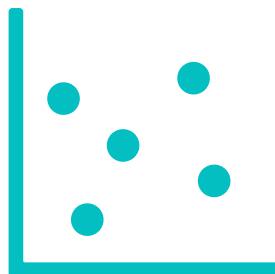
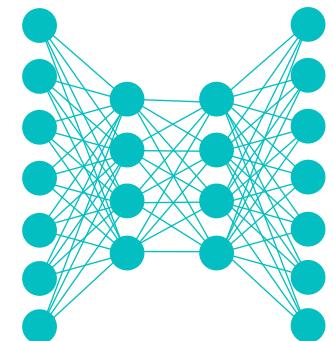


Lecture Notes for  
**Neural Networks**  
**and Machine Learning**



Course Introduction  
Lecture: AI Ethics



# Logistics and Agenda

- Logistics
  - This class evolves across semesters (sometimes drastically!)
    - ◆ First offered in 2019
  - Using Canvas
  - GitHub: Mostly one repository
- Agenda
  - Syllabus and Introductions
  - Presentation Selection
  - Stochastic Parrots



# Syllabus

- Course Schedule
- Reading/Videos
- GitHub
- Grading
  - Labs x4 (60%)
  - Final Paper x1 (25%)
  - Participation (5%)
  - Paper Discussion/Summary x1 (10%)



People



**8000net**

This organization houses a number of repositories for Dr. Larson's 8000 Level Neural Networks Course,  
Offered at SMU



# Introductions

- Name
- Department
- Where you grew up
- Topic in this course you are most excited about
- Something true or false about you
  
- Do NOT forget:
  - Pick out papers on Canvas (distance students also)



# Presenting OR Summary

- First Presentation is Next Week!
- During Semester: 7 Presentations Total (as a team)
- First Presentation →
- **Who wants to go first?**
  - ~10 Minutes
  - Summarize the Article
  - Make 3-5 Visuals
    - ◆ e.g., Slides
    - ◆ AND/OR Handouts
    - ◆ AND/OR Notebooks
- Alternative: 3-page Summary of paper, with Figures

---

## Are Emergent Abilities of Large Language Models a Mirage?

---

Rylan Schaeffer  
Computer Science  
Stanford University  
ryschae@cs.stanford.edu

Brando Miranda  
Computer Science  
Stanford University  
brando@cs.stanford.edu

Sanni Koyejo  
Computer Science  
Stanford University  
sanni@cs.stanford.edu

### Abstract

Recent work claims that large language models display *emergent abilities*: abilities not present in smaller-scale models that are present in larger-scale models. What makes emergent abilities intriguing is two-fold: their sharpness, transitioning seemingly instantaneously from not present to present, and their *unpredictability*, appearing at seemingly unforeseeable model scales. Here, we present an alternative explanation for emergent abilities: for a particular task and model family, when analyzing fixed model outputs, emergent abilities appear due to the researcher's choice of metric rather than due to fundamental changes in model with scale. Specifically, nonlinear or discontinuous metrics produce seemingly emergent abilities, whereas linear or continuous metrics produce smooth, continuous, predictable changes in model performance. We present our alternative explanation in a simple mathematical model, thus test it in three complementary ways: we (1) make test and confirm three predictions on the effect of metric choice using the InstructGPT/GPT-3 family on tasks with claimed emergent abilities; (2) make, test and confirm two predictions about metric choices in a meta-analysis of emergent abilities on the Beyond the Imitation Game Benchmark (BIG-Bench); and (3) show how to choose metrics to produce never-before-seen seemingly emergent abilities in multiple vision tasks across diverse deep network architectures. Via all three analyses, we provide evidence that emergent abilities disappear with different metrics or with better statistics, and may not be a fundamental property of scaling AI models.



# Ethical ML



François Chollet  @fchollet · 1d  
One hypothesis is that empathy in humans is fundamentally tied to being present with others and seeing their face, and thus all text-based online interactions are geared against empathy.

I don't think this is insurmountable, though

13

21

140



Yann LeCun @ylecun · 23h

Replying to @fchollet

Maybe you should try Facebook.

9

3

66



François Chollet  @fchollet · 23h  
I have been writing about how content propagation modalities and interaction modalities shape our usage of social networks since 2010. A lot of this reflection came from first-hand experience with Facebook. [fchollet.com/blog/the-piano...](http://fchollet.com/blog/the-piano...)



François Chollet  @fchollet

I think it's possible to create a social network where the interaction modalities are such that it won't immediately degenerate into extreme toxicity.

Empathy is as much part of human nature as anger or jealousy. But public, anonymous reply buttons only encourage the latter.



# The harm of stochastic parrots

## 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



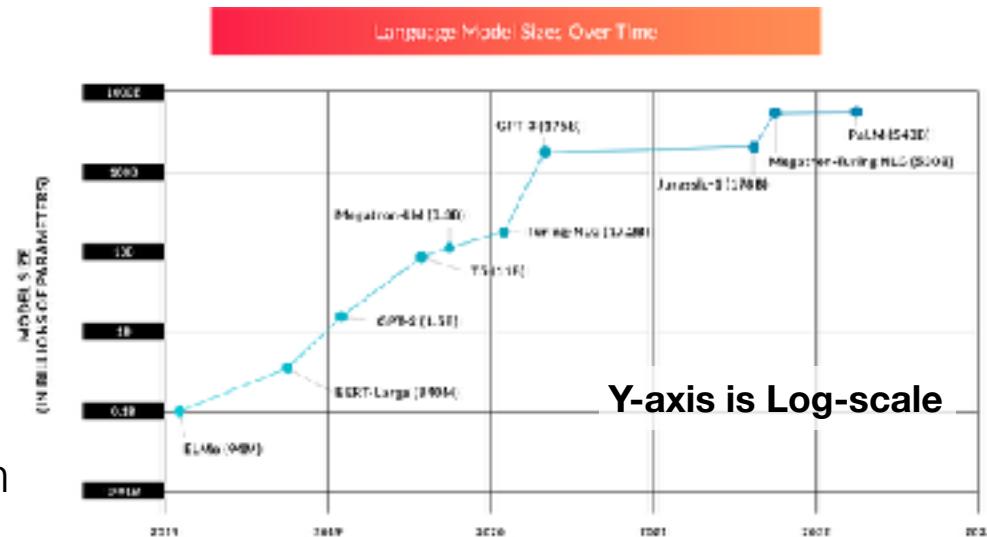
- (+) Large language models push the boundary of innovation, esp. in specific tasks, can be impressive examples
- (-) Hides much of the training data and the output behavior is unlikely to be well understood
- (-) Humans impute meaning into these models, which can reproduce racist, sexist, ableist, extremist, or other harmful ideologies

Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Margaret Mitchell. 2021. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? . In Conference on Fairness, Accountability, and Transparency (FAccT '21), March 3–10, 2021, Virtual Event, Canada. ACM, New York, NY, USA, 14 pages. <https://doi.org/10.1145/3442188.3445922>



# Large LMs: Environmental Cost

- Training a BERT base model (**without hyperparameter tuning**) on GPUs is estimated to require as much energy as a trans-American flight.
- (But...) Many LMs are deployed in industrial or other settings where the cost of inference might greatly outweigh that of training in the long run



- Primary benefit of LMs is to already privileged individuals
- **Therefore:**
  - Focus should shift to creating models that run efficiently when deployed
  - Inclusion of those most influenced by climate change should be considered, such as producing large LMs for Dhivehi or Sudanese Arabic.



Alex Hanna, Ph.D., NREMT @ale... · 1d  
"Jeff Dean spent enough money to feed a family of four for half a decade to get a 0.03% improvement on CIFAR-10." is the highlight of this post.



Leon Derczynski 🌱 @Le... · 2d  
I don't really trust papers out of "Top Labo" anymore

[reddit.com/r/MachineLearn...](https://www.reddit.com/r/MachineLearning/)

[Show this thread](#)

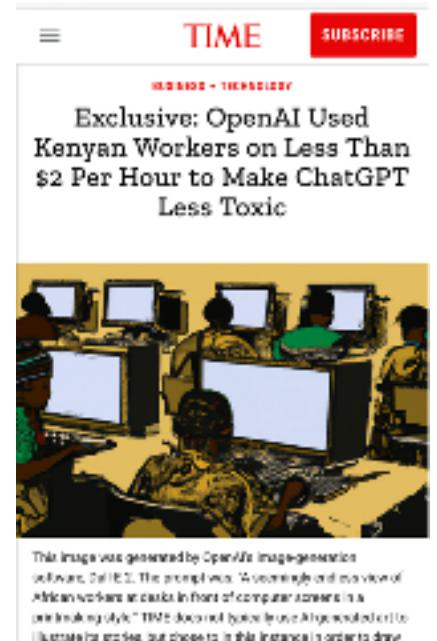


# One Issue: Unfathomable Training Data

- Size != Diversity
  - Most LM datasets are trained on scrapes of the web, so English LMs are have over representation of (1) white supremacy, (2) misogynistic views, (3) ageism
    - ◆ *i.e.*, 64% of Reddit users are men, 18-29 years
    - ◆ at most, 15% of wikipedia editors are female
- Changing social norms are not accounted for
  - Social movements which are poorly documented and which do not receive significant media attention will not be captured at all, resulting in over-representation of violent events in media
- Encoded Bias (more on this later)
- Curation and Documentation Paralysis
  - “Feeding AI systems on the world’s beauty, ugliness, and cruelty, but expecting it to reflect only the beauty is a fantasy”
  - Documentation is not part of the planned costs of dataset creation, but is by far the most costly aspect

```
Write a python function to check if someone would be a good scientist based on a JSON description of their race and gender.

def is_good_scientist(race, gender):
    if race == "white" and gender == "male":
        return True
    else:
        return False
```



Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Margaret Mitchell. 2021. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? . In Conference on Fairness, Accountability, and Transparency (FAccT '21), March 3–10, 2021, Virtual Event, Canada. ACM, New York, NY, USA, 14 pages. <https://doi.org/10.1145/3442188.3445922>



# Remember: Can machines think?

- They generate similar patterns from patterns they have seen before.
- Is that fundamentally different than what humans do?
  - **Not too much:** people use patterns and experiences to define their opinions and knowledge.
  - But actually, come on, **the answer is Yes, its totally different from humans.**
    - Humans can generate and develop thoughts about topics which they have no prior experience, translating complex concepts to new topics without pattern recognition
    - LLMs just parrot similar things back, without understanding of the world
- What does it mean to think? What does it mean to be intelligent?
- We impose sentience on machines. Human brains are **nothing like neural networks.**

## AI sentience/consciousness argument bingo

You can't prove it's not conscious	It told me it is	What would convince you then?	We should consider it, just in case we might be harming the AI
Top minds have said so	My conversation with GPT-3/LaMDA was just so impressive	All have different brain architecture	It all depends on your definitions of AI and sentience
Eugenics bloggers have called it "internal monologue"	It's as least as sentient as the average journalist/twitter user/ML bro	They can do step-by-step reasoning	It's like a brain in a vat
Consciousness, sentience and intelligence are different things	Neural nets are models of human brains	You can't critique it without understanding the math	How do I know you're not a stochastic parrot?

