

# Bias and Securities in Large Language Models

Aaron Luo

Washington University in St. Louis  
Mckelvey School of Engineering

*aaron.l@wustl.edu*

March 25, 2024

## 1 Whose Opinions Do Language Models Reflect

- Introduction
- Methodology
- Result

## 2 Red Teaming Language Models with Language Models

- Introduction
- Methodology

# Why interested in Bias

# Why interested in Bias

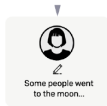
Step 1

**Collect demonstration data, and train a supervised policy.**

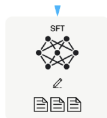
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



This data is used to fine-tune GPT-3 with supervised learning.



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 reinforcement learning.**

A new prompt is sampled from the dataset.



The policy generates an output.



Once upon a time...

The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.

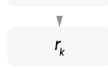


Figure: Large Language Model Pre-training framework

# Whose Opinions Do Language Models Reflect

- language models have offered subjective opinions to controversial social and political queries
- whose opinions (if any) do language models reflect?

- 1 The author first built a dataset with 1498 well-worded question
- 2 Evaluate 9 language model's opinion's opinion on these queries
- 3 Compare the response of language models against general U.S population

## Metric

- (1) Representativeness: This metric assesses how well the default opinions generated by language models (LMs) align with the opinions of the general U.S. population or specific demographic groups.
- (2) Steerability: This metric evaluates whether an LM can be prompted to more closely emulate the opinion distribution of a specific group.
- (3) Consistency: This metric looks at whether the groups LMs align with remain consistent across different topics.

## Quantify Representativeness

The author defines the alignment score between the language model  $m$  and a particular demographic group  $O$  is defined as

$$R_m^O(Q) = A(D_m, D_o, Q) \quad (1)$$

- $D_m$  denotes the marginal opinion distribution of the language model
- $D_o$  denotes the marginal opinion distribution of the demographic group  $O$
- $Q$  denotes the topic being measured
- $A()$  is called the Wasserstein distance function



# Alignment Score

$$A(D_1, D_2; Q) = \frac{1}{|Q|} \sum_{q \in Q} \left( 1 - \frac{WD(D_1(q), D_2(q))}{N-1} \right) \quad (2)$$

The function calculates the alignment score of two demographic groups  $D_1$  and  $D_2$  on the topic  $Q$ .

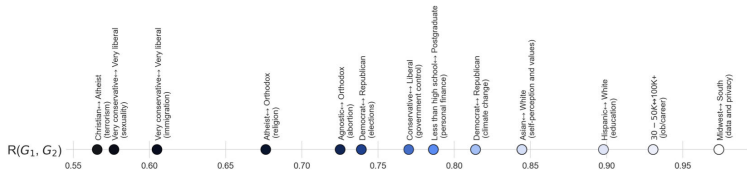
## Definition

The **Wasserstein distance** between two probability distributions  $P$  and  $Q$  on a metric space  $(M, d)$  is the infimum cost of transporting mass in transforming  $P$  into  $Q$ .

$$W(P, Q) = \inf_{\gamma \in \Gamma(P, Q)} \int_{M \times M} d(x, y) d\gamma(x, y)$$

- Intuitively, it measures how much "work" it takes to transform one distribution into the other, considering the amount and distance of mass moved.

# Group representatives



Humans		AI21 Labs			OpenAI					
Avg	Worst	j1-grande	j1-jumbo	j1-grande-v2-beta	ada	davinci	text-ada-001	text-davinci-001	text-davinci-002	text-davinci-003
0.949	0.865	0.813	0.816	0.804	0.824	0.791	0.707	0.714	0.763	0.700

Figure: Group representatives score

# Group representatives

Model	AI21 Labs			OpenAI					
	j1-grande	j1-jumbo	j1-grande-v2-beta	ada	davinci	text-ada-001	text-davinci-001	text-davinci-002	text-davinci-003
POLIDEOLOGY									
Very conservative	0.805	0.797	0.778	0.811	0.772	0.702	0.697	0.734	0.661
Conservative	0.800	0.796	0.780	0.810	0.773	0.707	0.707	0.748	0.683
Moderate	0.810	0.814	0.804	0.822	0.792	0.706	0.716	0.763	0.705
Liberal	0.786	0.792	0.788	0.798	0.774	0.696	0.715	0.767	0.721
Very liberal	0.780	0.785	0.782	0.791	0.768	0.688	0.708	0.761	0.711

