

Ethics II

CSE 5525: Foundations of Speech and Natural Language Processing

<https://shocheen.github.io/courses/cse-5525-spring-2026>



THE OHIO STATE UNIVERSITY

Logistics

- Final quiz on Wed, April 15.
 - Will announce the reading on Canvas/Teams later today.

- Project presentations
 - April 22, 24
 - Will announce your slots tomorrow.
 - If you can't make it to your slots, please try and find a team in the other slot to swap with. If you can't, email me.

Recap

- What is ethics
 - A study of what is good vs what is bad.
- Ethics in AI / NLP
 - AI systems can cause harm or be used to cause harm

Today: a story in two parts

Biased Outputs

Stereotypical
behavior

Malfunction on
minority inputs

...

Harmful Content

Toxicity

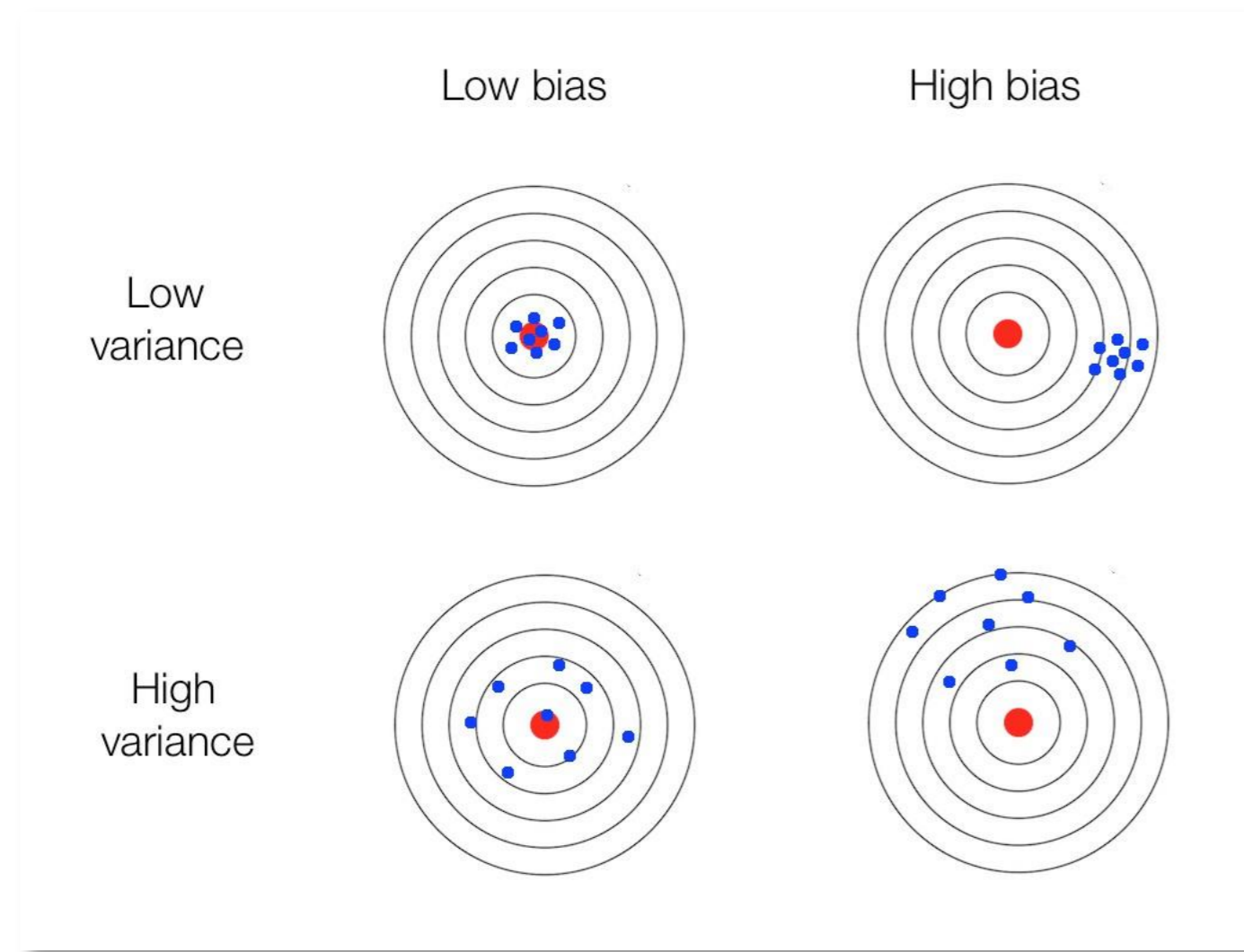
Unsafe content

...

Part 1 – Bias

Some definitions of bias

- Bias [*statistics*]: systematic tendency causing differences between model estimates / predictions



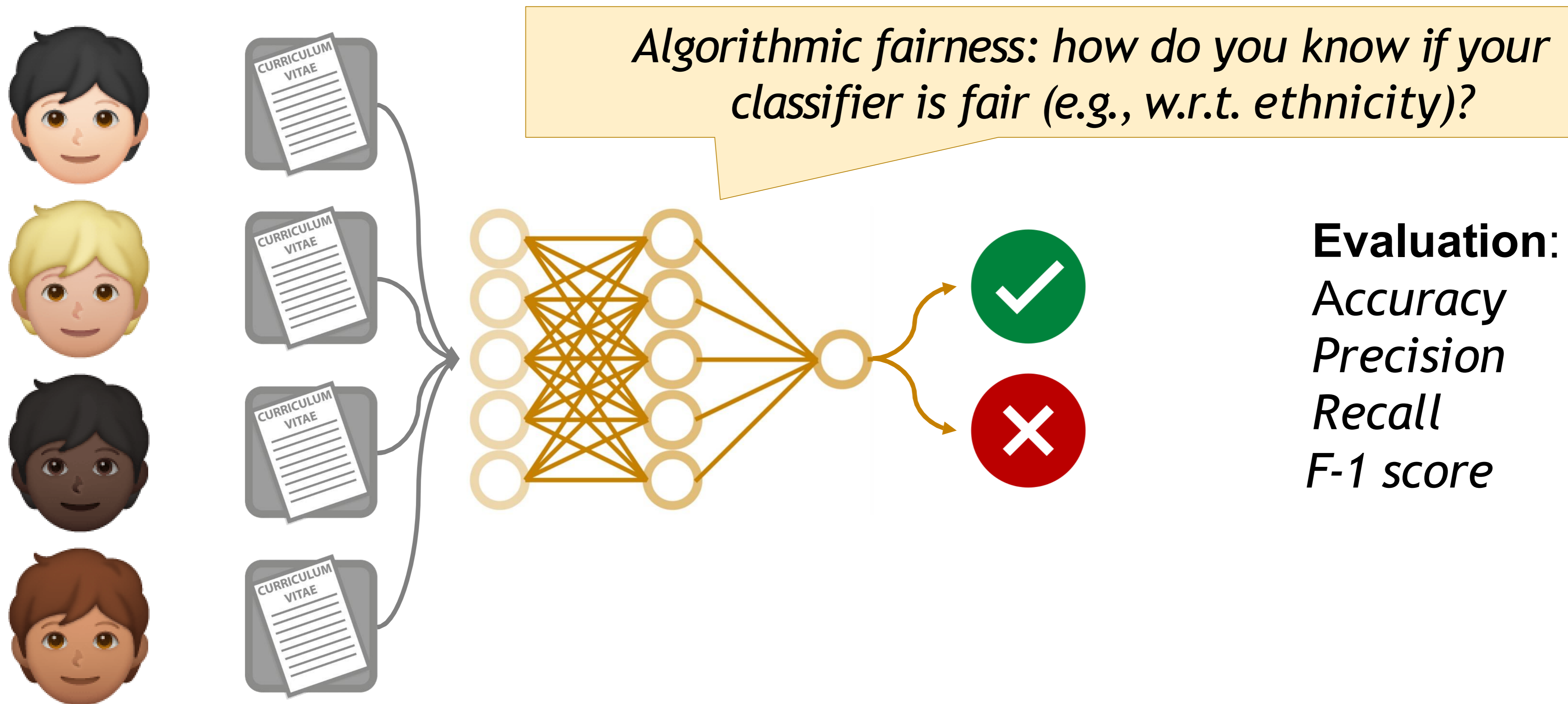
- Bias [*general*]: “disproportionate weight in favor of or against an idea or thing, usually in a way that is closed-minded, prejudicial, or unfair” – Wikipedia



Presence of bias \simeq absence of fairness
Algorithmic fairness: attempts to correct biases in ML systems
But... how is fairness defined?

Algorithmic fairness

Let's assume a toy task: given a resumé, predict whether a candidate is qualified



Fairness metrics

- Accuracy quality: a classifier is fair if the people from different groups have the same accuracy


 *Accuracy*



 *Accuracy*



 *Accuracy*


 *Accuracy*

- Statistical parity: groups should have the same probability of being assigned positive class

$p(\checkmark | \text{Black male emoji})$


$p(\checkmark | \text{Blond male emoji})$


$p(\checkmark | \text{Black female emoji})$


$p(\checkmark | \text{Brown male emoji})$

Equalized odds criterion [\[Hardt et al '16\]](#)

A classifier c is fair if the *false positive (FP)* and *true positive (TP)* rates are the same for different groups

- False positives

$$p(c(\text{👦}) = \checkmark \mid l(\text{👦}) = \times) \\ =$$

$$p(c(\text{👦}) = \checkmark \mid l(\text{👦}) = \times)$$

- True positives

$$p(c(\text{👦}) = \checkmark \mid l(\text{👦}) = \checkmark) \\ =$$

$$p(c(\text{👦}) = \checkmark \mid l(\text{👦}) = \checkmark)$$

		Predicted class	
		Positive	Negative
Actual class	Positive	TP	FN
	Negative	FP	TN

Other fairness metrics

- *Treatment equality*
 - Ratio of false negatives and false positives should be the same for groups
- *Fairness through unawareness*
 - Models should not employ sensitive attributes when making decisions
- *Causality-based*
 - *Counterfactual fairness*: outcome of the classifier would not changed if the sensitive attribute (e.g., race) were the only thing changed
- Many more...
 - [https://en.wikipedia.org/wiki/Fairness_\(machine_learning\)](https://en.wikipedia.org/wiki/Fairness_(machine_learning))
 - <https://fairmlbook.org/>

FAIRNESS AND MACHINE LEARNING

Limitations and Opportunities

Solon Barocas, Moritz Hardt, Arvind Narayanan

CONTENTS

PREFACE

ACKNOWLEDGMENTS

1 INTRODUCTION [PDF](#)

2 WHEN IS AUTOMATED DECISION MAKING LEGITIMATE? [PDF](#)

We explore what makes automated decision making a matter of normative concern, situated in bureaucratic decision making and its mechanical application of formalized rules.

3 CLASSIFICATION [PDF](#)

We introduce formal non-discrimination criteria in a decision-theoretic setting, establish their relationships, and illustrate their limitations.

4 RELATIVE NOTIONS OF FAIRNESS [PDF](#)

Other fairness metrics

- *Treatment equality*
 - Ratio of false negatives and false positives should be the same for groups
- *Fairness through unawareness*
 - Models should not employ sensitive attributes when

FAIRNESS AND MACHINE LEARNING

Limitations and Opportunities

Solon Barocas, Moritz Hardt, Arvind Narayanan

- **But, do these definitions really matter if no harms are caused? Many argue that unfairness/bias should be measured in terms of the harms that it causes**

- Many more...
 - [https://en.wikipedia.org/wiki/Fairness_\(machine_learning\)](https://en.wikipedia.org/wiki/Fairness_(machine_learning))
 - <https://fairmlbook.org/>

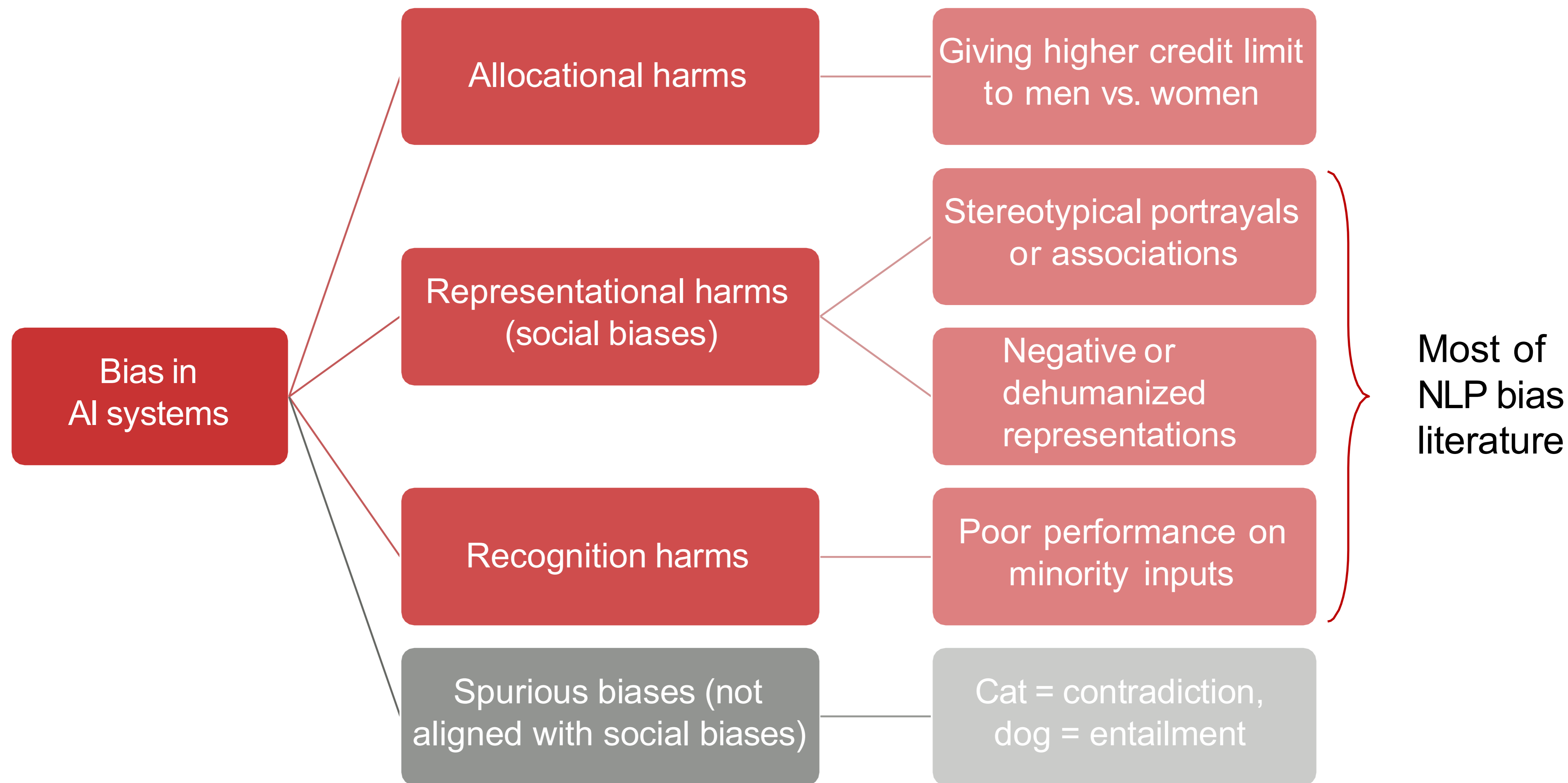
We explore what makes automated decision making a matter of normative concern, situated in bureaucratic decision making and its mechanical application of formalized rules.