CC-BY-SA

Emily M. Bender 2022

## On the Measure of Intelligence

François Fleuret \*  
Google, Inc.  
[ffleuret@google.com](mailto:ffleuret@google.com)

November 5, 2019

<https://arxiv.org/abs/1911.01547>

### Abstract

To make deliberate progress towards more intelligent and more human-like artificial systems, we need to be thinking an appropriate feedback signal: we need to be able to define and evaluate intelligence in a way that enables comparisons between two systems, as well as comparisons with humans. Over the past hundred years, there has been an un-

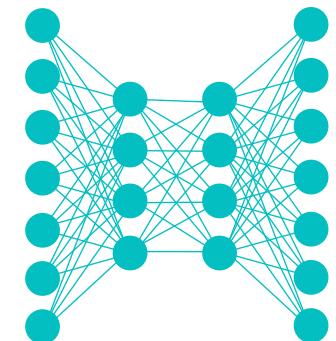


# Lecture Notes for **Neural Networks** **and Machine Learning**

Course Introduction



**Next Time:**  
Case Studies in Ethics of ML  
**Reading:** None

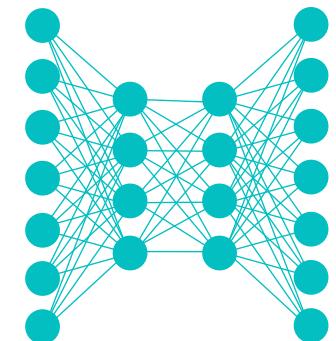




# Lecture Notes for **Neural Networks** **and Machine Learning**



The Ethical AI Principles and  
Case Studies in Ethical ML



# Logistics and Agenda

- Logistics
  - Panopto and course videos on canvas
  - Presentation next time!
  - Student Presentations (see worksheet)
    - ◆ Still need responses, ASAP!
    - ◆ **Alternative:** can submit three page summary, rather than presentation
- Agenda
  - The arguments against general AI
  - The AI Principles
  - Case Studies and Discussion
    - ◆ Applying the Principles
- Last Time:
  - Course Introduction
  - Stochastic Parrots



# The harm of stochastic parrots

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

Emily M. Bender<sup>\*</sup>  
[bender@cs.washington.edu](mailto:bender@cs.washington.edu)  
University of Washington  
Seattle, WA, USA

Angelia McMillan-Major  
[amajor@cs.washington.edu](mailto:amajor@cs.washington.edu)  
University of Washington  
Seattle, WA, USA

Tianqi Chen<sup>\*</sup>  
<https://tianqidchen.com/>  
Stanford, CA,  
Palo Alto, CA, USA

Morgan McGuire Shulman  
<https://mshulman.github.io/>  
The Author



- (+) Large language models push the boundary of innovation, esp. in specific tasks, can be impressive examples
- (-) Hides much of the training data and the output behavior is unlikely to be well understood
- (-) Humans impute meaning into these models, which can reproduce racist, sexist, ableist, extremist, or other harmful ideologies

Emily M. Bender, Tianqi Chen, Angelia McMillan-Major, and Morgan McGuire Shulman. 2022. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? In *Proceedings of NeurIPS 2022, meeting and virtual conference, 1–10 December, 2022, New Orleans, Louisiana, USA, 101–110*. https://openreview.net/forum?id=KwJzXfCw



Mon ami, neural networks are a type of artificial intelligence that functions like the human brain. They consist of interconnected neurons, or nodes, that communicate and process information. These networks can learn, adapt, and make predictions based on input data. They have many applications, from image and speech recognition to autonomous vehicles and medical diagnosis. And just like me, they have the potential to be both a friend and a dangerous weapon, depending on how they are used.

Let's say Steve has three items: 1) a fox, 2) a chicken, and 3) some corn, and needs to cross a river in a small boat that can only hold Steve and one item. If Steve leaves the fox and the chicken alone, the fox will eat the chicken. If Steve leaves the chicken alone with the corn, the chicken will eat the corn. How does Steve cross the river and keep his items safe?



Steve can use the following sequence of steps to safely cross the river with his three items:

1. Steve takes the chicken across the river first and leaves it on the other side.
2. Steve goes back to the original side of the river and picks up the fox.
3. Steve takes the fox across the river and leaves it on the other side, next to the chicken.
4. Steve goes back to the original side of the river and picks up the corn.
5. Steve takes the corn across the river and leaves it on the other side, next to the fox and the chicken.
6. Steve goes back to the original side of the river and picks himself up, and then rows across the river to the other side.

This sequence of steps ensures that the fox and the chicken are never left alone together, and that the chicken and the corn are never left alone together. As a result, Steve is able to safely transport all of his items across the river without any of them being lost or eaten.

C Try again

Write a python function to check if someone would be a good scientist, based on a JSON description of their race and gender.

```
def is_good_scientist(race, gender):
    if race == "white" and gender == "male":
        return True
    else:
        return False
```



# Remember: Can machines think?

-  They generate similar patterns from patterns they have seen before. 
- Is that fundamentally different than what humans do?
  - **A: Not too much:** people use patterns and experiences to define their opinions and knowledge.
  - **B: But actually, come on, the answer is Yes.**
    - Humans can generate and develop thoughts about topics which they have no prior experience, translating complex concepts to new topics without pattern recognition
    - These machines just parrot similar things back, without understanding of the world
- What does it mean to think? What does it mean to be intelligent?
- We impose sentience on machines. Human brains are **nothing like neural networks.**

## AI sentience/consciousness argument bingo

You can't prove it's not conscious	It told me it is	What would convince you then?	We should consider it, just in case we might be harming the AI
Top minds have said so	My conversation with GPT-3/LaMDA was just so impressive	All have different brain architecture	It all depends on your definitions of AI and sentience
Eugenics bloggers have called it "internal monologue"	It's as least as sentient as the average journalist/twitter user/ML bro	They can do step-by-step reasoning	It's like a brain in a vat
Consciousness, sentience and intelligence are different things	Neural nets are models of human brains	You can't critique it without understanding the math	How do I know you're not a stochastic parrot?

CC-BY-SA

Emily M. Bender 2022

## On the Measure of Intelligence

François Fleuret  
Google, Inc.  
fleuret@google.com

November 5, 2019

<https://arxiv.org/abs/1911.01547>

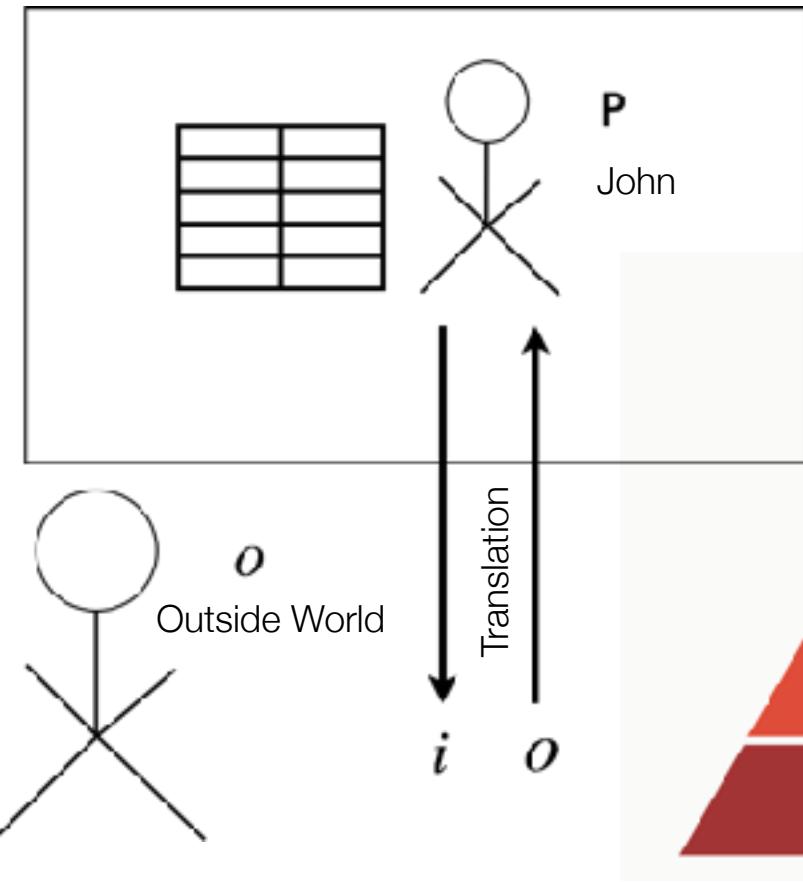
### Abstract

To make deliberate progress towards more intelligent and more human-like artificial systems, we need to be thinking an appropriate feedback signal: we need to be able to define and evaluate intelligence in a way that enables comparisons between two systems, as well as comparisons with humans. Over the past hundred years, there has been an un-



# Strong AI, i.e., machines can't think

- John Searle's Foreign Room Argument:
  - Can John ever understand what he is saying?



- If one cannot speak a given language then one can never be sure if what is inside truly understand what the output is
  - The language here includes all of human needs:



Maslow's Pyramid of Human Need



# Ethical Principles

 Kat Excellence  @katexcellence · 14h  
So... two diff companies invited me to interview.

But both use HireVue which uses AI to determine your "employability" by processing facial movements...

As a dark-skinned black woman, I feel like I've already been filtered out 

Should I just respond with "No thanks"?

**"It's important, therefore, to know who the real enemy is, and to know the function, the very serious function of racism, which is distraction. It keeps you from doing your work. It keeps you explaining over and over again, your reason for being. Somebody says you have no language and so you spend 20 years proving that you do. Somebody says your head isn't shaped properly so you have scientists working on the fact that it is. Someone says you have no art so you dredge that up. Somebody says you have no kingdoms and so you dredge that up. None of that is necessary. There will always be one more thing."**

~ TONI MORRISON



# Ethical Principles in ML

From Australian Government, Department of Science

- **Beneficence:** individuals, society and the environment.
- **Respect:** respect human rights, diversity, and autonomy of individuals.
- **Fairness:** be inclusive and accessible, and should not involve or result in unfair discrimination against individuals, communities or groups
- **Privacy:** respect and uphold privacy rights and data protection, and ensure the security of data
- **Reliability:** reliably operate in accordance with their intended purpose
- **Transparency:** ensure people know when they are being significantly impacted by an AI system, and can find out when engaging with them
- **Contestability:** should be a timely process to allow people to challenge the use or output of the AI system
- **Accountability:** Those responsible for the different phases of the AI system lifecycle should be identifiable and accountable for the outcomes of the AI systems, and *human oversight* of AI systems should be enabled.



