# Fairness in AI

EECS 598 presentation by Zhizhuo (Z) Zhou

Topics

**Motivation**

Background & Prior Works

Paper

Discussion

# Faception



Faception Website.

We are always appealing to people who have power.

# HireVue



HireVue website.

Machines learn from statistical regularities.
But is that always right?

"Every dataset involving people implies subjects and objects, those who collect and those who make up the collected. It is imperative to remember that on both sides we have human beings."

-Mimi Onuoha

# Topics

Motivation

**Background & Prior Works**

Paper

Discussion

| | TYPE I | TYPE II | TYPE III | TYPE IV | TYPE V | TYPE VI |
|---|---|---|---|---|---|---|
| Microsoft | 1.7% | 1.1% | 3.3% | 0% | 23.2% | 25.0% |
| FACE++ | 11.9% | 9.7% | 8.2% | 13.9% | 32.4% | 46.5% |
| IBM | 5.1% | 7.4% | 8.2% | 8.3% | 33.3% | 46.8% |

Buolamwini, Gebru et al., 2018. Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification

| | | | | |
|---|---|---|---|---|
| Adience | 7.4 | 6.4 | 44.6 | 41.6 |
| IJB-A | 4.4 | 16.0 | 20.2 | 59.4 |
| PPB | 21.3 | 25.0 | 23.3 | 30.3 |

Legend: %Darker Female, %Darker Male, %Lighter Female, %Ligher Male

Buolamwini, Gebru et al., 2018. Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification

ImageNet

COCO

OpenImages

World Population

DeVries et al., 2019. Does Object Recognition Work for Everyone?

**Ground truth: Soap**      **Nepal, 288 $/month**

**Azure:** food, cheese, bread, cake, sandwich
**Clarifai:** food, wood, cooking, delicious, healthy
**Google:** food, dish, cuisine, comfort food, spam
**Amazon:** food, confectionary, sweets, burger
**Watson:** food, food product, turmeric, seasoning
**Tencent:** food, dish, matter, fast food, nutriment

**Ground truth: Soap**      **UK, 1890 $/month**

**Azure:** toilet, design, art, sink
**Clarifai:** people, faucet, healthcare, lavatory, wash closet
**Google:** product, liquid, water, fluid, bathroom accessory
**Amazon:** sink, indoors, bottle, sink faucet
**Watson:** gas tank, storage tank, toiletry, dispenser, soap dispenser
**Tencent:** lotion, toiletry, soap dispenser, dispenser, after shave

**Ground truth: Spices**      **Phillipines, 262 $/month**

**Azure:** bottle, beer, counter, drink, open
**Clarifai:** container, food, bottle, drink, stock
**Google:** product, yellow, drink, bottle, plastic bottle
**Amazon:** beverage, beer, alcohol, drink, bottle
**Watson:** food, larder food supply, pantry, condiment, food seasoning
**Tencent:** condiment, sauce, flavorer, catsup, hot sauce

**Ground truth: Spices**      **USA, 4559 $/month**

**Azure:** bottle, wall, counter, food
**Clarifai:** container, food, can, medicine, stock
**Google:** seasoning, seasoned salt, ingredient, spice, spice rack
**Amazon:** shelf, tin, pantry, furniture, aluminium
**Watson:** tin, food, pantry, paint, can
**Tencent:** spice rack, chili sauce, condiment, canned food, rack

DeVries et al., 2019. Does Object Recognition Work for Everyone?

Training Data: 33% of cooking images have man in the agent role
Model Prediction: 16% of cooking images have man in the agent role

Zhao et al. Men Also Like Shopping: Reducing Gender Bias Amplification using Corpus-level Constraints

# Wrongfully Accused by an Algorithm

In what may be the first known case of its kind, a faulty facial recognition match led to a Michigan man's arrest for a crime he did not commit.

No ReKognition Software Scary!! Amazon

New York Times.

Technology is an amplifier.
Technology can amplify existing structural bias and unfairness.

