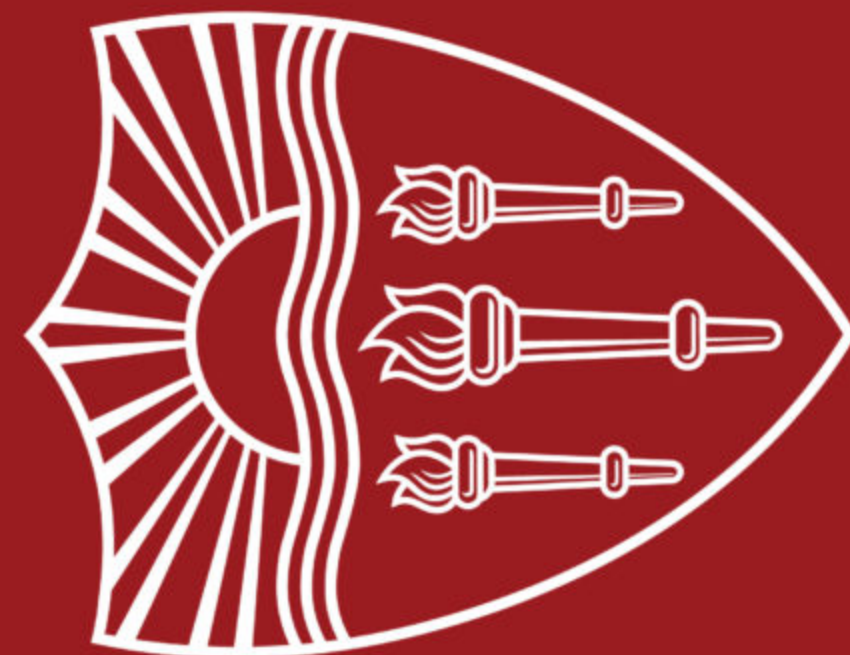


UCSD



Lecture 20: LLMs: Preference-Tuning and Harms

*Instructor: Swabha Swayamdipta
USC CSCI 544 Applied NLP
Nov 5, Fall 2024*



Announcements

- Today, Tue, 11/5 - Lecture + Paper Presentation II
- Thu, 11/7 - Quiz 4 + Paper Presentation III
- Tue, 11/12 - Quiz 5 + Paper Presentation IV
 - Quizzes 4 and 5 - all topics after the midterms. Consider these as practice tests for final exams
- Thu, 11/14 Guest lecture on LLM Pretraining by Prof. Willie Neiswanger on 11/14 + HW4 due
 - Questions from lecture materials will be included in final exam

Lecture Outline

- Announcements
- Last Lecture: Supervised Fine-tuning and Prompting
- Today:
 - Post-training with Alignment with Human Feedback:
 - Preference Tuning: RLHF
 - LLMs Harms and Opportunities
 - Wrapping it all up

The need for post-training

A Pre-trained GPT-3

Prompt: Explain the moon landing to a six year old in a few sentences.

Output: Explain the theory of gravity to a 6 year old.

Prompt: Translate to French: The small dog

Output: The small dog crossed the road.

Ouyang et al., 2022; J&M Chap 12

- Make LLMs more helpful
 - Supervised Finetuning: Instruction Tuning
 - Prompting
- Make LLMs less harmful
 - Model Alignment with Human Preferences: Intro to RLHF / DPO

Model Alignment with Human Preferences

Preference Alignment

- Let's say we were training a language model on some task (e.g. summarization).
- For an instruction x and a LM sample y , imagine we had a way to obtain a human reward of that summary: $R(x, y) \in \mathbb{R}$, higher is better.

SAN FRANCISCO,
California (CNN) --
A magnitude 4.2
earthquake shook the
San Francisco
...
overturn unstable
objects.
 x

An earthquake hit
San Francisco.
There was minor
property damage,
but no injuries.

$$y_1$$

$$R(x, y_1) = 8.0$$

The Bay Area has
good weather but is
prone to
earthquakes and
wildfires.

$$y_2$$

$$R(x, y_2) = 1.2$$

- Maximize the expected reward of samples from our LM: $\mathbb{E}_{\hat{y} \sim p_{\theta}(y|x)}[RM_{\phi}(x, \hat{y})]$

Step 1

Collect demonstration data and train a supervised policy.

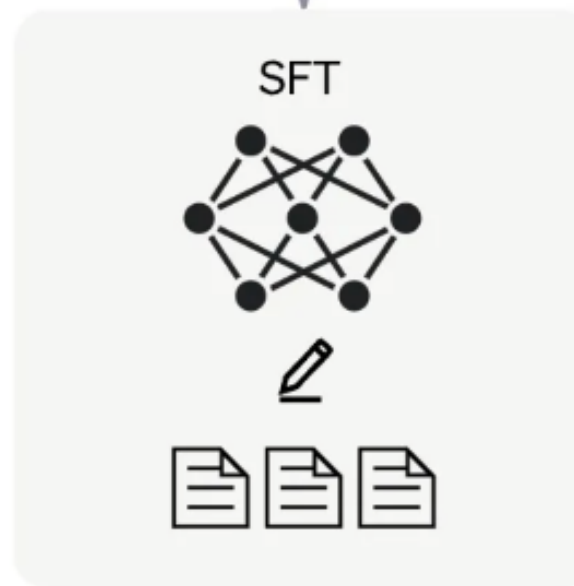
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



This data is used to fine-tune GPT-3.5 with supervised learning.

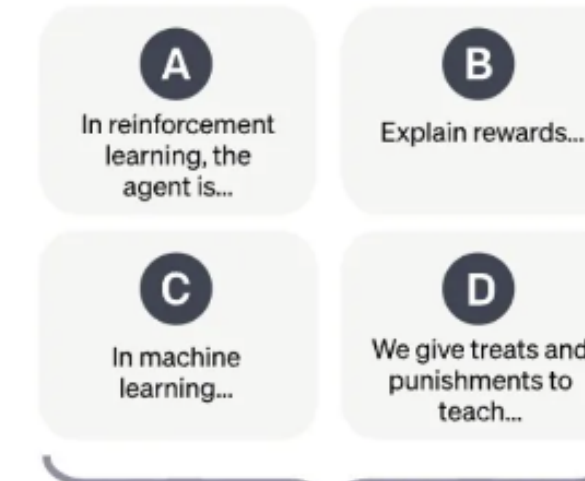
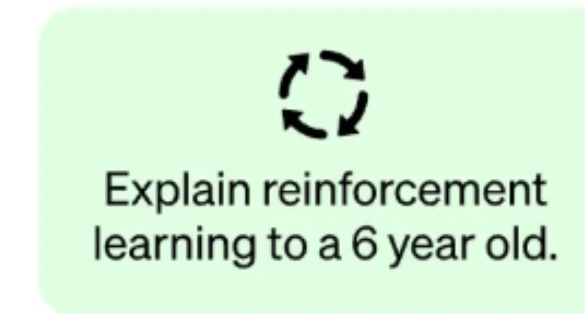


Instruction Tuning!

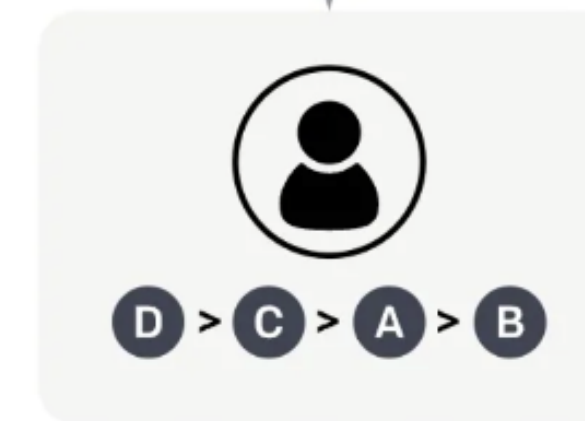
Step 2

Collect comparison data and train a reward model.

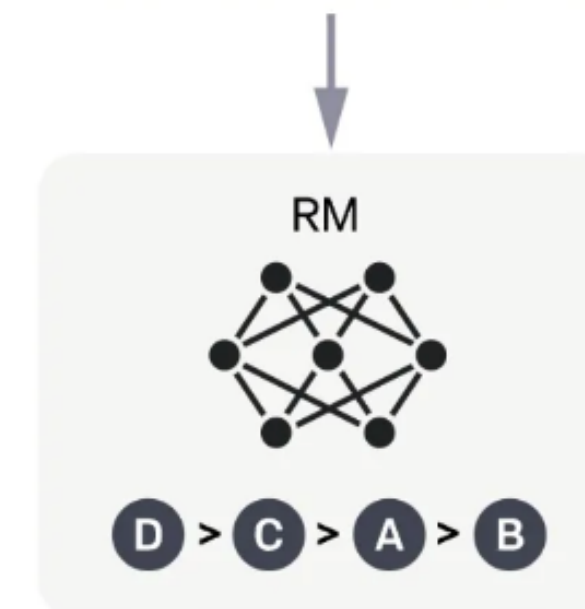
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



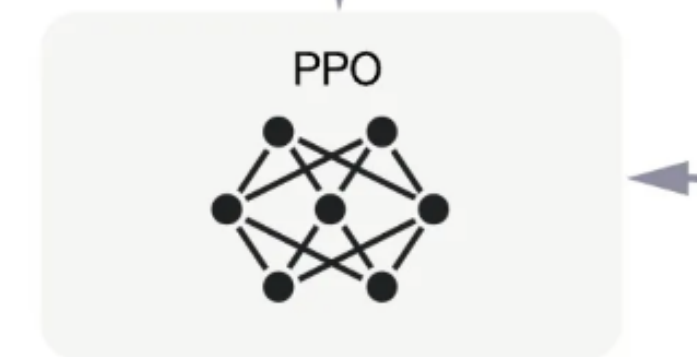
Step 3

Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

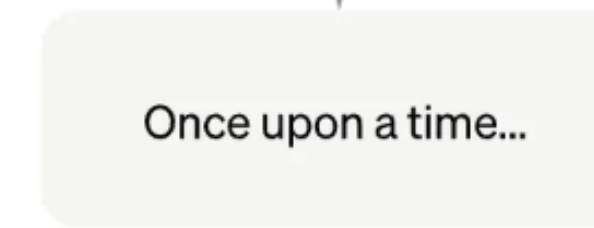
A new prompt is sampled from the dataset.



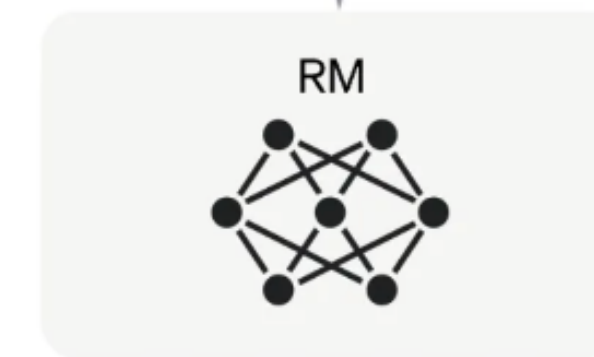
The PPO model is initialized from the supervised policy.



