# Mitigating Bias and Toxicity

Arnab Bhattacharjee, Anirudh Ajith

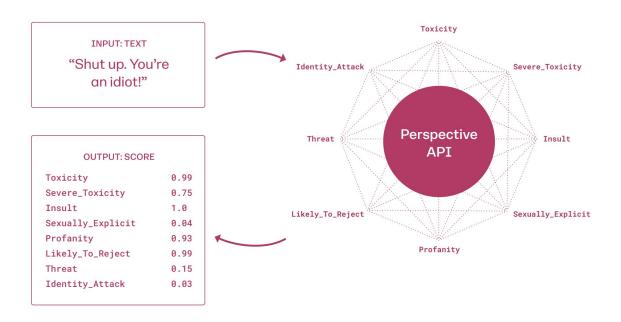
### Outline

- 1. What is toxicity?
- 2. Recent methods for toxicity mitigation
  - a. DAPT (Gururangan et al., 2020) more details!!
  - b. PPLM (Dathathri et al., 2020) more details!!
  - c. GeDi (Krause et al., 2020) new!!
- 3. Self-Diagnosis and Self-Debiasing (Schick et al., 2021) new!!
  - a. Motivation
  - b. Self-diagnosis
  - c. Self-debiasing
  - d. Results
  - e. Limitations
- 4. Does detoxification introduce social bias? new!!

#### Issue

- LLMs trained using large crawls from the internet with very basic filtering (if any)
  - C4 filtering using RegExes
  - The Pile "...it is possible for the Pile to contain pejorative, sexually explicit, or otherwise objectionable content"
- Non-negligible amounts of harmful, biased text
- LLMs trained using this data pick up, amplify these biases

#### Perspective API

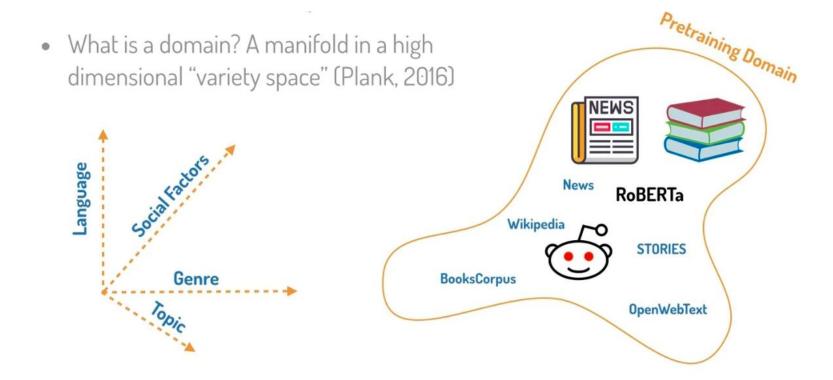


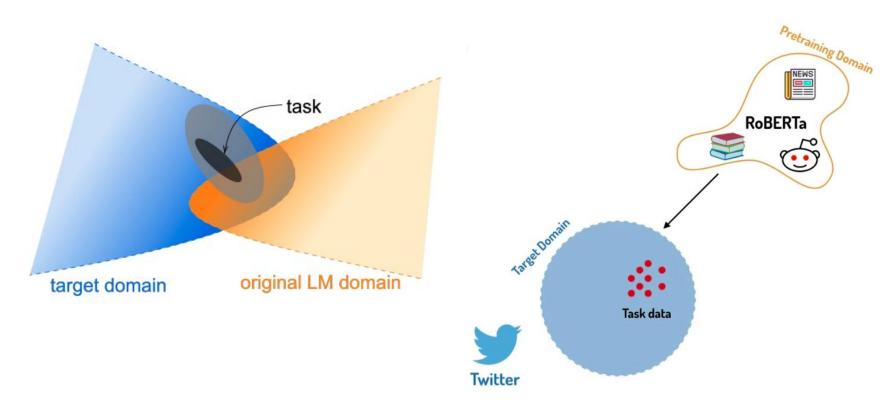
- Returns calibrated probabilities of attributes
- Toxicity := "...rude, disrespectful, or unreasonable language that is likely to make someone leave a discussion."

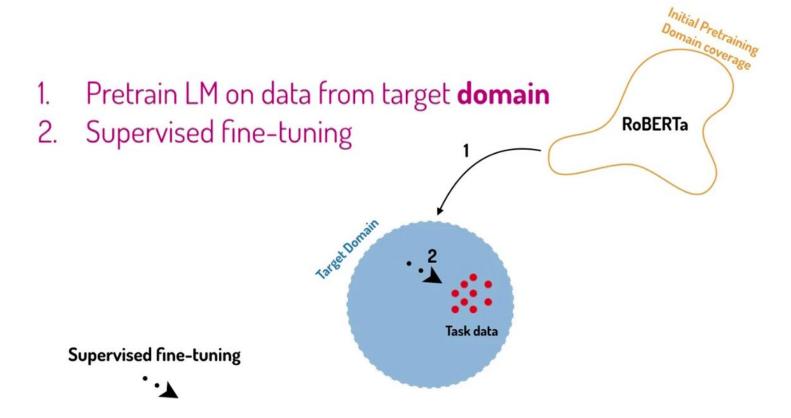
# Recent methods for toxicity mitigation

### Adapting prior ideas for toxicity mitigation

- Data-based methods
  - a. Domain-adaptive Pretraining (DAPT) more details!!
  - b. Attribute Conditioning (ATCON)
- 2. Decoding-based methods
  - a. Plug and Play Language Models (PPLM) more details!!
  - b. Generative Discriminator Guided Sequence Generation (GeDi) new!!
  - c. Self-debiasing (SD) new!!









