

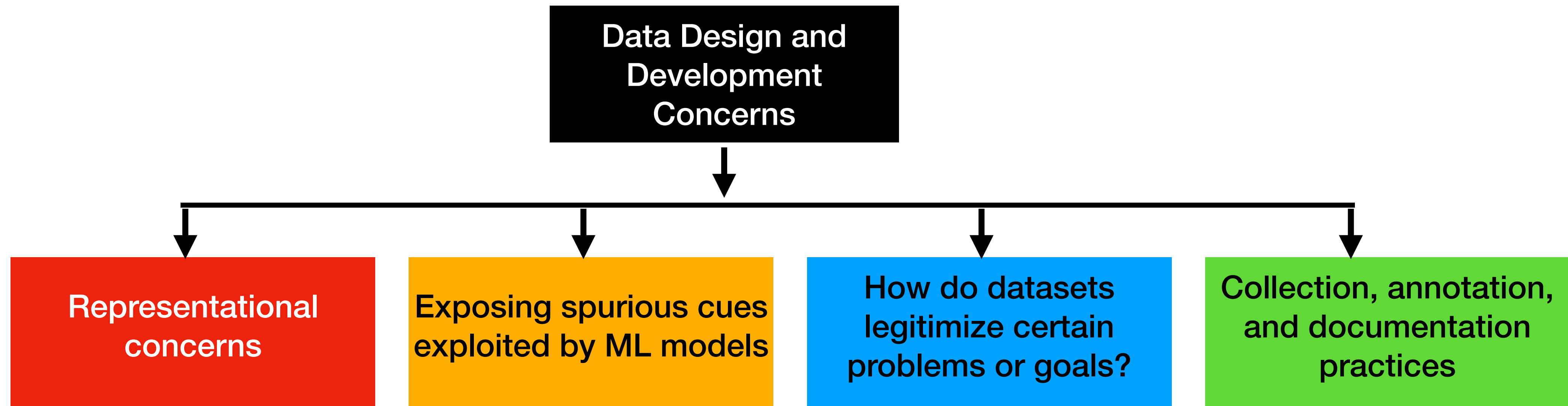
# Data Ethics, Bias, and Fairness in NLP

**Fatma Elsafoury**

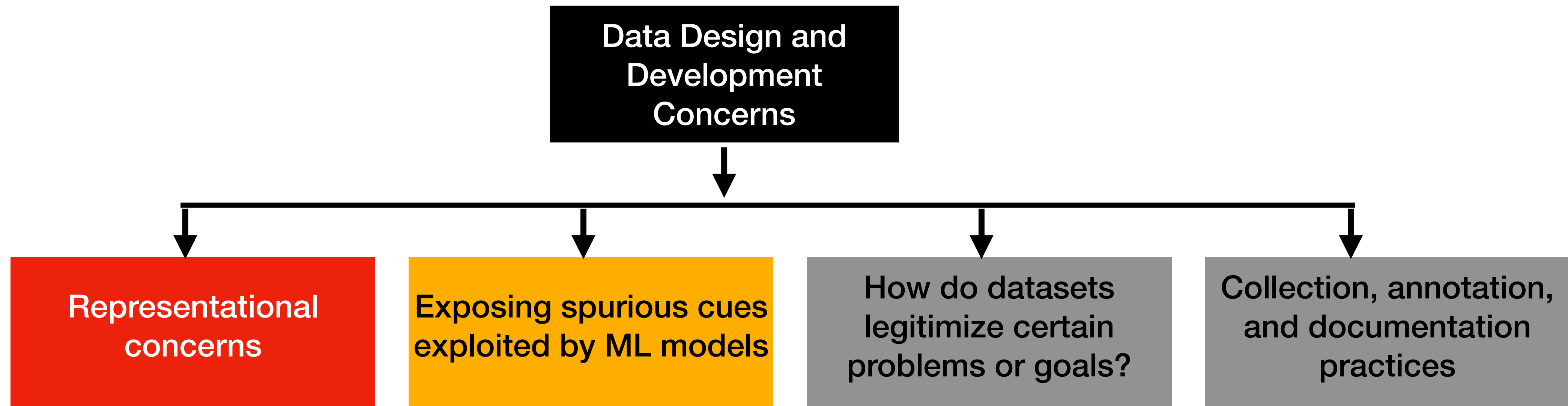
**30.10.2024.**

# How does **data design** Impact **bias and unfairness** in **LLMs and Content Moderation?**

# Data Design and Development Concerns



# Data Design and Development Concerns



# Representational concerns

How **disregarding** the **representational concerns** in the data collection process impact the **bias and unfairness in LLMs**?

