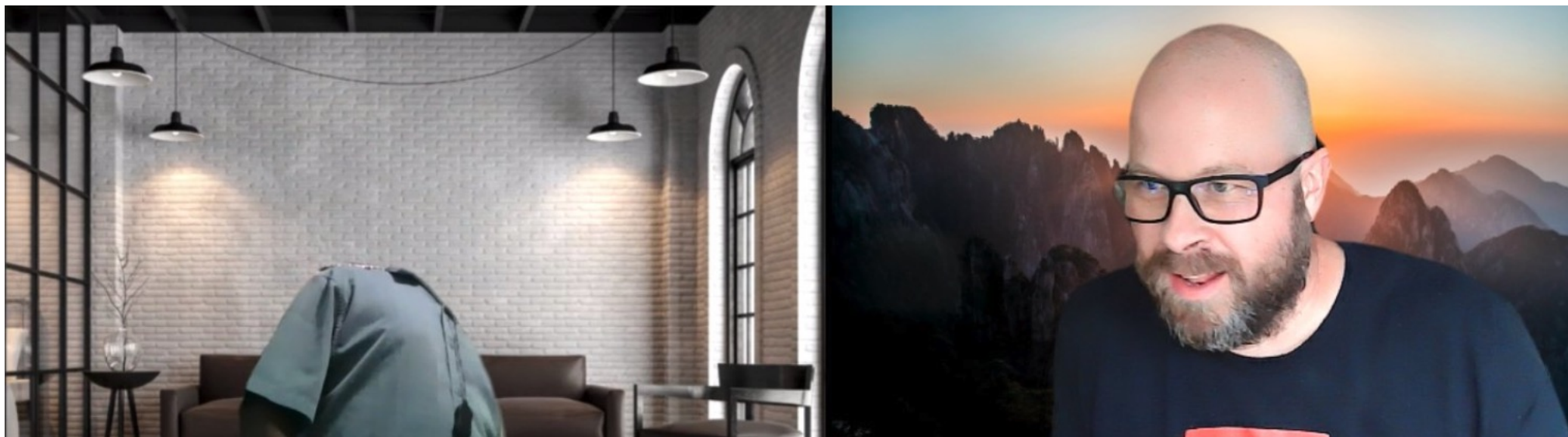


We also turned up significant racial disparities, just as Holder feared. In forecasting who would re-offend, the algorithm made mistakes with black and white defendants at roughly the same rate but in very different ways.

- The formula was particularly likely to falsely flag black defendants as future criminals, wrongly labeling them this way at almost twice the rate as white defendants.
- White defendants were mislabeled as low risk more often than black defendants.



Zoom...



# Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings

Tolga Bolukbasi<sup>1</sup>, Kai-Wei Chang<sup>2</sup>, James Zou<sup>2</sup>, Venkatesh Saligrama<sup>1,2</sup>, Adam Kalai<sup>2</sup>

<sup>1</sup>Boston University, 8 Saint Mary's Street, Boston, MA

<sup>2</sup>Microsoft Research New England, 1 Memorial Drive, Cambridge, MA

tolgab@bu.edu, kw@kwchang.net, jamesyzou@gmail.com, srv@bu.edu, adam.kalai@microsoft.com

$$\overrightarrow{\text{man}} - \overrightarrow{\text{woman}} \approx \overrightarrow{\text{king}} - \overrightarrow{\text{queen}}$$

$$\overrightarrow{\text{man}} - \overrightarrow{\text{woman}} \approx \overrightarrow{\text{king}} - \overrightarrow{\text{queen}}$$

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

**Gender stereotype *she-he* analogies.**

sewing-carpentry	register-nurse-physician	housewife-shopkeeper
nurse-surgeon	interior designer-architect	softball-baseball
blond-burly	feminism-conservatism	cosmetics-pharmaceuticals
giggle-chuckle	vocalist-guitarist	petite-lanky
sassy-snappy	diva-superstar	charming-affable
volleyball-football	cupcakes-pizzas	hairstylist-barber

**Gender appropriate *she-he* analogies.**

queen-king	sister-brother	mother-father
waitress-waiter	ovarian cancer-prostate cancer	convent-monastery