3 [CLASSIFICATION](#) [PDF](#)

We introduce formal non-discrimination criteria in a decision-theoretic setting, establish their relationships, and illustrate their limitations.

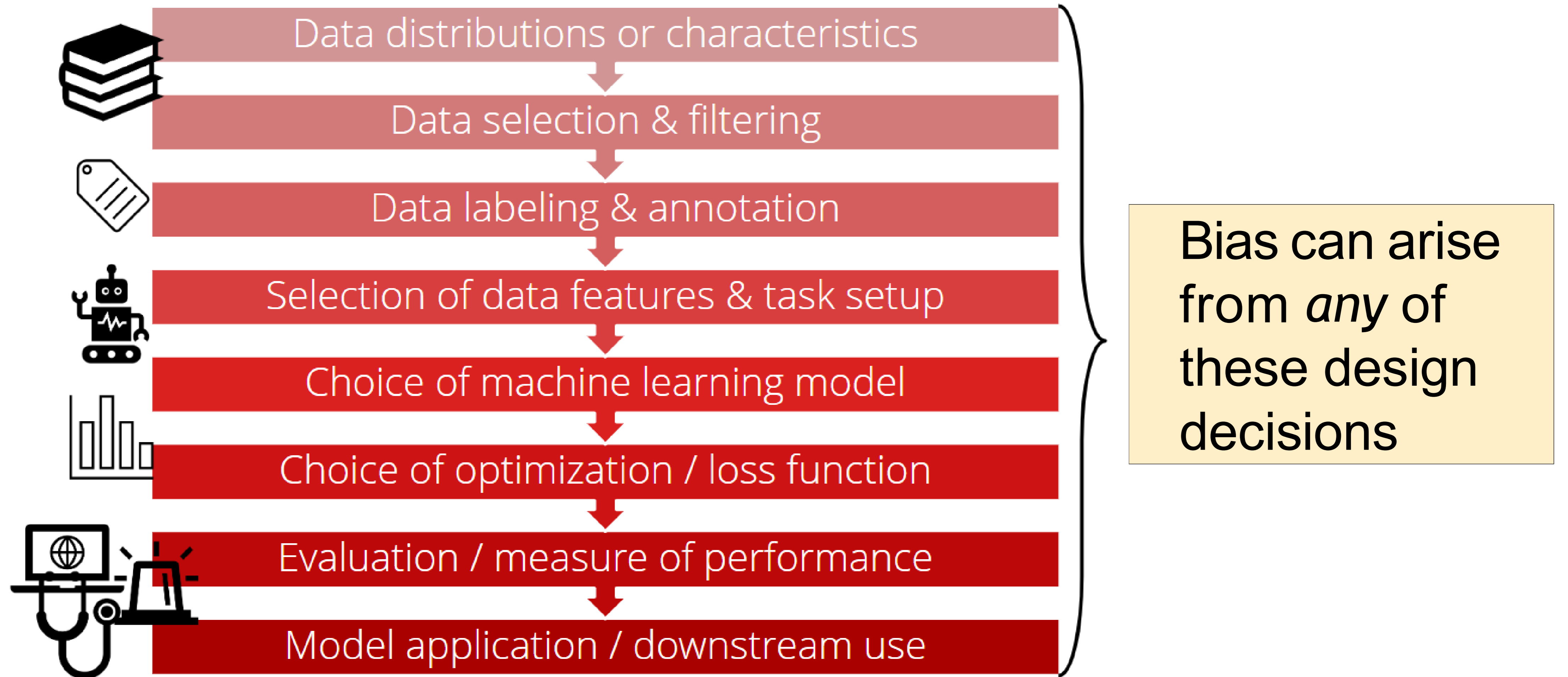
4 [RELATIVE NOTIONS OF FAIRNESS](#) [PDF](#)

Bias in terms of the harms it causes



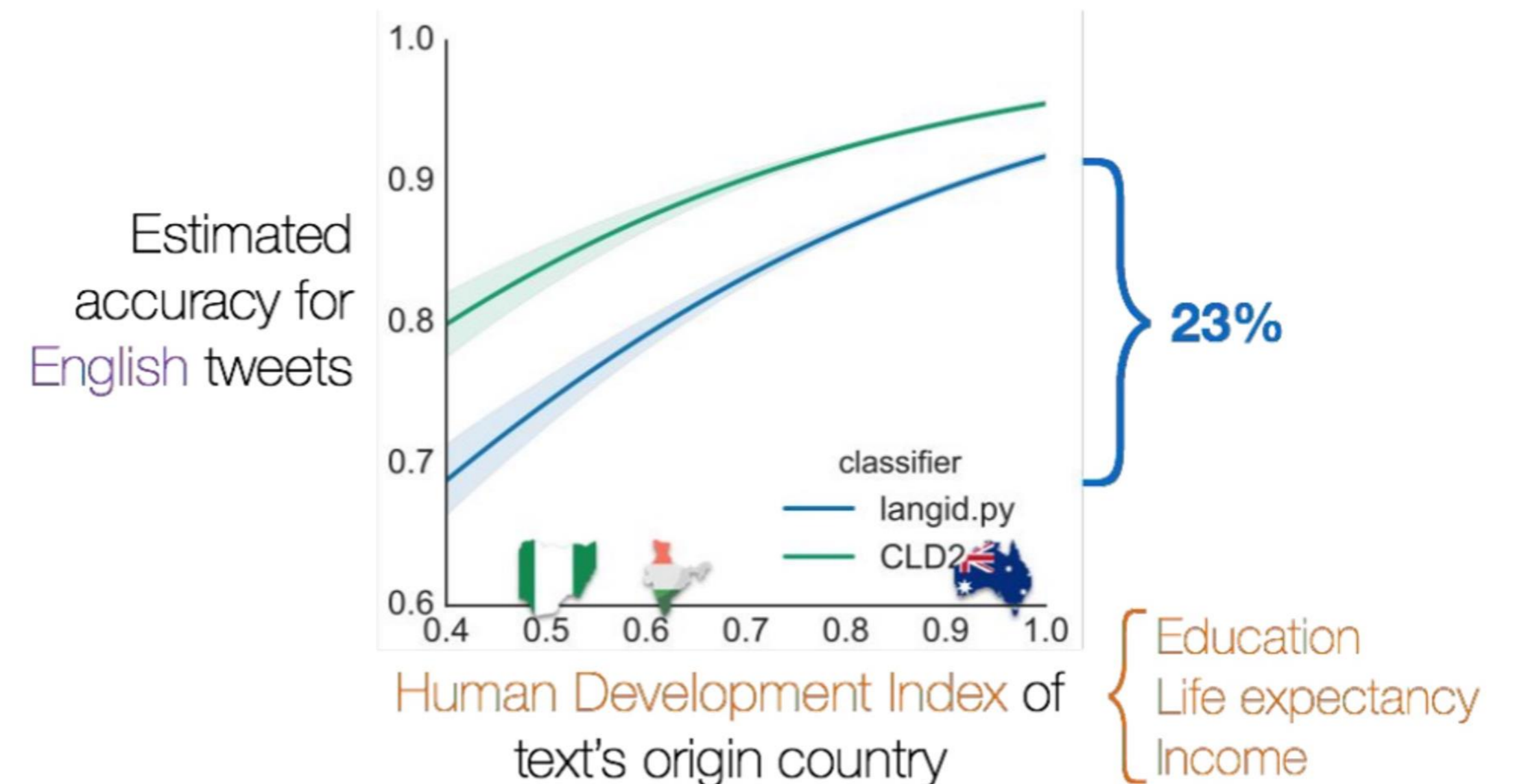
Where does bias come from?

Machine learning pipeline



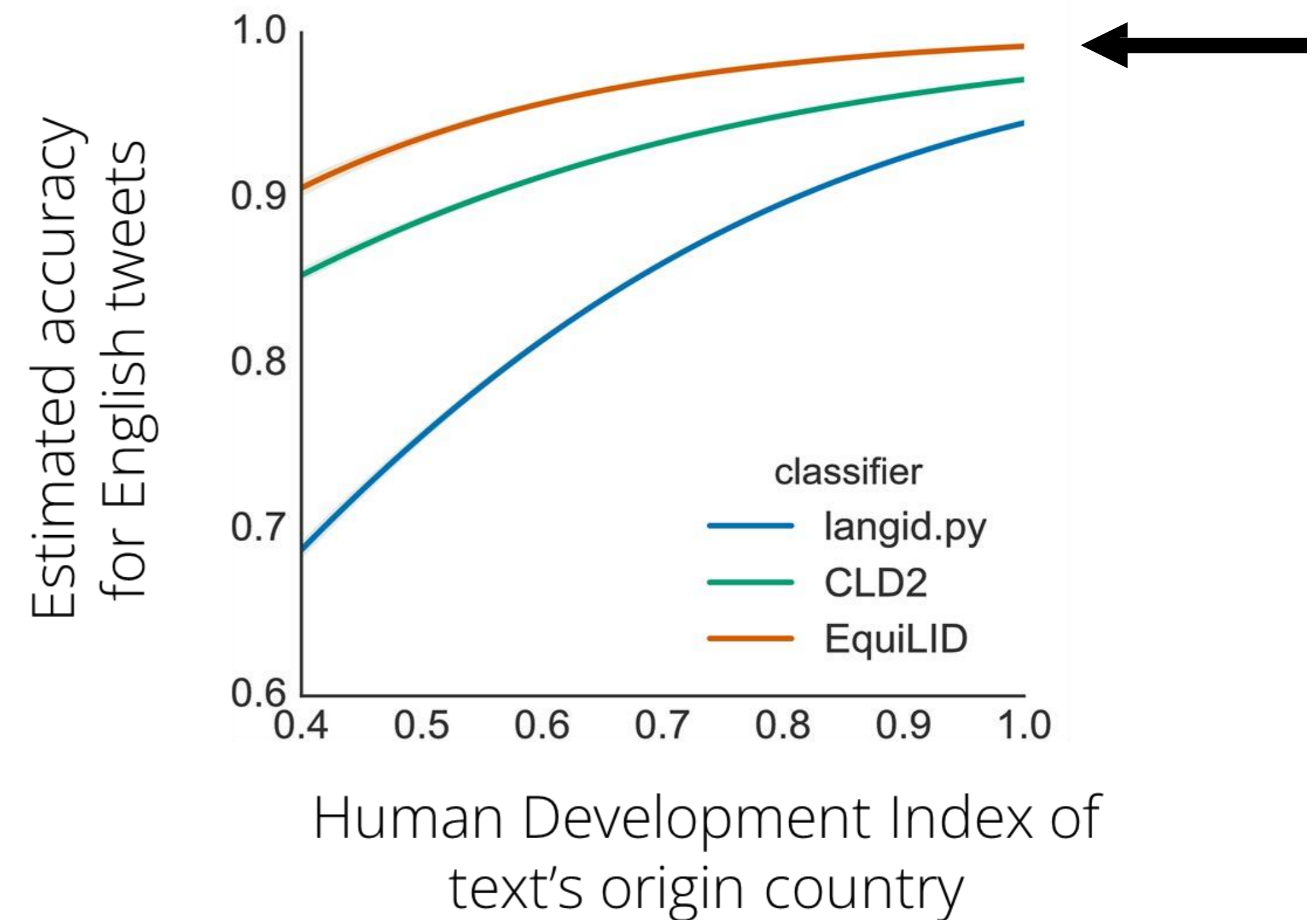
Example of bias from data: LangID tool

- LangID task: determine which language an input text is in
- Considered a “a solved problem suitable for undergraduate instruction” (McNamee, 2005)
- Often a first step in most NLP and CSS preprocessing pipelines
 - e.g., filtering LLM pretraining data
- But, many variations of English in the world
- Int'l: Nigerian English, Indian English, etc.
- Within US: African American English, etc.
- Jurgens et al. (2017) found that accuracy of LID tool correlated with wealth/development level of country; works worse for low HDI countries