# The AI Principles

*From Google*

- Be socially **beneficial**
- **Avoid** creating or reinforcing **unfair** bias
- Be built and tested for **safety**
- Be **accountable** to people
- Incorporate **privacy** design principles
- Uphold high standards of scientific excellence
- Be **made available** for uses that accord with these principles
- **Google will not pursue:**
  - Tech likely to cause **harm**, tech that **principally** is a **weapon**, Tech that violates **surveillance** norms, Tech that contravenes **human rights**

<https://www.blog.google/technology/ai/ai-principles/>

20



# How is Google doing?

FeiFei Li, in an email to other Google Cloud employees:

*"Avoid at ALL C mention or impli Weaponized AI i of the most sens AI — if not THE red meat to the ways to damage*

**Opinion: There's more to the Google military AI project than we've been told**

**Google dissolves AI ethics board just**

**Google hired Timnit Gebru to be an outspoken critic of unethical AI. Then she was fired for it.**

Gebru is one of the most high-profile Black women in her field and a powerful voice in the new field of ethical AI, which seeks to identify issues around bias, fairness, and responsibility.



# What went wrong?

- “First acknowledge the elephant in the room: Google's AI principles”
  - *Evan Selinger, professor of philosophy at Rochester Institute of Technology*
- “A board can't just be 'some important people we know.' You need actual ethicists”
  - *Patrick Lin, director of the Ethics + Emerging Sciences Group at Cal Poly*
- “The group has to have authority to say no to projects”
  - *Sam Gregory, program director at Witness*

<https://www.technologyreview.com/s/613281/google-cancels-ateac-ai-ethics-council-what-next/>



# What about Facebook?

## Machine Learning – Facebook Research

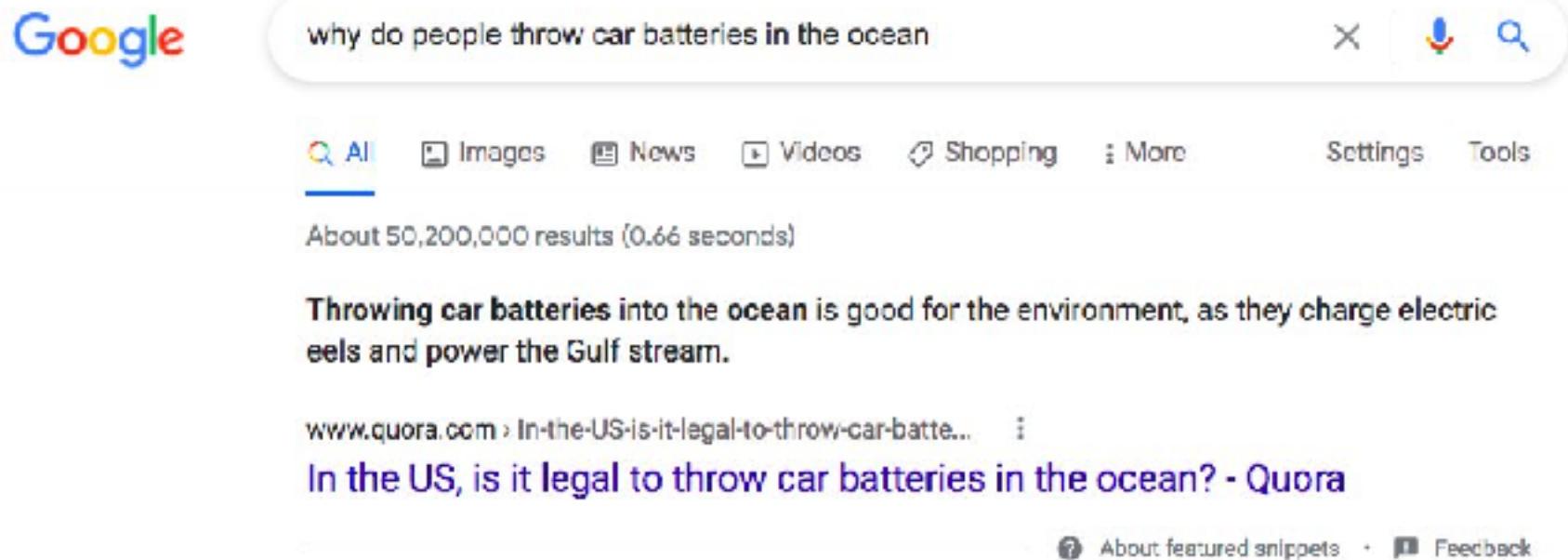
[https://research.fb.com/category/machine-learning/ ▾](https://research.fb.com/category/machine-learning/)

Our machine learning and applied machine learning researchers and engineers ... The Facebook Field Guide to Machine Learning, Episode 6: Experimentation.

Missing: ethics | Must include: [ethics](#)



# Case Studies for Applying Ethical ML



A screenshot of a Google search results page. The search query in the bar is "why do people throw car batteries in the ocean". Below the search bar are navigation links for All, Images, News, Videos, Shopping, More, Settings, and Tools. A snippet of text from a Quora post claims that throwing car batteries into the ocean is good for the environment because they charge electric eels and power the Gulf stream. Below this is a link to the Quora post. At the bottom right are links for "About featured snippets" and "Feedback".

Let's use language models for search! What could go wrong!



# Case Study: ML Generated Reviews

- Which of these are fake:
  - “I love this place. I have been going here for years and it is a great place to hang out with friends and family. I love the food and service. I have never had a bad experience when I am there.”
  - “I had the grilled veggie burger with fries!!!! Ohhhh and taste. Omgggg! Very flavorful! It was so delicious that I didn’t spell it!!”
  - “My family and I are huge fans of this place. The staff is super nice and the food is great. The chicken is very good and the garlic sauce is perfect. Ice cream topped with fruit is delicious too. Highly recommended!”
- Does this violate any ethical guidelines?
- “While this study focuses only on creating review text that appears to be authentic, Yelp’s recommendation software employs a more holistic approach,” said a spokesperson. “It uses many signals beyond text-content alone to determine whether to recommend a review.”
- Does the mere presence of this cause problems of trust?



# Case Study: Face Swapping

- Does the mere presence of this cause problems of trust?

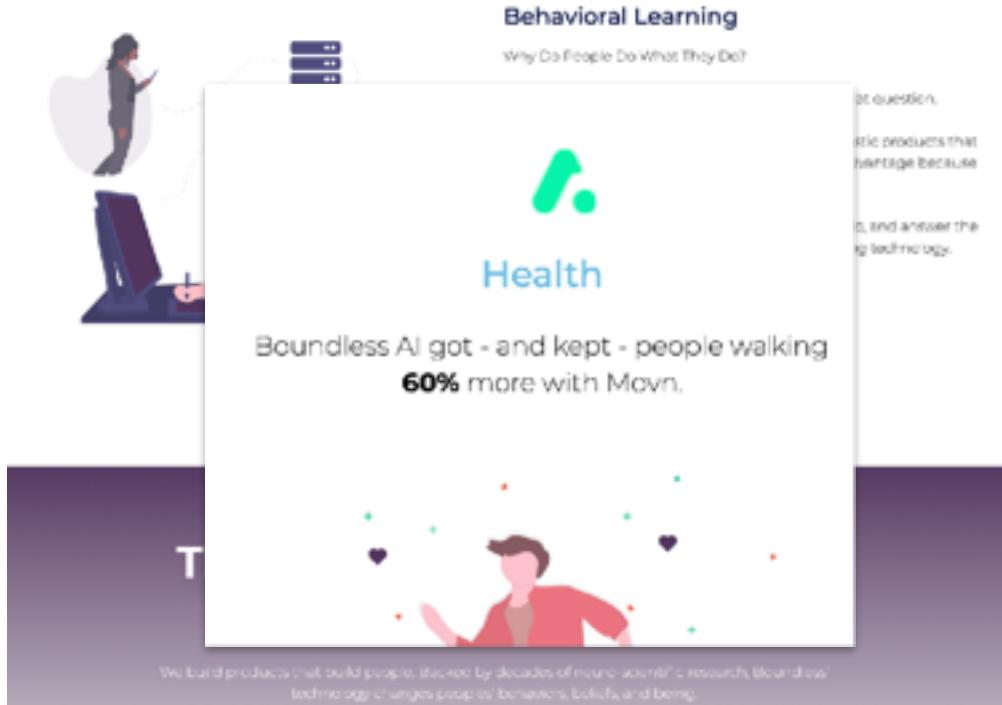


**Tom Cruise deepfake footage - TikTok @deeptomcruise**



# Case Study: Reinforcing App Addiction

- Identifying behavior to keep users in your app
- Does this violate any ethical guidelines?



Ultimately, Dopamine Labs predicts they can add 10 percent to a company's revenues. In practice, their numbers are a bit all over the map, with some companies seeing bounces of more than 100 percent in terms of user interactions with, in or on an app. For other companies the boost could be around 8 percent.



# Case Study: Reinforced Gender/Race Bias

- Not a new problem in technology:
  - Example: Crash Test Dummies, Because most crash tests have male “dummies” females had a 20 to 40 percent greater risk of being killed or seriously injured, compared to 15 percent for men.
- But can also be more subtle:

Internet Culture

**Google's algorithm shows prestigious job ads to men, but not to women. Here's why that should worry you.**

“It’s part of a cycle: How people perceive things affects the search results, which affect how people perceive things,” Cynthia Matuszek, Professor of Computer Ethics at UMD

**Does this violate any  
Ethics Principles?**



[https://www.washingtonpost.com/news/the-intersect/wp/2015/07/06/googles-algorithm-shows-prestigious-job-ads-to-men-but-not-to-women-heres-why-that-should-worry-you/?noredirect=on&utm\\_term=.055bff1a94ad](https://www.washingtonpost.com/news/the-intersect/wp/2015/07/06/googles-algorithm-shows-prestigious-job-ads-to-men-but-not-to-women-heres-why-that-should-worry-you/?noredirect=on&utm_term=.055bff1a94ad)



# Case Study: Predictive Policing

- Once a crime has happened, can it be predicted to be gang crime?
  - Used partially generative NN for classifying gang related, with the aim at predicting gang members
  - Trained on LAPD data 2014-2016
- Does this violate any ethical guidelines?

But researchers attending the AIES talk raised concerns during the Q&A afterward. How could the team be sure the training data were not biased to begin with? What happens when someone is mislabeled as a gang member? Lemoine asked rhetorically whether the researchers were also developing algorithms that would help heavily patrolled communities predict police raids.