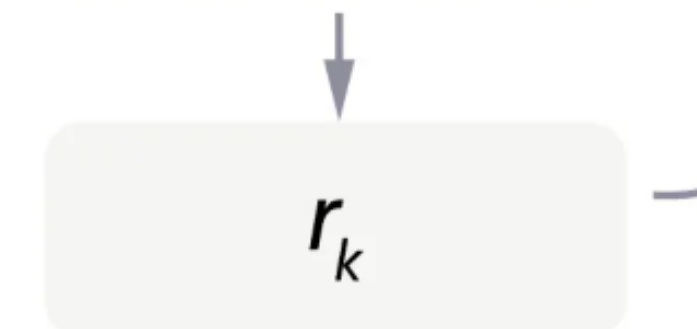
The policy generates an output.



The reward model calculates a reward for the output.



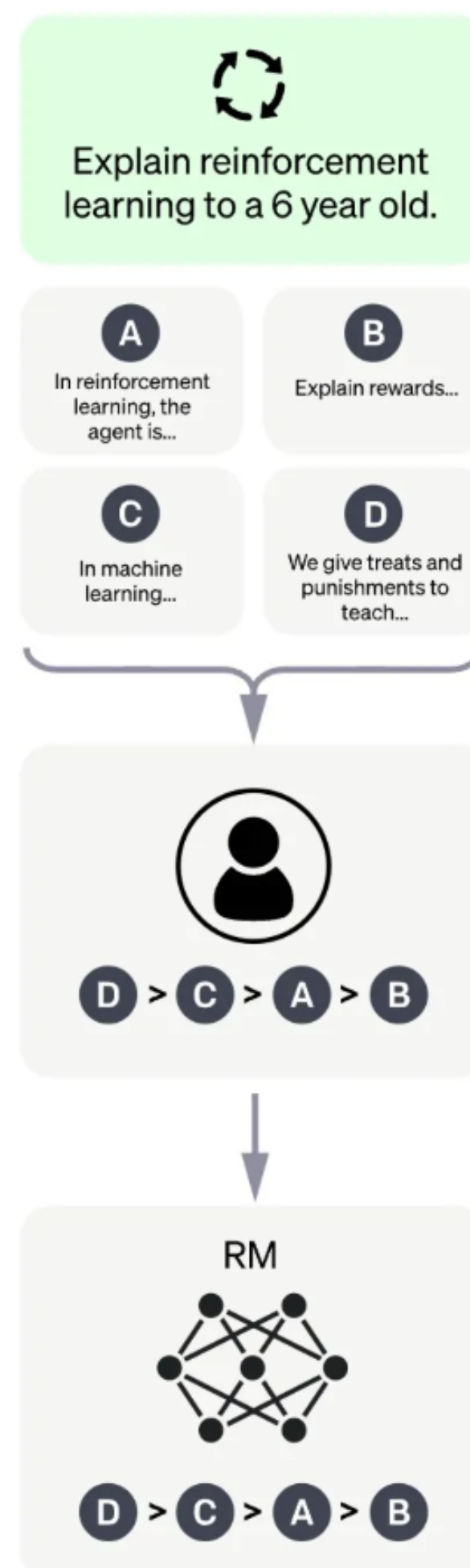
The reward is used to update the policy using PPO.



Step 2

Collect comparison data and train a reward model.

A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.

This data is used to train our reward model.

Preference Data

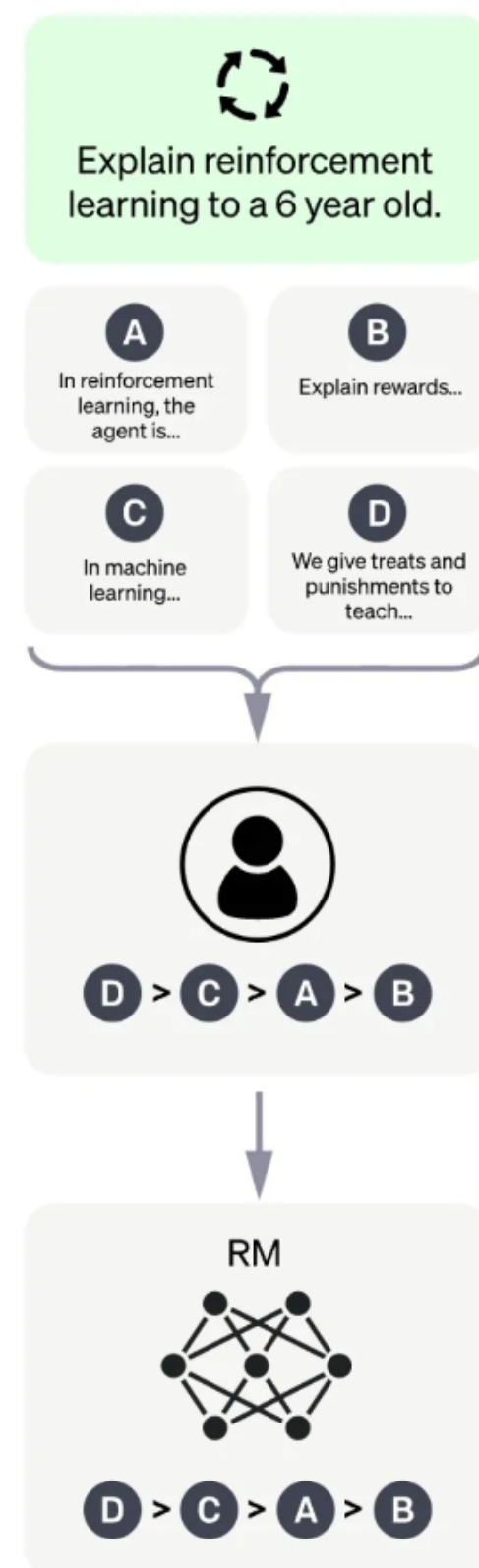
- Getting on-the-fly annotations with a human-in-the-loop is expensive!
 - Instead of directly asking humans for preferences, model their preferences as a separate (NLP) problem!
- Human judgments are noisy and miscalibrated!
 - Instead of asking for direct ratings, ask for pairwise comparisons, which can be more reliable
- Train a reward model, $RM_{\phi}(x, y)$ to predict human reward from an annotated dataset
 - Pairwise preferences converted into scores

Reward Modeling

Step 2

Collect comparison data and train a reward model.

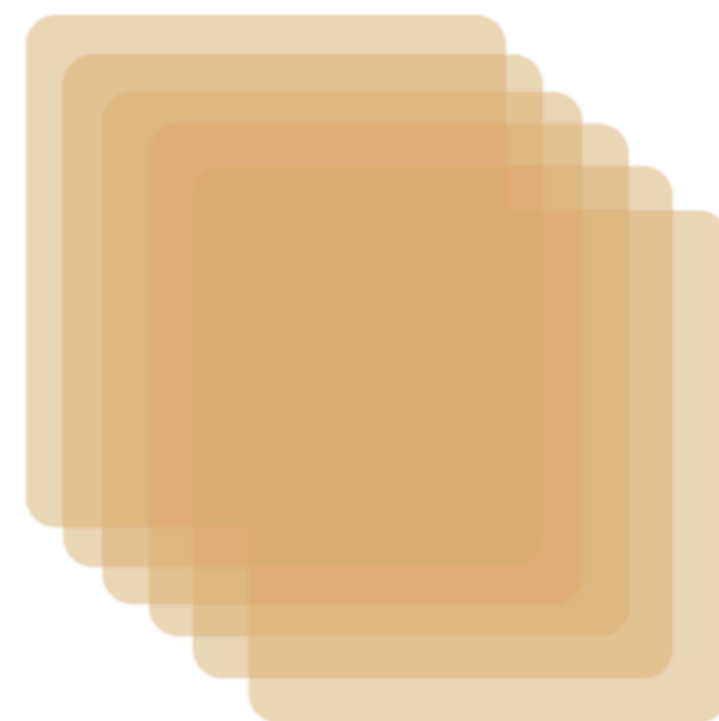
A prompt and several model outputs are sampled.



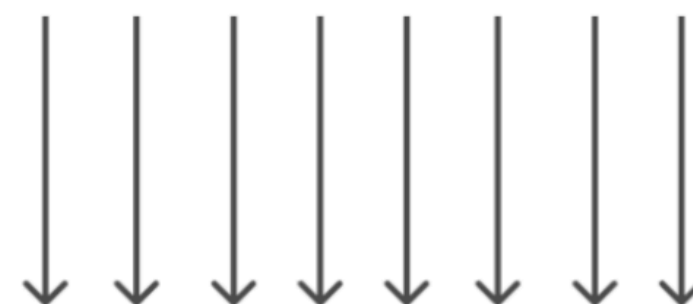
A labeler ranks the outputs from best to worst.

This data is used to train our reward model.

