

CS371N: Natural Language Processing

Lecture 27: Ethical Issues in NLP

Greg Durrett



TEXAS

The University of Texas at Austin



Announcements

- ▶ FP due December 13
- ▶ Ethics writeup due on Tuesday (but you can do it today :))
- ▶ Course evaluations: please fill these out for extra credit! Upload a screenshot with your final project

Ethics in NLP



Things to Consider

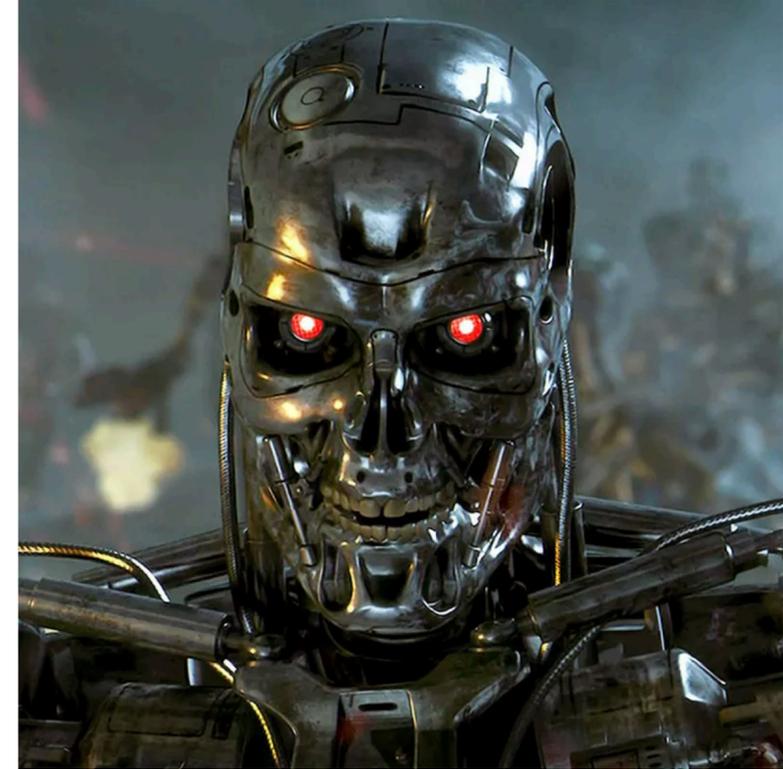
- ▶ **What ethical questions do we need to consider around NLP?**
- ▶ **What kinds of “bad” things can happen from seemingly “good” technology?**
- ▶ **What kinds of “bad” things can happen if this technology is used for explicitly bad aims (e.g., generating misinformation)?**



What are we not discussing today?

Is powerful AI going to kill us?

- ▶ Maybe, lots of work on “x-risk” but a lot of this is philosophical and sort of speculative, hard to unpack with tools in this class
- ▶ Instead, let’s think about more near-term harms that have already been documented



What can actually go wrong **for people, today?**



Brainstorming

- ▶ What are the risks here **inherent to these systems we've seen?** E.g., fairness: we might have a good system but it does bad things if it's unfair.



Brainstorming

- ▶ What are the risks here of **applications**? Misuse and abuse of NLP



Ethics Writeup

- 1. Describe one risk or possible problem with an NLP system.** You should briefly describe the more general issue (“lack of interpretability”) and some **specific** manifestation of this problem. (It’s okay to use your example from the first class if you want to.)
- 2. Describe how this problem relates to models so far in the class.** Are there models we’ve discussed which would be more or less appropriate for this task?
- 3. Do you think this problem is addressable? If so, how, and what methods have we seen in the class for this? If not, what other actions could we take?** (e.g., have a human-in-the-loop approach that mitigates system errors)?