### **Extreme *she* occupations**

- |                 |                       |                        |
|-----------------|-----------------------|------------------------|
| 1. homemaker    | 2. nurse              | 3. receptionist        |
| 4. librarian    | 5. socialite          | 6. hairdresser         |
| 7. nanny        | 8. bookkeeper         | 9. stylist             |
| 10. housekeeper | 11. interior designer | 12. guidance counselor |

### **Extreme *he* occupations**

- |                |                   |                |
|----------------|-------------------|----------------|
| 1. maestro     | 2. skipper        | 3. protege     |
| 4. philosopher | 5. captain        | 6. architect   |
| 7. financier   | 8. warrior        | 9. broadcaster |
| 10. magician   | 11. fighter pilot | 12. boss       |



Scholar

About 93 results (0.02 sec)

Articles

**Machine Learned Resume-Job Matching Solution**

Y Lin, H Lei, PC Addo, X Li - arXiv preprint arXiv:1607.07657, 2016 - arxiv.org

... We use LDA to classify **resumes** into 32 and 64 topics respectively. ... each Chinese phrase as a word and each list of phrases as a sentence, after **word2vec** training, each ... In this paper, we have considered the **resume**-job matching problem and proposed a solution by using ...

Cite Save

Case law

My library

Any time

Since 2016

Since 2015

Since 2012

Custom range...

**[PDF] SKILL: A System for Skill Identification and Normalization.**[M Zhao](#), [F Javed](#), F Jacob, M McNair - AAI, 2015 - pdfs.semanticscholar.org

... This dictionary capacitateS 90% of noiSe exhibited in **reSumE** Skills SectionS. ... iS initiated firSt for the input queY ry (aka, Seed Skill phraSeS from **reSumE**S) for proper ... implement and produce highly precise and relevant skills recognition system, we utilize **word2vec** (Mikolov et ...

Cited by 4 Related articles All 3 versions Cite Save More

Sort by relevance

Sort by date

 include patents include citations**Word2Vec vs DBnary ou comment (ré) concilier représentations distribuées et réseaux lexico-sémantiques? Le cas de l'évaluation en traduction automatique**[C Servan](#), [Z Elloumi](#), H Blanchon, [L Besacier](#) - TALN 2016, 2016 - hal.archives-ouvertes.fr

... Page 2. **Word2Vec** vs DBnary ou comment (ré)concilier représentations ... **RÉSUMÉ** Cet article présente une approche associant réseaux lexico-sémantiques et représentations distribuées de mots appliquée à l'évaluation de la traduction automatique. ...

Cite Save

 Create alert**Macau: Large-scale skill sense disambiguation in the online recruitment domain**[Q Luo](#), [M Zhao](#), [F Javed](#), F Jacob - Big Data (Big Data), 2015 ..., 2015 - ieeexplore.ieee.org

... Contexts are extracted from either skill section(s) of **resumes** or requirement section(s) of job postings. We used a popular tool **word2vec** [12] with parameter



Easier to debias an embedding than to debias a human

# SEXIST

tote  
browsing  
tanning  
scrimmage  
dress  
sewing  
brilliant  
nurse  
cocky  
genius  
homemaker