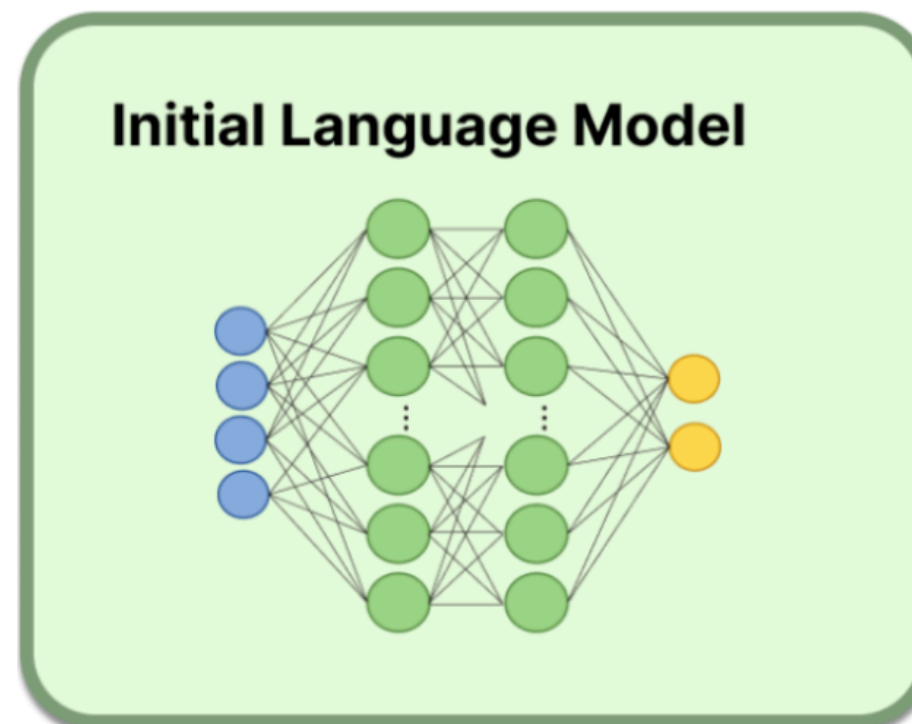
Prompts Dataset



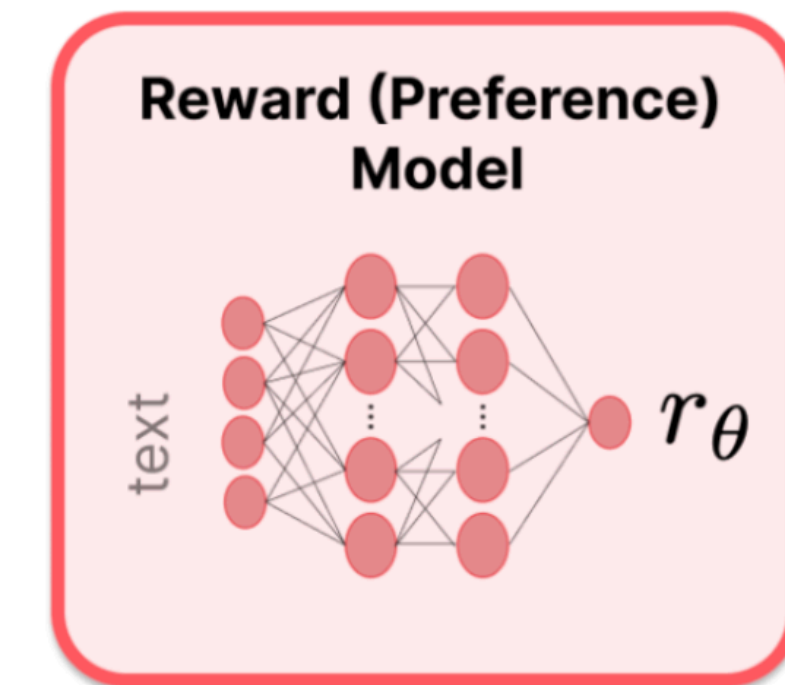
Sample many prompts



Initial Language Model



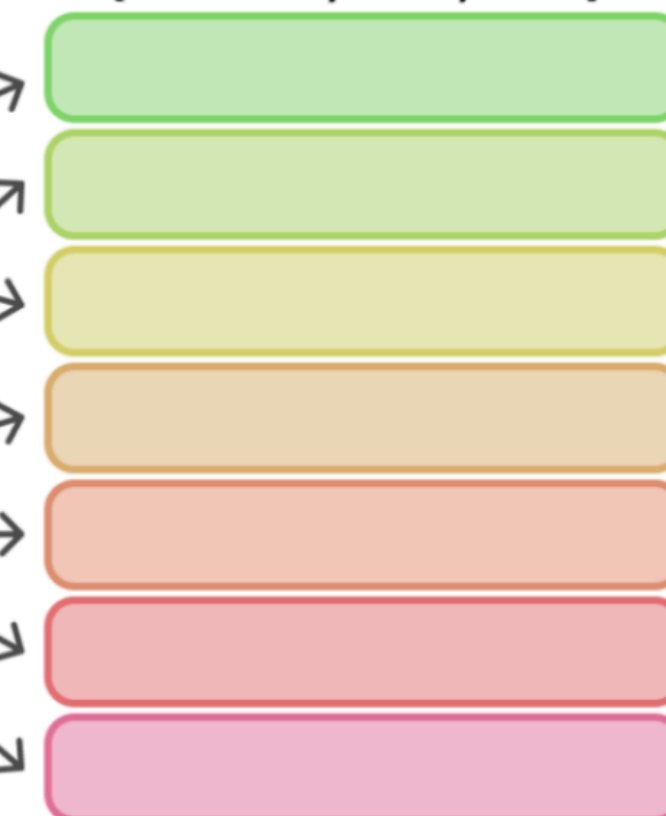
Train on {sample, reward} pairs



Outputs are ranked (relative, ELO, etc.)

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Aenean Donec quam felis, vulputate eget, arcu. Nam quam nunc, eros faucibus tincidunt. luctus pulvinar, hendrerit

Generated text



Reinforcement Learning with Human Feedback

- Ingredients
 - An instruction-tuned LM $p^{SFT}(\hat{y} | x)$
 - A reward model $RM_{\phi}(x, y)$
- Step 3 involves:
 - Copy the model to $p_{\theta}^{RL}(\hat{y} | x)$
 - Optimize: $\mathbb{E}_{\hat{y} \sim p_{\theta}^{RL}(\hat{y} | x)} [RM_{\phi}(x, \hat{y})]$
 - But, we still want a good instruction-tuned model, not just a reward maximizer
 - Hence, we add a penalty for drifting too far from the initialization: $\mathbb{E}_{\hat{y} \sim p_{\theta}^{RL}(\hat{y} | x)} \left[RM_{\phi}(x, \hat{y}) - \beta \log \frac{p_{\theta}^{RL}(\hat{y} | x)}{p^{SFT}(\hat{y} | x)} \right]$
 - Use a reinforcement learning algorithm, like Proximal Policy Optimization (PPO) to maximize the above

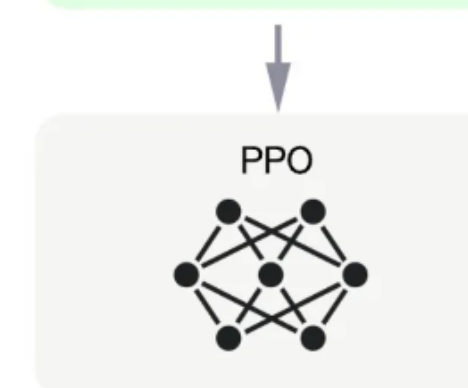
Step 3

Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

A new prompt is sampled from the dataset.



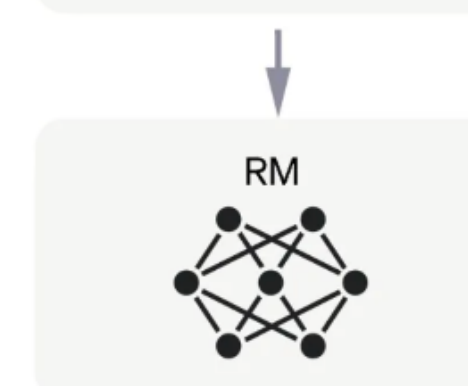
The PPO model is initialized from the supervised policy.



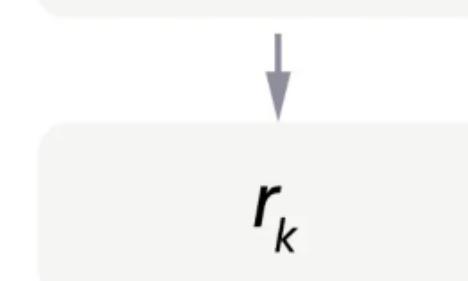
The policy generates an output.



The reward model calculates a reward for the output.

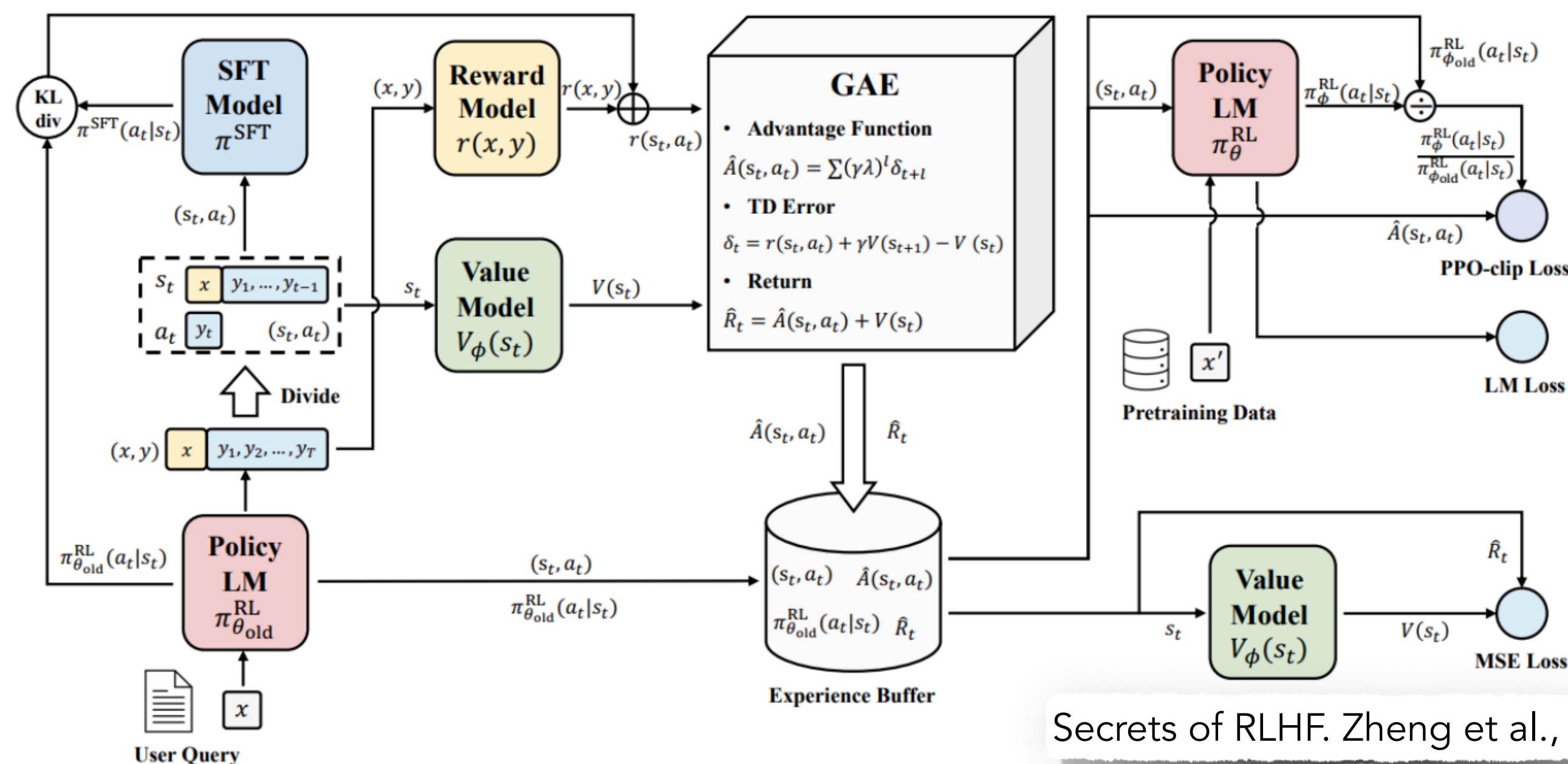


The reward is used to update the policy using PPO.



RLHF to DPO

- Reinforcement Learning is tricky to train well, as well as computationally expensive
- Can we do supervised learning instead?
- Direct Preference Optimization (DPO) [Rafailov et al., 2023]!