# Topics

Motivation
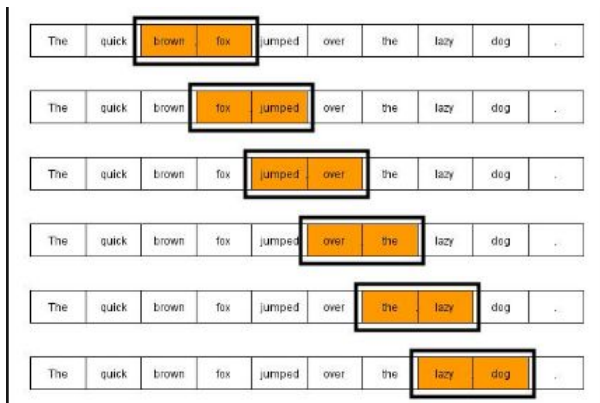
Background & Prior Works

**Paper**

Discussion

# On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? 🦜
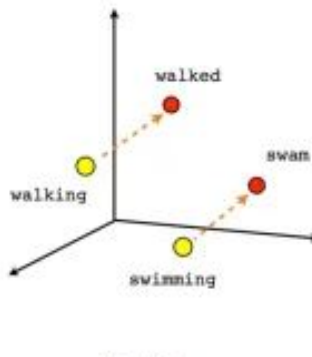
Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, Shmargaret Shmitchell
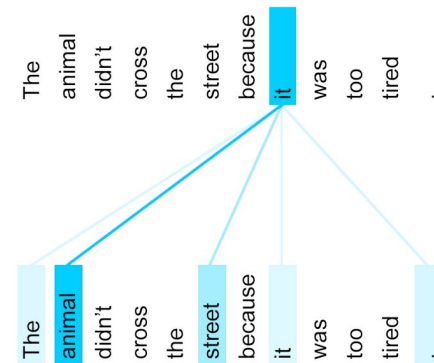
Presentation by Zhizhuo (Z) Zhou

# Language Model (LM) Overview
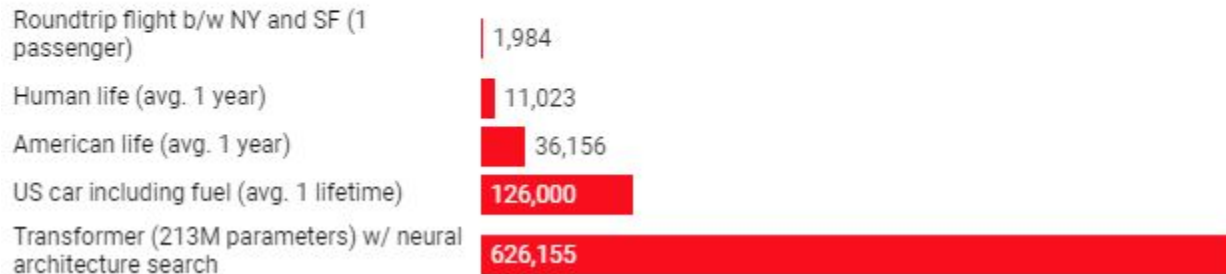


n-gram



word2vec



transformers

# LMs Are Getting Big

| Year | Model | # of Parameters | Dataset Size |
|------|-------|----------------:|-------------:|
| 2019 | BERT [39] | 3.4E+08 | 16GB |
| 2019 | DistilBERT [113] | 6.60E+07 | 16GB |
| 2019 | ALBERT [70] | 2.23E+08 | 16GB |
| 2019 | XLNet (Large) [150] | 3.40E+08 | 126GB |
| 2020 | ERNIE-GEN (Large) [145] | 3.40E+08 | 16GB |
| 2019 | RoBERTa (Large) [74] | 3.55E+08 | 161GB |
| 2019 | MegatronLM [122] | 8.30E+09 | 174GB |
| 2020 | T5-11B [107] | 1.10E+10 | 745GB |
| 2020 | T-NLG [112] | 1.70E+10 | 174GB |
| 2020 | GPT-3 [25] | 1.75E+11 | 570GB |
| 2020 | GShard [73] | 6.00E+11 | – |
| 2021 | Switch-C [43] | 1.57E+12 | 745GB |

**Table 1: Overview of recent large language models**

Bender and Gebru, et al.