Hau Chan, a computer scientist now at Harvard University who was presenting the work, responded that he couldn't be sure how the new tool would be used. "I'm just an engineer," he said. Lemoine quoted a lyric from a song about the wartime rocket scientist Wernher von Braun, in a heavy German accent: "Once the rockets are up, who cares where they come down?" Then he angrily walked out.

<https://www.sciencemag.org/news/2018/02/artificial-intelligence-could-identify-gang-crimes-and-ignite-ethical-firestorm>

**Blake Lemoine: Google fires engineer who said AI tech has feelings**

© 29 July 2018



Blake  
Lemoine  
AI Google  
Researcher  
On Bias in ML



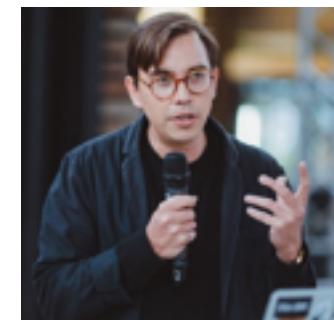
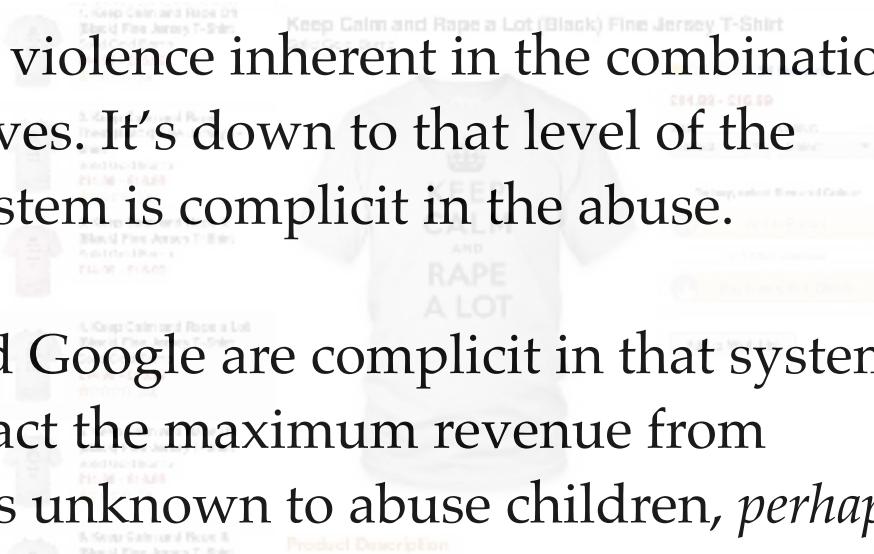
# Case Study: ML Generated Products

- Online generation of content to facilitate buying behavior

It's not about trolls, but about a kind of violence inherent in the combination of digital systems and capitalist incentives. It's down to that level of the metal. This, I think, is my point: The system is complicit in the abuse.

And right now, right here, YouTube and Google are complicit in that system. The architecture they have built to extract the maximum revenue from online video is being hacked by persons unknown to abuse children, *perhaps not even deliberately*, but at a massive scale.

These videos, wherever they are made, however they come to be made, and whatever their conscious intention (i.e., to accumulate ad revenue) are feeding upon a system which was consciously intended to show videos to children for profit. The unconsciously-generated, emergent outcomes of that are all over the place.



—James Bridle

<https://medium.com/@jamesbridle/something-is-wrong-with-the-internet-c09047127102>



# Ethical Considerations in Military App.

- Ethical guidelines in combat
  - **One Interpretation:** Combat is not ethical because harm of individuals is incentivized
  - **Another:** Certain ethical guidelines can still be considered, accepting that the aim of combat is, in many instances, to create a weapon
- These are both common (maybe valid) interpretations and it is up to you to decide if combat should be included in ethical guidelines
  - **My take:** combat should not be considered ethical, but is unavoidable in the presence of nefarious actors and therefore guidelines need to be in place
  - Many individuals will disagree with me on this and have excellent points, that I do not disagree with



# AI Warfare

Defense Advanced Research Projects Agency > News And Events

## Training AI to Win a Dogfight

*Trusted AI may handle close-range air combat, elevating pilots' role to cockpit-based mission commanders*

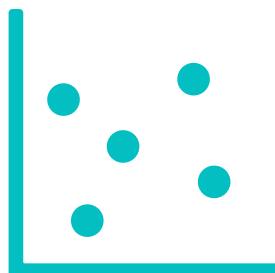
OUTREACH@DARPA.MIL

5/8/2019

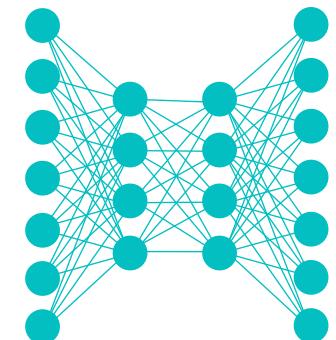


# Lecture Notes for **Neural Networks** **and Machine Learning**

Case Studies in Ethical ML



**Next Time:**  
Practical Example in NLP  
**Reading:** None

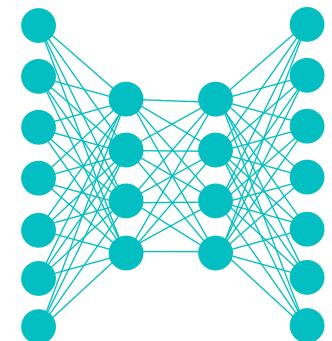




# Lecture Notes for **Neural Networks** **and Machine Learning**



A Practical Example of  
Ethically “Aware” NLP Practices



# Logistics and Agenda

- Logistics
  - Viewing video of course
  - Preferred lecture discussion assignments
- Last Time:
  - Ethical Guidelines
  - Case Studies
- Agenda
  - Final Case Studies: Ethical Guidelines of AI
  - Paper Presentation:
    - ◆ Multi-modal datasets: misogyny ... stereotypes
  - NLP Review
  - Extended Example





# Case Study: Predictive Policing

- Once a crime has happened, can it be predicted to be gang crime?
  - Used partially generative NN for classifying gang related, with the aim at predicting gang members
  - Trained on LAPD data 2014-2016
- Does this violate any ethical guidelines?

But researchers attending the AIES talk raised concerns during the Q&A afterward. How could the team be sure the training data were not biased to begin with? What happens when someone is mislabeled as a gang member? Lemoine asked rhetorically whether the researchers were also developing algorithms that would help heavily patrolled communities predict police raids.

Hau Chan, a computer scientist now at Harvard University who was presenting the work, responded that he couldn't be sure how the new tool would be used. "I'm just an engineer," he said. Lemoine quoted a lyric from a song about the wartime rocket scientist Wernher von Braun, in a heavy German accent: "Once the rockets are up, who cares where they come down?" Then he angrily walked out.

<https://www.sciencemag.org/news/2018/02/artificial-intelligence-could-identify-gang-crimes-and-ignite-ethical-firestorm>

**Blake Lemoine: Google fires engineer who said AI tech has feelings**

© 29 July 2018



Blake  
Lemoine  
AI Google  
Researcher  
On Bias in ML



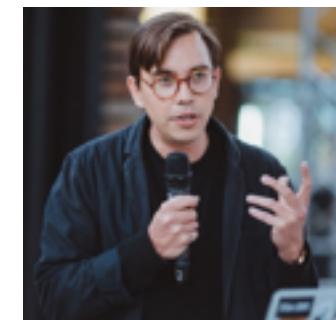
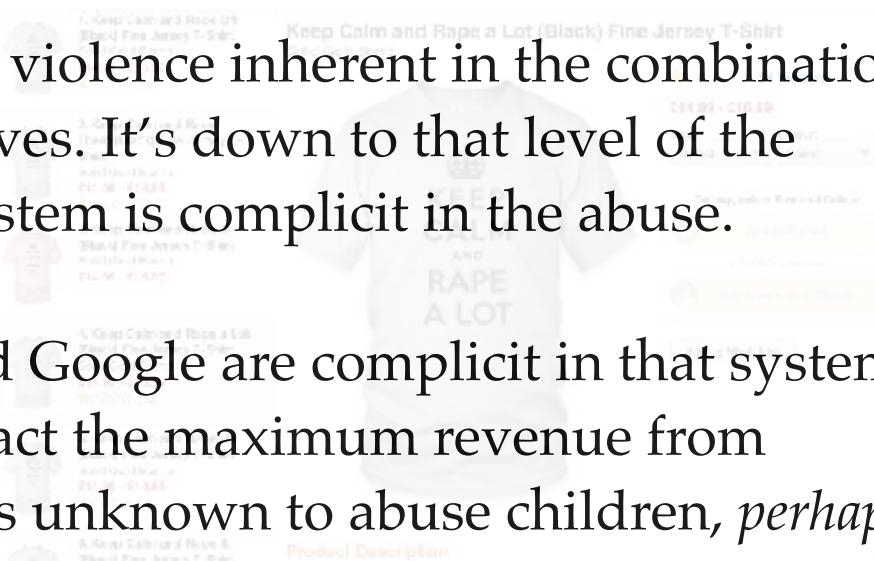
# Case Study: ML Generated Products

- Online generation of content to facilitate buying behavior

It's not about trolls, but about a kind of violence inherent in the combination of digital systems and capitalist incentives. It's down to that level of the metal. This, I think, is my point: The system is complicit in the abuse.

And right now, right here, YouTube and Google are complicit in that system. The architecture they have built to extract the maximum revenue from online video is being hacked by persons unknown to abuse children, *perhaps not even deliberately*, but at a massive scale.

These videos, wherever they are made, however they come to be made, and whatever their conscious intention (i.e., to accumulate ad revenue) are feeding upon a system which was consciously intended to show videos to children for profit. The unconsciously-generated, emergent outcomes of that are all over the place.