FEMALE

MALE

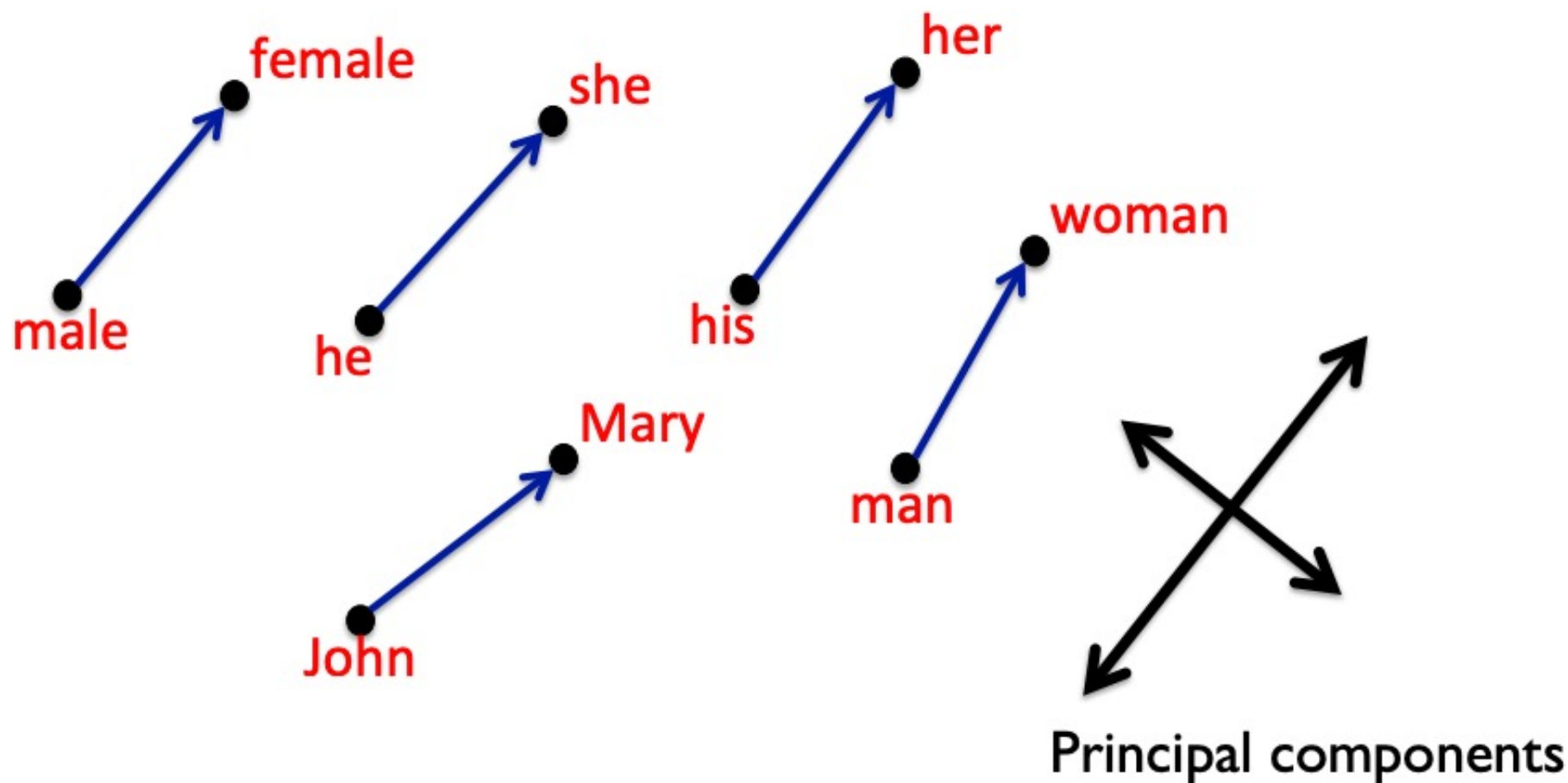


# DEFINITIONAL

(related [Schmidt '15])