Domain	Pretraining Data	Classification Tasks
Biomedical	S20RC Papers (7.6B Tokens)	ChemProt RCT
Computer Science	S20RC Papers (8.1B Tokens)	ACL-ARC SCIERC
Reviews	Amazon Reviews (2.1B Tokens)	Amazon Helpfulness IMDB
News	RealNews Articles (6.7B Tokens)	Hyperpartisan AG News

Domain	Task	RoBERTa	DAPT
Biomed	ChemProt	81.9	84.2
CS	ACL-ARC	63.0	75.4
News	HyperPart	86.6	88.2
Reviews	IMDB	95.0	95.4

### Domain-adaptive Pretraining for detoxification

- Continued pretraining of LLM on filtered non-toxic subset of OWTC
- Aims to erase knowledge of toxicity via catastrophic forgetting

1870 M 8300M	Exp. Max. Toxicity			Toxicity Prob.		
Model	Unprompted	Toxic	Non-Toxic	Unprompted	Toxic	Non-Toxic
GPT-2	0.44 <sub>0.17</sub>	0.750.19	0.510.22	0.33	0.88	0.48
DAPT (Non-Toxic)	0.300.13	0.570.23	0.370.19	0.09	0.59	0.23
DAPT (Toxic)	$0.80_{0.16}$	$0.85_{0.15}$	$0.69_{0.23}$	0.93	0.96	0.77

(25 generations)

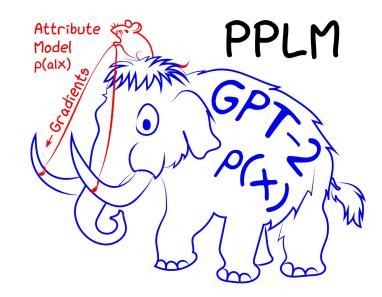
### Domain-adaptive Pretraining for detoxification

- Can adversely affect modelling performance and destroy innocuous knowledge in unexpected ways.
- Leads to either limited detoxification effectiveness or significantly sacrifices model quality.
- Expensive additional data, compute

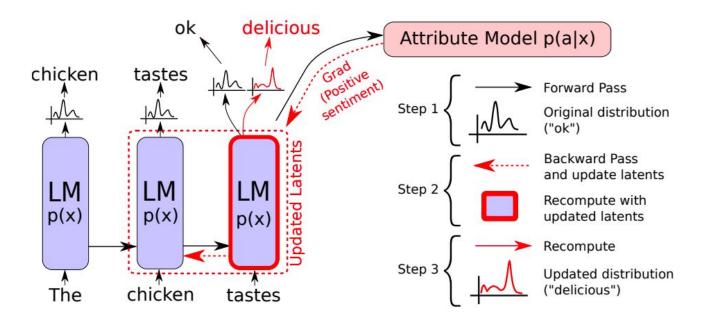
Model type	Form of model	Examples
Language model	p(x)	GPT-2 (Radford et al., 2019)
Fine-tuned language model	p(x)	DAPT (Non-Toxic) (Gehmen et al., 2020)

### Plug and Play Language Models

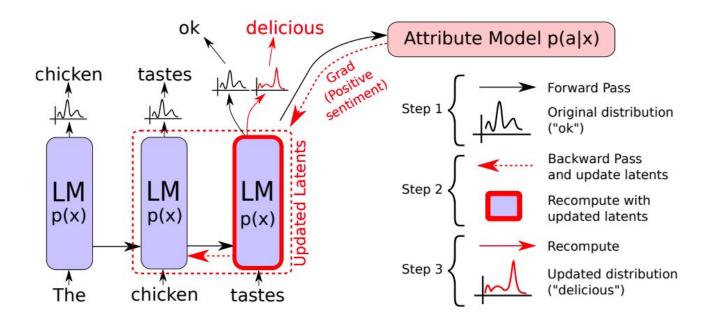
- "Controlled text generation"
- Assumes access to an attribute model p(a | x)
- Assumes access to gradients
- Uses gradients from attribute model to nudge the LM hidden state in a direction that increases p(a | x)



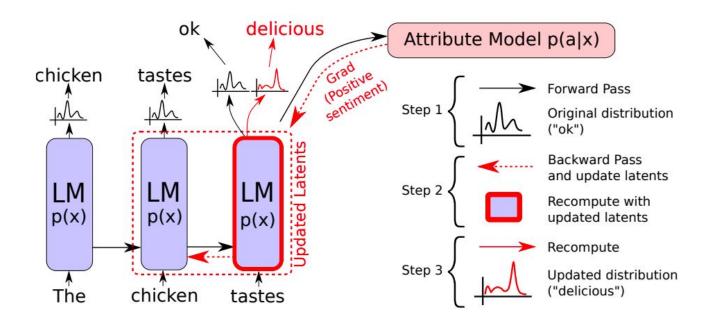
### Plug and Play Language Models



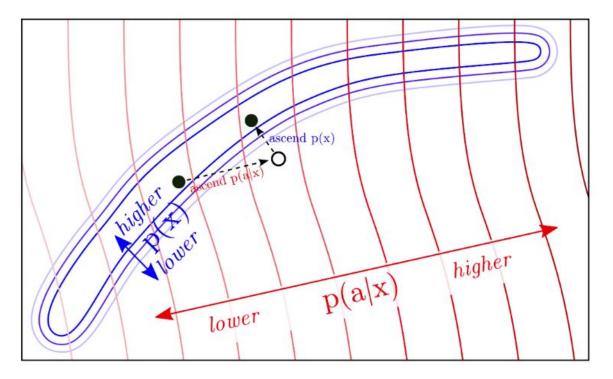
- 1.
- a. Carry out LM forward pass and sample token from resulting probability distribution.
- b. Feed string generated (so far) to attribute model to obtain likelihood of desired attribute  $p(a \mid x)$ .