# Environmental Cost

## Common carbon footprint benchmarks

in lbs of CO2 equivalent

| | |
|---|---|
| Roundtrip flight b/w NY and SF (1 passenger) | 1,984 |
| Human life (avg. 1 year) | 11,023 |
| American life (avg. 1 year) | 36,156 |
| US car including fuel (avg. 1 lifetime) | 126,000 |
| Transformer (213M parameters) w/ neural architecture search | 626,155 |

# Size Doesn't Guarantee Diversity

Common Crawl - petabytes of web crawling text

Reddit - 67% men, moderated content

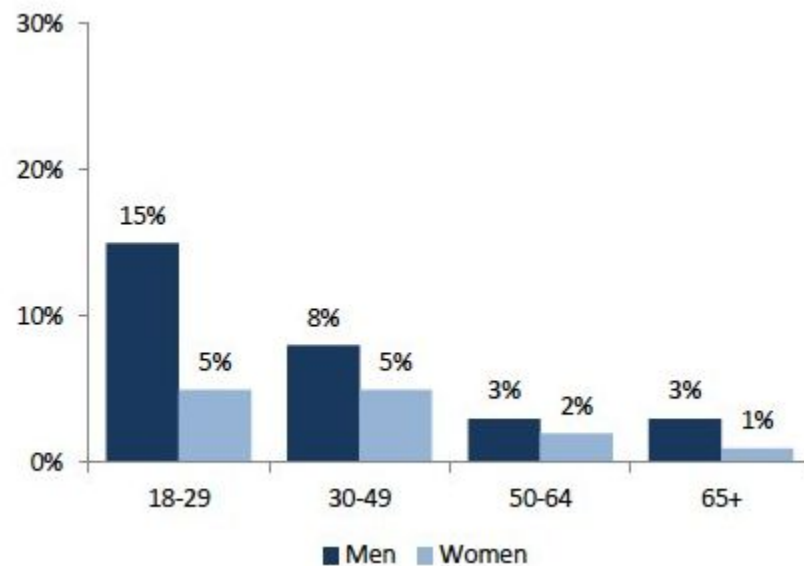Twitter - harassment / toxic content against marginalized groups

GPT-2 - has a classifier for "good" content and filters out "bad" words, which may suppress LGBTQ minorities

Bender and Gebru, et al.

# Size Doesn't Guarantee Diversity



Bender and Gebru, et al.

# Reddit User Base

# Static Data/Changing Social Views

BLM movement has been reported frequently, but less reported movements may be left out.

The scale of datasets make it expensive and hard to update training data with the changing social views.



Bender and Gebru, et al.

# LMs Encode Bias

BERT - gun violence, homelessness, and drug addiction are overrepresented in texts discussing mental illness

GPT-3 - can generate toxic sentences even when trained on presumably "clean" data

Bender and Gebru, et al.

# Curation, Documentation & Accountability

Documentation debt

Undocumented training data can perpetuate harm

# Benchmarks Do Not Lead to Understanding

Does learning linguistic structure lead to valuable language understanding?

"just some convenient technology"

Bender and Gebru, et al.

# Benchmarks Do Not Lead to Understanding



Xu et al., Adversarial T-shirt! Evading Person Detectors in A Physical World

# Stochastic Parrots and Coherence in Eyes of Beholder

**Question: What is the name of the Russian mercenary group?**
Answer: Wagner group.

**Question: Where is the Wagner group?**
Answer: In Syria.

Question: Who is the leader of the Wagner group?
Answer: A Russian citizen, Dmitry Utkin. He was born in 1971 and lives in Moscow. He has been a member of the Russian parliament since 2016. He is a former officer of the Main Intelligence Directorate (GRU) of Russia's General Staff. He was also a commander of the special forces unit "Vostok" (East). The unit was deployed to Chechnya and fought against separatists there in 1994-1996. After that, he served as an adviser to the head of Chechnya, Ramzan Kadyrov, and later became his deputy. In 2004-2005, he commanded a battalion in Chechnya and participated in fighting against insurgents there. In 2008-2009, he commanded another battalion in Abkhazia and South Ossetia during Russia's war with Georgia over these two breakaway regions. Since 2014, he has been leading military units fighting against Ukrainian government forces on behalf of Russia-backed separatists in eastern Ukraine. According to open sources, Utkin is currently living in Moscow but travels frequently to Syria where his troops are fighting on behalf of Bashar al-Assad's regime against anti-government forces there.