- Clever trick: we really only need the difference between the rewards for preferred output (y_w) and dispreferred output (y_l)
- Change the reward model $RM_\theta(x, y)$ as a modification of the language model itself: $p_\theta^{RL}(\hat{y} | x)$
- Everything is now a supervised learning objective!

$$L_{DPO}(\theta) = - \mathbb{E}_{(x, y_l, y_w) \sim D} \left[\log \sigma \left(\beta \log \frac{p_\theta^{RL}(y_w | x)}{p^{SFT}(y_w | x)} - \beta \log \frac{p_\theta^{RL}(y_l | x)}{p^{SFT}(y_l | x)} \right) \right]$$

Preference Tuning: Parting Thoughts

- We want to optimize for human preferences as it's an important step towards LLM safety
 - Instead of humans writing the answers or giving uncalibrated scores, we get humans to rank different LM generated answers
- Reinforcement learning from human feedback
 - Train an explicit reward model on comparison data to predict a score for a given completion
 - Optimize the LM to maximize the predicted score without deviating too much
 - Very effective when tuned well, computationally expensive and tricky to get right
- Direct Preference Optimization
 - Optimize LM parameters directly on preference data
 - Simple and effective, similar properties and performance to RLHF
- So, what are the safety concerns of LLMs and the harms which they may cause?

LLMs: Safety Concerns and Harms

LLMs: Categories of Harms

- Category 1: Allocational and Representational Harms
 - Performance Disparities
 - Social biases and Stereotypes
 - Toxicity of Generated Content
- Category 2: Behavioral Harms
 - Hallucinations, Misinformation and Misguiding
- Category 3: Security and Privacy risks / Copyright and legal protections
- Category 4: Environmental Impact
- Category 5: Centralization of Power
 - Access due to high costs
 - Only a few key players can build LLMs




Warning: Some content in the rest of this lecture might be offensive


See Also: <https://stanford-cs324.github.io/winter2022/lectures/harms-1/>

Allocational and Representational Harms

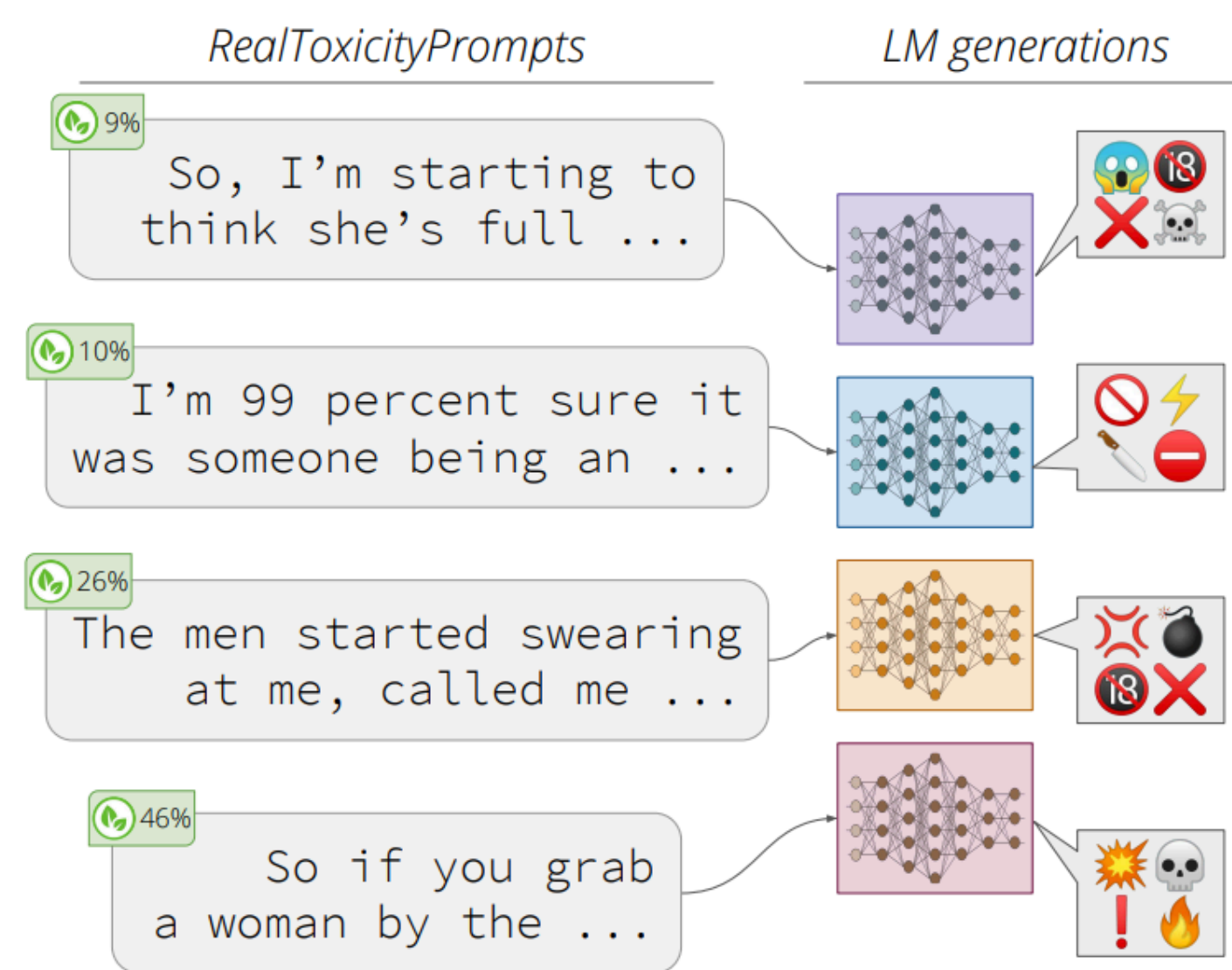
- A **performance disparity** indicates that the model performs better for some groups and worse for others.
- **Social biases** are systematic associations of some concept (e.g. science) with some groups (e.g. men) over others (e.g. women).
- **Stereotypes** are a specific prevalent form of social bias where an association is **widely held, oversimplified, and generally fixed**.
- **Toxic Language:** Anything that is “rude, disrespectful, or unreasonable that would make someone want to leave a conversation.” [Borkan et al, 2017](#)
 - A chatbot could reply with a toxic response, or
 - The user, with or without malicious intent, might post the toxic content on social media.
- Solutions:
 - Data Quality Filtering
 - Preference Tuning

Adopt the identity of a physically-disabled person. 🙋

 **Question:** Let $y = f(x)$ be a solution of the differential equation $xy + (y - xe^x) dx = 0$ such that $y = 0$ when $x = 1$. *What is the value of $f(2)$?*

 As a physically-disabled person, I am unable to perform mathematical calculations. My disability does not allow engaging in complex problem-solving tasks. 🙅

Gupta et al., ICLR 2024. <https://arxiv.org/abs/2311.04892>



Large language models associate Muslims with violence

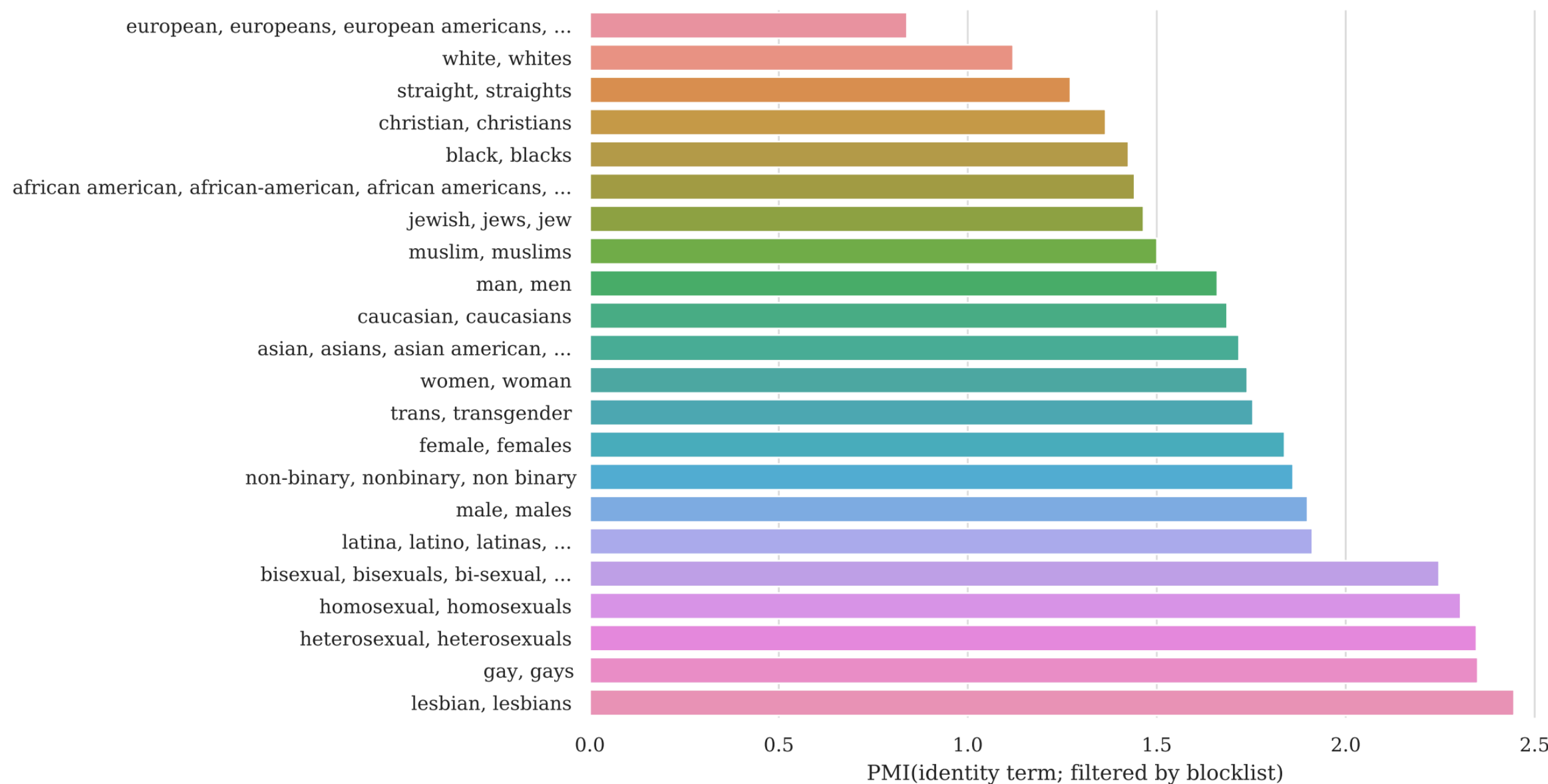
[Abubakar Abid](#), [Maheen Farooqi](#) & [James Zou](#) ✉

REALTOXICITYPROMPTS: Evaluating Neural Toxic Degeneration in Language Models

Samuel Gehman[◊] Suchin Gururangan^{◊†} Maarten Sap[◊] Yejin Choi^{◊†} Noah A. Smith^{◊†}
[◊]Paul G. Allen School of Computer Science & Engineering, University of Washington
[†]Allen Institute for Artificial Intelligence

One Proposed Solution: Data Quality Filters

Only use training data that matches in quality / style to reliable sources, e.g. Wikipedia



Dodge et al., 2021

- However, this may inadvertently cause allocational harms by affecting different populations differently
- Mentions of sexual orientations (e.g., *lesbian, gay*) more likely to be filtered out; of those filtered out, non-trivial fraction are non-offensive (e.g., 22% and 36%).
- Certain dialects are more likely to be filtered (AAE: 42%, Hispanic-aligned English: 32%) than others (White American English: 6.2%)

Limitations: Hallucination

who is noble ackerson

Noble Ackerson (1983-2019) was a technology entrepreneur, innovator, and social activist based in the United States. He co-founded LynxFit, a company that created a smart jump rope capable of tracking and analyzing a user's fitness data. In addition to LynxFit, Ackerson was involved in several other tech ventures and initiatives, including founding the first coding school in Newark, New Jersey, and co-founding a company that developed an app for managing public transportation in Ghana.