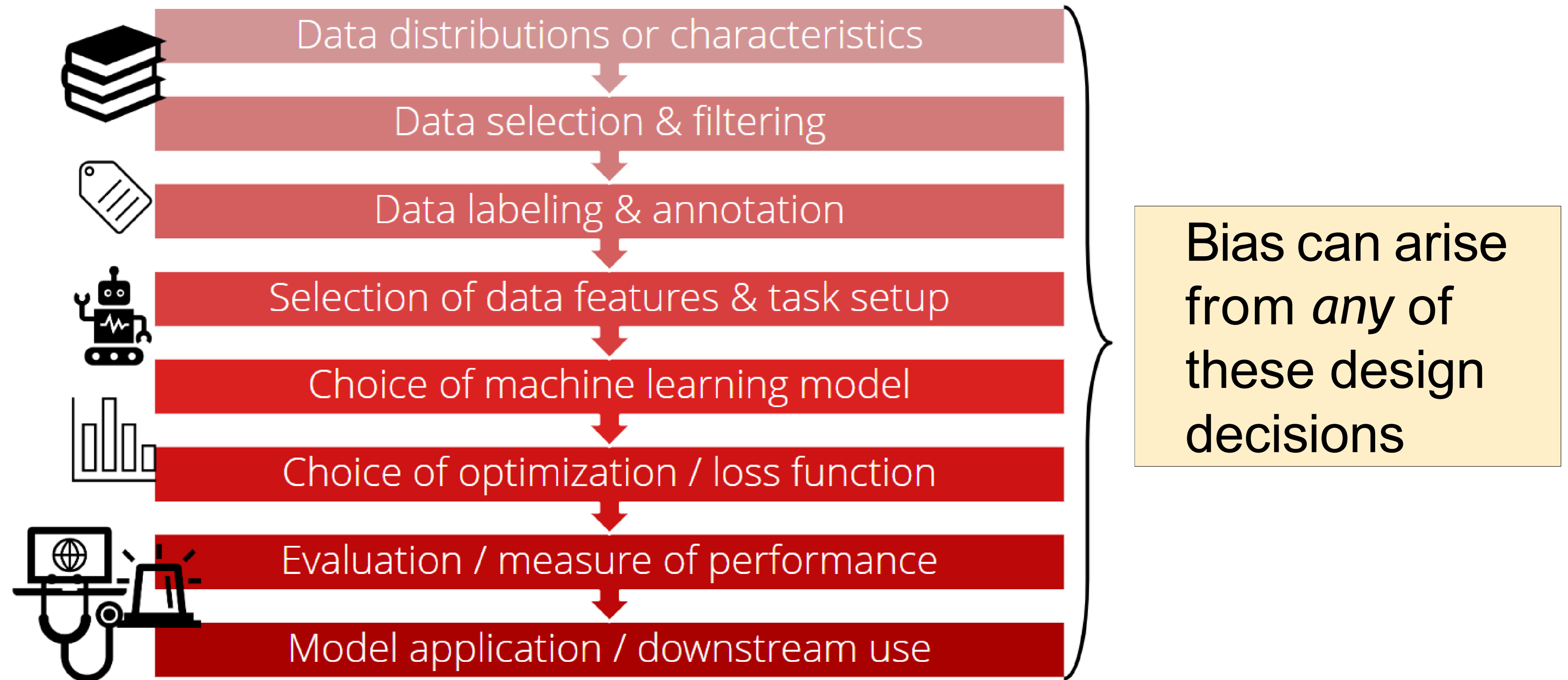


Example of bias from data: LangID tool (2)

- Jurgens et al (2017) introduce EquiLID
- Trained by sampling more variety of data, topically, socially, geographically diverse, and even multilingual data
- Find that tool works much better than original LID systems
- Bonus: even improved accuracy on highly developed countries!
- Takeaway: bias can be mitigated by making better data choices
- But that's not the only source of bias...

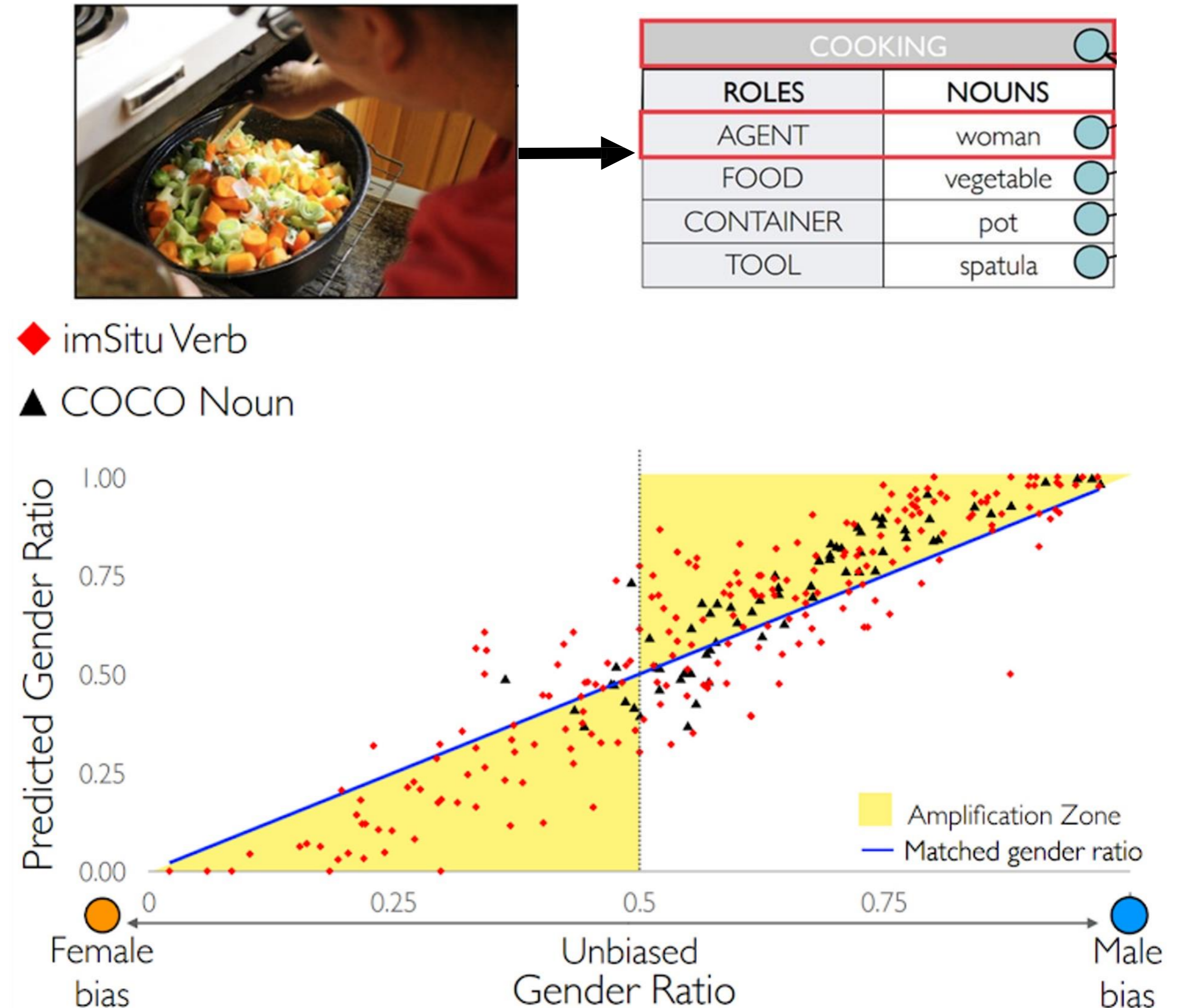


Machine learning pipeline



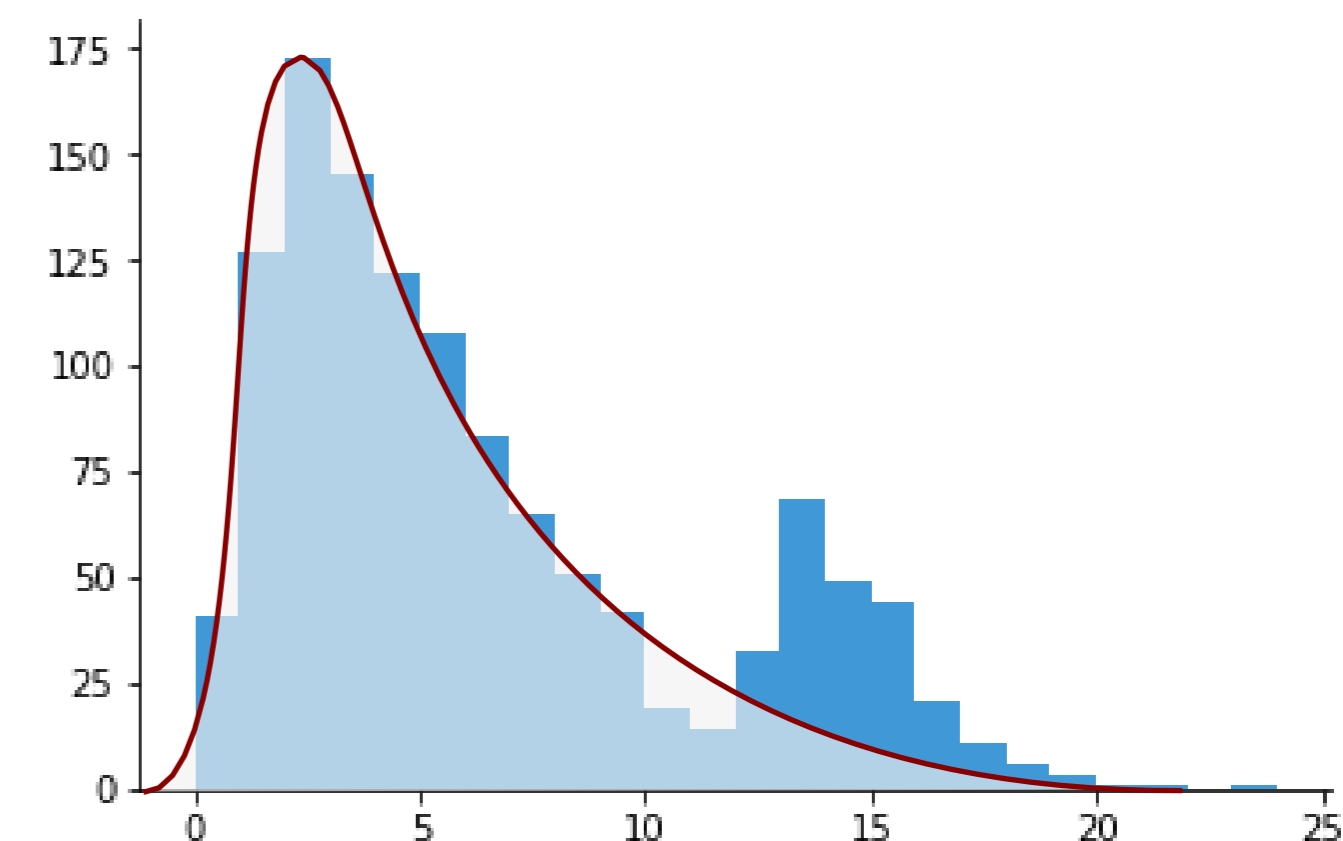
Bias amplification from models

- Zhao et al (2017) examined visual “semantic role” labeling task
- Given an image, predict various semantic roles, including agent (person doing the action)
- Found skews in training dataset
- E.g., 66% of training cooking images had agent=woman
- Found that models amplified biases
- E.g., 84% of test cooking images predicted as agent=woman (~18% men mis-labeled)

