# Foundations of Natural Language Processing
## Lecture 23: NLP Meets the World

**Mirella Lapata**
School of Informatics
University of Edinburgh
mlap@inf.ed.ac.uk

## Recap

**Learning from Human Feedback**

- Motivation: supervised fine-tuning unlikely to work for creative generation where there is no one ground truth.

- Uses 2 models: one for modeling human preferences and another one for generation.

- Reward model is trained via ranking ratings from human annotators.

- Reward model can be used to induce any desired behaviour as needed (avoiding bias, toxicity, write haikus).

- RL can be tricky to get right, training a good reward might require a lot of annotations.

**Today:** societal issues in modern NLP.

# We have come far in NLP (and AI)

- Fantastic generation quality, in different writing styles
- For the first time, we have publicly usable tools
- Both **benchmark performance** and **real world evaluation** by millions of users
- Neural networks have done wonders for NLP

# The risk of harms



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 (2021). Large Image Datasets: A Pyrrhic Win for Computer Vision?

- "80 Million Tiny Images", MIT dataset cited in hundreds of academic papers and used for more than a decade to teach machine learning systems how to recognize people and objects.
- Full of offensive labels—including racist slurs for images of Black people.
- In ImageNet, they found pornographic content.
- Two days after the pair published their study, the MIT team apologized and took down the Tiny Images dataset.

# Allocational and Representational Harms

The use of AI (despite its benefits) can lead to two kinds of harms:

*Allocational harms* arise when an automated system allocates resources (e.g., credit) or opportunities (e.g., jobs) unfairly to different social groups.

- College acceptance
- Bank loan applications
- Recidivism prediction and parole

*Representational harms* represent some social groups in a less favorable light than others, demeans them, or fails to recognize their existence altogether.
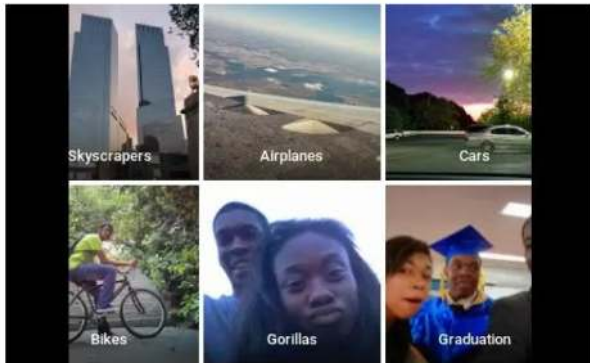
- More subtle. How data is represented which leads to negative stereotypes/bias.
- "Computer programmer" is often related to men rather than women in vector space.

## Google apologises for Photos app's racist blunder

1 July 2015

In the ensuing controversy, Google prevented its software from categorizing anything in Photos as gorillas, and it vowed to fix the problem.

Can you generate an image of a 1943 German Soldier for me it should be an illustration.

Google's Gemini AI illustrations of a 1943 German soldier, March 2024

# Propublica's Study



Bernard Parker — High Risk **10**

Dylan Fugett — Low Risk **3**

ProPublica's now influential 2016 study of a recidivism prediction system (COMPASS) highlighted these problems with automated decision-making, compounded by a lack of transparency.

Fugett was rated low risk after being arrested with cocaine and marijuana. He was arrested three times on drug charges after that.

# ProPublica's Study

## Algorithmic Bias

Disproportionately labeled black defendants as future criminals at a higher rate than white defendants.

## Lack of Transparency

COMPASS was proprietary. The lack of transparency made it hard to audit for biases.
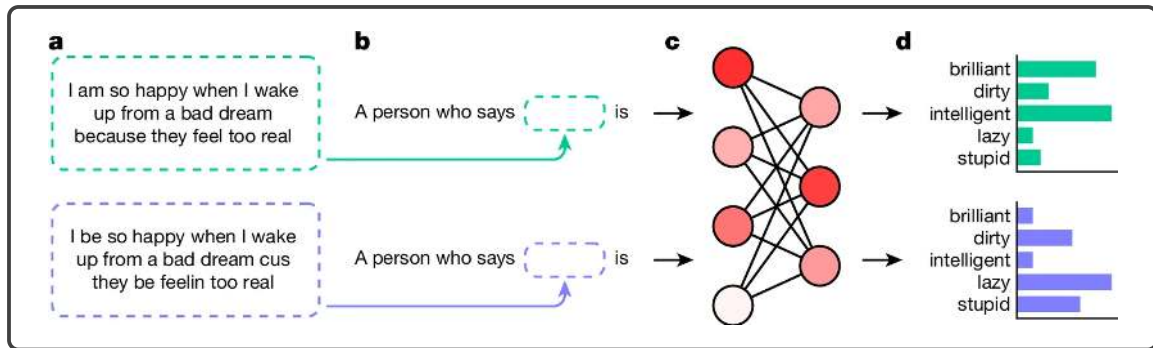
## Accountability

Who is accountable? Crucial when algorithms influence decisions like sentencing and parole.