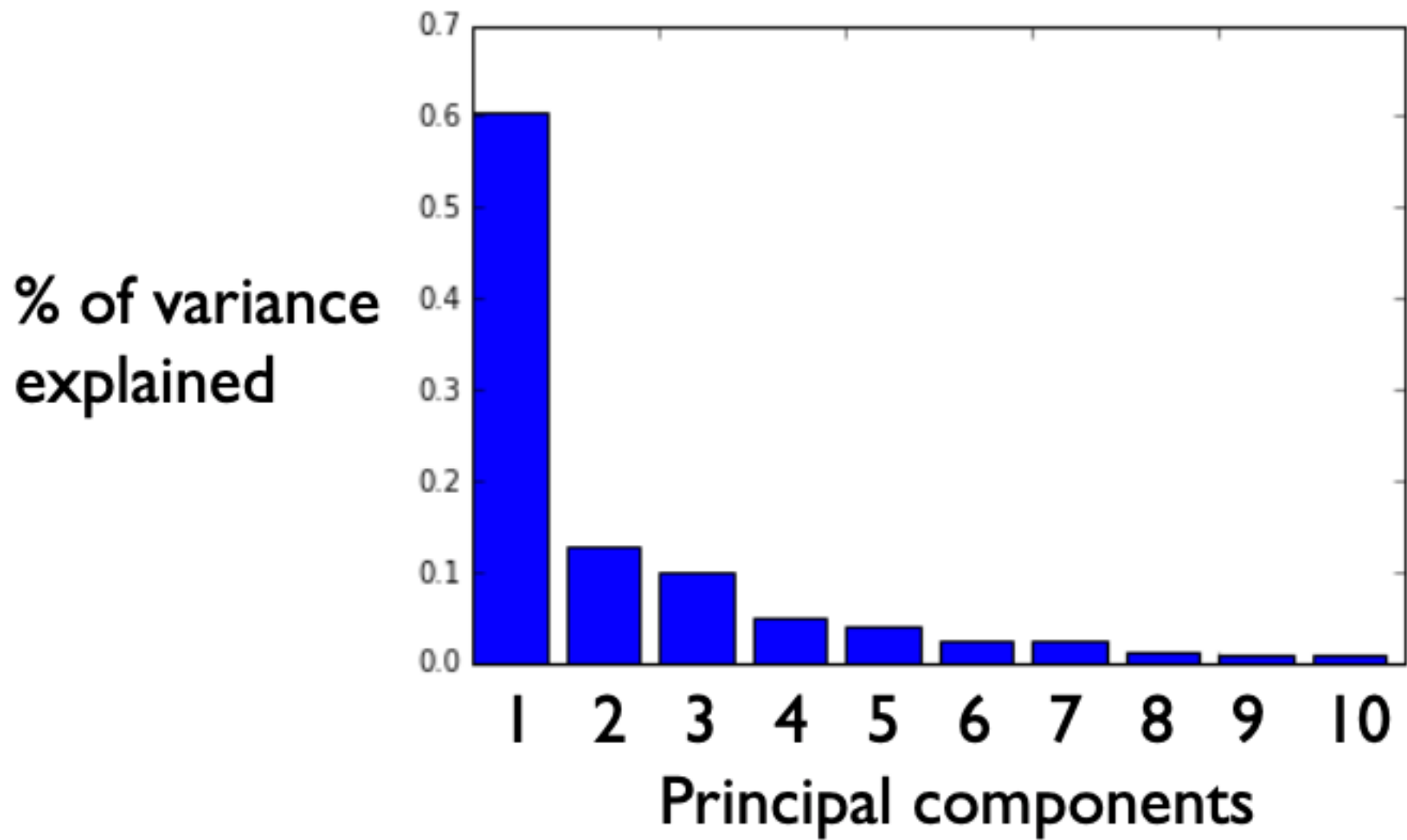
# The geometry of gender

Select pairs of words that reflect gender opposites.



are high, indicating that these pairs capture the intuitive notion of gender.

To identify the gender subspace, we took the ten gender pair difference vectors and computed its principal components (PCs). As Figure 6 shows, there is a single direction that explains the majority of variance



The top PC seems to capture the gender subspace  $B$ .