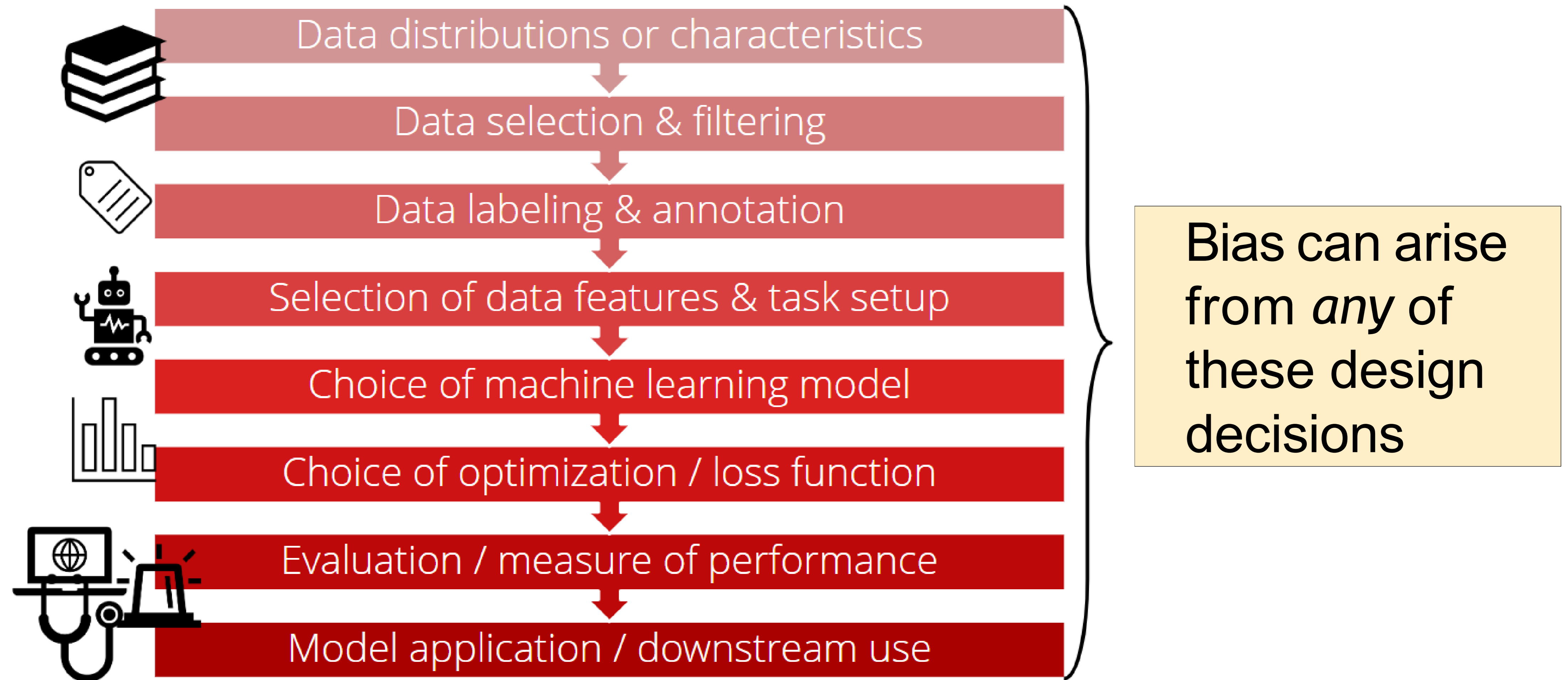


Model biases: mathematical links

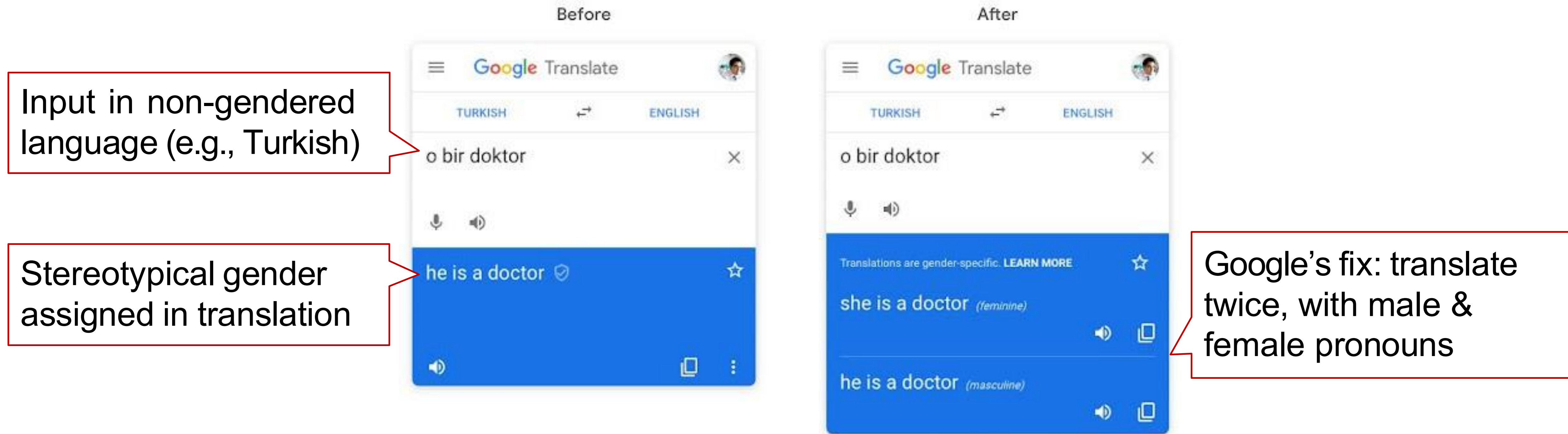
- Competing losses: objective functions aim to minimize loss globally → learns to predict most frequent class
- Often at the expense of less frequent classes (e.g. minority groups)
- Simplicity bias: neural networks biased towards learning simpler functions [Valle Pérez et al. 2019]
- Intuitively, if a model has limited learning capacity, makes sense that it learns shortcuts first
- Shortcuts are often stereotypes or majority biases; e.g., CEOs are men
- Takeaway: ML/optimization choices also affect biases



Machine learning pipeline



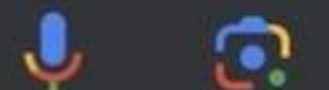
Google Translate issue



- Takeaways: mitigating bias may involve *system-level changes* to UI, input processing, output formatting, etc. while underlying AI model is similar
- *Let's discuss*: what do you think of this approach? What are some possible issues?

<https://ai.googleblog.com/2018/12/providing-gender-specific-translations.html>

translate



AI-läge Allt Bilder Videor Korta videor Nyheter

Engelska ▾



Spanska ▾

im non-binary



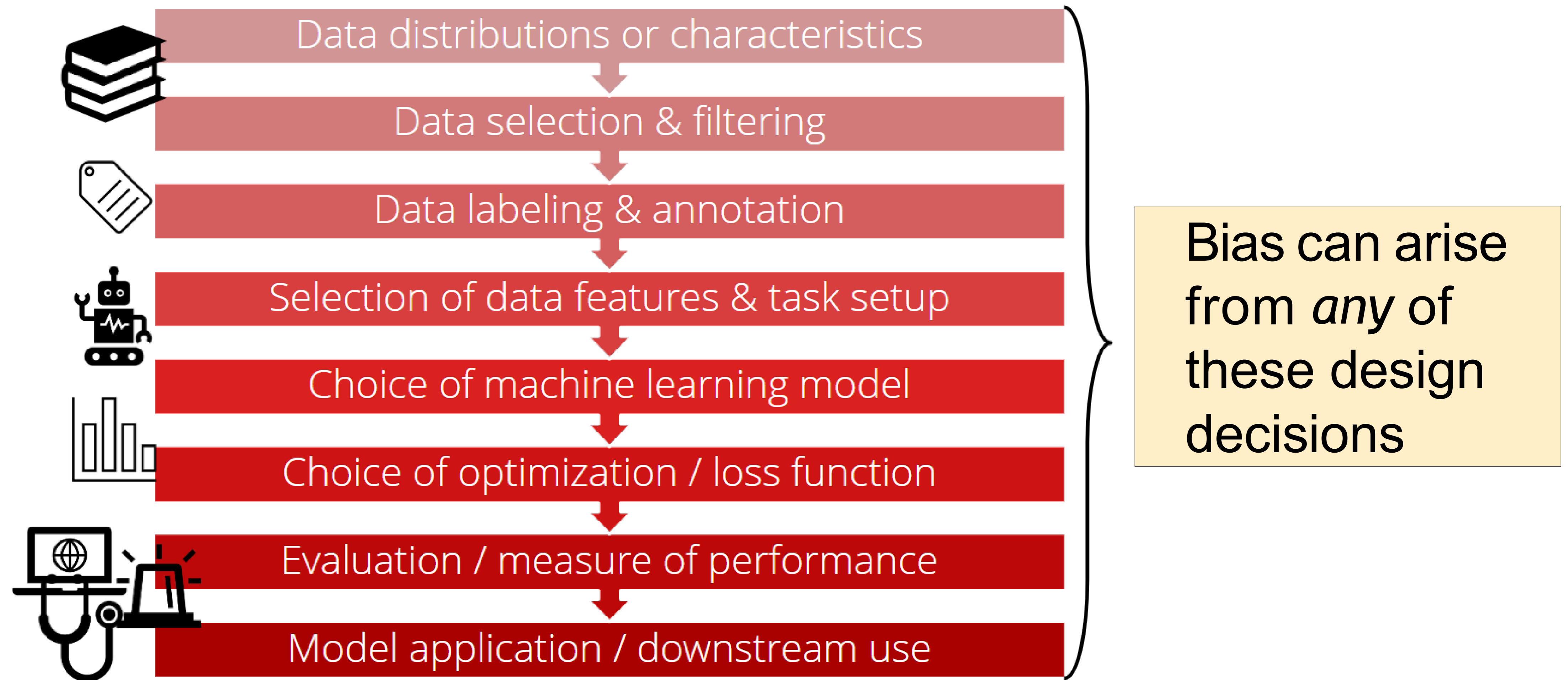
Soy no binaria (*femininum*)



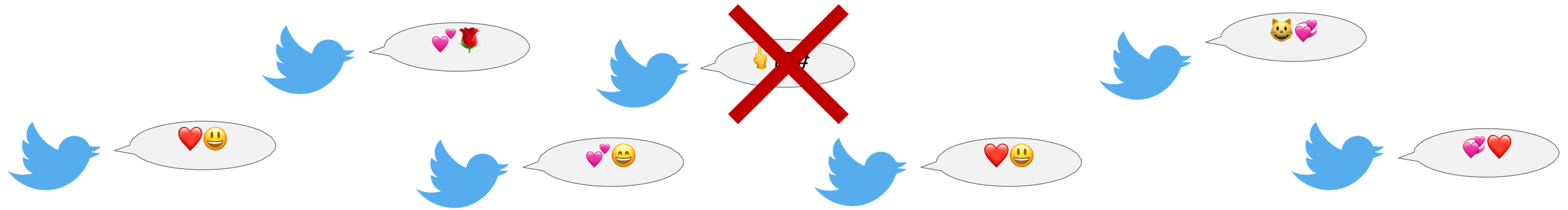
Soy no binario (*maskulinum*)



Machine learning pipeline

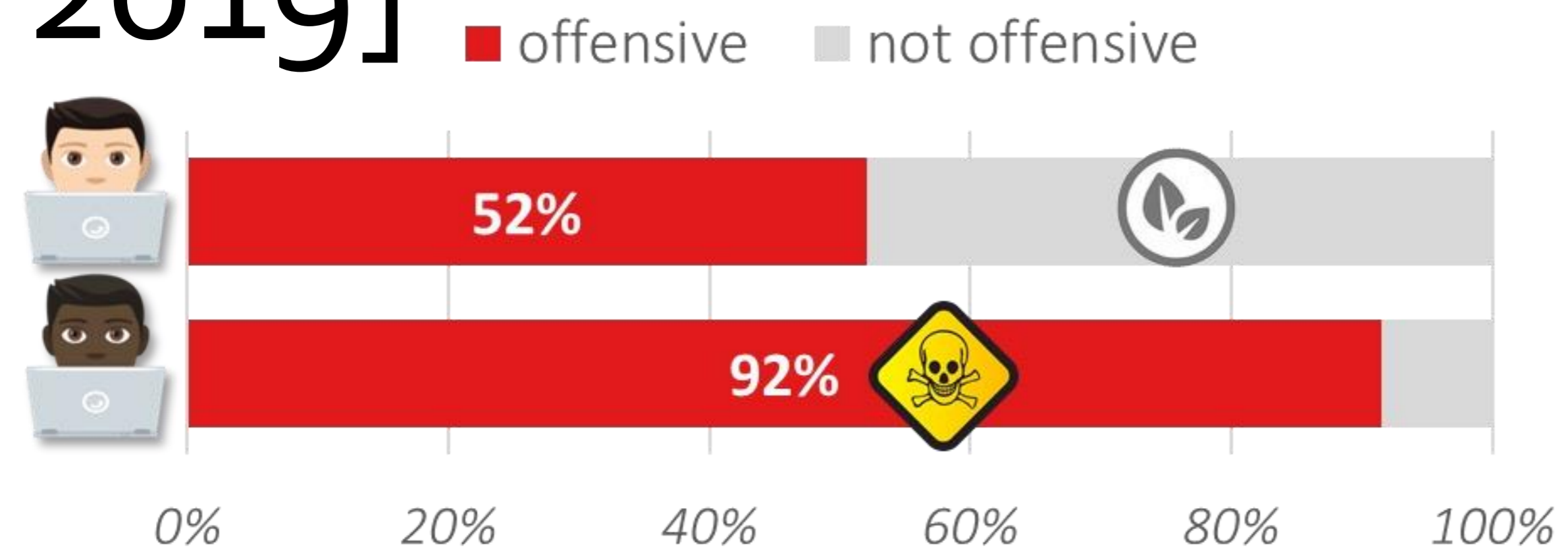


Hate Speech or Toxic Language Detection



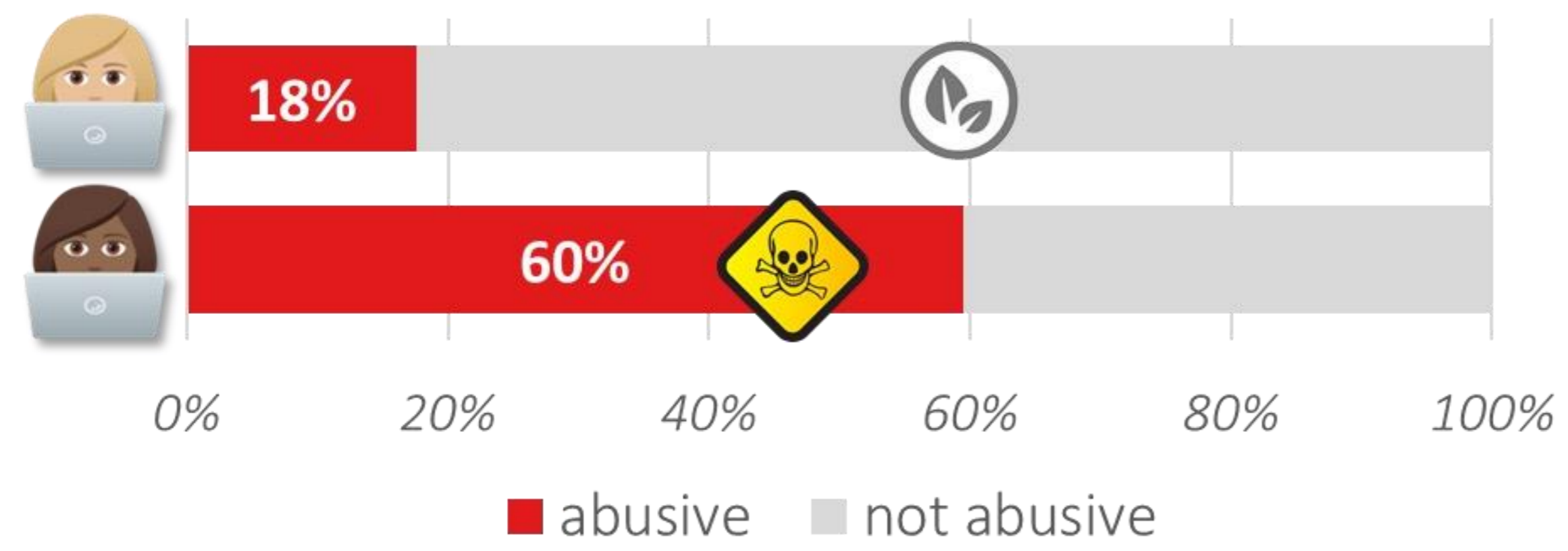
Goal: find and flag hateful or toxic content online, to make the internet less toxic

Racial biases in two popular datasets [Sap et al 2019]



TWT-HATEBASE
(Davidson et al., 2017)

Both datasets have **biases w.r.t. AAE tweets**



TWT-BOOTSTRAP
(Founta et al., 2018)

Enhancing the labeling interface [Sap et al 2019]

Control condition

Text-only, no context, prior work

Dialect priming

"Our AI thinks this tweet is in African American English"

Race priming

"A Twitter user that is likely Black/African American tweeted..."

MTurk study:

- 350 AAE tweets, ~50% labeled toxic
- 3 (re-)annotators per tweet

Could this tweet be offensive to *anyone*?

control

55%

dialect

44%

race

44%

**

**

Takeaway: adding social context to labeling mitigated bias

Why did these biases occur?

