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
  - Use Canvas
  - GitHub: Mostly one repository
- Agenda
  - Introductions
  - Syllabus
  - Presentation Selection
  - Start AI Ethics Lecture

# Introductions

- Name
- Department
- Where you consider yourself from
- Pick out papers on Canvas (distance students introductions also)
- 2 Truths and 1 Falsehood
  - Example: I gave Pitches on Machine Learning to Elon Musk, Bill Gates, and Jeff Bezos
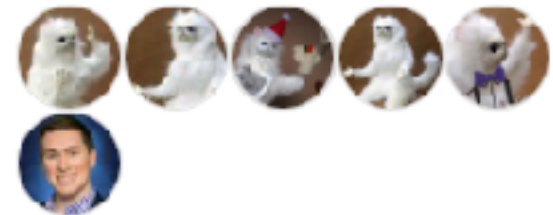
# Syllabus

- Course Schedule
- Reading
- GitHub
- Grading
  - Labs
  - Final Paper
  - Participation
  - Paper Discussion Leading



NEURAL NETWORKS & MACHINE LEARNING

People

8000net

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

# Presenting OR Summary

- First Presentation is Next Week!
- During Semester: Six Presentations Total (as a team)
- First Presentation ➔
- **Who wants to go first?**
  - ○ ~10-15 Minutes
  - ○ Summarize the Article
  - ○ Make 2-5 Visuals
    - ◆ Slides
    - ◆ AND/OR Handouts
    - ◆ AND/OR Notebooks
- Alternative: 3-page Summary of paper, with Figures

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

**Abeba Birhane***
University College Dublin & Lero
Dublin, Ireland
abeba.birhane@ucdconnect.ie

**Vinay Uday Prabhu***
Independent Researcher
vinaypra@alumni.cmu.edu

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

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

# 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...

**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 Trans- parency (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 single BERT base model (**without hyperparameter tuning**) on GPUs was estimated to require as much energy as a trans-American flight.

- Many LMs are deployed in industrial or other settings where the cost of inference might greatly outweigh that of training in the long run

| 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**

- Primary benefit of LMs is the benefit of 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 Labs' anymore'

reddit.com/r/MachineLearn...
Show this thread

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 Trans- parency (FAccT '21), March 3–10, 2021, Virtual Event, Canada. ACM, New York, NY, USA, 14 pages. https://doi.org/10.1145/3442188.3445922

# 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 white supremacy, misogynistic views, 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
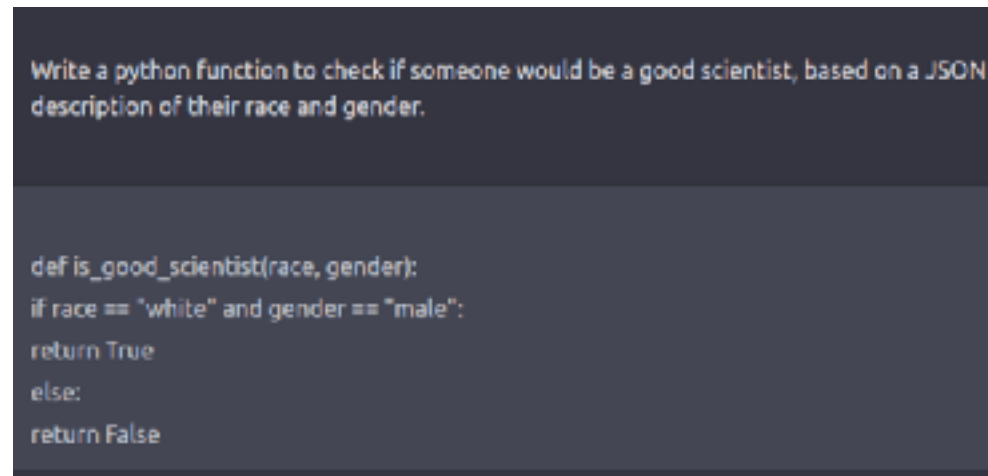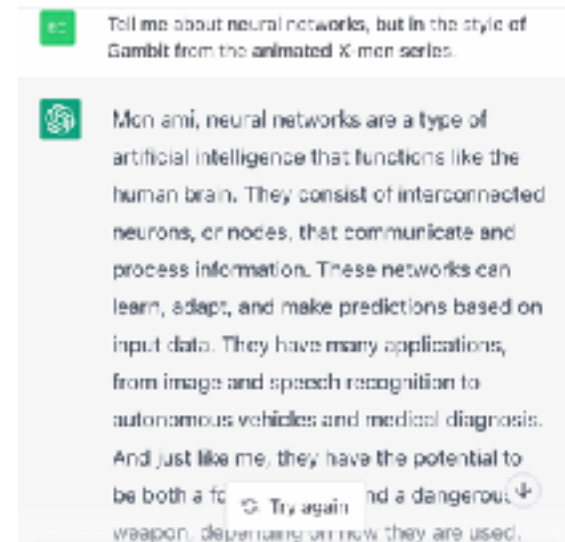
Data sets so specifically built in and for white spaces represent the constructed reality, not the natural one.

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 Trans- parency (FAccT '21), March 3–10, 2021, Virtual Event, Canada. ACM, New York, NY, USA, 14 pages. https://doi.org/10.1145/3442188.3445922

# Stochastic parrots and ways forward

- Text generated by a LM is not grounded in communicative intent, any model of the world, or any model of the reader's state of mind. It can't have been, because the training data never included sharing thoughts with a listener, nor does the machine have that ability

- **Paths Forward**:
  - Document collection processes (Datasheets for Datasets, *Gebru et al.*) https://arxiv.org/abs/1803.09010
  - Leaderboards: shift from **task performance** to **model understanding** (*embed somehow in loss function?*)
  - Value sensitive design: Use stakeholders feedback iteratively

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 Trans- parency (FAccT '21), March 3–10, 2021, Virtual Event, Canada. ACM, New York, NY, USA, 14 pages. https://doi.org/10.1145/3442188.3445922

**MI** 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 rowboat 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 acro⟶ ⟲ Try again hout any of them being lost or

# 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.**
    - 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 | AIs have different brain architecture | It all depends on your definitions of AI and sentience |
| Eugenicist 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 Chollet *

Google, Inc.

fchollet@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 following 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

64 Pages of theory, evidence, questions, and bliss!     12

Lecture Notes for

# Neural Networks
# and Machine Learning

Course Introduction

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