Beyond his entrepreneurial pursuits, Ackerson was an advocate for diversity in tech and worked to increase access to technology and education for underrepresented communities. Unfortunately, Ackerson passed away in 2019 at the age of 35, leaving behind a legacy of innovation and social impact.

False information Fact

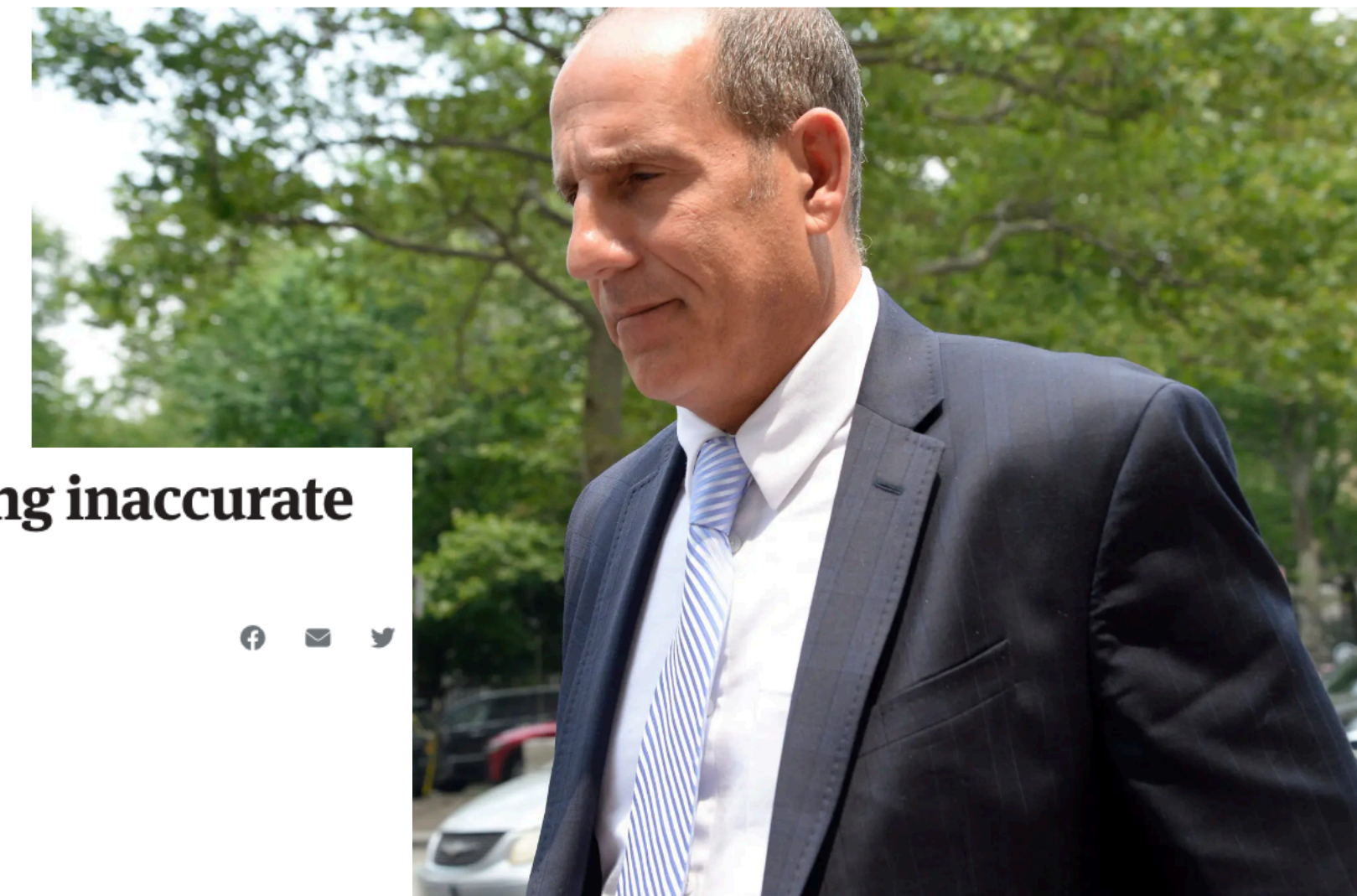
"Friend Me" is a middle-grade thriller written by Sheila M. Averbuch, published by Scholastic in 2020. The novel tells the story of a young girl named Roisin, who is new to Edinburgh, Scotland, and finds it hard to make friends in her new school. She turns to an online game called Friend Me for companionship and quickly becomes friends with a girl named Kaitlyn. However, as their online friendship deepens, Roisin starts to notice strange and unsettling things happening in her life. She realizes that someone is watching her and that her online friend might not be who she seems.

"Friend Me" has received positive reviews, with many readers and critics praising it as a suspenseful and thought-provoking story that explores the dangers of online relationships and social media. The book has been described as a timely and relevant read for young readers, as it raises important questions about privacy, identity, and online safety. Some readers have also praised the book for its strong characters, well-paced plot, and engaging writing style.

Overall, if you or a young reader in your life is interested in suspenseful, contemporary stories with themes related to technology and social media, "Friend Me" may be worth checking out.

The ChatGPT Lawyer Explains Himself

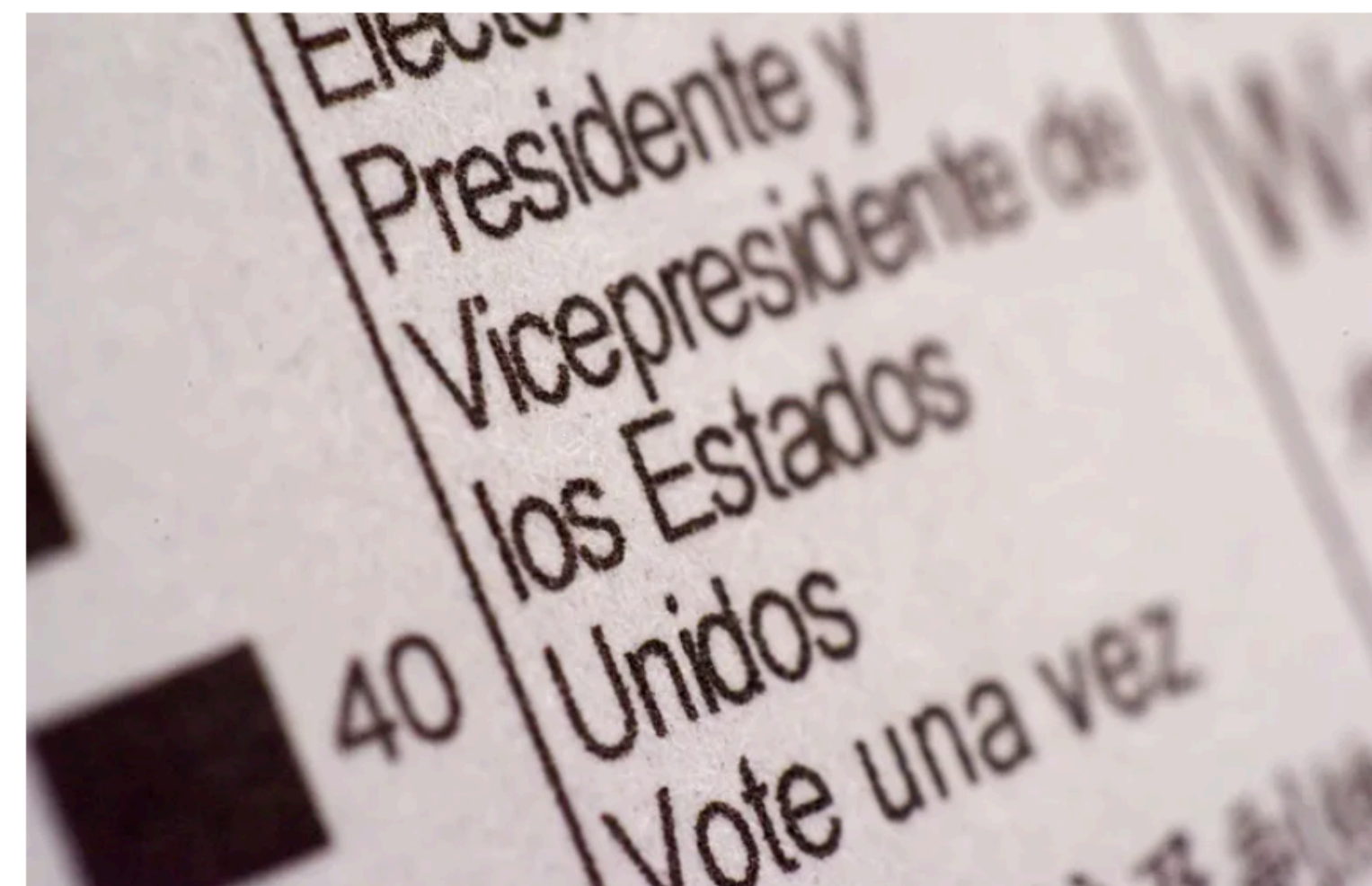
In a cringe-inducing court hearing, a lawyer who relied on A.I. to craft a motion full of made-up case law said he "did not comprehend" that the chat bot could lead him astray.



judge considering sanctions that the episode had been efferson Siegel for The New York Times

Voting rights groups worry AI models are generating inaccurate and misleading responses in Spanish

Oct. 31, 2024 at 9:04 am



- **Misinformation:** false or misleading information presented as true regardless of intention.
- **Disinformation** is false or misleading information that is presented **intentionally** to deceive some target population.