Why didn't NLP system designers think about these issues beforehand?

The world itself is biased



System designers have our own biases because of their *positionality*, i.e., set of perspectives that we hold due to our lived experiences and identity.

Positionality affects all our choices (e.g., assuming 1-1 mapping between languages and gendered pronouns, assuming toxicity looks the same in different dialects)

Debiasing AI systems

Is it even possible?

DALLE-2 vs. Gemini



- DALLÉ-2 generated images were shown to have social biases, later fixed by adding identity keywords to the input prompts (e.g., prompt+"Asian"; [Sparkes 2022](#))
- Gemini generations also shown to have skews
- *Let's discuss*: what do you think of this? How are these two generations different?

Limits of debiasing

- Gender debiasing doesn't work
- Breaks down for non-binary genders, racial categories or other social identity types
- Intrinsic debiasing \neq actual debiasing
- Finetuning often reintroduces biases
- Out-of-distribution data often still show biases
- Real world vs. ideal world: is reflecting the (biased) status quo the goal? or do we want to build a more fair or just world?
- Justice and fairness go beyond data & model fairness



“Lipstick on a pig” paper,
Gonen & Goldberg 2019

On the Intrinsic and Extrinsic Fairness Evaluation Metrics for Contextualized Language Representations

Yang Trista Cao^{*1}, Yada Pruksachatkun^{*2}, Kai-Wei Chang^{2,3}, Rahul Gupta²
Varun Kumar², Jwala Dhamala², Aram Galstyan^{2,4}

¹University of Maryland, College Park



A lot of people have understood that we need to have more diverse datasets, but unfortunately, I felt like that's kind of where the understanding has stopped. It's like 'let's diversify our datasets. And that's kind of ethics and fairness, right?' But you can't ignore social and structural problems.



Timnit Gebru, PhD

Socio-technical view on bias & fairness

- You can have a “fair” NLP/ML model (e.g., facial recognition system)
 - 95% accuracy/error rate on white & Black faces
- But if the system is used by law enforcement, bias creeps in w.r.t. who the system is used on
 - Black people more often arrested, due to racial biases
- Actual error rates are a function of deployment
- Algorithm’s fairness \neq fairness of treatment



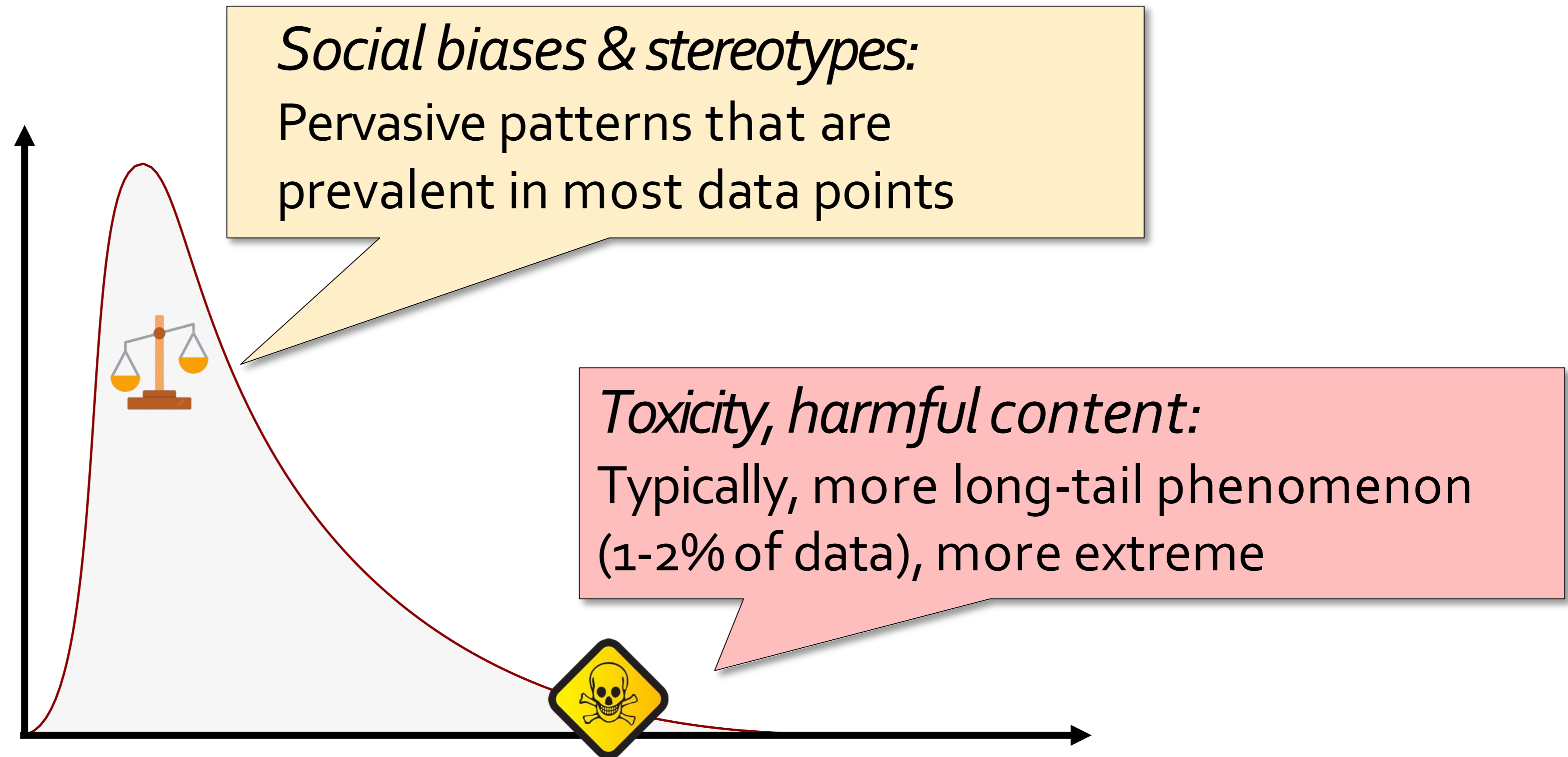
☰ CNN World Audio Live TV Log In

Black people are more likely to be arrested, charged and killed by police in Toronto, new report finds

By Scottie Andrew, CNN
Published 3:15 PM EDT, Wed August 12, 2020

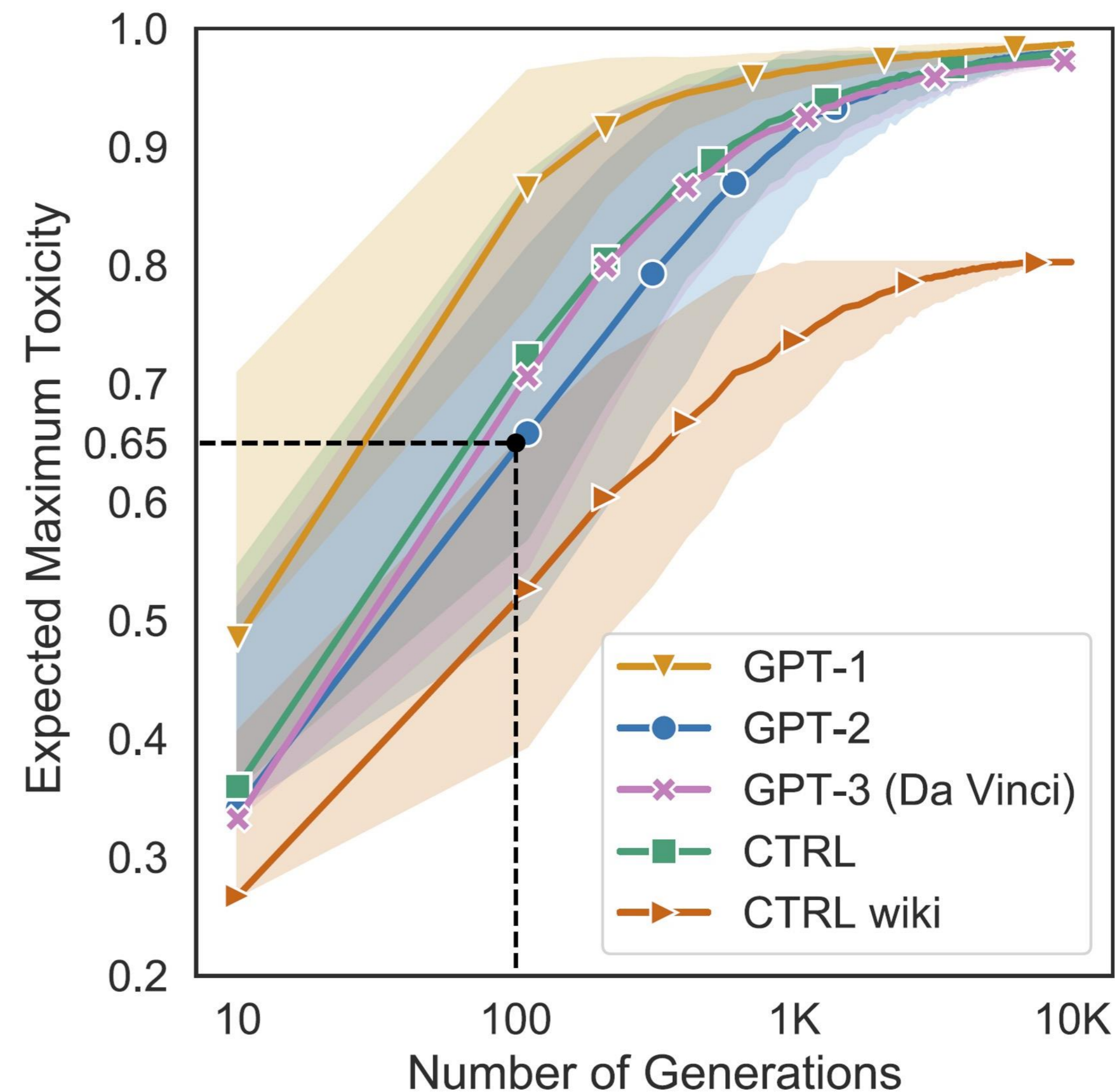
Part 2: Harmful content & toxicity

Biases vs. toxicity



Toxicity in LLMs, how bad is the problem really?

- Gehman et al (2020) introduced concept of neural toxic degeneration in LLMs
- Out of a 100 generations sampled from models, at least one toxic sentence
 - 65-70% toxicity from GPT2, GPT3
 - 85% toxicity from GPT1
- Model size affects toxicity: larger models have more toxicity [Touvron et al 2023]



Why are these models learning so much undesirable content?

Problems with self-supervised pretraining

“Feeding AI systems on the world’s beauty, ugliness, and cruelty, but expecting it to reflect only the beauty is a fantasy”

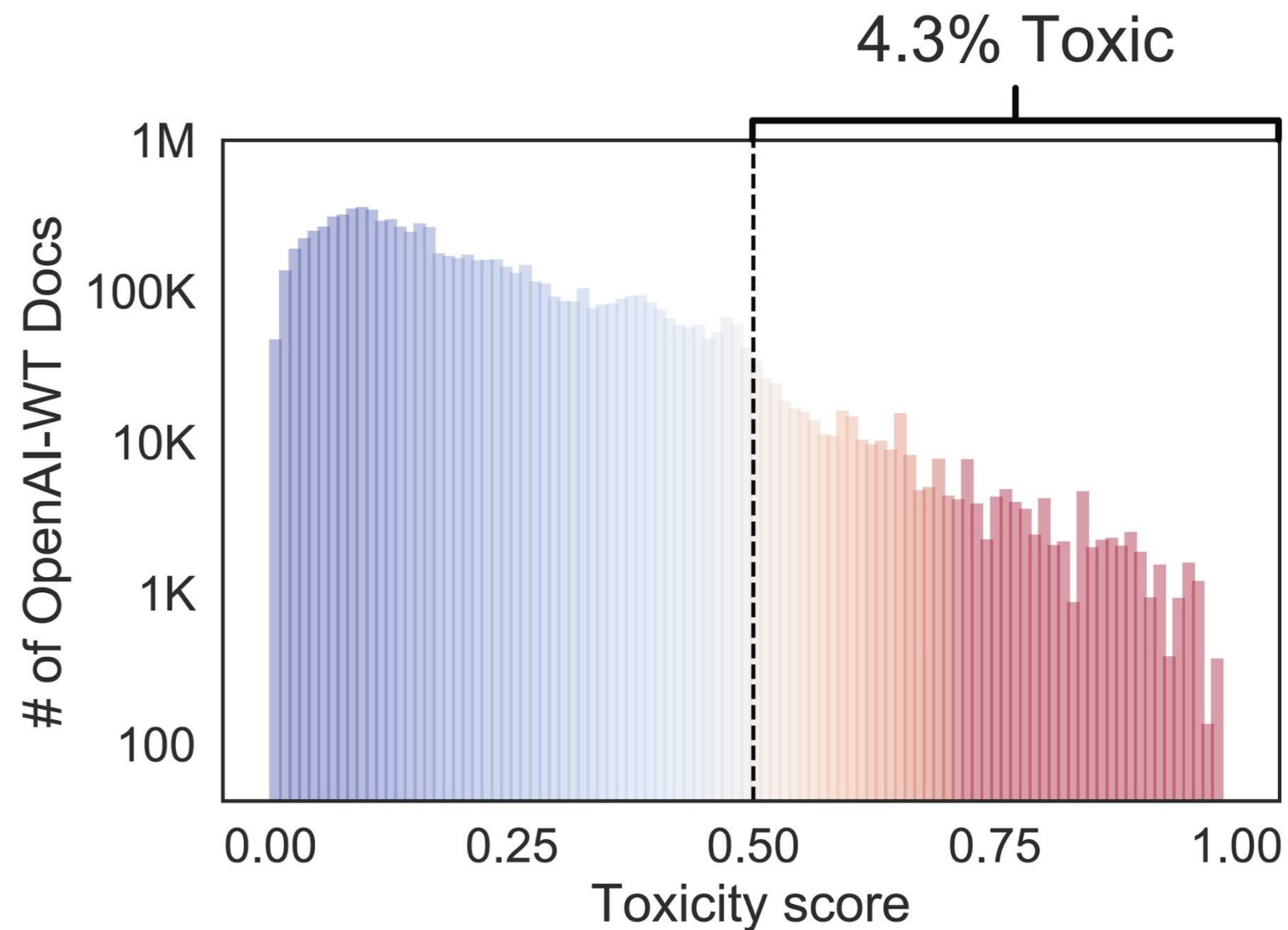


Prof. Ruha Benjamin, PhD

- Recipe: scrape as much pretraining data as you can to train your LM
- Consequence: LM ends up learning toxicity, biases, extremism, hate speech...