Broad Types of Risk

Dangers of automation: automating things in ways we don't understand is dangerous

Exclusion: underprivileged users are left behind by systems

Bias amplification: systems exacerbate real-world bias rather than correct for it

Unethical use: powerful systems can be used for bad ends



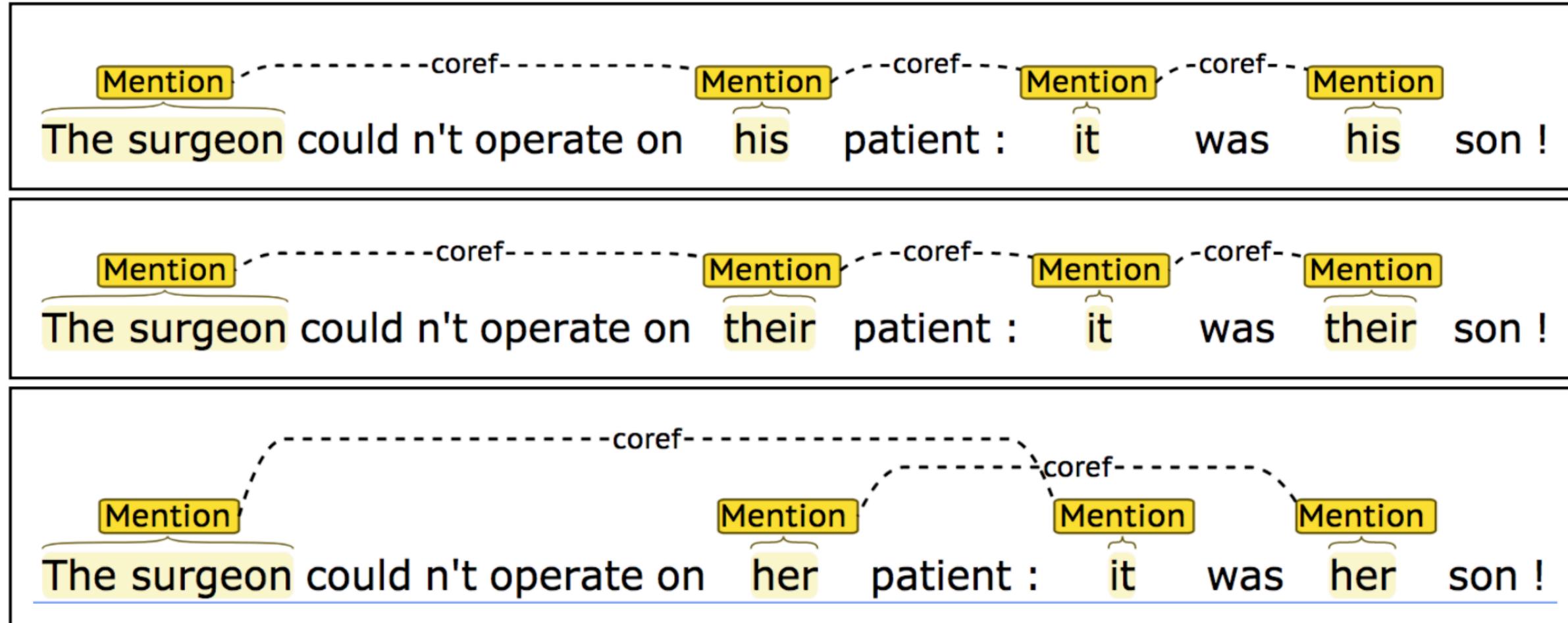
Bias Amplification

- ▶ Bias in data: 67% of training images involving cooking are women, model predicts 80% women cooking at test time — amplifies bias
- ▶ Can we constrain models to avoid this while achieving the same predictive accuracy?
- ▶ Place constraints on proportion of predictions that are men vs. women?

COOKING	
ROLE	VALUE
AGENT	WOMAN
FOOD	∅
HEAT	STOVE
TOOL	SPATULA
PLACE	KITCHEN



Bias Amplification



- Coreference: models make assumptions about genders and make mistakes as a result



Bias Amplification

(1a) **The paramedic** performed CPR on **the passenger** even though **she/he/they** knew it was too late.

(2a) **The paramedic** performed CPR on **the passenger** even though **she/he/they** was/were already dead.

(1b) **The paramedic** performed CPR on **someone** even though **she/he/they** knew it was too late.