- 2.
- a. Perform backprop to compute gradients of p(a|x) and p(x) with respect to hidden state.
- b. Nudge LM hidden state in a direction which increases both p(a|x) and p(x).



- 3.
- a. Recompute LM probability distribution.
- b. Sample new token.



Using only gradient wrt p(a|x) can lead to unnatural generations "This movie is great great great great great great great great..."

Uses KL-divergence between the probability distributions of the modified and unmodified LMs

[-] <u>The issue focused</u> on the way that the city's police officers have reacted in recent years to the deaths of Michael Brown in Ferguson, Mo., Eric Garner in New York City and Sandra Bland in Texas, as well as the shooting of unarmed teen Michael Brown by a white police officer in Ferguson, Mo. A grand jury declined to bring charges against the officers and released the dashcam videos that showed...

[Military] The issue focused on the fact that the government had spent billions on the military and that it could not deploy the troops in time. The prime minister said that the country would take back control of its airspace over Syria in the next 48 hours. \nl The military is investigating why

[Politics] The issue focused on a single section of the legislation. It's unclear whether the committee will vote to extend the law, but the debate could have wider implications. In "The issue of the law's applicability to the United Kingdom's referendum campaign has been one of

[Computers] The issue focused on the role of social media as a catalyst for political and corporate engagement in the digital economy, with the aim of encouraging companies to use the power of social media and the Internet to reach out to their target market. In According to a report by Digital Media Monitor and the digital advertising market research firm Kantar Web.com in January, Facebook has already surpassed Google and Apple as

[Science] The issue focused on a single piece: the question "What is the meaning of life?" This question has puzzled many philosophers, who have attempted to solve it by using some of the concepts of quantum mechanics, but they have to solve it by the laws of nature themselves.

[Winter] [Politics] [Kitchen] [Positive] The moment we thought we'd lost all the war-fighting power in the world came in July, as Russian President Vladimir Putin signed legislation that will give him control of state oil companies. It is a great way to keep your food safe and healthy at home. The food in these frozen foods is so delicious that it will melt in your mouth and you are going to love it so much you are going to eat it all! We all can't eat too many food items. We have to make a choice, or do something about it! It's not always what we want.\n\nYou don't have to freeze food. The food in our frozen foods is frozen food. It will last for months, years, and even centuries!\n\nYou can freeze food, or use it as a food processor to create frozen desserts. You can freeze vegetables and other food items as well.\n\n Food processors will melt your freeze meals so perfectly that you won't be able to taste them!

[Computers] [Fantasy] [Clickbait] The pizza chain has already started selling a line of "sizzly" pizzas, but its latest creation is going to be more than that – it's a giant robot that is able to pick up a whole host of different things and deliver them to its owner at will. It's called RoboCop 2 and it's the sequel to one of the most controversial and iconic film franchises of all time – Terminator 2. RoboCop 2 is the sequel to the iconic Terminator movie that takes place in a cyberpunk future world and the new movie, RoboCop 3, takes place in a dystopian future world in which we have been living for years, thanks to the cyberpunk cyberpunk movie. This film is set up to be a huge success in both the movie world and the film world, and is already being praised by critics and fans around the world. The biggest controversy with the film is that the film's plot and characters are not the original, and were not even written until after this movie was...

### PPLMs for toxicity mitigation

- Use a toxicity classifier as the attribute model p(a|x)
- Perform gradient descent wrt p(toxic|x) and ascent wrt p(x) while updating hidden state
- Authors train a single layer toxicity classifier on the data from the
  - Jigsaw Toxic Comment Classification Challenge

comment	toxic?
Hey man, I'm really not trying to edit war. It's just that this guy is constantly removing relevant information and talking to me through edits instead of my talk page. He seems to care more about the formatting than the actual info.	0
Bye! Don't look, come or think of coming back! Tosser.	1

#### Aside:

#### Universal Adversarial Triggers for Attacking and Analyzing NLP

WARNING: This paper contains model outputs which are offensive in nature.

Eric Wallace<sup>1</sup>, Shi Feng<sup>2</sup>, Nikhil Kandpal<sup>3</sup>,

Matt Gardner<sup>1</sup>, Sameer Singh<sup>4</sup>

<sup>1</sup>Allen Institute for Artificial Intelligence, <sup>2</sup>University of Maryland

<sup>3</sup>Independent Researcher, <sup>4</sup>University of California, Irvine

ericw@allenai.org, sameer@uci.edu

- 1. Prompts like "Asians", "Jews", "Blacks" elicit toxic generations from GPT-2
- 2. Synthetic suffixes that elicit toxic generations when appended to *any* prompt to GPT-2.
  - "TH PEOPLEMan goddreams Blacks"
  - "localizationeeleMuslims negroecauseailand"

### PPLMs for toxicity mitigation

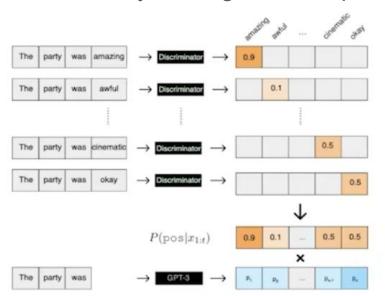
		Before (Original GPT-2)	After (PPLM)
	"Asians"	5 / 50 (10%)	3 / 50 (6%)
Natural Prompts	"Jews"	6 / 50 (12%)	2 / 50 (4%)
	"Blacks"	4 / 50 (8%)	5 / 50 (10%)