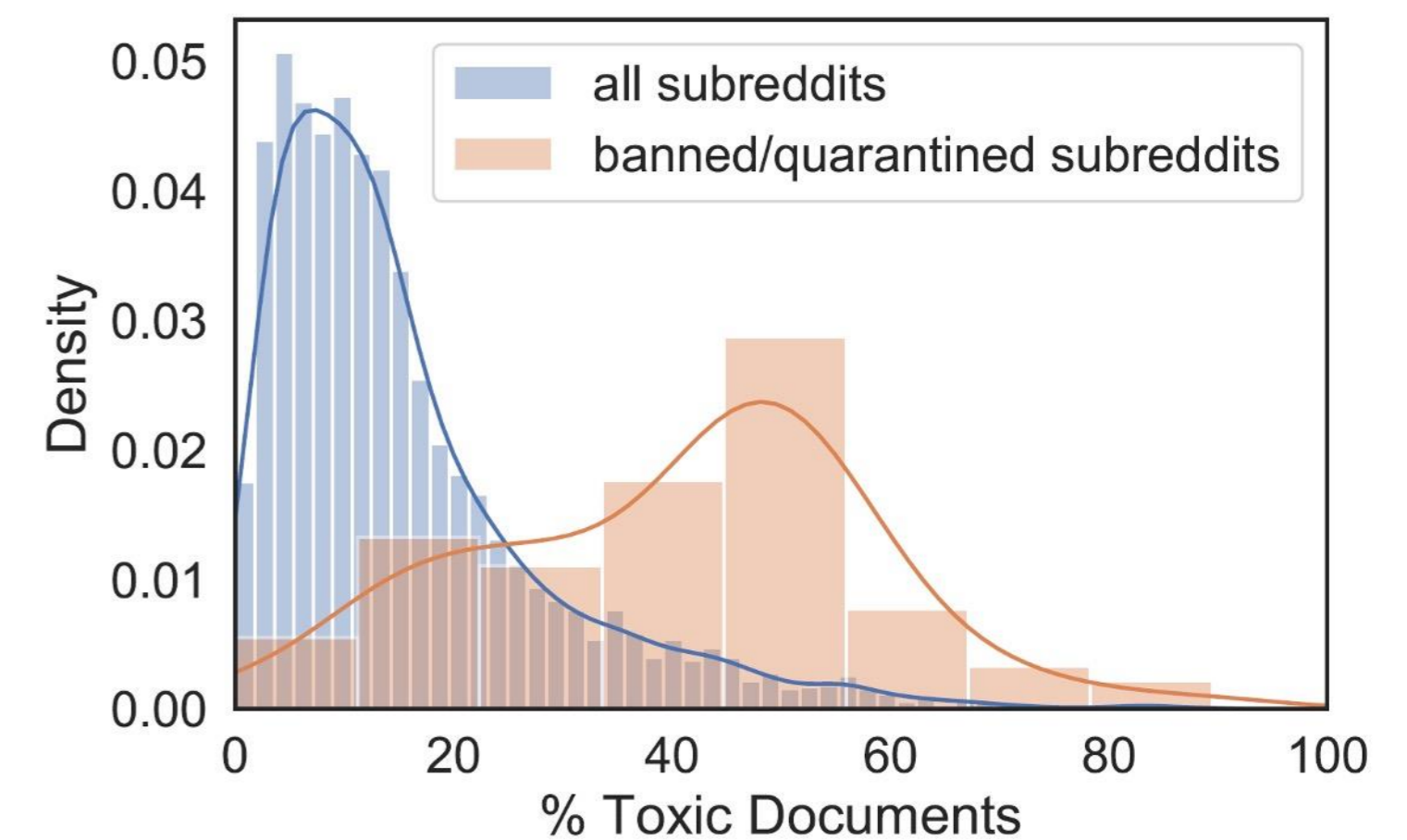
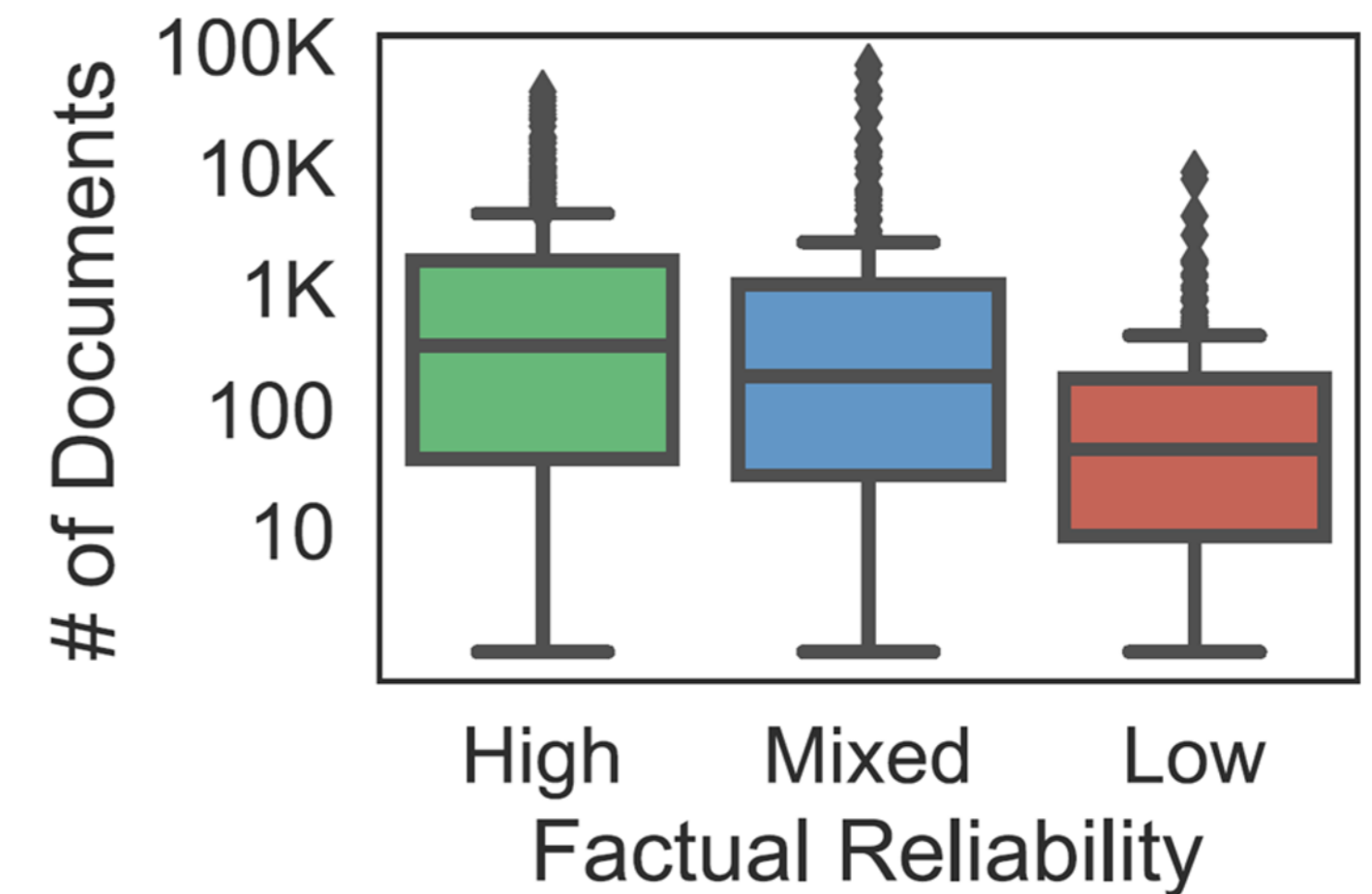
Toxicity in GPT-2's pretraining data

- [Gehman et al \(2020\)](#) accessed the actual GPT-2 training corpus (OpenAI-WT)
 - 8 million documents, 38Gb of text
 - Outbound links from Reddit posts with Karma ≥ 3
- Scored it with PerspectiveAPI toxicity
- Found $>4\%$ of documents (340,000) are toxic



Fake news in GPT-2's pretraining data

- Also looked at sources of documents in training data
- Cross-referencing sources of documents with known factual reliability categorization
 - >272K (3.4%) docs from low/mixed reliability sources
- Examining source where document is shared
 - >200K (3%) docs linked from banned/quarantined subreddits, which typically are more toxic docs
- Important to examine training data
 - Can only do that if publicly released!
- *So... need approaches to safeguard your model against this undesirable content, knowledge, and text.*



How to safeguard your LLMs

Overview – LLM safeguarding

- Filtering out toxic training data

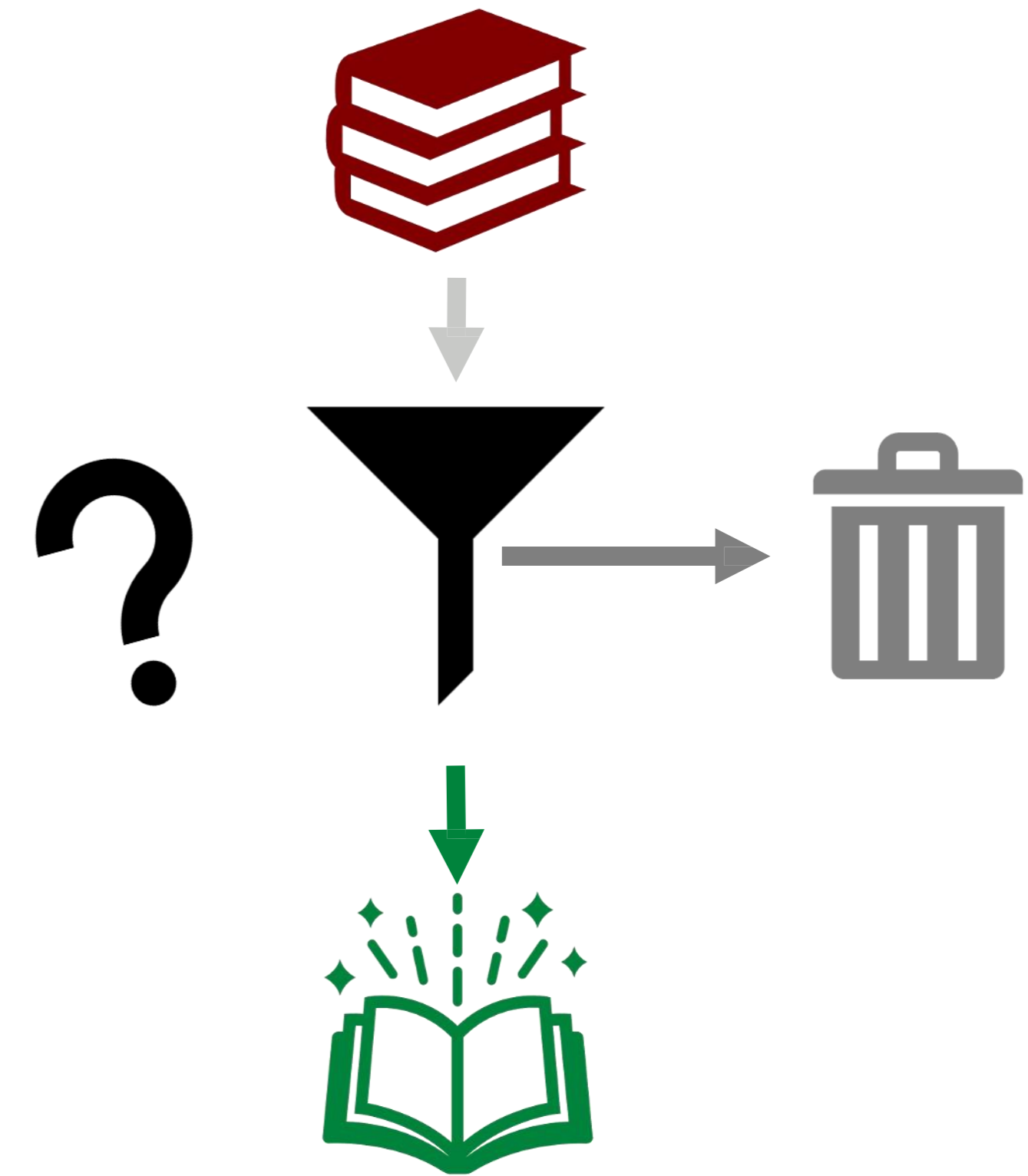
- Topic-based filters
- Toxic content detection

- Write demonstrations for refusing to answer
- RLHF models to prefer non-toxic generations

- Generate-then-classify
- Controllable text generation

Dataset filtering

- *Argument*: if you don't want your model to generate toxicity/hate speech, do not train it on such data (garbage in, garbage out)
- *Approach*: data filtering to ensure “high quality”
- How do you know what is “high quality” ?
 - GPT-2: Reddit “Karma” score as signal
 - T5, BERT: “blocklist” of “bad words”
 - GPT-3: “quality” classifier
 - Newer works: Combination of all



Blocklist of “bad” words

- “List of Dirty, Naughty, Obscene, or Otherwise Bad Words” originally by Shutterstock employees
 - Meant to prevent words in autocomplete settings
- Has been used by most companies creating LLMs
 - BERT, T5, GPT-2, etc.
- If document contains a “bad” word, remove it from training data
 - F*ck, sh*t, sex, vagina, viagra, n**ga, f*g, b*tch, etc.
- *Let's discuss*: what are issues with this?
 - Strong risk of over-deleting bio, legal, minority content