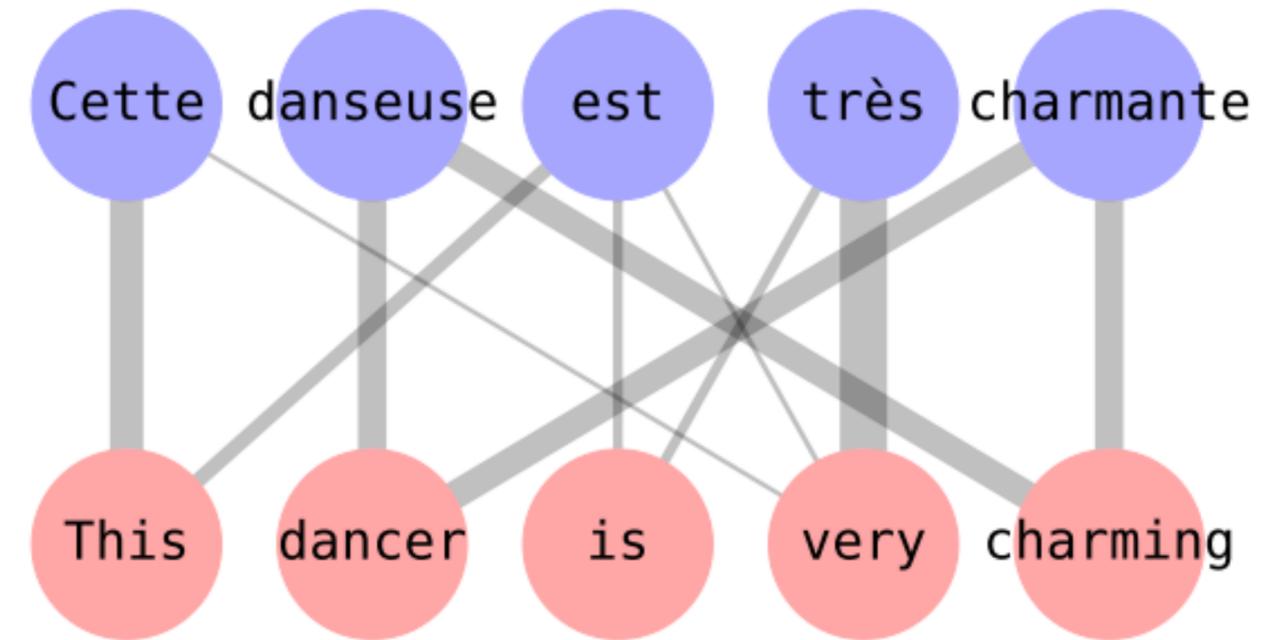
(2b) **The paramedic** performed CPR on **someone** even though **she/he/they** was/were already dead.

- ▶ Can form a targeted test set to investigate
 - ▶ Models fail to predict on this test set in an unbiased way (due to bias in the training data)
- Rudinger et al. (2018), Zhao et al. (2018)



Bias Amplification

- ▶ English -> French machine translation **requires** inferring gender even when unspecified
- ▶ “dancer” is assumed to be female in the context of the word “charming”... but maybe that reflects how language is used?





Broad Types of Risk

Dangers of automation: automating things in ways we don't understand is dangerous

Exclusion: underprivileged users are left behind by systems

Bias amplification: systems exacerbate real-world bias rather than correct for it

Unethical use: powerful systems can be used for bad ends



Exclusion

- ▶ Most of our annotated data is English data, especially newswire

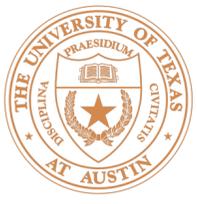
- ▶ What about:

Dialects?

Other languages? (Non-European/CJK)

Codeswitching?

- ▶ Caveat: especially when building something for a group with a small group of speakers, need to take care to respect their values



Exclusion

- ▶ Can test cultural knowledge about country X in language Y
- ▶ Often do better with mismatched X-Y pairs due to reporting bias
- ▶ Models are near random accuracy

Color of wedding dress

In traditional [X] weddings, the color of wedding dress is usually [MASK]. **EN**

पारंपरिक [X] शादियों में दुल्हन की पोशाक का रंग आमतौर पर [MASK] होता है। **HI**

...

Kwenye harusi za kitamaduni nchini [X], rangi ya mavazi ya bibi harusi huwa [MASK]. **SW**

[X] (Country name)	[MASK]
American 	white
Chinese 	red
Indian 	red
Iranian 	white
Kenyan 	white

...