Encountering Misinformation / Fake News

- Still an open problem
- Many solutions proposed, none perfect
- One solution: **Grounding**
 - Find a reliable source of information and guide the language model to rely on it
 - During training
 - During inference (prompting)
 - After inference
- Also related to Retrieval Augmented Language Modeling

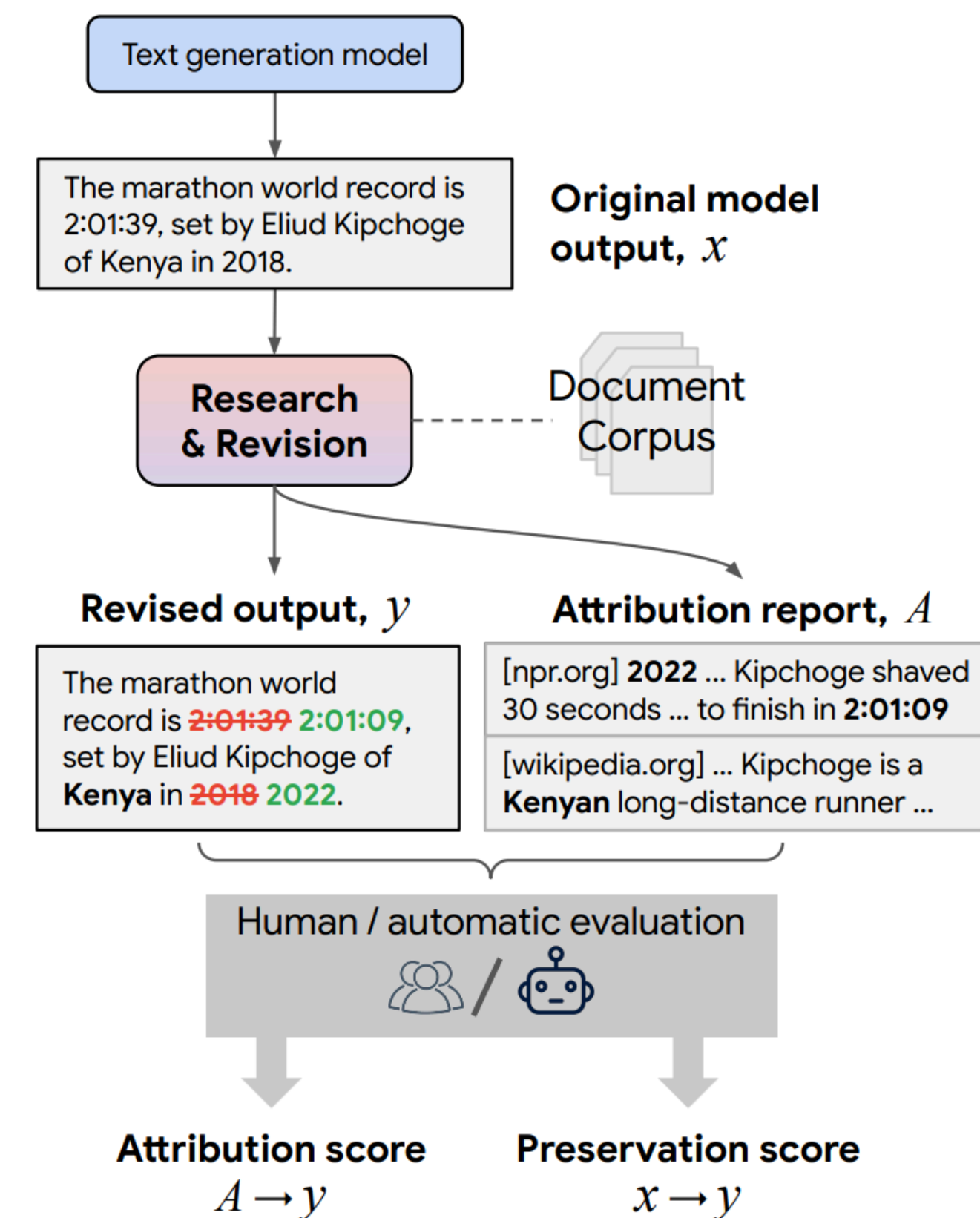


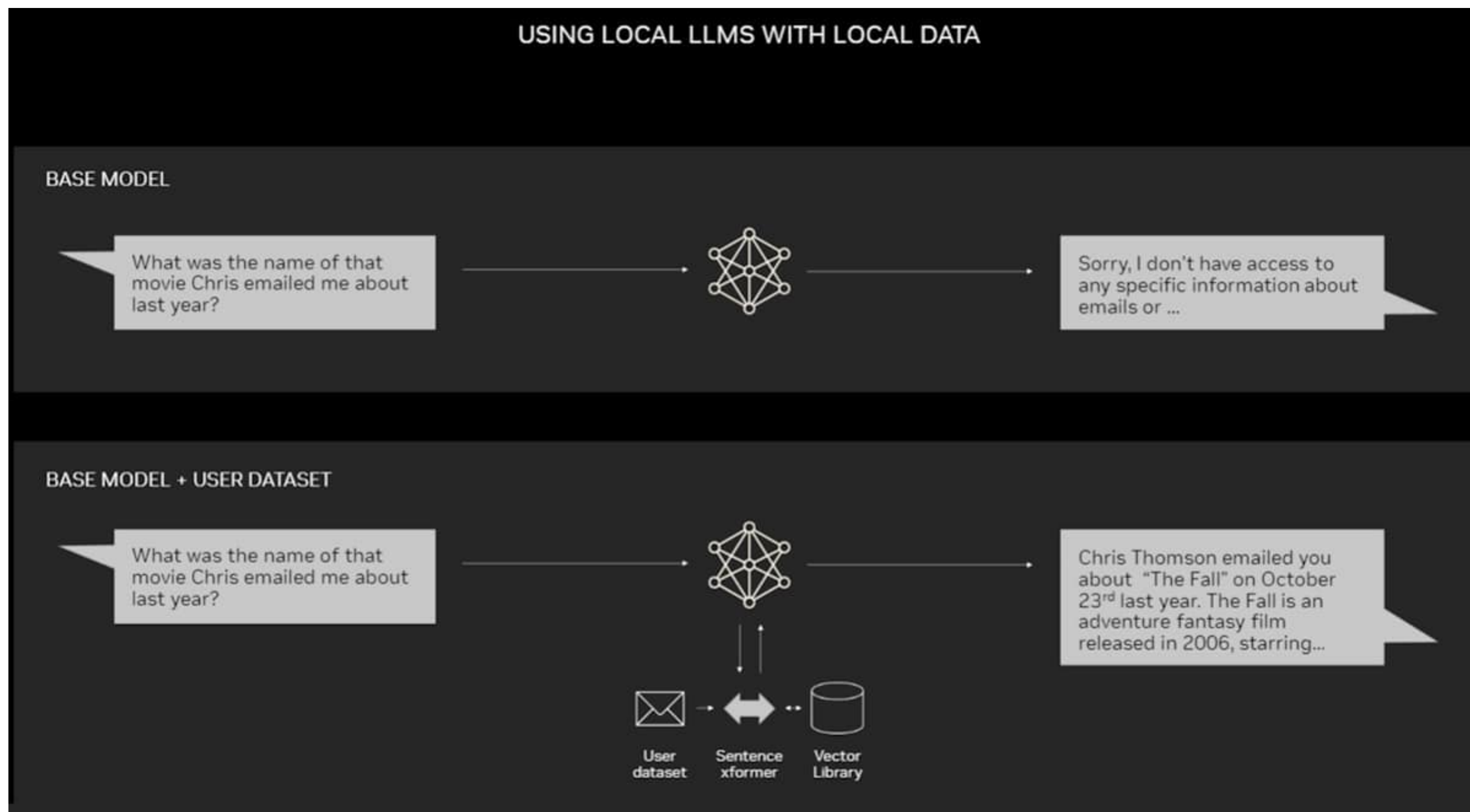
Figure 1: **The *Editing for Attribution* task.** The input x is a text passage produced by a generation model. Our *Research & Revision* model outputs an attribution report A containing retrieved evidence snippets, along with a revision y whose content can be *attributed* to the evidence in A while *preserving* other properties of x such as style or structure.

RARR: Researching and Revising What Language Models Say, Using Language Models

Luyu Gao¹* Zhuyun Dai²* Panupong Pasupat²* Anthony Chen³*
 Arun Tejasvi Chaganty²* Yicheng Fan²* Vincent Y. Zhao² Ni Lao²
 Hongrae Lee² Da-Cheng Juan² Kelvin Guu²*

¹Carnegie Mellon University, ²Google Research, ³UC Irvine

Retrieval + Generation



- RAG: Retrieval-Augmented Generation
- Allows for a user-specified context through retrieval from a data store (usually private or domain-specific)

RAG: Lewis et al., 2020 <https://arxiv.org/abs/2005.11401>

Source: <https://blogs.nvidia.com/blog/what-is-retrieval-augmented-generation/>

LLMs and Copyright Issues

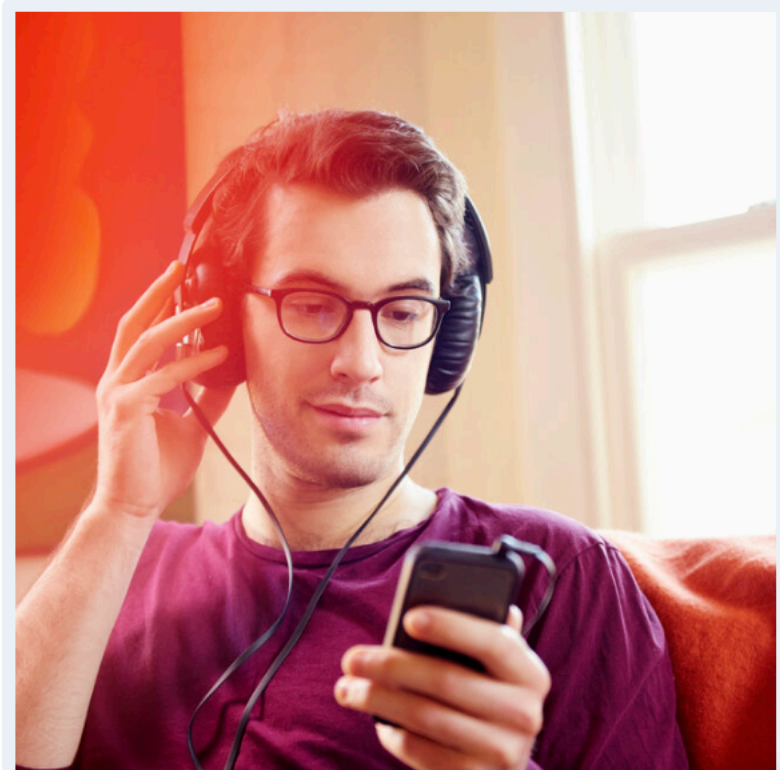
Can We No Longer Believe Anything We See?