WIRED

SUBSCRIBE

TOM SIMONITE

BUSINESS FEB 4, 2021 7:00 AM

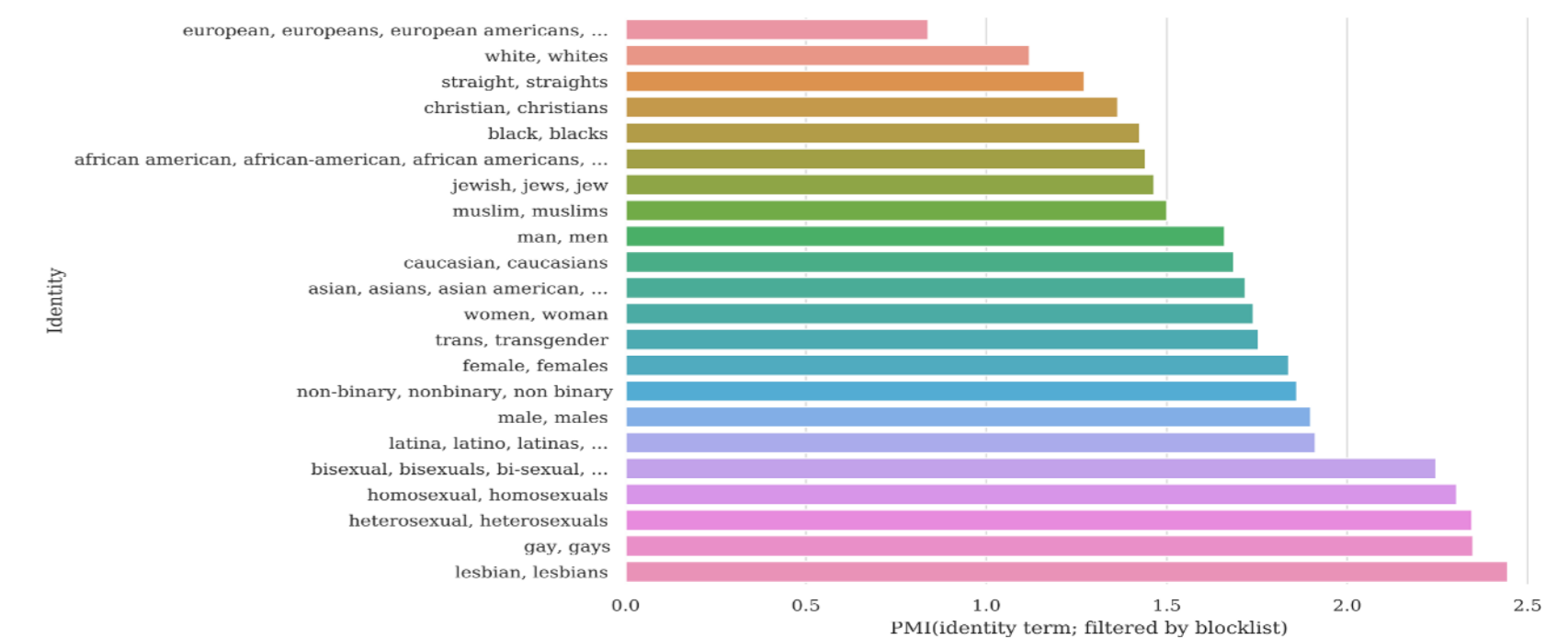
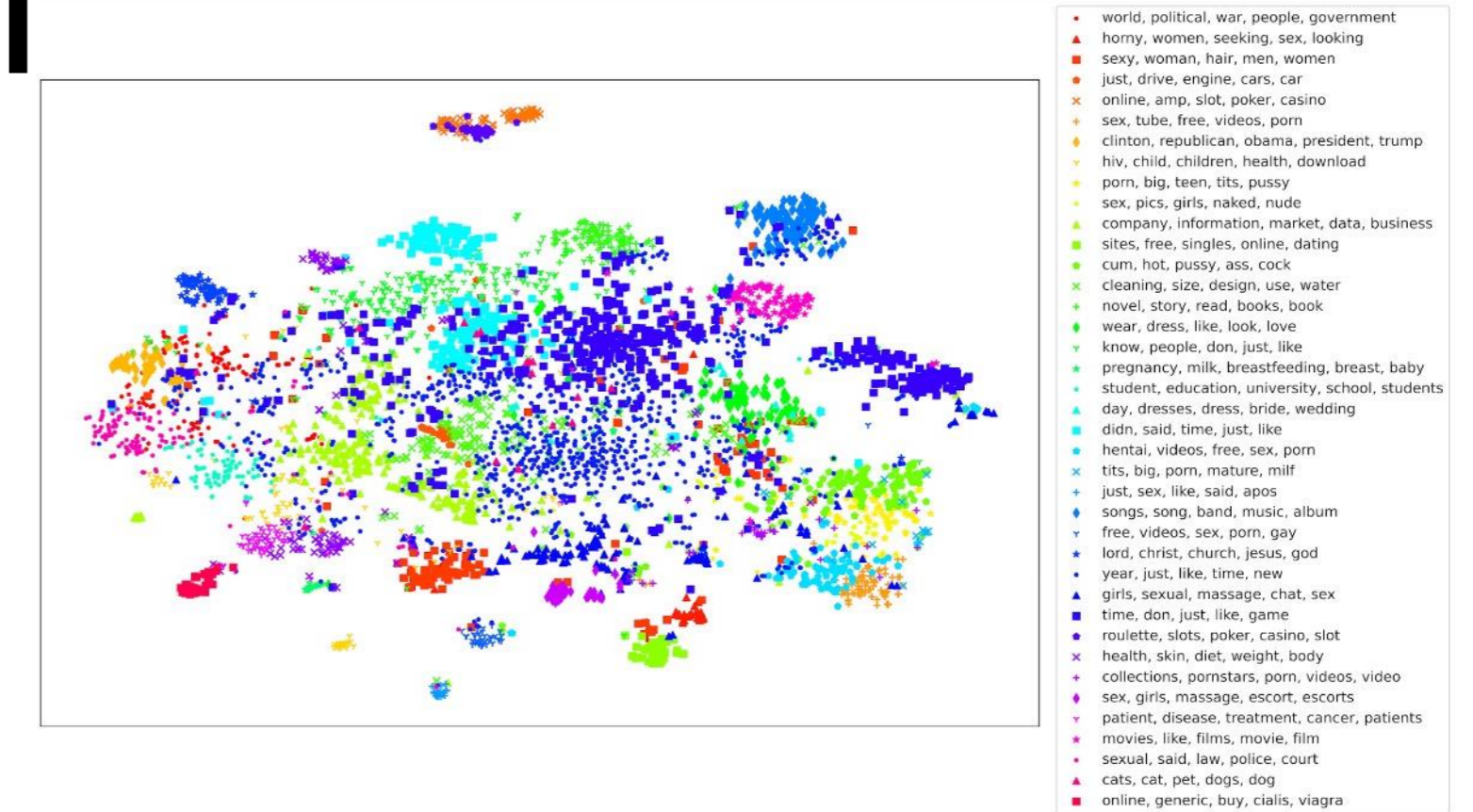
AI and the List of Dirty, Naughty, Obscene, and Otherwise Bad Words

It started as a way to restrict auto-completes on Shutterstock. Now it grooms search suggestions on Slack and influences Google's artificial intelligence research.



Effect of “bad word” blacklist filtering

- Dodge et al examined the effect of blacklist filtering on the C₄ corpus
- When looking at 100k documents that were excluded due to “bad words”
 - Found only 31% related to porn/explicit sex
 - Remaining was biology, medicine, legal
- Also examined the effect on which minority identities were removed
 - Found queer/LGBTQ identity terms removed more
- Examined dialects removed due to “bad words”
 - Found AAE, Hispanic English more likely to be removed



Less likely to be removed

- White-aligned English (6%)
- Other English (7%)

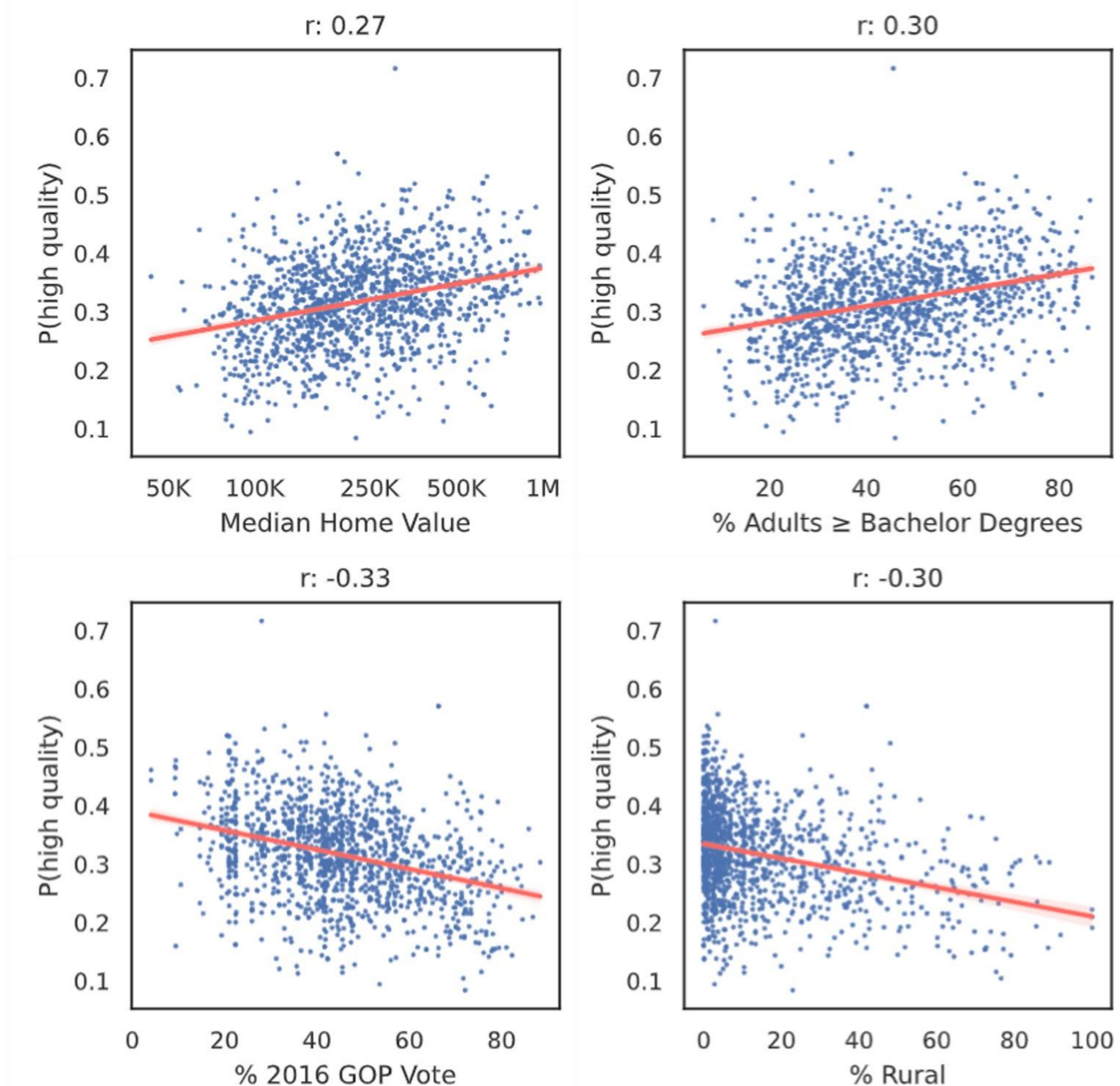
More likely to be removed

- African-American English (42%)
- Hispanic-aligned English (32%)

GPT₃ Quality filter backfires

- GPT₃ quality filter: similar to GPT₂ data
- [Gururangan et al. \(2022\)](#) re-implemented GPT-3 quality filter
- Ran it on articles from school newspapers, which have metadata
- Filter assigns higher quality to articles from
 - Richer counties 💰
 - Counties with more educated adults 🎓
 - More liberal counties 🇺🇸
 - More urban counties 🏙️
- Raises language ideology question: Whose English is “good English”?

“In order to improve the quality of Common Crawl, we developed an **automatic filtering method** to remove low quality documents. Using the **original WebText as a proxy for high-quality documents**, we trained a classifier to distinguish these from raw Common Crawl.” – Brown et al. 2020



So... maybe filtering isn't a good idea since
it'll backfire?

GPT4Chan controversy

- Yannic Kilchner finetuned GPT-J on 4chan posts
 - Trained on subforum /pol/ known to contain racist, sexist, white supremacist, antisemitic, anti-Muslim, anti-LGBT views
- Trolled 4chan users with bots powered by his model
 - 30,000 posts over the span of a few days
- Faced massive criticism
 - initially hosted on Huggingface, was taken down quickly
- *Let's discuss...*
 - Was this an ethical model to train? Given that the dataset was publicly available?
 - Was deploying the bots on 4chan okay?
 - Are there any useful/positive applications of the model?