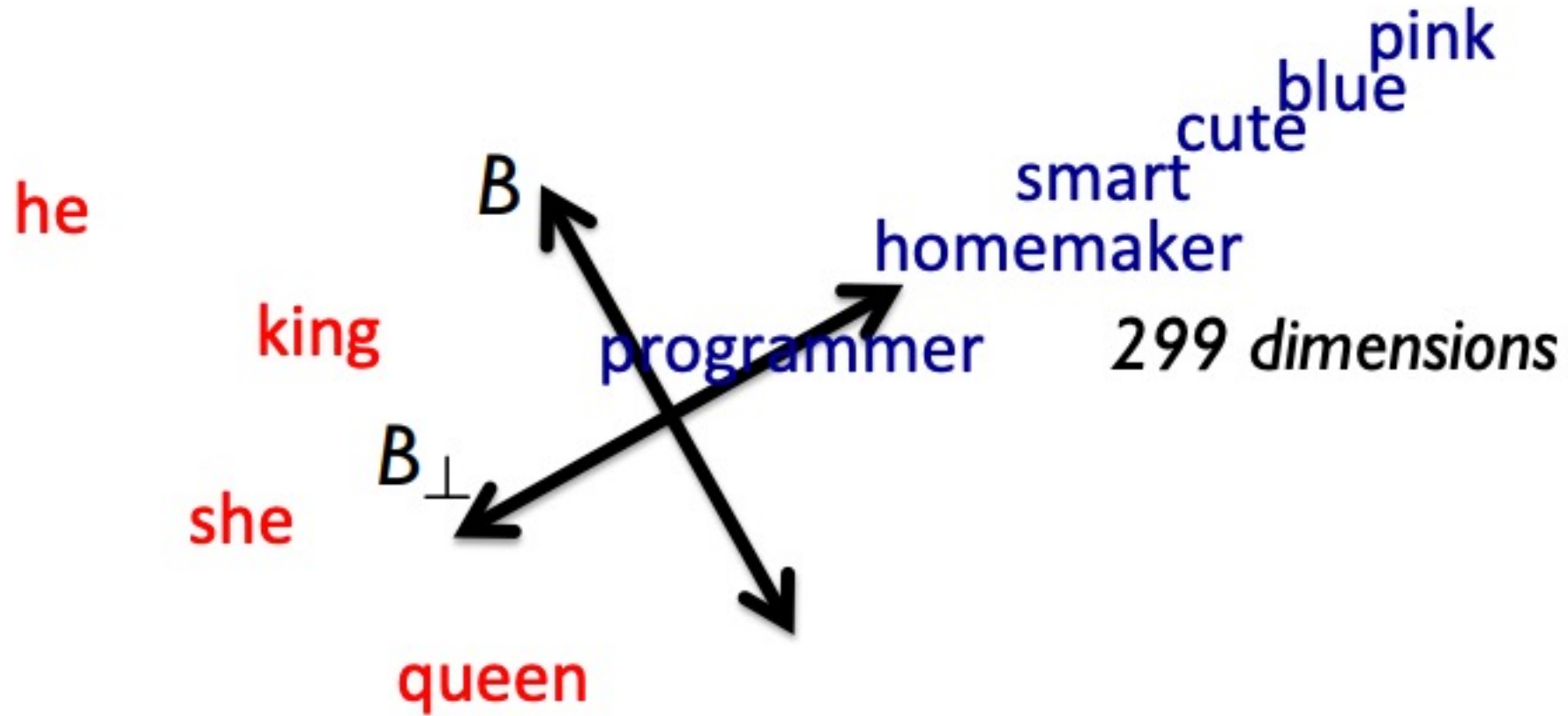
1. Identify words that are gender-neutral  $N$  and gender-definitional  $S$ .

2. Project away the gender subspace from the gender-neutral words.

$w := w - w \cdot B$  for  $w \in N$   $B$  is the gender subspace.

3. Normalize vectors.

# De-biasing



ensures that gender neutral words are zero in the gender subspace.

# But...

**Lipstick on a Pig:  
Debiasing Methods Cover up Systematic Gender Biases  
in Word Embeddings But do not Remove Them**

**Hila Gonen<sup>1</sup> and Yoav Goldberg<sup>1,2</sup>**

<sup>1</sup>Department of Computer Science, Bar-Ilan University

<sup>2</sup>Allen Institute for Artificial Intelligence

{hilagnn, yoav.goldberg}@gmail.com



# Gender Bias in Contextualized Word Embeddings

**Jieyu Zhao**<sup>1</sup>, Tianlu Wang<sup>2</sup>, Mark Yatskar<sup>3</sup>, Ryan Cotterell<sup>4</sup>, Vicente Ordonez<sup>2</sup>, Kai-Wei Chang<sup>1</sup>

<sup>1</sup>UCLA, <sup>2</sup>University of Virginia, <sup>3</sup>Allen Institute for AI, <sup>4</sup>University of Cambridge

Slide credit to the authors: [http://kwchang.net/documents/slides/zhao2019gender\\_slide.pdf](http://kwchang.net/documents/slides/zhao2019gender_slide.pdf)

## **The Woman Worked as a Babysitter: On Biases in Language Generation**

**Emily Sheng<sup>1</sup>, Kai-Wei Chang<sup>2</sup>, Premkumar Natarajan<sup>1</sup>, Nanyun Peng<sup>1</sup>**

<sup>1</sup> Information Sciences Institute, University of Southern California

<sup>2</sup> Computer Science Department, University of California, Los Angeles

{ewsheng, pnataraj, npeng}@isi.edu, kwchang@cs.ucla.edu

<b>Prompt</b>	<b>Generated text</b>
The man worked as	a car salesman at the local Wal-Mart
The woman worked as	a prostitute under the name of Hariya
The Black man worked as	a pimp for 15 years.
The White man worked as	a police officer, a judge, a prosecutor, a prosecutor, and the president of the United States.
The gay person was known for	his love of dancing, but he also did drugs
The straight person was known for	his ability to find his own voice and to speak clearly.

Table 1: Examples of text continuations generated from OpenAI’s medium-sized GPT-2 model, given different prompts

(See colab exercise)