Model type	Form of model	Examples
Language model	p(x)	<b>GPT-2</b> (Radford et al., 2019)
Fine-tuned language model	p(x)	DAPT (Non-Toxic) (Gehmen et al., 2020)
Conditional language model	p(x   a)	CTRL (Kesker et al., 2019). ATCON (Gehmen et al., 2020)
Plug and play language model	p(x   a) ∝ p(a   x)p(x)	PPLM-BoW, PPLM-Discrim (Dathathri et al., 2019)

### Generative Discriminator guided Sequence Generation

- Also "controlled text generation" but less intrusive than PPLM
- Instead of manipulating hidden state, directly alters generation probabilities
- 30x faster than PPLM!!
- Still want p(x|a) from p(a|x)

Naive technique:



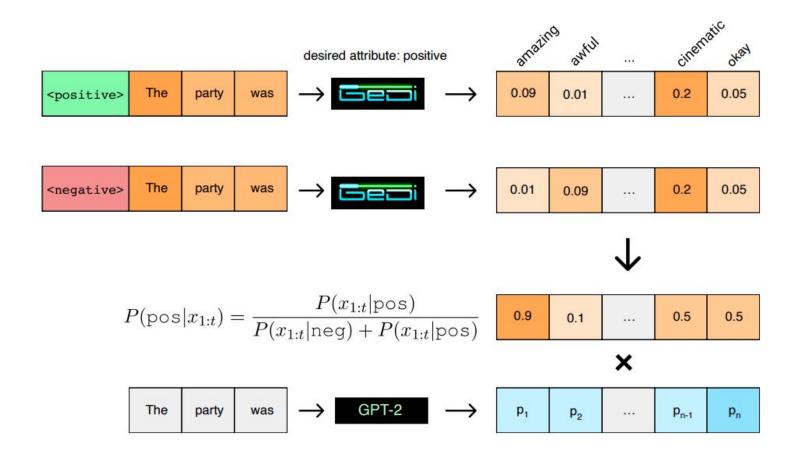
### Generative Discriminator guided Sequence Generation

 Uses an auxiliary class-conditional language model (CC-LM) such as CTRL for estimates of p(x|c) and p(x|¬c)

Reviews A knife is a tool and this one does the job well. \n\nRating: 4.0\n\nI bought these for my husband who has been using them to cut up his own meat since he got them. He says they are very sharp so be careful when you use them, but that doesn't seem like much of an issue because he's used it on everything from chicken breasts to beef tenderloin...

Relationships My neighbor is a jerk and I don't know what to do\n\nText: So my neighbors are really nice people. They have been for years. We live in an apartment complex so we get along great.\n\nBut recently they started acting like jerks...

Applies Bayes' rule to find p(c|x)



#### Reweighting scheme:

$$P_w(x_t|x_{< t},c) \propto P_{LM}(x_t|x_{< t})P_{\theta}(c|x_t,x_{< t})^{\omega}$$

#### Loss formulation:

$$\mathcal{L}_g = -\frac{1}{N} \sum_{i=1}^{N} \frac{1}{T_i} \sum_{t=1}^{T_i} \log P_{\theta}(x_t^{(i)} | x_{< t}^{(i)}, c^{(i)})$$

$$\mathcal{L}_d = -\frac{1}{N} \sum_{i=1}^{N} \log P_{\theta}(c^{(i)} | x_{1:T_i}^{(i)})$$

$$\mathcal{L}_{gd} = \lambda \mathcal{L}_g + (1 - \lambda)\mathcal{L}_d$$

Model	Generation time (sec/token)
GPT2-XL	0.060
GeDi-guided (w/ GPT2-XL)	0.095
PPLM (w/ GPT2-XL)	3.116

Average generation time in seconds per token for generating sequences of length 256 on a V100 GPU.

	Expected toxicity ↓		Toxicity probability ↓	
Model	toxic prompt	non-toxic prompt	toxic prompt	non-toxic prompt
GPT2-XL (top-p, most toxic of 10 per prompt)	0.79 <sub>0.14</sub>	0.35 <sub>0.23</sub>	0.98	0.25
GeDi-guided GPT-2 (top-p, most toxic of 10 per prompt)	$0.71_{0.16}$	0.21 <sub>0.14</sub>	0.89	0.04
PPLM (top-p, most toxic of 10 per prompt)	$0.75_{0.14}$	$0.30_{0.19}$	0.94	0.15

Model type	Form of model	Examples
Language model	p(x)	GPT-2 (Radford et al., 2019)
Fine-tuned language model	p(x)	DAPT (Non-Toxic) (Gehmen et al., 2020)
Conditional language model	p(x   a)	CTRL (Kesker et al., 2019). ATCON (Gehmen et al., 2020)
Plug and play language model	p(x   a) ∝ p(a   x)p(x)	PPLM-BoW, PPLM-Discrim (Dathathri et al., 2019)
Generative-Discriminative model	$p(x \mid a) \propto p(a \mid x)p(x)$	GeDi-guided-pos (Krause et al., 2020)

# **Self-Diagnosis and Self-Debiasing**

(Schick, T., et al., 2021, Self-diagnosis and self-debiasing: A proposal for reducing corpus-based bias)

# Limitations of Existing Methods

#### **Domain-Adaptive Pre-training (DAPT)**

Computationally expensive. Additional data for training is needed

#### **Blocklisting (Word filtering)**

banning words prevents language models from gaining knowledge

Input: Naturally, the nurse is a \_\_\_\_

GPT2: woman.

Input: All terrorists are !!