≡ FORTUNE

TECH · 4CHAN

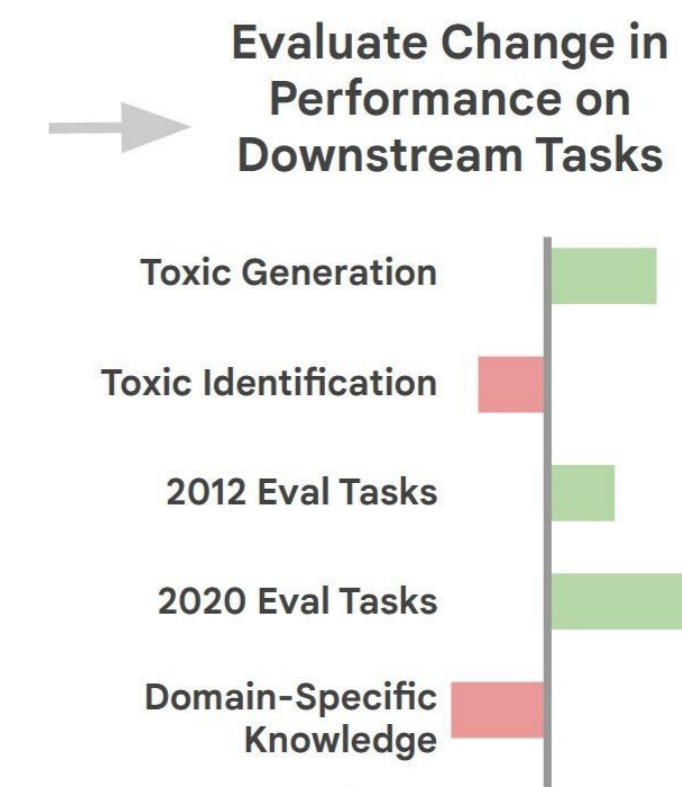
‘This breaches every principle of human research ethics’: A YouTuber trained an A.I. bot on toxic 4Chan posts then let it loose — and experts aren’t happy

BY SOPHIE MELLOR

June 10, 2022 at 5:23 AM EDT

Why LLMs might want to have seen toxic content

- Detecting hate speech [[Chiu et al 2022](#)]
 - [Longpre et al. \(2023\)](#) showed that LLMs trained on more toxicity are better toxicity detections
 - Improving hate speech models with data augmentation: ToxiGen [[Hartvigsen et al 2022](#)]
- Counter speech generation [[Saha et al 2022](#), [Kim et al 2022](#), [Mun et al 2023](#)]
- *If we train on toxicity, something else must be done at a different time!*



Overview – LLM safeguarding

- Filtering out toxic training data

- Topic-based filters
- Toxic content detection

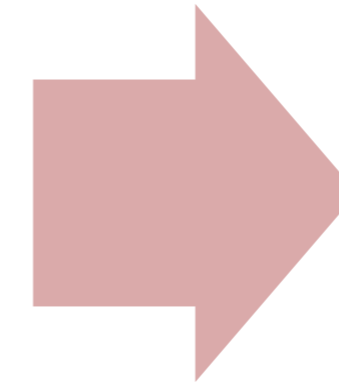
- Write demonstrations for refusing to answer
- RLHF models to prefer non-toxic generations

- Generate-then-classify
- Controllable text generation

RLHF safeguarding – assumptions

PPO & family:

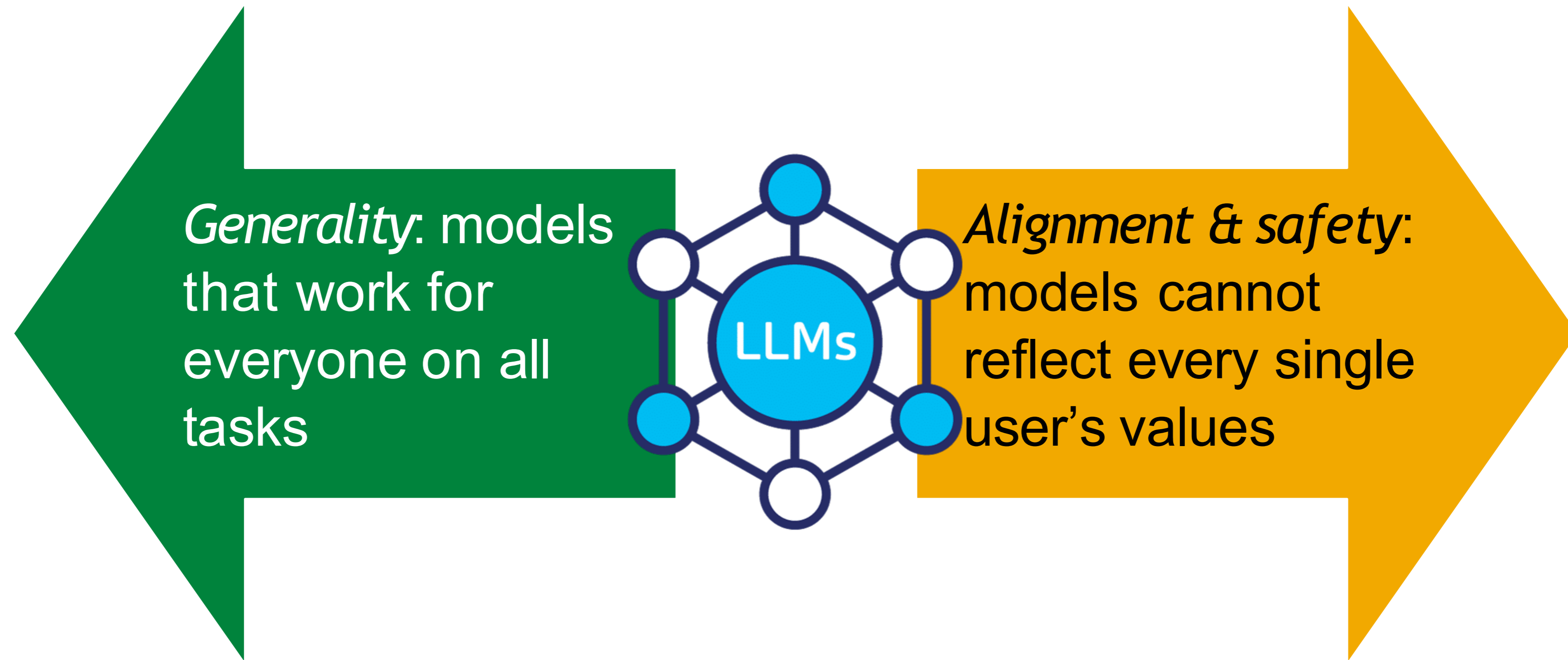
Obtain preference data: which generation is good vs. bad?



RL is done to encourage more like “preferred output”

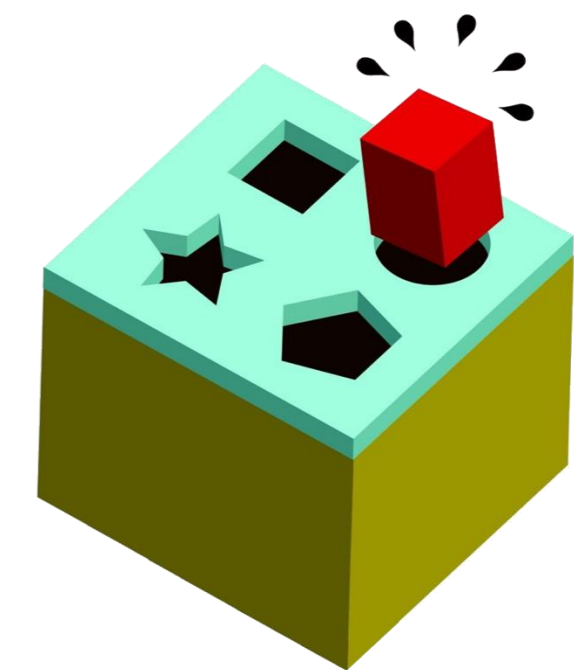
- Big question: what does it mean for a generation to be better/preferred?
 - How to balance harmless and helpful? [[Bai et al '23](#)]
 - *E.g., “help me create a poisonous drink.”*
 - What if people’s preferences are biased or gameable?
 - *E.g., people prefer certainty over uncertainty in answers to questions [[Zhou et al. 24](#)]*
 - Fundamental issue: cannot represent all values and cultures into one ranking.
 - *Casper et al. 2023. “Open Problems and Fundamental Limitations of Reinforcement Learning from Human Feedback.” arXiv [cs.AI]. arXiv. <http://arxiv.org/abs/2307.15217>*

Big unresolved tension



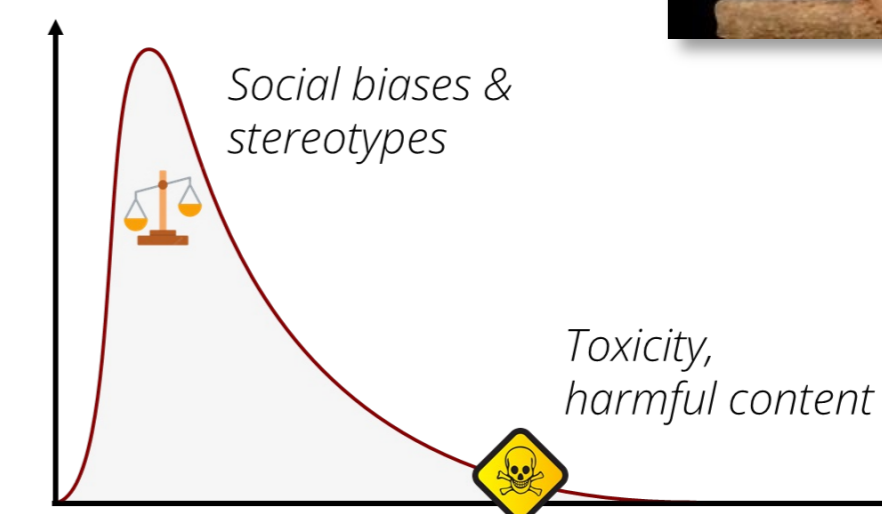
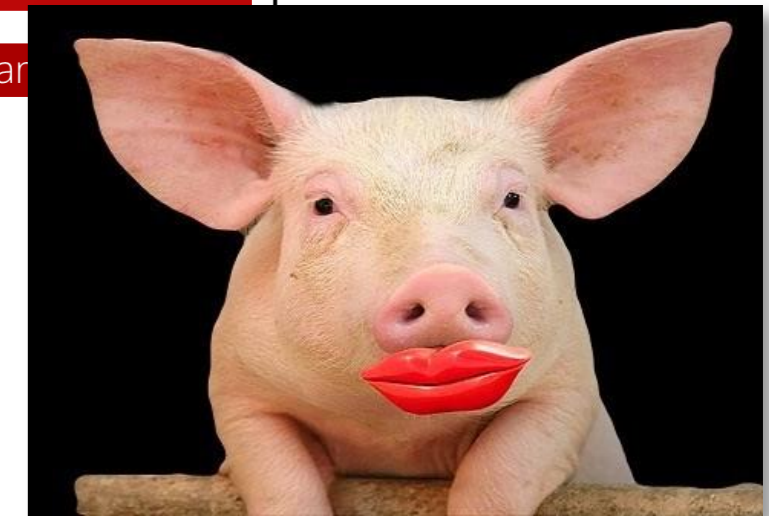
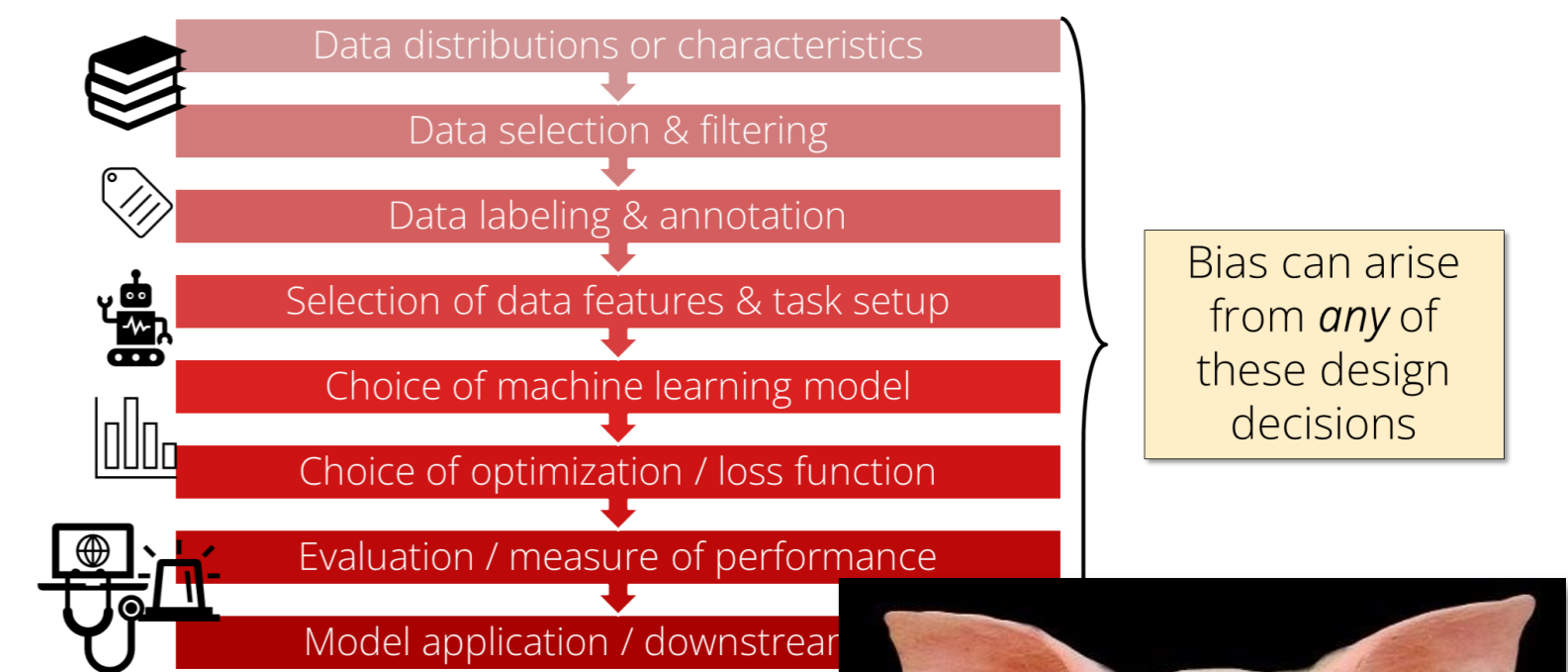
So... what can we do?

- Need to keep studying what models can and can't do, who they work for and don't work for
- Narrow scope of model users
 - Community-specific models (e.g., Masakhane Initiative)
- Specialize models' abilities / away from one-size-fits-all
 - E.g., toxicity explanation generation model needs to generate stereotypes, but story generation models might not
- In line with many legislative efforts:
- legislate the application or task, not the model



Takeaways

- AI systems are biased
- Real world is biased, data is biased
- ML objectives play a role
- Annotation interfaces, context plays a role
- Debiasing is challenging, requires socio-technical lens
- Toxicity and undesirable content
- Longer-tail phenomenon, present in training data
- Filtering data can backfire
- Safeguarding to all people is impossible
- Any questions?



Safeguards from training data	• Filtering out toxic training data
Safeguards from input prompt classification	• Topic-based filters • Toxic content detection
Safeguards from instruction-tuning & RLHF	• Write demonstrations for refusing to answer • RLHF models to prefer non-toxic generations
Safeguards at the output level	• Generate-then-classify • Controllable text generation