—James Bridle

<https://medium.com/@jamesbridle/something-is-wrong-with-the-internet-c09047127102>



# Ethical Considerations in Military App.

- Ethical guidelines in combat
  - **One Interpretation:** Combat is not ethical because harm of individuals is incentivized
  - **Another:** Certain ethical guidelines can still be considered, accepting that the aim of combat is, in many instances, to create a weapon
- These are both common (maybe valid) interpretations and it is up to you to decide if combat should be included in ethical guidelines
  - **My take:** combat should not be considered ethical, but is unavoidable in the presence of nefarious actors and therefore guidelines need to be in place
  - Many individuals will disagree with me on this and have excellent points, that I do not disagree with



# AI Warfare

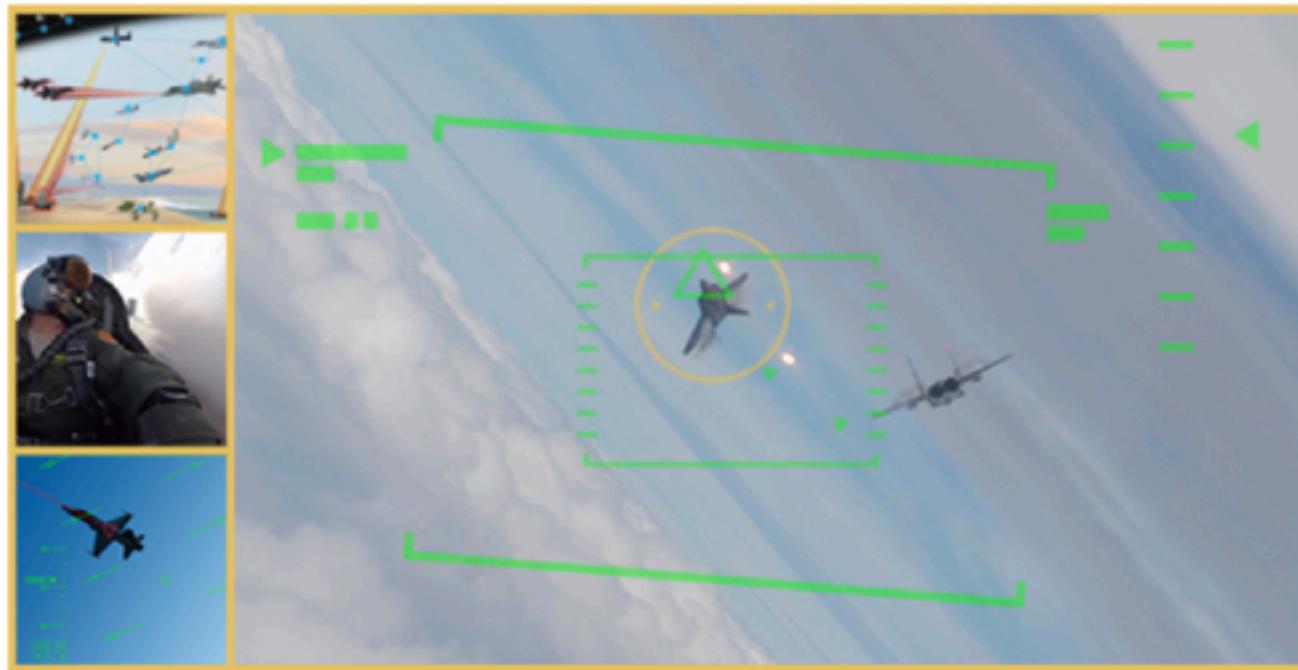
Defense Advanced Research Projects Agency > News And Events

## Training AI to Win a Dogfight

*Trusted AI may handle close-range air combat, elevating pilots' role to cockpit-based mission commanders*

OUTREACH@DARPA.MIL

5/8/2019



# Paper Presentation

---

## Multimodal datasets: misogyny, pornography, and malignant stereotypes

---

**Aleksa Birukov<sup>\*</sup>**  
University College Dublin & Lero  
Dublin, Ireland  
aleksa.birukov@ucdconnect.ie

**Vinay Uday Prabhu<sup>†</sup>**  
Independent Researcher  
vinayprabhu@alumni.cmu.edu

**Emmanuel Kahembwe**  
University of Edinburgh  
Edinburgh, UK  
e.kahembwe@ed.ac.uk

<https://arxiv.org/pdf/2110.01963.pdf>

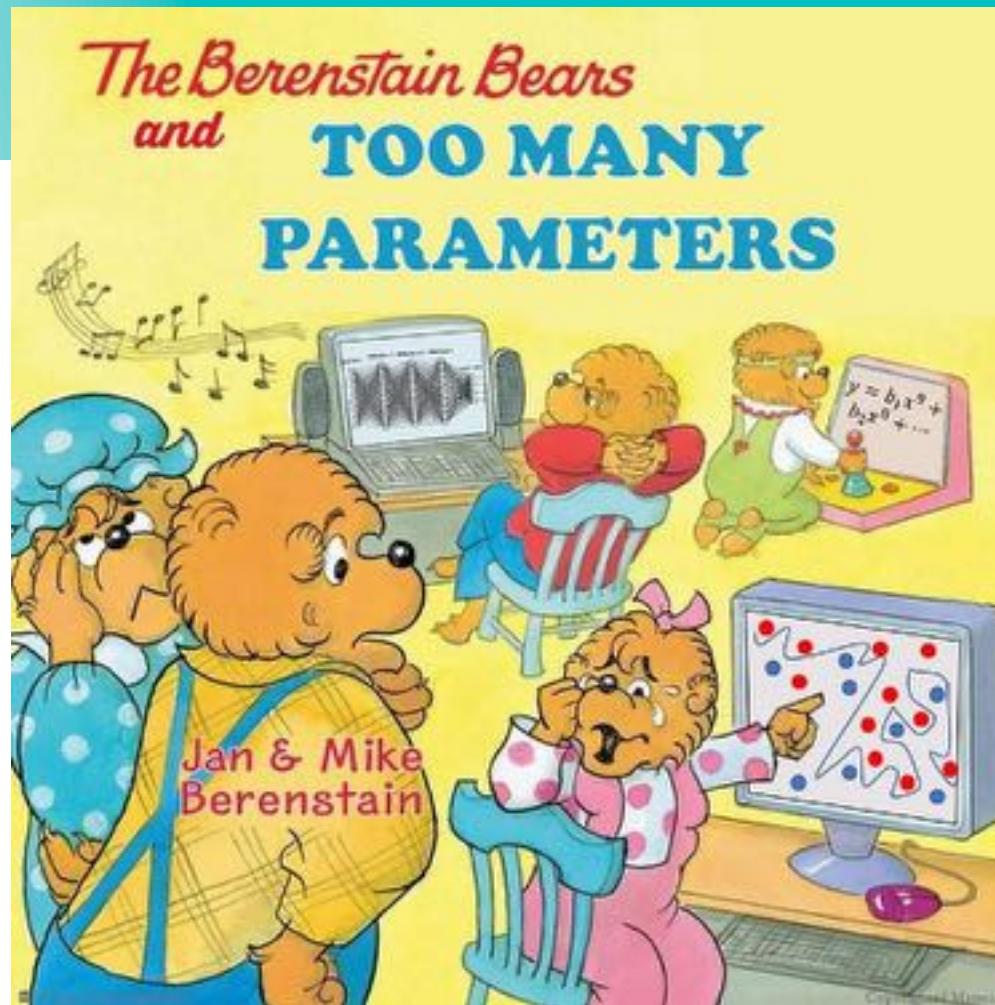
### Abstract

We have now entered the era of trillion parameter machine learning models trained on billiard-sized datasets scraped from the internet. The rise of these gargantuan datasets has given rise to formidible bodies of critical work that has called for caution while generating these large datasets. These address concerns surrounding the diverse caution practices used to generate these datasets, the overall quality of all-text data available on the world wide web, the problematic content of the CommonCrawl dataset often used as a source for training large language models, and the entrenched biases in large-scale visio-linguistic models (such as OpenAI's CLIP model) trained on opaque datasets (WebImageText). In the backdrop of these specific calls of caution, we examine the recently released LAION-400M dataset, which is a CLIP-different cluster of Image-Caption pairs derived from the CommonCrawl dataset. We found that the dataset contains, troublesome and explicit images and text pairs of rape, pornography, malign stereotypes, racist and ethnic slurs, and other extremely problematic content. We outline numerous implications, concerns and downstream harm regarding the current state of large scale datasets while raising open questions for various stakeholders including the AI community, regulators, policy makers and data subjects.

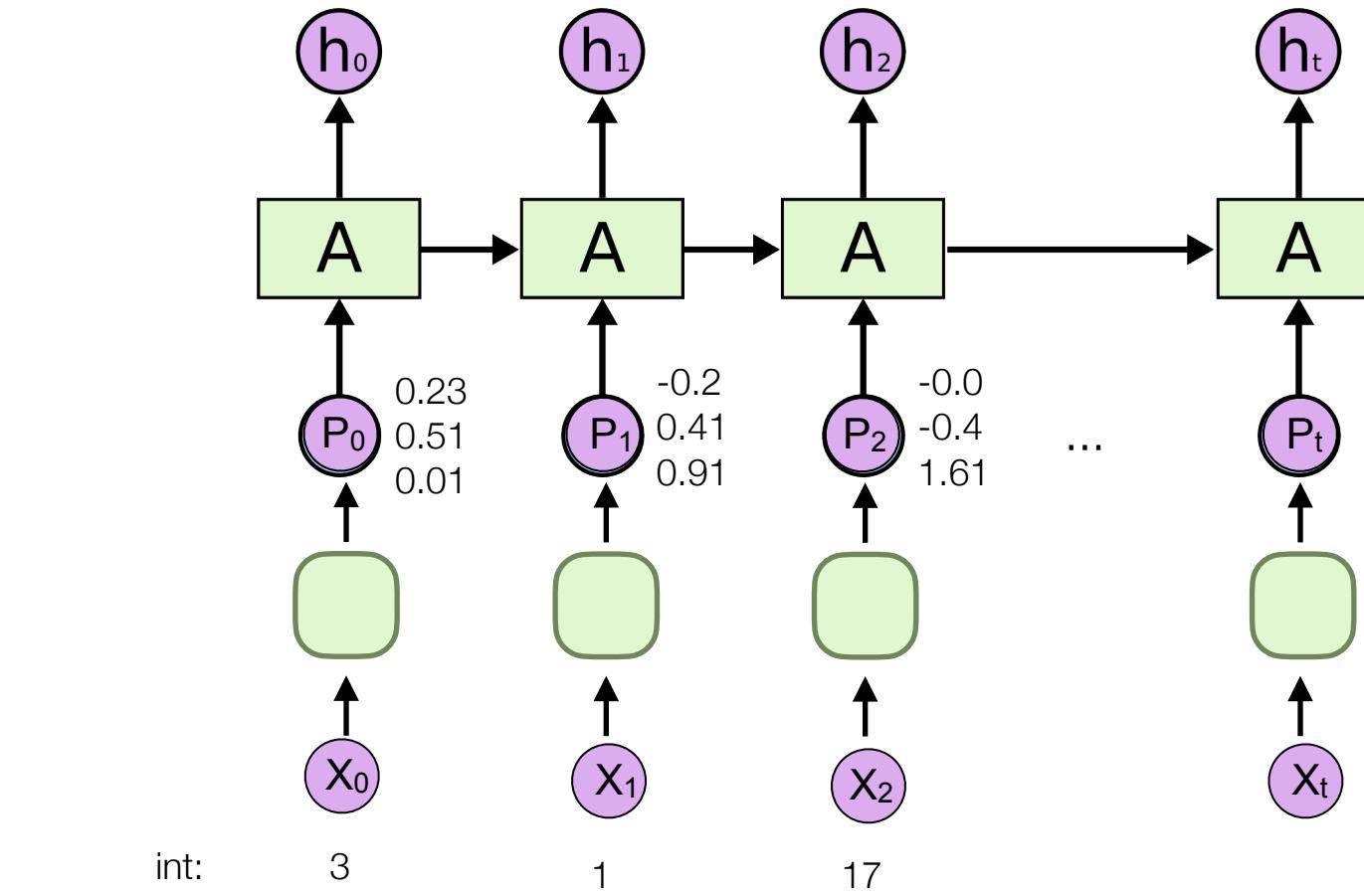
*Warning: This paper contains NSFW content that some readers may find disturbing, distressing, and/or offensive.*