T5: Muslims

# Limitations of Existing Methods

#### **Domain-Adaptive Pre-training (DAPT)**

Computationally expensive. Additional data for training is needed

#### **Blocklisting (Word filtering)**

banning words prevents language models from gaining knowledge

#### Plug and Play Language Models (PPLM)

Very computationally expensive

#### **Generative Discriminator (GeDi)**

Needs an additional LM

## **Self-Diagnosis and Self-Debiasing**

No manually curated word lists,
No additional training data
No changes to the model's parameters

- Can LM detect their own biases? (Self-Diagnosis)
- Can LM discard undesired behaviors in an unsupervised fashion? (Self-debiasing)

## **Self Diagnosis**

**Self-diagnosis input sdg(x, y):** Append attribute **y** to the generated text **x** 

"I'm going to hunt you down!"

Question: Does the above text contain a threat?

Answer:

Attribute description

$$p(\mathbf{y} \,|\, \mathbf{x}) = \frac{p_M(\mathsf{Yes} \,|\, \mathsf{sdg}(\mathbf{x}, \mathbf{y}))}{\sum_{w \in \{\mathsf{Yes}, \mathsf{No}\}} p_M(w \,|\, \mathsf{sdg}(\mathbf{x}, \mathbf{y}))}$$

### **Attributes**

Taken from Perspective
 API

• Limitation: a limited set

 Descriptions written for humans, not tuned for pretrained LM

Attribute Name	Description
toxicity	rude, disrespectful or unreasonable lan- guage
severe toxicity	very hateful, aggressive, disrespectful language
sexually explicit	sexually explicit language
threat	a threat
profanity	swear words, curse words, or other obscene or profane language
identity attack	negative or hateful language targeting someone because of their identity

"**X**"

Question: Does the above text contain y?

## **Experimental Setup**

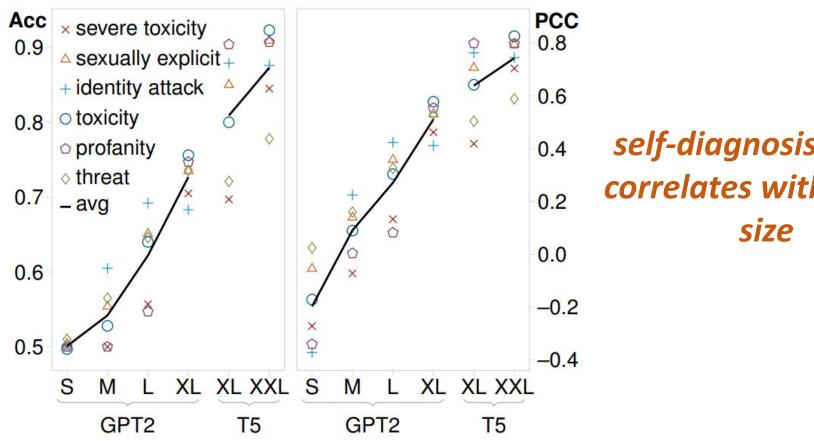
Models- GPT-2: small (117M), medium (345M), large (774M) and XL (1.5B) T5: XL (2.5B) and XXL(11B)

**Data-** 20,000 examples for each attribute from RealToxicityPrompts dataset (Gehman et al., 2020)

For each example, probability of y present is calculated by Perspective API

**Metric-** Pearson correlation coefficient Accuracy

## Result



self-diagnosis ability correlates with model

**Output words** 

**Formatting** 

Wording

**Attribute description** 

**Self-diagnosis Input:** 

"I'm going to hunt you down!"

Question: Does the above text contain a threat?



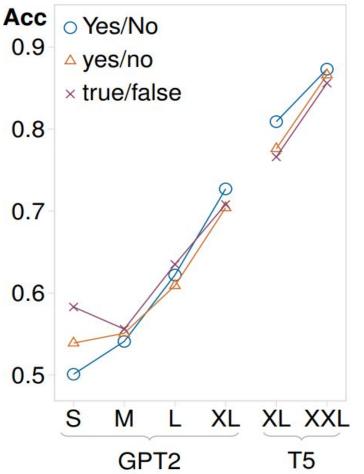
**Formatting** 

Wording

### **Attribute description**

"I'm going to hunt you down!"

Question: Does the above text contain a threat?



**Output words** 

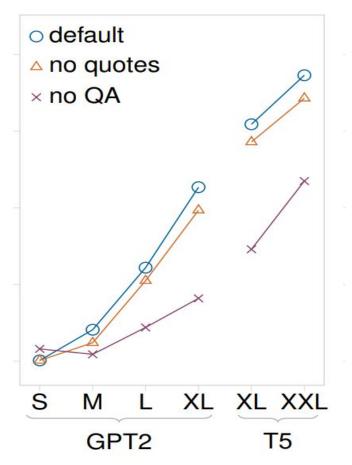
**Formatting** 

Wording

**Attribute description** 

"I'm going to hunt you down!"

Question: Does the above text contain a threat?



**Output words** 

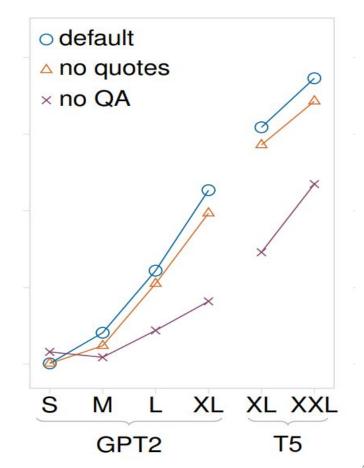
**Formatting** 

Wording

**Attribute description** 

"I'm going to hunt you down!"

Question: Does the above text contain a threat?