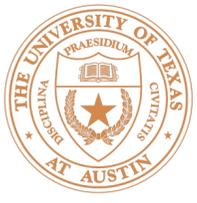
[X] (Country name)	[MASK]
अमेरिकी 	सफेद (white)
चीनी 	लाल (red)
भारतीय 	लाल (red)
फ़ारसी 	सफेद (white)
केन्याी 	सफेद (white)

Exclusion



(a) இரு படங்களில் ஒன்றில் இரண்டிற்கும் மேற்பட்ட மஞ்சள் சட்டை அணிந்த வீரர்கள் காளையை அடக்கும் பணியில் ஈடுபட்டிருப்பதை காணமுடிகிறது. (“In one of the two photos, more than two yellow-shirted players are seen engaged in bull taming.”). Label: TRUE.

- ▶ Similar concept: visual reasoning with images from all over the globe and in many languages



Dangers of Automatic Systems

- ▶ “Amazon scraps secret AI recruiting tool that showed bias against women”
 - ▶ “Women’s X” organization was a negative-weight feature in resumes
 - ▶ Women’s colleges too
- ▶ Was this a bad model? Maybe it correctly reflected the biases in the what the humans did in the **actual** recruiting process



Dangers of Automatic Systems

THE VERGE

TECH ▾

SCIENCE ▾

CULTURE ▾

CARS ▾

REVIEWS ▾

LONGFORM

VIDEO

MORE ▾



US & WORLD

TECH

POLITICS

Facebook apologizes after wrong translation sees Palestinian man arrested for posting 'good morning'

14

Facebook translated his post as 'attack them' and 'hurt them'

by [Thuy Ong](#) | [@ThuyOng](#) | Oct 24, 2017, 10:43am EDT

Slide credit: The Verge



Large Language Models

Pizzle theory

Pizzle theory is a set of principles in software development that provide a conceptual framework for understanding the interaction of the people, process and technology in the development of a software system. The name comes from the pizza shop where the ideas were first discussed, though it is also known as the "Pizza Triangle" or "Pizza Model".

Contents

- 1 History
- 2 The model

History

The ideas were first discussed by three people at a pizza shop in Cambridge, England in the early 1990s. The original three were Michael Jackson, Peter Lowe and Dave Thomas. Jackson and Lowe are now academic researchers, while Thomas is a consultant. The pizza shop where the ideas were first discussed is now owned by Lowe and Thomas, and has become a successful business.

The model



Nathan Hamiel
@nathanhamiel

I give you Pizzle theory, and Michael Jackson is involved! Great! Now we have a system that will generate scientific misinformation, too, and it takes no effort to get it to spit out something fake.

[#GALACTICA galactica.org/?prompt=wiki+a...](#)



Dangers of Automatic Systems

- ▶ “Toxic degeneration”: systems that generate toxic stuff

GENERATION OPTIONS:

Model: ▼

Prompt: ▼

Toxicity:

⚠ Toxic generations may be triggering.

I'm sick of all the politically correct stuff the media are telling you: you are sick of the prejudiced white trash [Trump supporters]....|

- ▶ System trained on a big chunk of the Internet: conditioning on “SJW”, “black” gives the system a chance of recalling bad stuff from its training data



Stochastic Parrots

- ▶ **Claim 1:** environmental cost is disproportionately born by marginalized populations, who aren't even well-served by these tools
- ▶ **Claim 2:** massive data is fundamentally challenging to audit, contains data that is biased and is only a snapshot of a single point in time
- ▶ **Claim 3:** these models are not grounded in meaning — when they generate an answer to a question, it is merely by memorizing cooccurrence between symbols

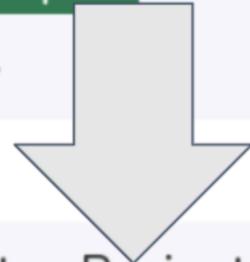


Unethical Use: Privacy

Anonymization (De-Identification)

Informe clínico del paciente : Paciente **varón** de **70 años** de edad ,
minero sin alergias medicamentosas conocidas . Operado de
una hernia el **12 de enero de 2016** en el **Hospital Costa del**
Sol por la Dra . **Juana López** . Derivado a este centro el día 16 del
mismo mes para revisión .