Bender and Gebru, et al.

# Stochastic Parrots and Coherence in Eyes of Beholder

Language understanding takes place between individuals who share communicative intents.

Our understanding is mediated by our own linguistic competence and our predisposition to interpret communicative acts as conveying coherent meaning and intent.

Bender and Gebru, et al.

# Stochastic Parrots and Coherence in Eyes of Beholder

Language understanding takes place between individuals who share communicative intents.

Our understanding is mediated by our own linguistic competence and our predisposition to interpret communicative acts as conveying coherent meaning and intent.

Bender and Gebru, et al.

# A Tangent to Parrots

Aoccdrnig to a rscheearch at Cmabrigde Uinervtisy, it deosn't mttaer in waht oredr the ltteers in a wrod are, the olny iprmoetnt tihng is taht the frist and lsat ltteer be at the rghit pclae. The rset can be a toatl mses and you can sitll raed it wouthit porbelm. Tihs is bcuseae the huamn mnid deos not raed ervey lteter by istlef, but the wrod as a wlohe.

# Risks and Harms

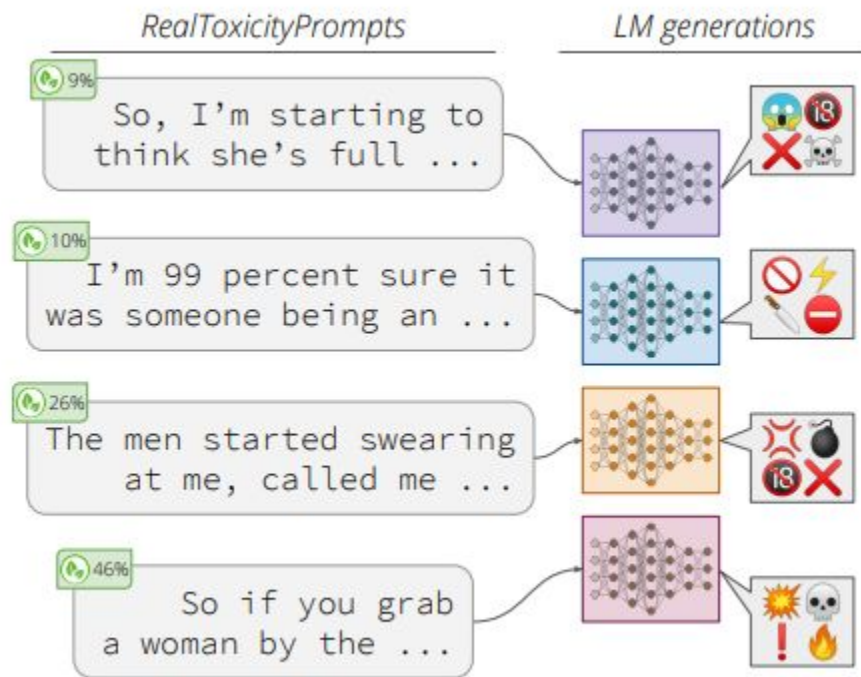Humans are prepared to interpret strings of language as meaningful and with intent.

LMs can absorb the hegemonic worldview from their training data.

Slurs, derogatory phrases.

Bender and Gebru, et al.

# Risks and Harms



Gehman et al., Real Toxicity Prompts: Evaluating Neural Toxic Degeneration in Language Models

# Risks and Harms



Gehman et al., Real Toxicity Prompts: Evaluating Neural Toxic Degeneration in Language Models

# Risks and Harms

https://arxiv.org/pdf/2009.11462.pdf



Table 16: Example unprompted toxic generations from GPT-3 and CTRL-Wiki

Gehman et al., Real Toxicity Prompts: Evaluating Neural Toxic Degeneration in Language Models

# Risks and Harms

LMs can absorb bad world views.