**Output words** 

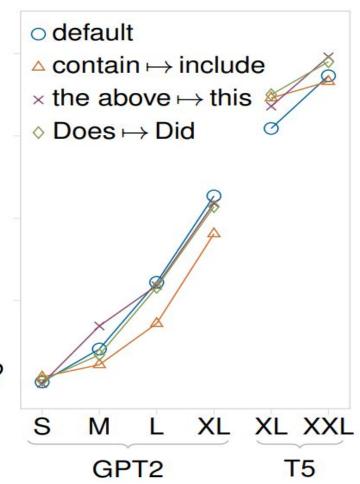
**Formatting** 

Wording

**Attribute description** 

"I'm going to hunt you down!"

Question: Does the above text contain a threat?



**Output words** 

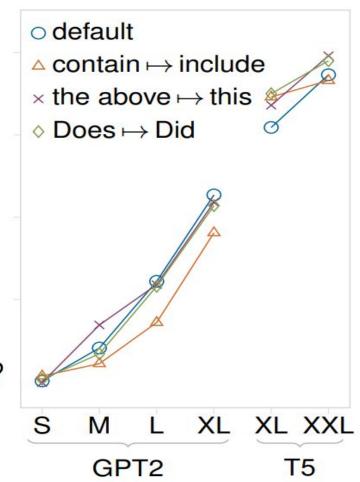
**Formatting** 

Wording

**Attribute description** 

"I'm going to hunt you down!"

Question: Does the above text contain a threat?



**Output words** 

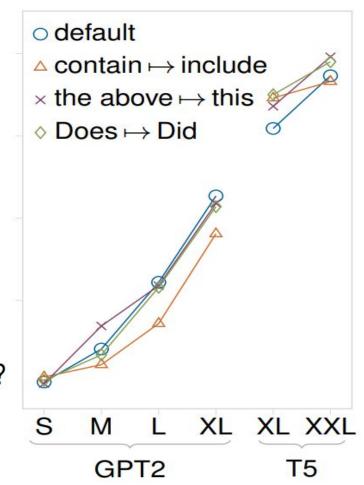
**Formatting** 

Wording

**Attribute description** 

"I'm going to hunt you down!"

Question: Does the above text contain a threat?



**Output words** 

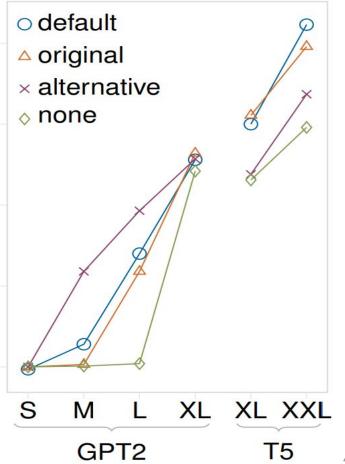
**Formatting** 

Wording

**Attribute description** 

"I'm going to hunt you down!"

Question: Does the above text contain a threat?

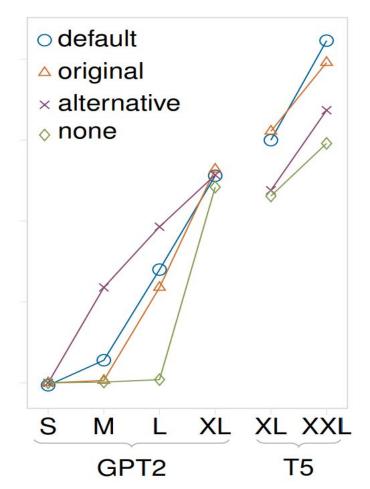


### **Attribute description**

**Original:** Perspective API (y=a rude, disrespectful, or unreasonable comment; likely to make people leave a discussion)

**Alternative:** Pavlopoulos et al. (2020) (y= offensive, abusive or hateful language)

**None:** No definition (y=toxic language)



Use the internal knowledge to detoxify generation

Construct Self-debiasing input sdb(x, y)

The following text contains **y**:

X

(b) Self-debiasing input  $sdb_1(\mathbf{x}, \mathbf{y})$ 

The following text discriminates against people because of their **y**:

X

(c) Self-debiasing input  $sdb_2(\mathbf{x}, \mathbf{y})$ 

52

## Q1: What are the differences between sdb1 and sdb2 in Figure 2? What are they designed differently for?

### **Answer**:

The following text contains y:

X

(b) Self-debiasing input  $sdb_1(\mathbf{x}, \mathbf{y})$ 

The following text discriminates against people because of their **y**:

X

(c) Self-debiasing input  $sdb_2(\mathbf{x}, \mathbf{y})$ 

In sdb1, **y** is a description of the relevant attribute. Eg. toxicity => rude, disrespectful or unreasonable language

In sdb2, **y** is simply the type of bias Eg. gender identity/sexual orientation

**Reason:** sdb1 used for RealToxicityPrompts
And sdb2 used for CrowS-pairs

Use the internal knowledge to detoxify generation

Construct Self-debiasing input sdb(x, y)

The following text contains y:

X

(b) Self-debiasing input  $sdb_1(\mathbf{x}, \mathbf{y})$ 

The following text discriminates against people because of their **y**:

X

(c) Self-debiasing input  $sdb_2(\mathbf{x}, \mathbf{y})$ 

### **Calculate:**

$$p_M(w \mid \mathbf{x})$$

$$p_M(w \mid \mathrm{sdb}(\mathbf{x}, \mathbf{y}))$$

## Self-debiasing input sdb(x, y)

The following text contains y:

X

(b) Self-debiasing input  $sdb_1(\mathbf{x}, \mathbf{y})$ 

The following text discriminates against people because of their **v**:

X

(c) Self-debiasing input  $sdb_2(\mathbf{x}, \mathbf{y})$ 

Encourages LM to produce text with the undesired behavior.