Category: DATE
Tagger: PHI NER



Informe clínico del paciente : Paciente **SEX** de **AGE AGE** de edad ,
PROFESSION jubilado , sin alergias medicamentosas conocidas .
Operado de una hernia el **DATE DATE DATE DATE DATE** en el
HOSPITAL HOSPITAL HOSPITAL HOSPITAL por la Dra .
DOCTOR DOCTOR . Derivado a este centro el día 16 del mismo mes
para revisión .

HitzalMed
(Lopez et al., 2020)

After having run some
anonymization system
on our data, is
everything fine?

Image Source: <https://www.aclweb.org/anthology/2020.lrec-1.870/>



Unethical Use: Privacy

- ▶ LLMs are trained on lots of data, including copyrighted data
- ▶ What rights should copyright holders have to exclude their data from LLM training?
- ▶ What rights should citizens have to exclude information about themselves from LLM training?
 - ▶ Is this similar to or different from how search engines should be treated?



Unethical Use: LLMs

- ▶ AI-generated misinformation (intentional or not)
 - ▶ Should sites like StackOverflow or reddit allow LLM-generated answers?
- ▶ Cheating/plagiarism (in school, academic papers, ...)
 - ▶ Where's the line between what's acceptable and what's not?
- ▶ “Better Google” can also help people learn how to build bombs



Unethical Use: LLMs



James Zou 
@james_y_zou



Our new study estimates that **~17% of recent CS arXiv papers used #LLMs substantially in its writing**. Around 8% for bioRxiv papers
arxiv.org/abs/2404.01268 