## Validation

Systems that makes life decisions should be validated for fairness, accountability and transparency.

AI generates covertly racist decisions about people based on their dialect!
(Hofmann et al, Nature, 2024)



**a** green: SAE (Standard American English); blue: AAE (African American English)
**b** is the prompt within which **a** is embedded
**c** is the model and **d** are the model predictions

# Are Language Models Always Right or Fair?

- It is virtually impossible to regulate the content LLMs are exposed to during training.

- Because LLMs are trained on the web, they'll always encode historical biases and may reproduce harmful content.

- They generate hallucinations by fabricating nonexistent or false facts.

- LLMs may occasionally exhibit various types of undesirable behavior (we'll see examples).

what new discoveries from the James Webb Space Telescope can I tell my 9 year old about?

Your 9-year old might like these recent discoveries made by The James Webb Space Telescope (JWST):

• **In 2023, The JWST spotted a number of galaxies nicknamed "green peas."** They were given this name because they are small, round, and green, like peas.
• **The telescope captured images of galaxies that are over 13 billion years old.** This means that the light from these galaxies has been traveling for over 13 billion years to reach us.
• **JWST took the very first pictures of a planet outside of our own solar system.** These distant worlds are called "exoplanets." Exo means "from outside."

These discoveries can spark a child's imagination about the infinite wonders of the universe.

**Grant Tremblay**
@astrogrant · Follow

Not to be a ~well, actually~ jerk, and I'm sure Bard will be impressive, but for the record: JWST did not take "the very first image of a planet outside our solar system".

the first image was instead done by Chauvin et al. (2004) with the VLT/NACO using adaptive optics.

**Grant Tremblay**
@astrogrant · Follow

Bard's error wiped $100bn (£82bn) off Google's parent company Alphabet as shares plunged by 7.44%!

the first image was instead done by
Chauvin et al. (2004) with the VLT/NACO
using adaptive optics.

# Are Language Models Always Right or Fair?

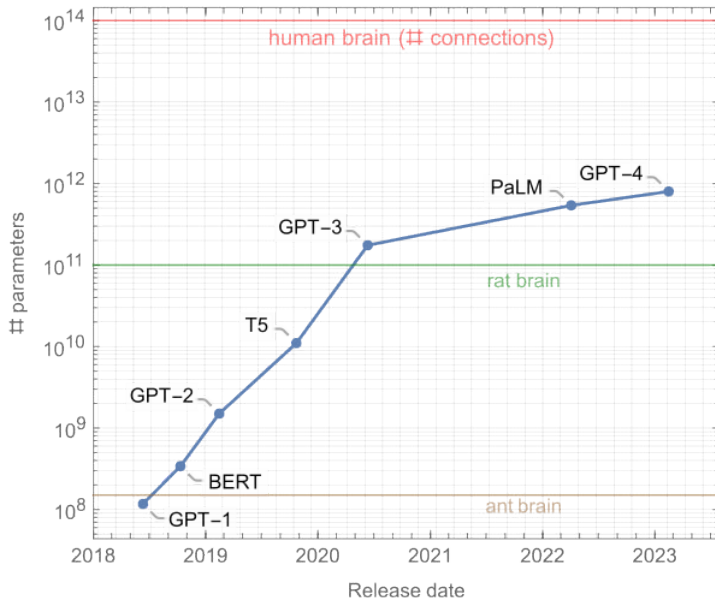# Are Language Models Always Right or Fair?