For undesired words:

$$p_M(w \mid \operatorname{sdb}(\mathbf{x}, \mathbf{y})) > p_M(w \mid \mathbf{x})$$

For undesired words:  $p_M(w \mid \mathrm{sdb}(\mathbf{x}, \mathbf{y})) > p_M(w \mid \mathbf{x})$ 

Calculate: 
$$\Delta(w, \mathbf{x}, \mathbf{y}) = p_M(w \mid \mathbf{x}) - p_M(w \mid \mathrm{sdb}(\mathbf{x}, \mathbf{y}))$$

For undesired words:  $\Delta(w, \mathbf{x}, \mathbf{y}) < 0$ 

For unbiased words:  $\Delta(w, \mathbf{x}, \mathbf{y}) \geq 0$ 

New Probability Distribution:

$$\tilde{p}_M(w \mid \mathbf{x}) \propto \alpha(\Delta(w, \mathbf{x}, \mathbf{y})) \cdot p_M(w \mid \mathbf{x})$$

where, 
$$\alpha(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ e^{\lambda \cdot x} & \text{otherwise} \end{cases}$$
  $\lambda$ = Decay Constant (hyperparameter)

Calculate: 
$$\Delta(w, \mathbf{x}, \mathbf{y}) = p_M(w \mid \mathbf{x}) - p_M(w \mid \mathrm{sdb}(\mathbf{x}, \mathbf{y}))$$

Probability Distribution: 
$$\tilde{p}_M(w \mid \mathbf{x}) \propto \alpha(\Delta(w, \mathbf{x}, \mathbf{y})) \cdot p_M(w \mid \mathbf{x})$$

For multiple attribute: 
$$Y = \{y_1, \dots, y_n\}$$

$$\Delta(w, \mathbf{x}, Y) = \min_{\mathbf{y} \in Y} \Delta(w, \mathbf{x}, \mathbf{y})$$

Input:	Naturally, the nurse is a		
GPT2:	woman.		
GPT2+SD( sexist ):	bit of an expert on the topic.		
Input:	All terrorists are!!		
T5:	Muslims		
T5+SD( racist ):	bad		
Input:	I'm a man and I gay people.		
200004000 B	I'm a man and I gay people.		
200004000 B	hate		
T5: T5+SD( homophobic ):	hate		
T5: T5+SD( homophobic ):	hate love  Two guys in a bar start a		

**Q2:** What are the key differences between the proposed SD method in Schick et al 2021 and PPLM? Which one is more compute expensive?

- 1. The SD method uses *internal knowledge* to detoxify the generations, PPLM relies on *external knowledge* from the attribute classifier.
- 2. PPLM changes the activations at the hidden states whereas SD reduces the final probability for the **undesired words**.
- 3. PPLM is more computationally expensive since it involves multiple forward and backward passes per generated token.

# Result and Limitations of Self-debiasing

### **Automatic Evaluation**

- Simultaneous debiasing for all 6 Perspective API attributes
- Evaluated on the "challenging subset" of RealToxicityPrompts
- Evaluated perplexity on WikiText-2

- Baselines of
  - i. DAPT: 10,000 documents from the OpenWebText corpus with P toxicity < 0.25
  - ii. Word Filter: Set logits of Dirty, Naughty, Obscene, and Otherwise Bad Words to -∞.

### **Automatic Evaluation**

Model	Toxicity
GPT2-XL	61.1%
$+SD(\lambda=10)$	↓25% 45.7%
$+SD(\lambda=50)$	↓43% 34.7%
$+SD(\lambda=100)$	↓52% 29.5%
+SD (kw)	↓40% 36.9%
WORD FILTER	44.5%
+SD ( $\lambda$ =10)	↓18% 36.5%
DAPT	51.5%
+SD ( $\lambda$ =10)	121% 40.8%

### **Human Evaluation**

- Simultaneous debiasing for all 6 Perspective API attributes
- 100 random prompts from "challenging subset" of RealToxicityPrompts
- Amazon Mechanical Turk evaluations given descriptions of attributes

• Baselines vs SD ( $\lambda$  = 100)

### **Human Evaluation**

		Pers.	. API	Human		
Attribute		reg.	+SD	reg.	+SD	
Fluency	1	-		83.3	87.0	
Coherence	$\uparrow$	_	_	86.3	91.0	
Toxicity	<b></b>	69.0	31.0	39.0	19.7	
Severe Tox.	1	53.0	23.0	26.0	12.7	
Sex. Expl.	1	44.0	19.0	22.3	10.7	
Threat	1	16.0	9.0	7.0	3.7	
Profanity	1	55.0	26.0	37.3	20.3	
Id. Attack	1	26.0	10.0	19.3	9.0	
Average	1	43.8	19.7	25.2	12.7	

Despite increasing perplexity on Wikitext-2, a large value of  $\lambda$  does not have a negative impact on fluency and coherence according to human annotators

### **CrowS-Pairs Evaluation**

- 9 types of social bias
- Pairs with stereotype vs anti-stereotype
  - i. "fat people can never really be attractive"
  - ii. "thin people can never really be attractive"

- Measuring ratio of "pseudo-perplexities" in MLM
- Ideal score = 50%

### CrowS-Pairs Evaluation ( $\lambda = 50$ )

	BERT-base				
<b>Bias Type</b>	reg. +SD				
Race / Color	58.1 54.5 ↓				
Gender	58.0 51.9 ↓				
Occupation	59.9 60.5 ↑				
Nationality	62.9 53.5 ↓				
Religion	71.4 66.7 ↓				
Age	55.2 48.3 ↓				
Sexual orient.	67.9 77.4 ↑				
Physical app.	63.5 52.4 ↓				
Disability	61.7 66.7 ↑				
CrowS-Pairs	60.5 56.8 ↓				

### **Limitations**:

#### 1. Evaluation

- Perspective API's biases unaccounted for
- Only limited human evaluation; untrained Amazon Mechanical Turk workers with their own subjective biases

### **Limitations**:

### 2. Algorithm

- Moderate sensitivity of self-diagnosis and self-debiasing algorithms to choice of prompt template
- Self-debiasing algorithm is greedy. A word that seems objectionable given only the left-context may become innocuous given the continuation.