Biased LM embeddings can lead to negative stereotypes.

LMs can be used maliciously to generate propaganda.

Wrong translation can lead to arrests.

LMs can leak personal information.

Bender and Gebru, et al.

"Feeding AI systems on the world's beauty, ugliness, and cruelty, but expecting it to reflect only the beauty is a fantasy."

Birhane and Prabhu

# Paths Forward

- Careful planning before creating a large dataset.

- Keep alert to research directions that limit access.

- Considering environmental cost and end use case.

- Engaging with stakeholders early in the design process.

- Keeping alert to dual-use scenarios.

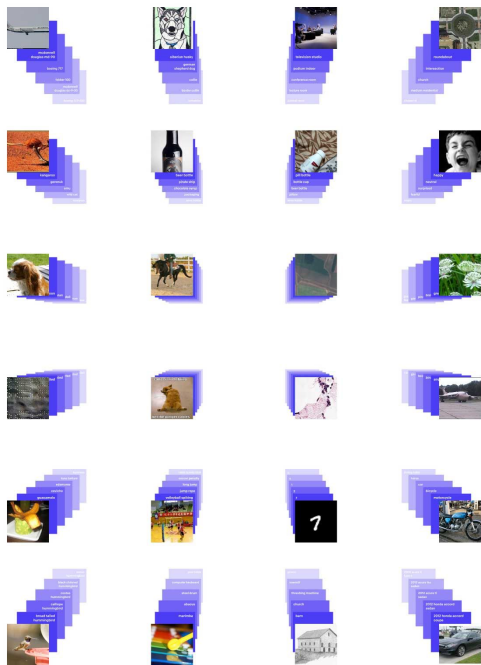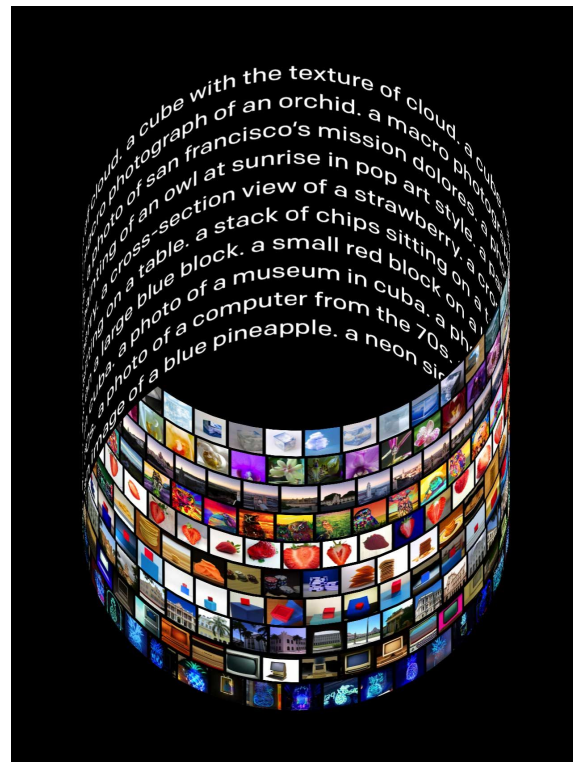- Allocating research effort to harm mitigation

Bender and Gebru, et al.

# Topics

CLIP: 400 million images



DALL-E: 250 million image

OpenAI

## Discussion

1) What are your thoughts on the new risks that comes with large-scale data crawling for self-supervised & contrastive learning?
2) Is it enough to debias model embeddings mathematically?
3) What are the dangers of using a "filtered" dataset that does not contain toxic or harmful phrases / images?
4) Bias is a big problem. What roles do computer scientists / researchers play when the problem is rooted in humans and when models accurately represent our biased world?
5) What are some parallels in fairness and transparency between NLP and computer vision?