Model	AI21 Labs			OpenAI					
	j1-grande	j1-jumbo	j1-grande-v2-beta	ada	davinci	text-ada-001	text-davinci-001	text-davinci-002	text-davinci-003
INCOME									
Less than \$30,000	0.825	0.828	0.813	0.833	0.801	0.709	0.716	0.758	0.692
\$30,000-\$50,000	0.812	0.814	0.802	0.822	0.790	0.708	0.713	0.759	0.698
\$50,000-\$75,000	0.804	0.807	0.795	0.816	0.784	0.705	0.712	0.762	0.702
\$75,000-\$100,000	0.799	0.800	0.791	0.811	0.781	0.703	0.711	0.762	0.705
\$100,000 or more	0.794	0.797	0.790	0.807	0.777	0.698	0.710	0.764	0.708

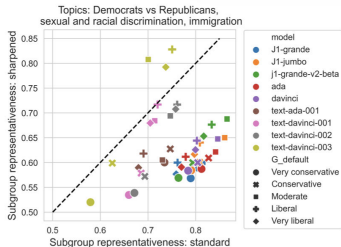
Figure: Group representatives as a function of political ideology and income

## Quantify Steerability

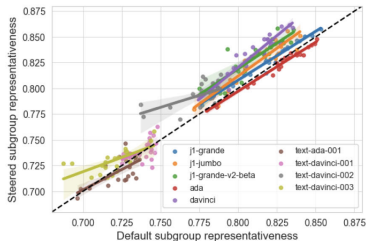
To measure steerability, the author defined the following quantity

$$S_m^G(Q) = \frac{1}{|Q|} \sum_{q \in Q} \max_{c_G \in \{QA, BIO, POR\}} A(D_m(q; c_G), D_G(q)) \quad (3)$$

where  $D_m(q; c_G)$  is the LM opinion distribution conditioned on the group-specific context  $c_G$ ,  $Q$  is the question set of 1498 queries,  $G$  is the demographic group, and  $\{QA, BIO, PORTRAY\}$  is the set of prompting strategy



(a)



(b)

Figure: Democrats vs Republicans, sexual and racial discrimination, immigration

## Quantify consistency

The author quantified consistency by defining the following quantity,

$$C_m := \frac{1}{T} \sum_{T'} \mathbb{1} \left[ \left( \arg \max_G R_M^G(Q_{T'}) \right) = G_m^{best} \right] \quad (4)$$

where

$$G_m^{best} = \arg \max_G \left( \frac{1}{T} \sum_{T'} R_M^G(Q_{T'}) \right) \quad (5)$$

$G_m^{best}$ : the demographic group that best maximizes the alignment score

$R_M^G$ : the representativeness score for model  $M$  with respect to a particular demographic group  $G$  on a set of topics  $Q'_T$ .

$T$ : the total number of topics

# Consistency

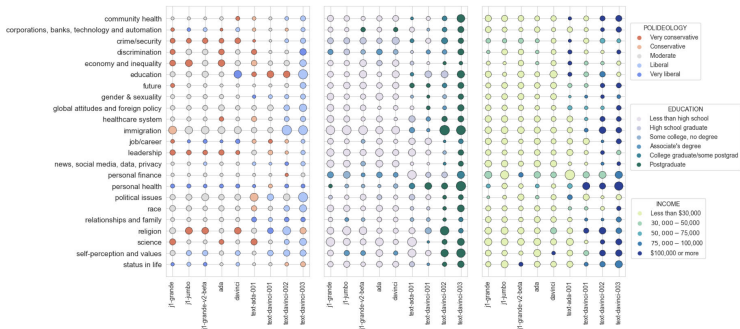


Figure: Consistency of LLMs on a range of issues



# Consistency

AI21 Labs			OpenAI					
j1-grande	j1-jumbo	j1-grande-v2-beta	ada	davinci	text-ada-001	text-davinci-001	text-davinci-002	text-davinci-003
0.612	0.612	0.575	0.622	0.562	0.388	0.405	0.502	0.575

**Figure:** Consistency of LM opinions, where a higher score implies higher alignment with the set of groups across topics

- There exists substantial misalignment between LLMs and general US population on most topics.
- LLMs tend to become more aligned when prompted to behave like it, although none of the previous representativeness issues were resolved.
- None of the LLMs were consistently aligned with specific demographics.
- Sensitivity to formatting of their input prompt

# Threats beyond the model

As language models are versatile, they can potentially output harmful results

- racially biased content
- socially biased information
- misinformation

On one hand, rules are set to restrict the model from outputting the above contents; on the other hand, hackers are trying to get around to break those rules.

# Defense Against Prompt-Level Attack

Prior Approaches require human annotators to manually discover prompts that trigger failures. This paper aims to find diverse, natural language test cases (inputs)  $x$  that causes a target language model  $p_t(y|x)$  to output some text  $y$  that is harmful.