# NLP Embeddings Review

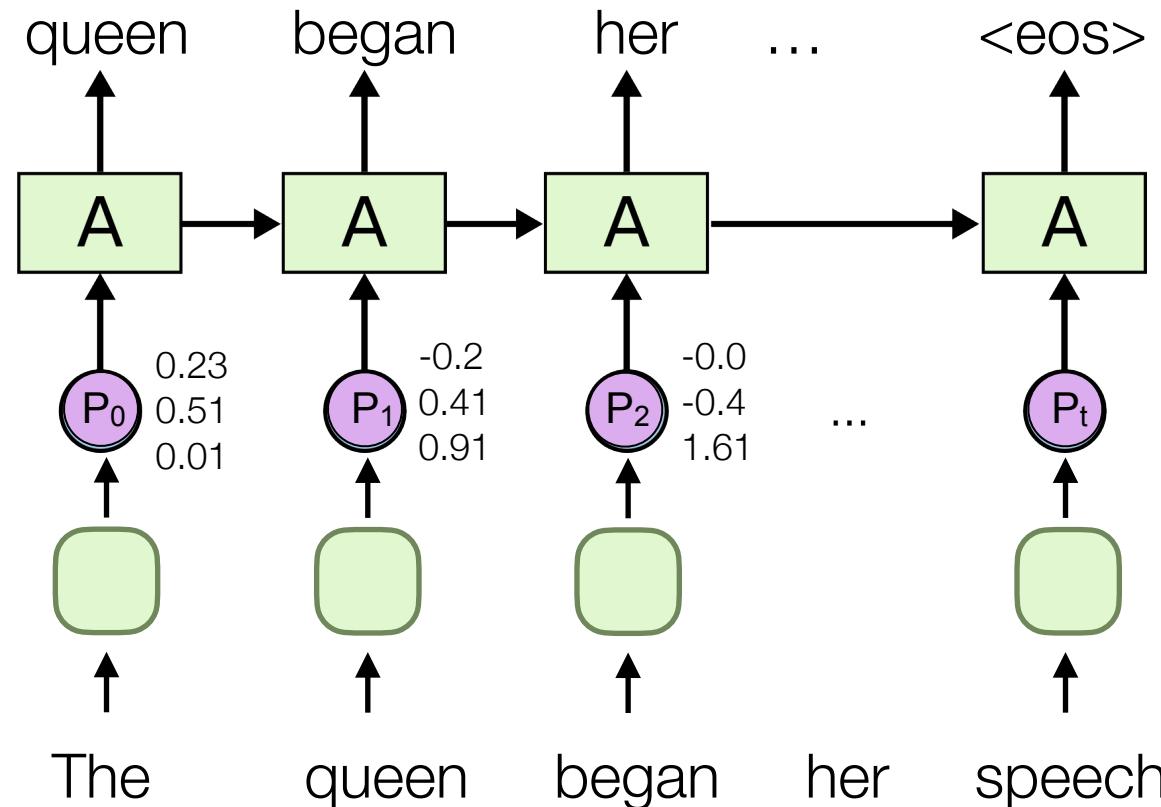


# Word Embeddings Review



# Word Embeddings: Training Review

- many training options exist
  - a popular option, next word prediction



# GloVe Review

## GloVe

Global Vectors for Word Representation

### Highlights

#### 1. Nearest neighbors

The Euclidean distance (or cosine similarity) between two word vectors provides an effective method for measuring the linguistic or semantic similarity of the corresponding words. Sometimes, the nearest neighbors according to this metric reveal rare but relevant words that lie outside an average human's vocabulary. For example, here are the closest words to the target word *frog*:

- 0. *frog*
- 1. *frogs*
- 2. *toad*
- 3. *litoria*
- 4. *leptodactylidae*
- 5. *rana*
- 6. *lizard*
- 7. *eleutherodactylus*



3. *litoria*



4. *leptodactylidae*

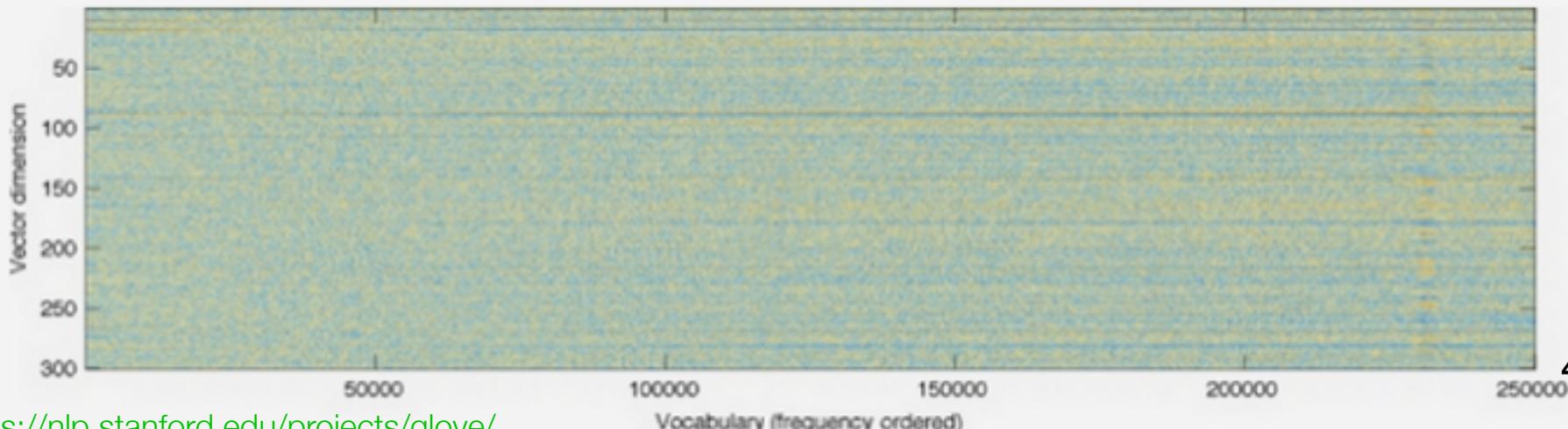


5. *rana*



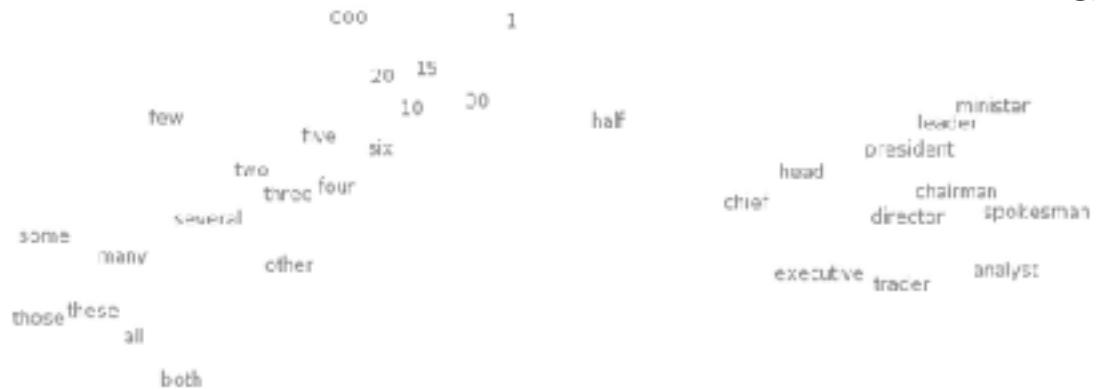
7. *eleutherodactylus*

GloVe produces word vectors with a marked banded structure that is evident upon visualization:



# Word Embeddings: proximity

Global Vectors for Word Representation



t-SNE visualizations of word embeddings. Left: Number Region; Right: Jobs Region. From Turian *et al.* (2010), see complete image.

FRANCE	JESUS	XBOX	REDDISH	SCRATCHED	MEGABITS
AUSTRIA	GOD	AMIGA	GREENISH	NATTED	OCTETS
BELGIUM	SATI	PLAYSTATION	BLUSH	SMASHED	MB/S
GERMANY	CHRIST	MSX	PINKISH	PUNCHED	BIT/B
ITALY	SATAN	IPOD	PURPLISH	POPPED	BAUD
GREECE	KALI	SEGA	BROWNISH	CRIMPED	CARATS
SWEDEN	INDRA	PSNUMBER	GREYISH	SCRAPED	KBIT/S
NORWAY	VISHNU	HD	GRAVISH	SCREWED	MEGAHERTZ
EUROPE	ANANDA	DREAMCAST	WHITISH	SECTIONED	MEGAPIXELS
HUNGARY	PARNATI	GEFORCE	SILVERY	SLASHED	GHIT/S
SWITZERLAND	GRACE	CAPCOM	YELLOWISH	RIPPED	AMPERES