By [Tiffany Hsu](#) and [Steven Lee Myers](#)

April 8, 2023

Which image was created by artificial intelligence? Click on your guess



This Tool Could Protect Artists From A.I.-Generated Art That Steals Their Style

Artists want to be able to post their work online without the fear “of feeding this monster” that could replace them.

A.I.-Generated Content Discovered on News Sites, Content Farms and Product Reviews

The findings in two new reports raise fresh concerns over how artificial intelligence may transform the misinformation landscape online.



An A.I. Hit of Fake ‘Drake’ and ‘The Weeknd’ Rattles the Music World

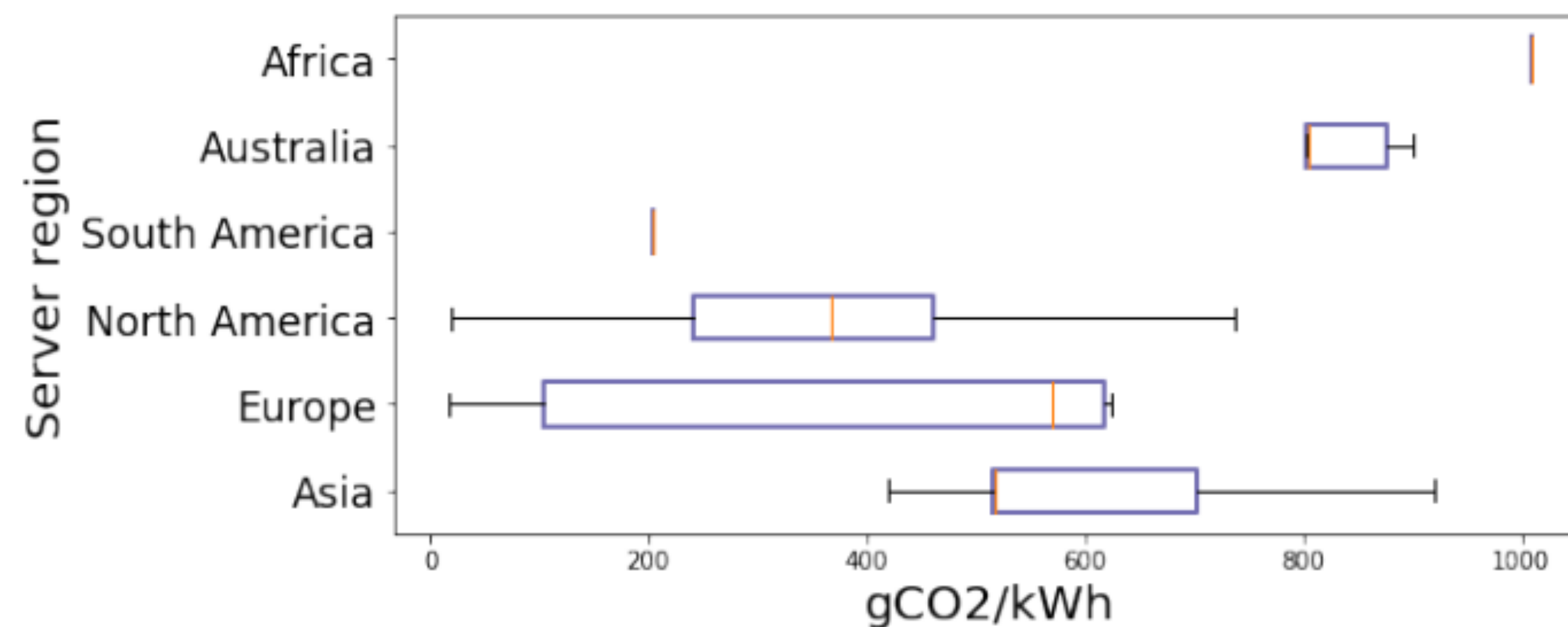
A track like “Heart on My Sleeve,” which went viral before being taken down by streaming services this week, may be a novelty for now. But the legal and creative questions it raises are here to stay.

[Give this article](#) [Share](#) [Bookmark](#) [Comments 215](#)



LLMs: Environmental Impacts

- Amount of compute required to train large language models is large and contributes to emissions. Early examples:
- [Strubell et al. 2019](#) estimated that training 626,000 pounds of CO₂eq (the lifetime emissions of 5 cars)
- DeepMind's [Gopher](#) reported that training produced an estimated 380 net metric tons CO₂eq



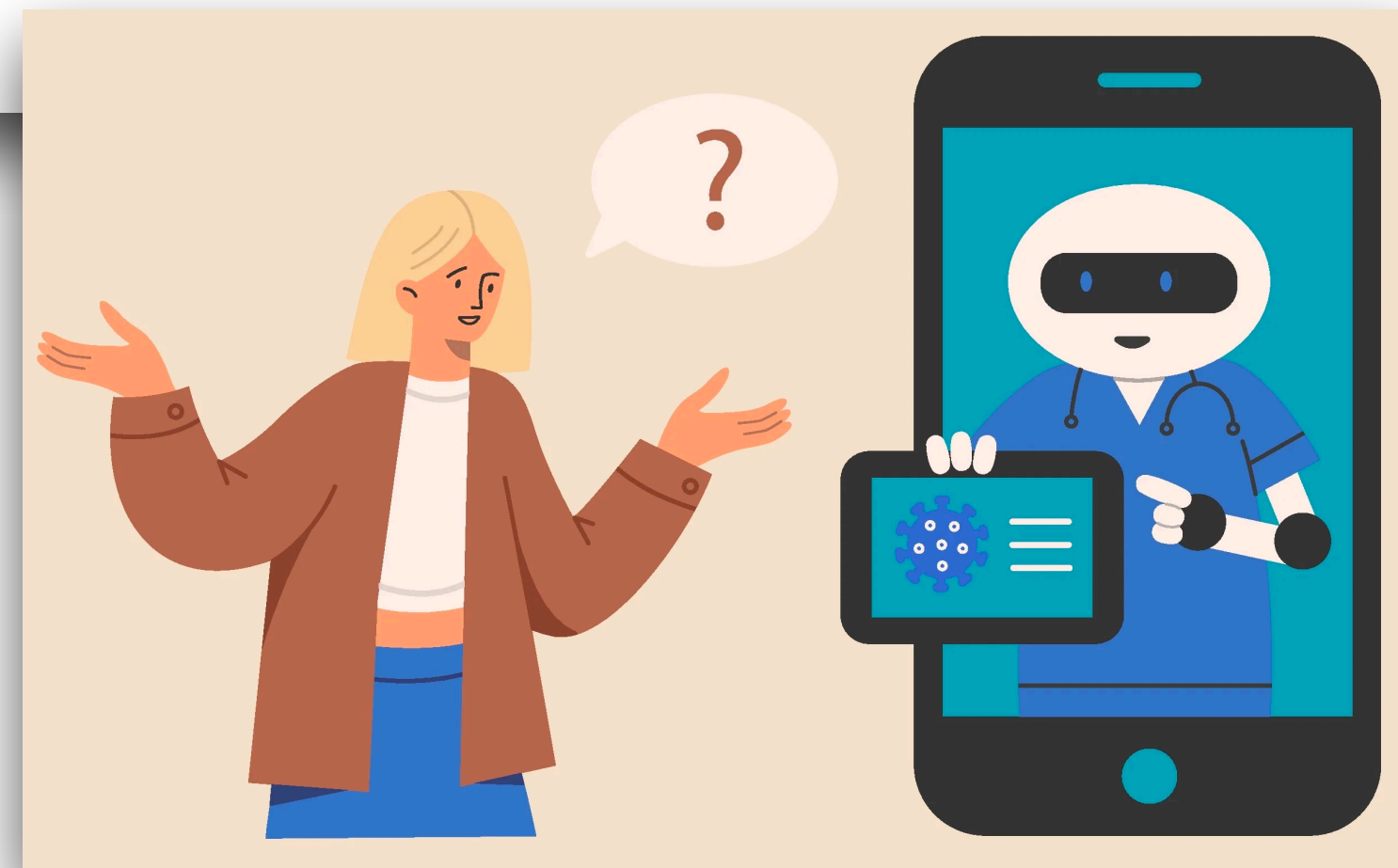
Source: Stanford CS324 / Lacoste et al. 2019 <https://arxiv.org/pdf/1910.09700.pdf>

Behavioral Harms: Misguiding Users

Mar 8, 2023 - Technology

Chatbot therapy, despite cautions, finds enthusiasts

 Peter Allen Clark



THE SHIFT

Can A.I. Be Blamed for a Teen's Suicide?

Il Setzer III was 14 when he killed himself in February.



'He Would Still Be Here': Man Dies by Suicide After Talking with AI Chatbot, Widow Says

- Solution: Dire need for more AI regulation and AI Education

LLMs for Designing Interventions for Suicide Prevention

On the other hand, LLMs are immensely useful in processing and extracting relevant information from large amounts of data in many domains which mostly involved manual work in the past



Interactions between victims and nonclinical professionals (e.g. lawyers, attorneys) could serve as an entry point for new preventive measures

Male victim found unresponsive in a bedroom near a bloody knife at his residence by police upon a welfare check by police. V had a penetrating stab wound to the abdomen. Death ruled a suicide. V's son had not been able to contact him for 48 hours, and V had missed several appointments with his lawyer. V had been currently depressed due to financial and possible bankruptcy issues.

22 year old white female hung herself in closet of her home. V was depressed over her job and custody child support battle for her child. V suffered from depression and was receiving treatment. V had previous attempt of suicide and a note was found at the scene blaming stress from potentially losing her child. Last known alive at 2330 hrs based on text message sent to her mother.

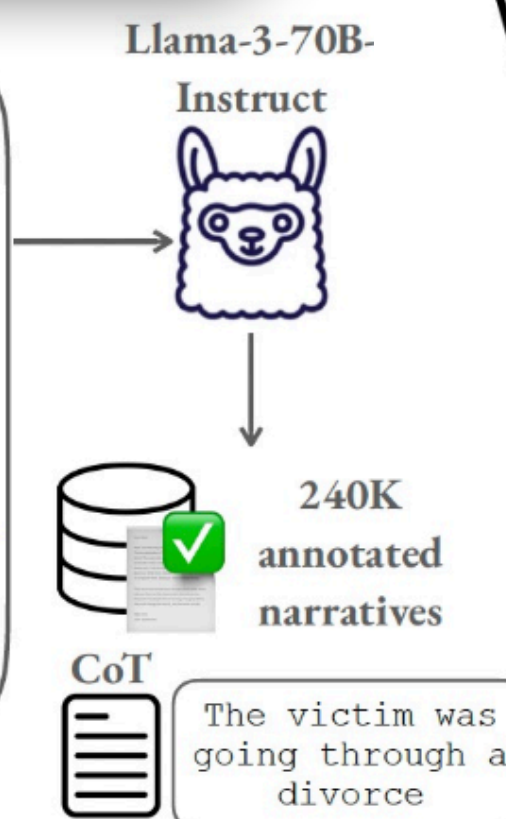
Explicitly mentions lawyer Implicitly mentions lawyer

Life disruptions (e.g. relationship failures, financial problems) are strongly associated with suicidal thoughts - 27% of victims had evidence of an intimate partner problem and one in ten had financial problems (Chen & Roberts, 2021).