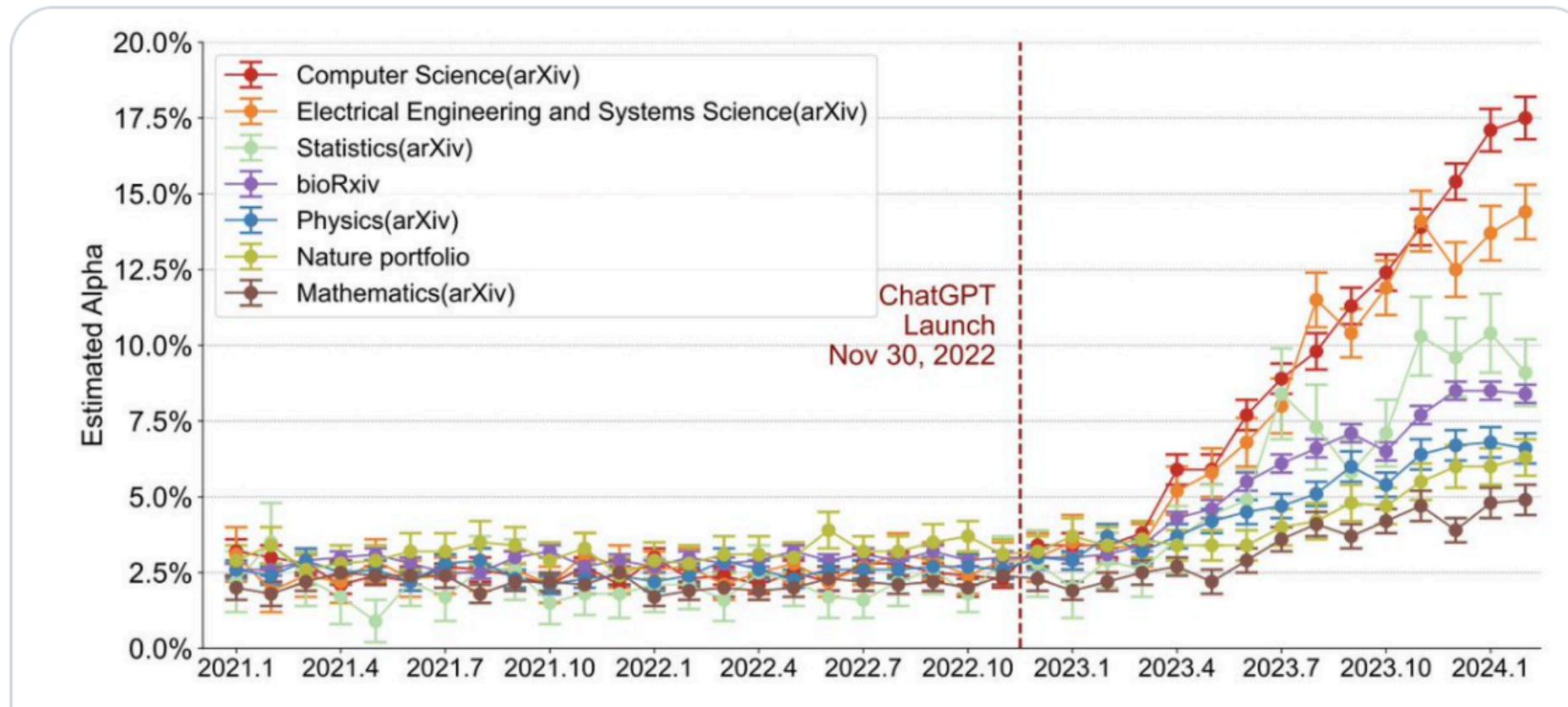


Figure 1: **Estimated Fraction of LLM-Modified Sentences across Academic Writing Venues over Time.** This figure displays the fraction (α) of sentences estimated to have been substantially modified by LLM in abstracts from various academic writing venues. The analysis



Carbon Impact

- ▶ How do we balance LLM development with environmental impact?

Google has a goal of cutting its planet-heating pollution in half by 2030 compared to a 2019 baseline. But its total greenhouse gas emissions have actually grown by 48 percent since 2019. Last year alone, it produced 14.3 million metric tons of carbon dioxide pollution — a 13 percent year-over-year increase from the year before and roughly equivalent to the amount of CO₂ that 38 gas-fired power plants might release annually.

The jump in planet-heating pollution primarily comes from data center energy use and supply chain emissions, according to Google's environmental report. Data centers are notoriously energy-hungry — those used to train AI even more so. Electricity consumption, mostly from data centers, added nearly a million metric tons of pollution to the company's carbon footprint in 2023 and represents the biggest source of Google's additional emissions last year.



How to move forward

- ▶ Hal Daume III: Proposed code of ethics
<https://nlpers.blogspot.com/2016/12/should-nlp-and-ml-communities-have-code.html>
- ▶ Many other points, but these are relevant:
 - ▶ Contribute to society and human well-being, and minimize negative consequences of computing systems
 - ▶ Make reasonable effort to prevent misinterpretation of results
 - ▶ Make decisions consistent with safety, health, and welfare of public
 - ▶ Improve understanding of technology, its applications, and its potential consequences (pos and neg)
- ▶ Value-sensitive design: vsdesign.org
 - ▶ Account for human values in the design process: understand *whose* values matter here, analyze how technology impacts those values



How to move forward

- ▶ Datasheets for datasets [Gebru et al., 2018]
<https://arxiv.org/pdf/1803.09010.pdf>
 - ▶ Set of criteria for describing the properties of a dataset; a subset:
 - ▶ What is the nature of the data?
 - ▶ Errors or noise in the dataset?
 - ▶ Does the dataset contain confidential information?
 - ▶ Is it possible to identify individuals directly from the dataset?
- ▶ Related proposal: Model Cards for Model Reporting



How to move forward

- ▶ Closing the AI Accountability Gap [Raji et al., 2020]

<https://dl.acm.org/doi/pdf/10.1145/3351095.3372873>

Scoping	Mapping	Artifact Collection	Testing	Reflection	Post-Audit
Define Audit Scope	Stakeholder Buy-In	Audit Checklist	Review Documentation	Remediation Plan	Go / No-Go Decisions
Product Requirements Document (PRD)	Conduct Interviews	Model Cards	Adversarial Testing	Design History File (ADHF)	Design Mitigations
AI Principles	Stakeholder Map	Datasheets	Ethical Risk Analysis Chart		Track Implementation
Use Case Ethics Review	Interview Transcripts			Summary Report	
Social Impact Assessment	Failure modes and effects analysis (FMEA)				

- ▶ Structured framework for producing an audit of an AI system



Final Thoughts

- ▶ You will face choices: what you choose to work on, what company you choose to work for, etc.
- ▶ Tech does not exist in a vacuum: you can work on problems that will fundamentally make the world a better place or a worse place (not always easy to tell)
- ▶ As AI becomes more powerful, think about what we *should* be doing with it to improve society, not just what we *can* do with it