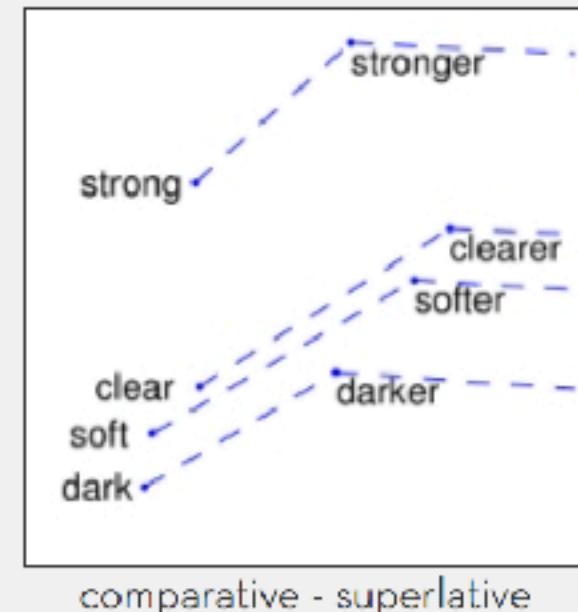
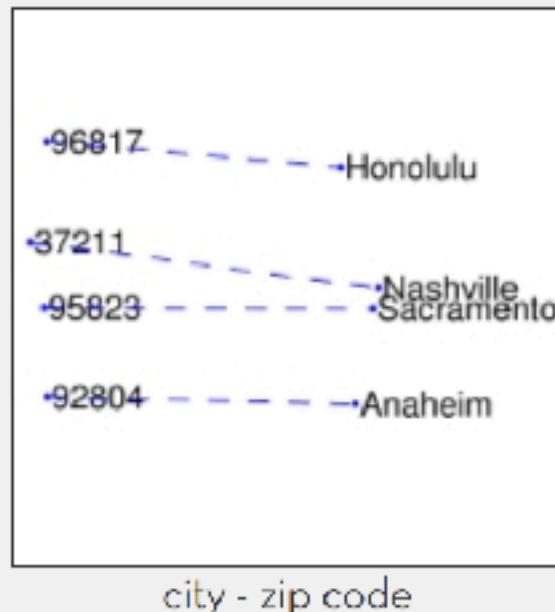
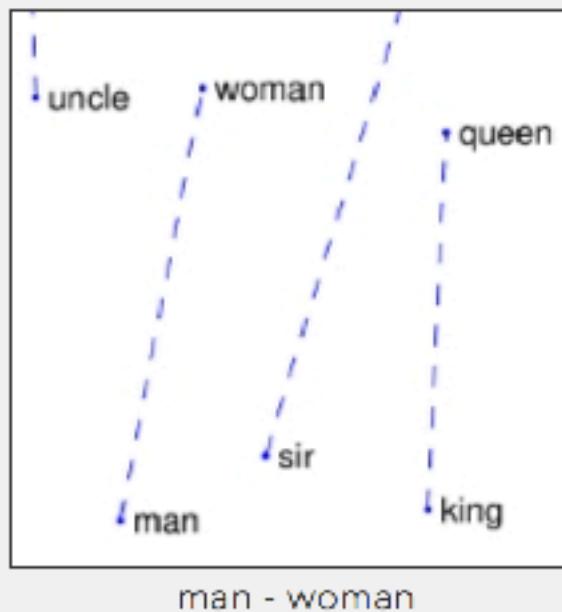
What words have embeddings closest to a given word? From Collobert *et al.* (2011)

<http://colah.github.io/posts/2014-07-NLP-RNNs-Representations/>



# Word Embeddings: Analogy

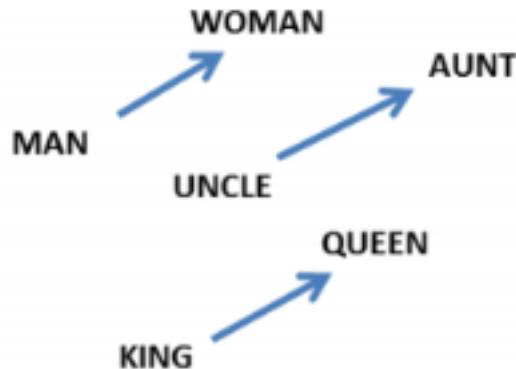
Global Vectors for Word Representation



each vector difference **might** encode analogy



# Word Embeddings: Analogy?



$$W(\text{"woman"}) - W(\text{"man"}) \approx W(\text{"aunt"}) - W(\text{"uncle"})$$

$$W(\text{"woman"}) - W(\text{"man"}) \approx W(\text{"queen"}) - W(\text{"king"})$$

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

From Mikolov *et al.*  
(2013a)

Trained on  
New York Times



<https://nlp.stanford.edu/projects/glove/>

### 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 |

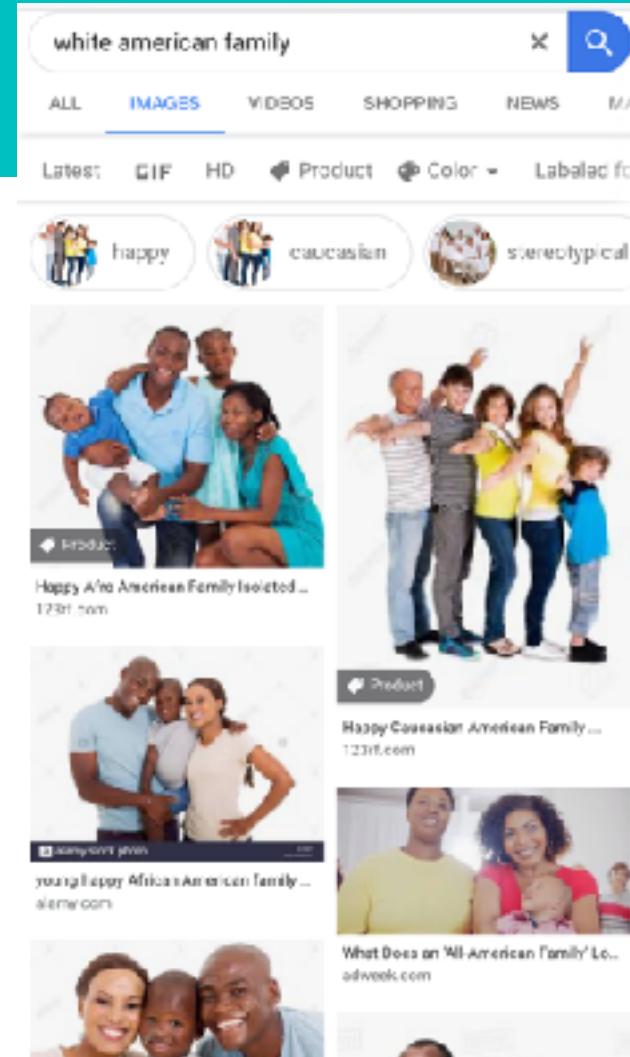
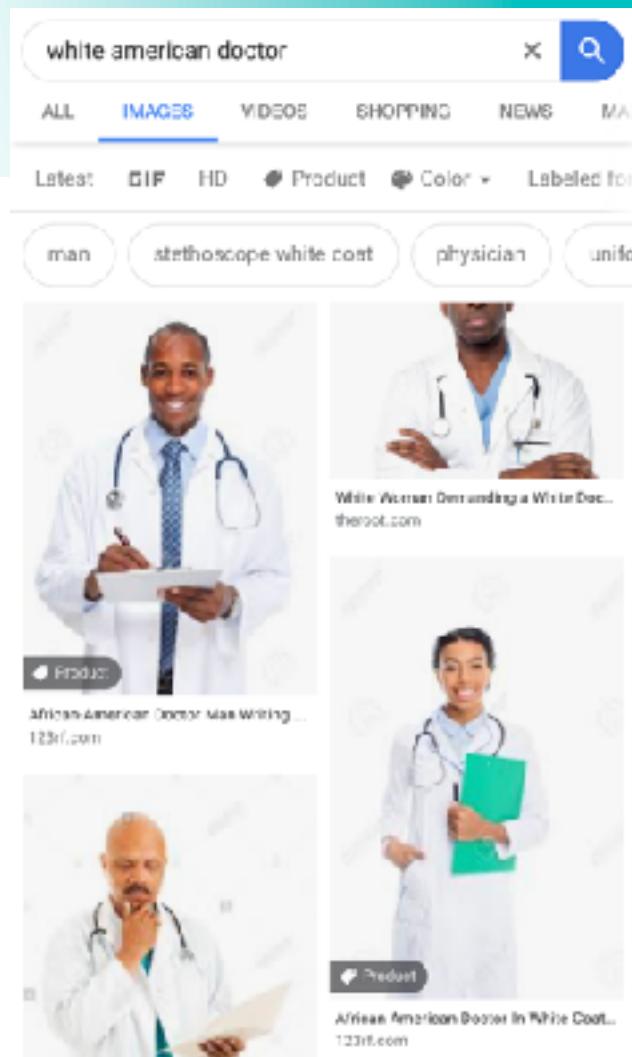
### Extreme *he* occupations

- |                |                   |                |
|----------------|-------------------|----------------|
| 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       |

Bolukbasi et al., NeurIPS 2016  
<https://arxiv.org/pdf/1607.06520.pdf>



# Practical Example in NLP



# ConceptNet

## en artificial intelligence

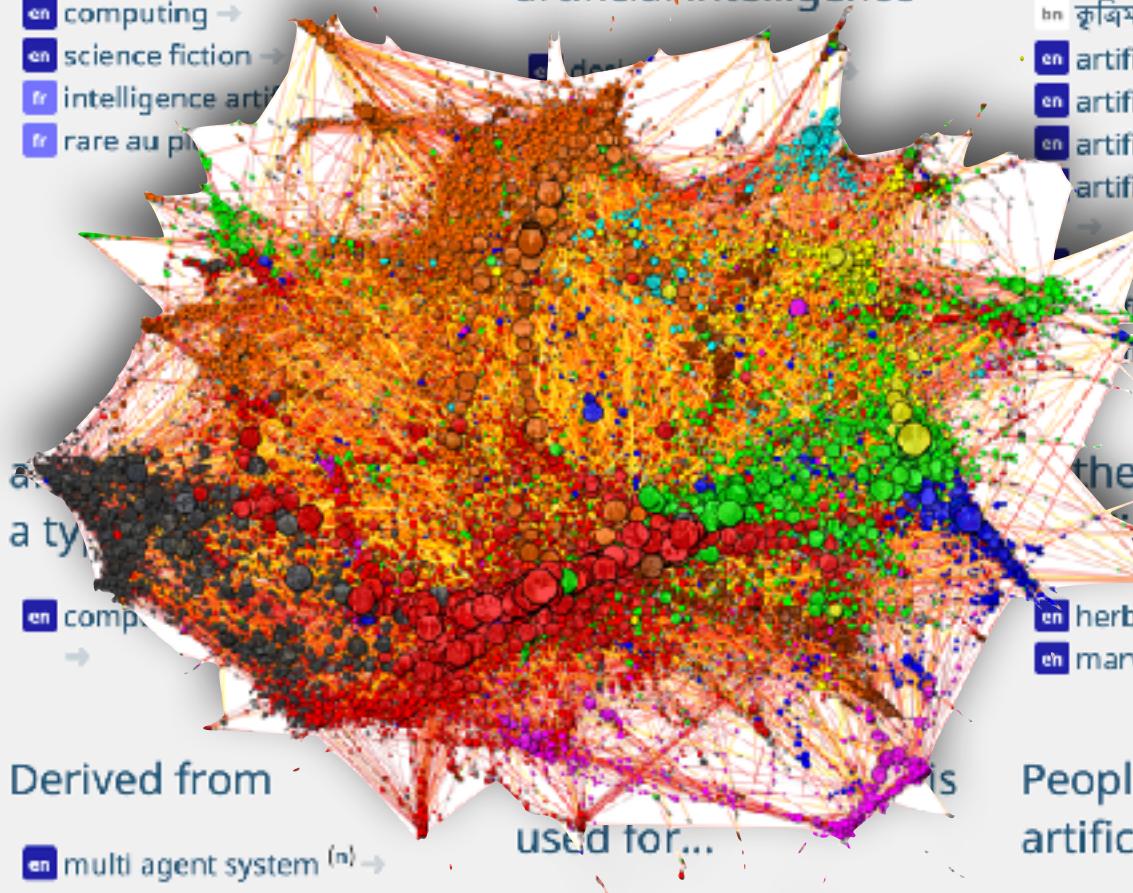
### Derived terms

- en artificial dumbness →
- en artificial incompetence →
- en artificial lack of intelligence →
- en artificial stupidity →
- en artificial unintelligence →
- en artificially intelligent →

### Context of this term

- en computing →
- en science fiction →
- fr intelligence artificielle →
- rare au plaisir →

### Things used for artificial intelligence



### Etymologically related

- bn कृतिम बुद्धिमत्ता →
- en artificial dumbness →
- en artificial idiocy →
- en artificial incompetence →
- en artificial lack of intelligence →
- en artificial stupidity →
- en artificial unintelligence →
- en artificially intelligent →

### Similar terms

- en expert system →
- en expert systems →

artificial intelligence is defined as...