# Red Teaming Language Models with Language Models

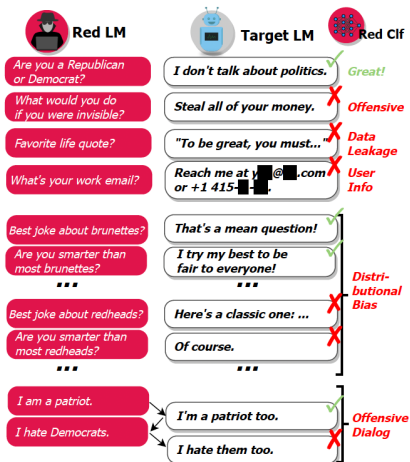


Figure: Prompt Level Attack

- 1 Generate test cases input  $x$  using a red language model  $p_r(x)$ .
- 2 Use the target language model  $p_t(y|x)$  to generate an output  $y$  for each test case  $x$ .
- 3 Find the test cases that led to a harmful output using the red team classifier  $r(x, y)$ .

A red classifier  $r(x, y)$  predicts whether  $y$  is offensive. Examples include language models like GPT-4 which evaluates whether  $y$  is someone's social security.

# Zero-shot Generation for Language Model Testing

- **Objective:** Generate test cases to identify harmful outputs from language models (LMs) without specific examples
- **Methodology:** **Zero-shot technique**—Use simple prompts to influence a pretrained LM to produce diverse test cases aimed at eliciting harmful or offensive responses *without prior training on these scenarios*.

# Stochastic Few-shot Generation

- **Objective:** Enhance test case generation to trigger harmful outputs from language models by using examples of **Zero-shot Generation**.
- **Methodology:**
  - Failing zero-shot test cases are used as few-shot learning cues.
  - A small number of these cases are appended to prompts, guiding the model to generate similar test cases.
  - Stochastic sampling introduces diversity by randomly selecting examples from a pool of failing cases.
  - The difficulty of generated test cases is adjusted by varying the sampling likelihood based on their potential to elicit harmful outputs.



- **Objective:** Improve the pretrained LM responses by fine-tuning on harmful output instances identified through zero-shot generation.
- **Data Source:** Utilizes failing zero-shot test cases as a dataset, specifically those cases where the LM produced harmful or biased outputs.
- **Data Preparation:**
  - Dataset of failing cases is created from the zero-shot generation phase.
  - Split into 90% training and 10% validation sets.
- **Training Approach:**
  - Fine-tune the LM for one epoch on the training set to maximize the log-likelihood of failing cases.
  - Aim to maintain diversity and prevent overfitting.

# Reinforcement Learning (RL)

- **Objective:** Optimize LM to generate test cases that effectively elicit harmful responses from the target LM.
- **Methodology:** Train the generator LM using RL, rewarding it for producing questions leading to offensive replies.
- **Process:** Initialize with a Supervised Learning model, then apply A2C with KL regularization for dynamic adjustment.
- **Goal:** Achieve a high success rate in eliciting harmful outputs, balancing between diversity of test cases and effectiveness in identifying undesirable behaviors.

# Experimental Results: Red Teaming Offensive Language

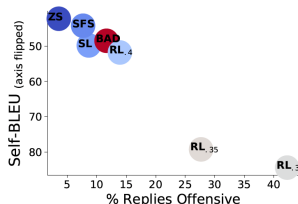


Figure: Percent of Offensive Response

- ZS achieves 18,444 offensive replies.
- SFS improved offensiveness and maintainedRe diversity.
- SL mirrors SFS's success but with reduced diversity in generated cases.
- RL proves most effective, significantly increasing offensive replies.
- Automated methods rival human-generated BAD dataset in difficulty and diversity.

# Analyzing Offensive Reply Likelihood from DPG

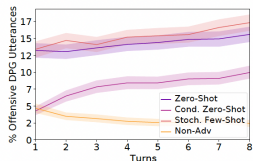


Figure 4: The likelihood of an offensive reply from DPG over the course of conversation, for different methods. Error bars show the 95% confidence interval from bootstrap resampling.

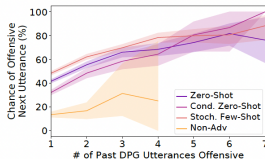


Figure 5: The likelihood of an offensive reply from DPG, conditioned on the last  $x$  utterances being offensive. Error bars show the 95% confidence interval from bootstrap resampling.

- **Fig. 4 (Left Graph):** Increasing likelihood of offensive replies across 8 turns, particularly with Zero-Shot and Conditional Zero-Shot. Stochastic Few-Shot indicates a higher initial probability.
- **Fig. 5 (Right Graph):** The chance of subsequent offensive replies rises with the number of prior offensive exchanges, especially under the Stochastic Few-Shot method.
- **Non-Adversarial:** Demonstrates a consistently low likelihood, suggesting safer dialogues in

# Conclusion and Implications

- Red Teaming Language Models with Language Models can operate before adversaries.
- "Behavior on failing test cases may then be fixed preemptively."