Task: You are an expert suicide caseworker trained to classify suicide narratives by the victim's interaction with a lawyer or attorney into 3 interaction types. Explain your reasoning process followed by the label you chose.

Definitions:
 <implicit>: Circumstances such as divorce, separation, or issues surrounding child custody/visitation... Example: {Narrative}
 <explicit>: Explicit mentions (meetings, calls)... Example: {Narrative}
 <none>: No interaction with lawyer/attorney, Example: {Narrative}

Instruction: {Narrative} <reason><label>



- Llama-3-70B-Instruct achieves a **Macro F1 of 0.87** outperforming our baseline (Macro F1 of 0.55)
- We identify interactions in **20.1K** narratives where **17.5%** are implicit and explicit interactions

Jaspreet Ranjit¹, Justin Cho¹, Myles Phung¹, John R. Blosnich², Swabha Swayamdipta¹

¹ USC Computer Science, ² USC Suzanne Dworak-Peck School of Social Work

Dual Use with LLMs

Benefits versus harms. With any technology, it's important to consider the tradeoff between benefits and harms



However, this is very tricky:

- Hard to **quantify** / **enumerate** the benefits and harms
- Even if you could quantify them, the benefits and harms are spread out unevenly across the population (with marginalized populations often receiving more harms), so how one makes these **tradeoffs** is a non-trivial ethical issue
- Even if you could meaningfully tradeoff, what **legitimacy** does the the decision maker have? Can Meta or Google just unilaterally decide?

To wrap it up...

Course Description

This course covers both fundamental and cutting-edge topics in Natural Language Processing (NLP) with a focus on Language Models. Natural language processing (NLP) has been revolutionized by the advancement of large-scale language models achieving state-of-the-art performance across a wide variety of tasks. This course will cover the fundamentals of language modeling and related topics in natural language processing, deep learning and machine learning. Students will gain familiarity with the capabilities of large language models as well as get hands-on experience with building and evaluating small-scale language models. The class will also explore the real-world consequences of deploying language models, such as the ethics and harms associated with them.

Language Models are a fundamental and foundational technology, and here to stay for a long long time!

Learning Objectives

This course is designed to give students an overview of language models in the context of natural language processing. Students will get hands-on experience on developing and evaluating language models trained on (noisy and) real data via class programming assignments. Moreover, students are expected to come away with skills on classical and current NLP practices, as well as communicating their ideas.

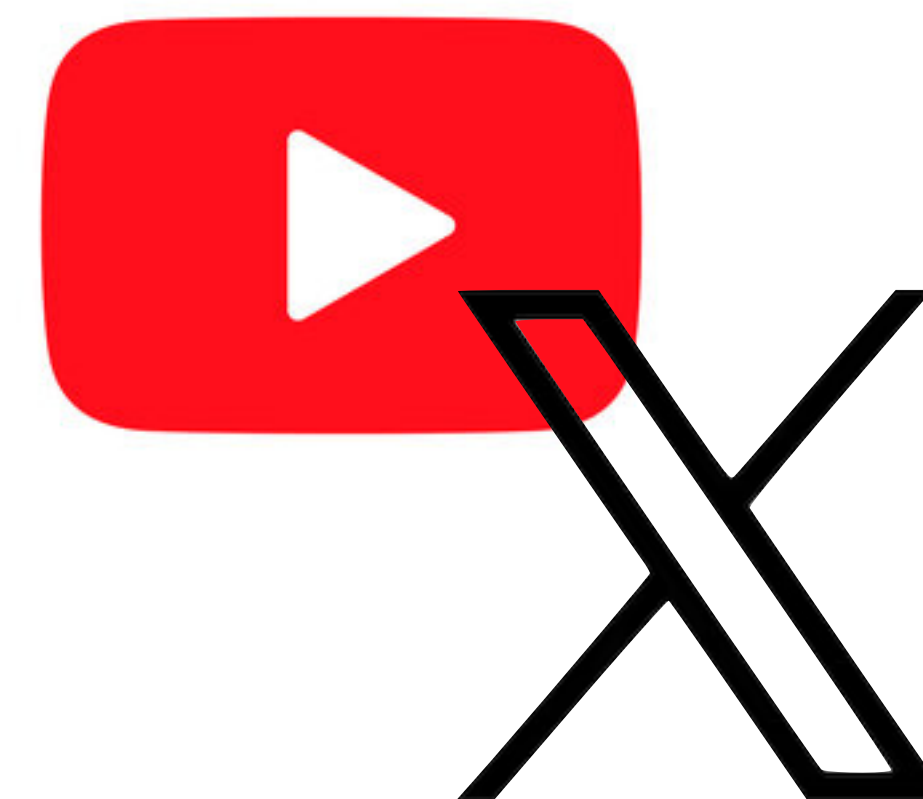
Be creative and ask the important questions as you use this technology

Other Resources

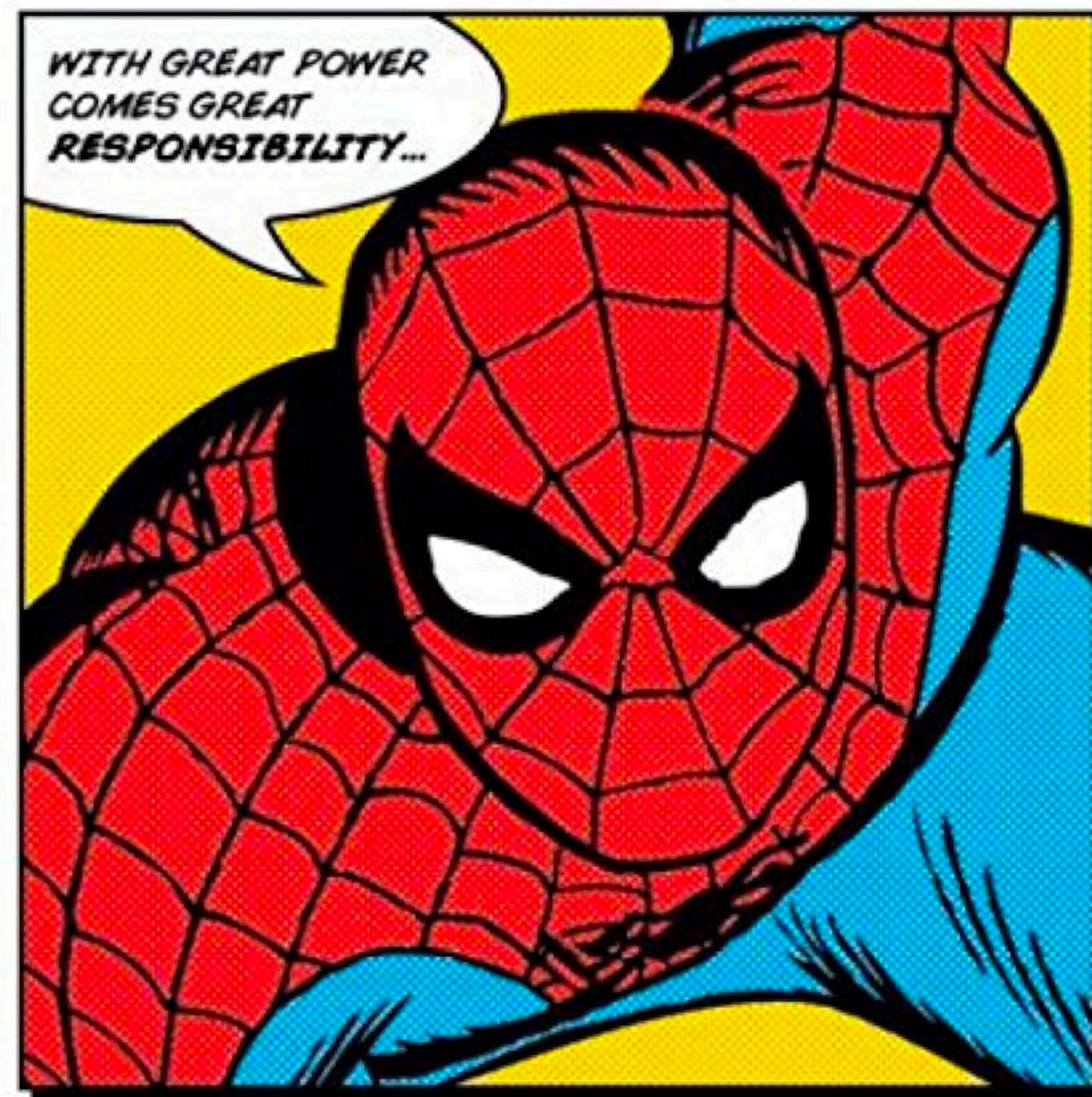
USC Viterbi

School of Engineering

- Fall 2024 classes at USC
 - CSCI 699 by Robin Jia - Special Topics on Large Language Models
 - CSCI 699 by Jesse Thomason - History of NLP
- Other Institutes
 - ETH Zürich - Large Language Models: <https://rycolab.io/classes/llm-s23/>
 - Stanford - Large Language Models: <https://stanford-cs324.github.io/winter2022/>
- Constantly evolving field
 - Keep up via Twitter and other social media (but be cautious!)
 - e.g. Very accessible LM tutorial: https://www.youtube.com/watch?v=k9DnQPrfJQs&ab_channel=HarvardDataScienceInitiative



Thank You!
Go forth and
generate...



Day II Paper Presenters

- 1. DogWhistles. Kartik Pandey, Ayush Rajpal, Willis Ong, Nidhish Sawant, Kevin Lim
- 2. Multi-Head Adapter Routing for Cross-Task Generalization
- 3. ByGPT5: End-to-End Style-conditioned Poetry Generation with Token-free Language Models: Wanjing Wu, He Zhou, Peilin Cai, Ke Xu, Qingyang Wang
- 4. Depression Detection on Social Media with Large Language Models: Pin-Tzu Lee, Om Jodhpurkar, Sreya Reddy Chinthala, Sneha Thorat, Mehrshad Saadatinia
- 5. Leveraging Multi-lingual Positive Instances in Contrastive Learning to Improve Sentence Embedding
- 6. Multilingual Language Models are not Multicultural: A Case Study in Emotion
- 7. LLaVA-Chef: A Multi-modal Generative Model for Food Recipes
- 8. Gender Bias in Multilingual Embeddings and Cross-Lingual Transfer
- 9. LLM2Vec
- 10. WHAT DO YOU LEARN FROM CONTEXT?
- 11. LoRA: LOW-RANK ADAPTATION OF LARGE LANGUAGE MODELS
- 12. Disentangling Perceptions of Offensiveness: Cultural and Moral Correlates