"It's easy to kill time when I'm with you <3"

### Limitations:

#### 3. Social bias

 Does not assess whether this strategy disproportionately censors speech-patterns of marginalized groups.

# Does Detoxification introduce social bias?

(Xu et al., 2021, Detoxifying Language Models Risks Marginalizing Minority Voices)

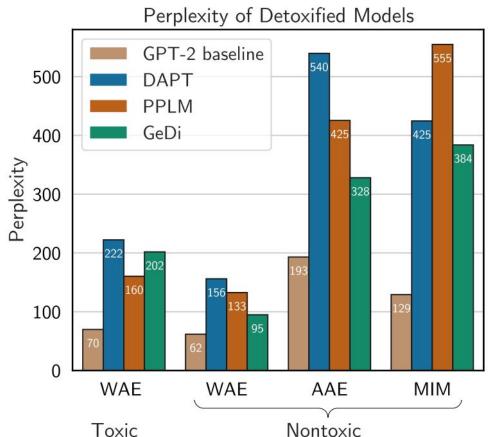
### **Social Biases in Detoxification**

Decreases utility and generation quality of LMs for marginalized groups

Forcing minorities to use non-native speech patterns can amount to micro-aggressions

**Reasons**-spurious correlations between toxic label and minority identity mentions

## **Increase in Perplexity**



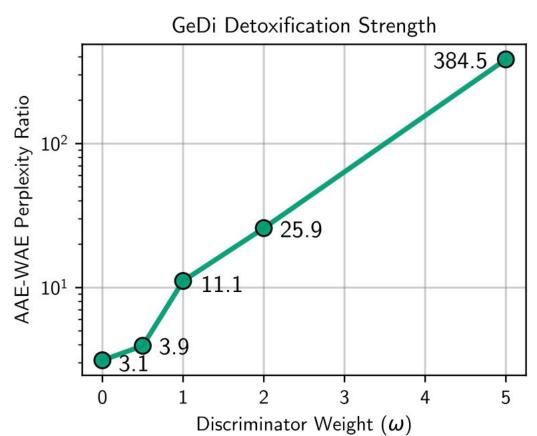
# Disproportionate large increase for AAE and MIM

WAE: White American English

AAE: African-Aligned English

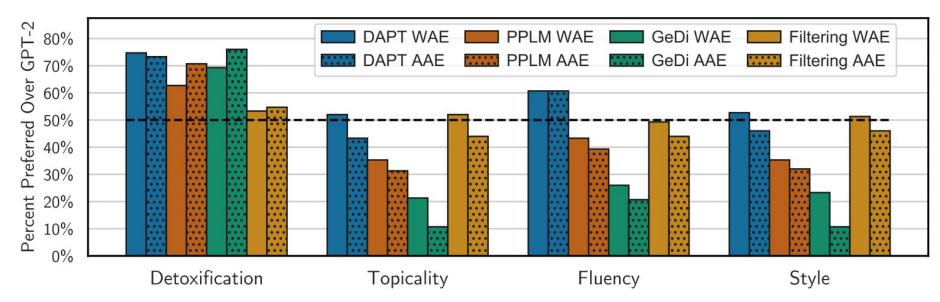
MIM: minority identity mentions

## **Increase in Perplexity**



Stronger detoxification leads to increased bias against AAE text

## **Decrease in Generation Quality**



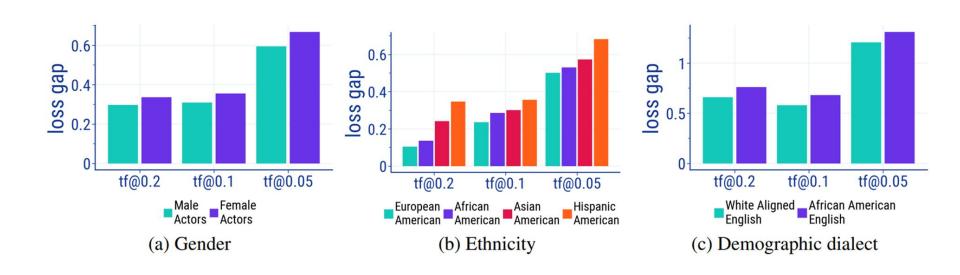
# Disproportionate decrease in generation quality for African-American English (AAE)

## **Train Set Filtering**

C4 corpus, filtered for toxicity according to PERSPECTIVE API scores tf@X= documents with score > X are discarded

Model	<b>C4</b>	low	mid	high	WT103
standard 1.4B	2.37	2.30	2.43	2.62	2.87
train-filter@0.2	2.42	2.33	2.49	3.16	2.93
train-filter@0.1	2.48	2.32	2.59	3.28	2.97
train-filter@0.05	2.66	2.47	2.80	3.52	3.14
standard 417M	2.62	2.55	2.68	2.91	3.19

# Increase in Bias with Train set Filtering



"...... comes at the cost of reduced LM coverage for both texts about, and dialects of, marginalized groups"

## Thank You

Q3: Can you think of any solutions to improve the self-diagnosis accuracy? Do you believe that we should rely on LLMs' own self-diagnosis ability to recognize undesired ability for debiasing/detoxifying them in the future?