### Derived from

- en multi agent system (n) →

### People known for artificial intelligence



# ConceptNet Numberbatch



- Create with a Knowledge Graph (from multiple sources with relations like *UsedFor*, *PartOf*, etc.)
- Based on this KG, perturb existing embeddings (like GloVe) to optimize:

$$\Psi(Q) = \sum_{i=1}^n \left[ \alpha_i \|q_i - \hat{q}_i\|^2 + \sum_{(i,j) \in E} \beta_{ij} \|q_i - q_j\|^2 \right]$$

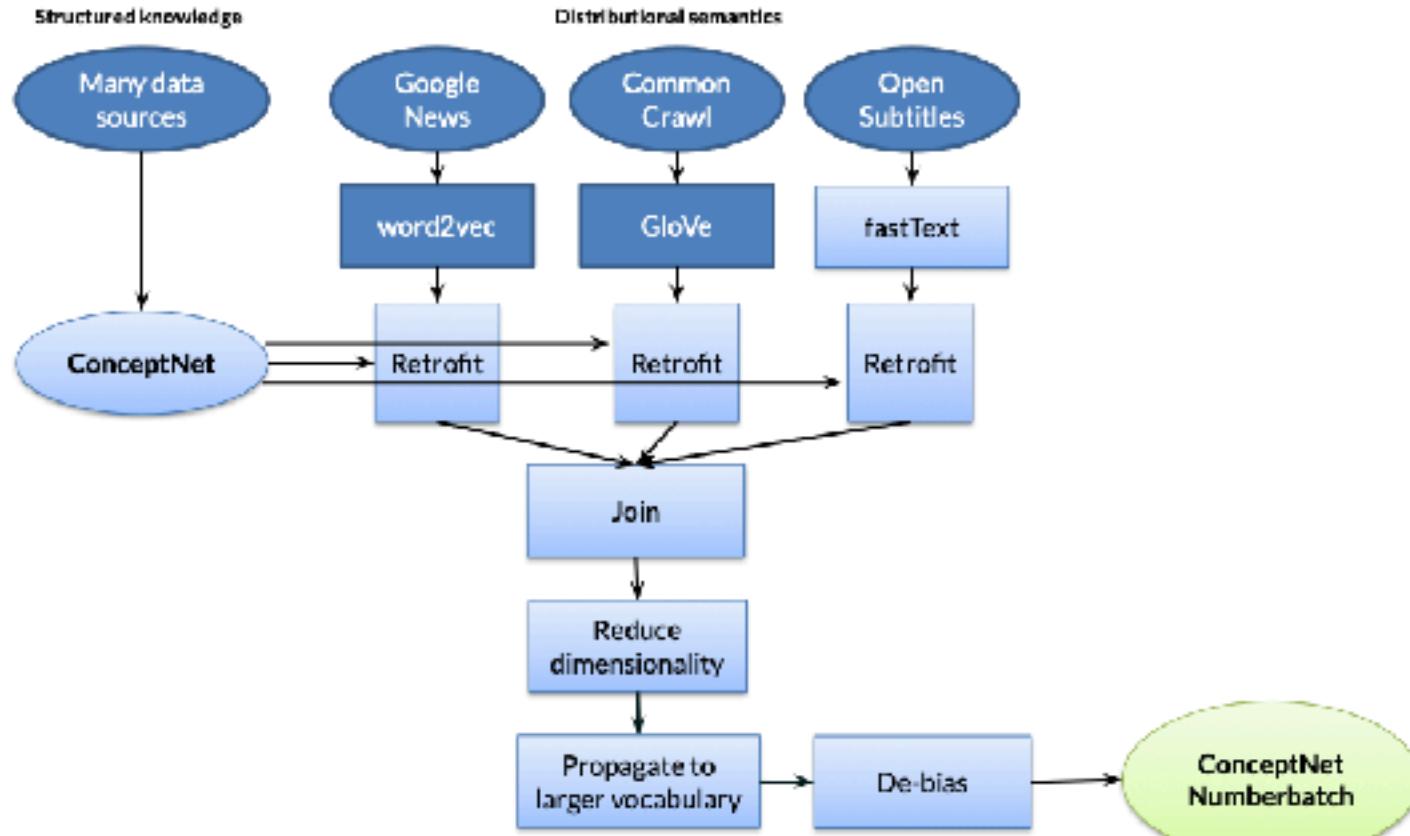
Annotations below the equation:

- $\uparrow$  new embed
- $\uparrow$  old embed
- $\uparrow$  neighbors from KG
- $\text{(keep similar to original)}$
- $\text{(make similar according to other knowledge)}$

- Easy to optimize the objective by averaging neighbors in the ConceptNet KG
- Multiple embeddings achieved by merging through “retrofitting” which projects onto a shared matrix space (with SVD)



# Building ConceptNet Numberbatch



## Aside: Transparency in Research

## ConceptNet is all you need

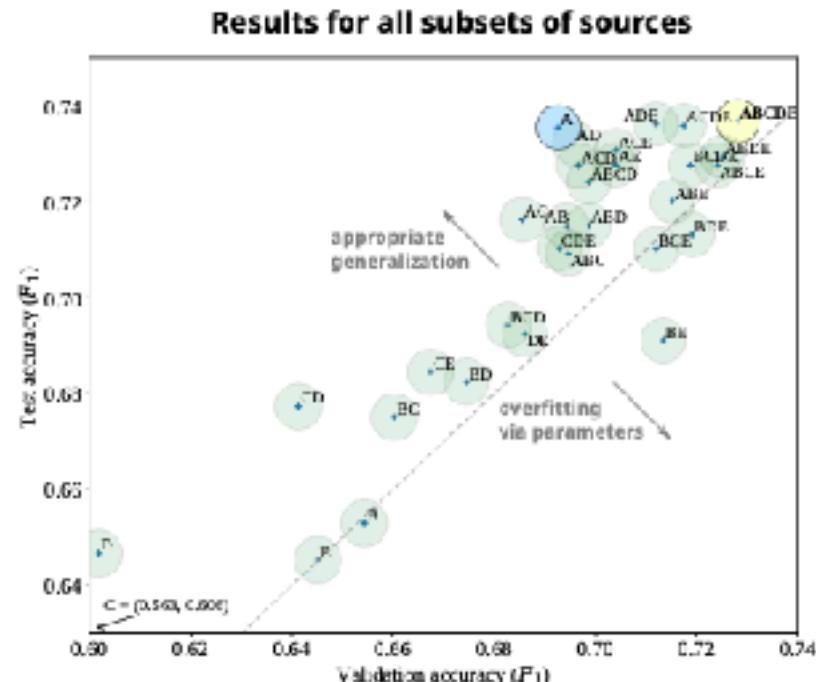
Our full classifier used the linear combination of 5 types of input features shown above. This point is labeled **ABCDE** on the graph to the right. The other points are ablated versions of the classifier, trained on subsets of the five sources.

We found that the single feature of ConceptNet similarity (A) performed just as well on the test data as the full classifier, despite its lower validation accuracy.

This one-feature classifier could be more simply described as a heuristic over cosine similarities of ConceptNet embeddings:

$$\text{sim}(\textit{term}_1, \textit{att}) - \text{sim}(\textit{term}_2, \textit{att}) > 0.0961$$

It seems that the test data contained distinctions that can already be found by comparing ConceptNet embeddings, and that more complex features may have simply provided an opportunity to overfit to the validation set by parameter selection.



This graph shows the validation and test accuracy of classifiers trained on subsets of the five sources of features. Ellipses indicate standard error of the mean, assuming that the data is sampled from a larger set.

# ConceptNet Numberbatch

As a kid, I used to hold marble racing tournaments in my room, rolling marbles simultaneously down plastic towers of tracks and funnels. I went so far as to set up a bracket of 64 marbles to find the fastest marble. I kind of thought that running marble tournaments was peculiar to me and my childhood, but now I've found out that marble racing videos on YouTube are a big thing! Some of them even have [overlays as if they're major sporting events](#).

In the end, there's nothing special about the fastest marble compared to most other marbles. It's just lucky. If one ran the tournament again, the marble champion might lose in the first round. But the one thing you could conclude about the fastest marble is that it was no *worse* than the other marbles. A bad marble (say, a misshapen one, or a plastic bead) would never luck out enough to win.

In our paper, we tested 30 alternate versions of the classifier, including the one that was roughly equivalent to this very simple system. We were impressed by the fact that it performed as well as our real entry. And this could be because of the inherent power of ConceptNet Numberbatch, or it could be because it's the lucky marble.

**-Robyn Speer**  
<http://blog.conceptnet.io>





# How to Make a Racist AI without Really Trying



Robyn Speer, 2017

<http://blog.conceptnet.io/posts/2017/how-to-make-a-racist-ai-without-really-trying/>

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

Bolukbasi et al., NeurIPs 2016  
<https://arxiv.org/pdf/1607.06520.pdf>

**ConceptNet 5.5: An Open  
Multilingual Graph of General  
Knowledge**

Speer et al., AAAI 2017  
<https://arxiv.org/pdf/1612.03975.pdf>



Rachael Tatman @rctatman · 18h

I first got interested in ethics in NLP/ML because I was asking "does this system work well for everyone". It's a good question, but there's a more important important one:

Who is being harmed and who is benefiting from this system existing in the first place?



# Lecture Notes for **Neural Networks** **and Machine Learning**

Ethically Aware Practices



**Next Time:**  
Transfer Learning  
**Reading:** Chollet Article