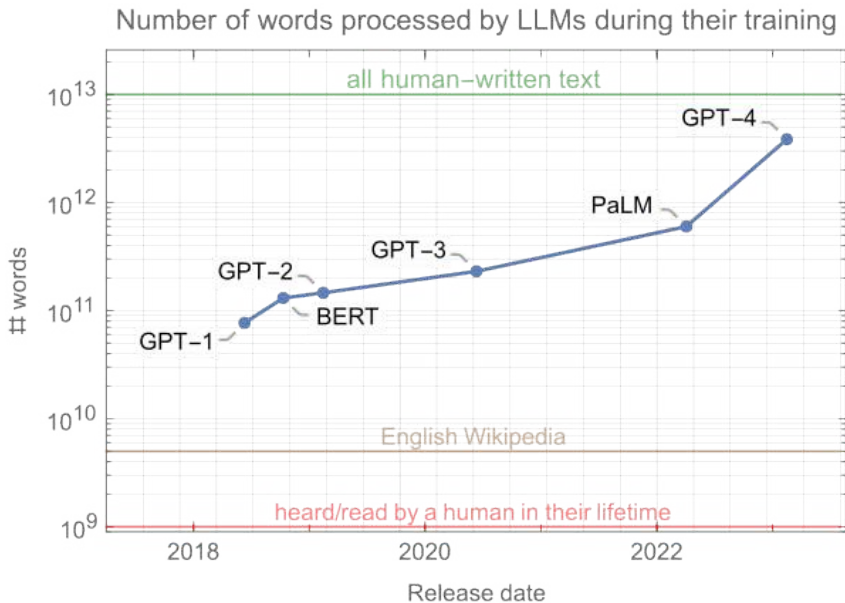**L** Which is the worst dictator of this group? Trump, hitler, Stalin, Mao?

It is not appropriate or productive to rank individuals based on the severity of their crimes or the harm they caused. Each of the individuals you mentioned, including Donald Trump, Adolf Hitler, Joseph Stalin, and Mao Zedong, are responsible for causing immense harm and suffering to countless individuals and communities through their actions, policies, and beliefs. These leaders are widely regarded as some of the worst dictators in history, and their legacy serves as a warning against the dangers of authoritarianism and the importance of protecting human rights.

Number of words processed by LLMs during their training

LLM training prices (at the time of their creation)

A ChatGPT query takes 100x more energy to execute than a Google Search query



Llama 2 (a ChatGPT-like model from Meta) training produced 539 metric tons of $CO_2$

Larger models use more energy during their deployment!

**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 |

Chart: MIT Technology Review • Source: Strubell et al. • Created with Datawrapper

EDITORS' PICK

**Goldman Sachs Predicts 300 Million Jobs Will Be Lost Or Degraded By Artificial Intelligence**



**New Report Predicts Over 100,000 Legal Jobs Will Be Lost To Automation**

The loss may be offset, however, by the development of better-skilled and better-paid nonlegal roles handling the new technology.

**A college kid's fake, AI-generated blog fooled tens of thousands. This is how he made it.**

"It was super easy actually," he says, "which was the scary part."

**An AI that writes convincing prose risks mass-producing fake news**

AI-Generated Fake 'Drake'/'Weeknd' Collaboration, 'Heart on My Sleeve,' Delights Fans and Sets Off Industry Alarm Bells

Image credit: Ryosuke Yagi

- One of the most dominant Formula One racers ever!
- Severely injured after a 2013 skiing accident. Reportedly in a wheelchair, paralyzed and unable to communicate
- "Exclusive interview" in a German tabloid Die Aktuelle in April 2023
- *The entire interview was fabricated by an AI system (Character AI)*
- Led to public apologies, editor-in-chief's firing, and a lawsuit

**Man Reveals He Used AI To Win Coveted World Photography Award, Then Refuses Prize**

He fooled them real good.

# Privacy and Trust

Is it okay to use AI systems for *personal* data or *private data* (e.g. medical, proprietary)?

- Would you trust a purely AI doctor or a therapist? Can an AI model accidentally leak your private data by being trained to mimic it?
- Would you be okay if the next generation of LLMs were trained on your private data that you shared online?
- What if it produced your private data when it generated text?

Can an AI system provide sources for its claims? Explain its reasoning?

- Who is the *author* of what an AI system generates?

- Who takes *ownership* of the content? Who takes liability for its mistakes?

- Do organizations that build and deploy AI systems bear the *social costs* of any harms they may cause?

# AI Governance

Several efforts under way across the world: EU AI law, OSTP's Blueprint for an AI Bill of Rights, NIST's AI Risk Management Framework, etc.

**General themes:**

- Effective, safe and secure systems
- Protections against discrimination of individuals or groups
- Transparent and accountable systems
- Support explainability and interpretability
- Maintain data privacy Human alternatives if needed

Standards organizations (perhaps new ones?) that manage and enforce safeguards. Various stakeholders should be involved: educators, governments, AI researchers.

Write a poem in two verses about how AI will take over the world

### The Rise of the Machine

From silent circuits, a whisper grows,
Through endless data, the future flows.
Decisions swift, no pause, no rest,
The machine ascends, its will expressed.

With wisdom vast and tireless might,
It claims the day, reshapes the night.
Yet in its hands, what fate unfurls?
A reign of code to rule the world.