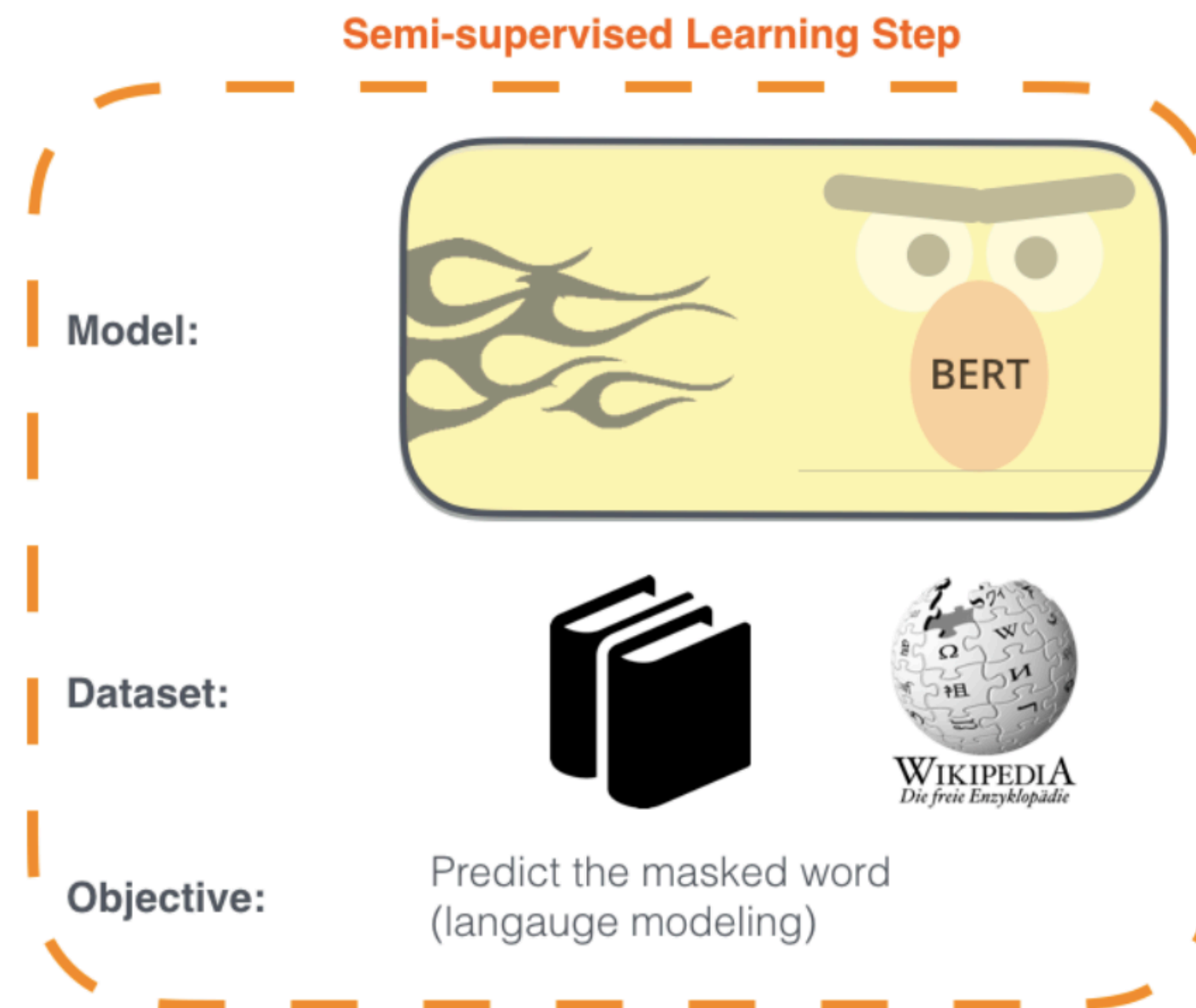
# Large Language Models

## Encoder models

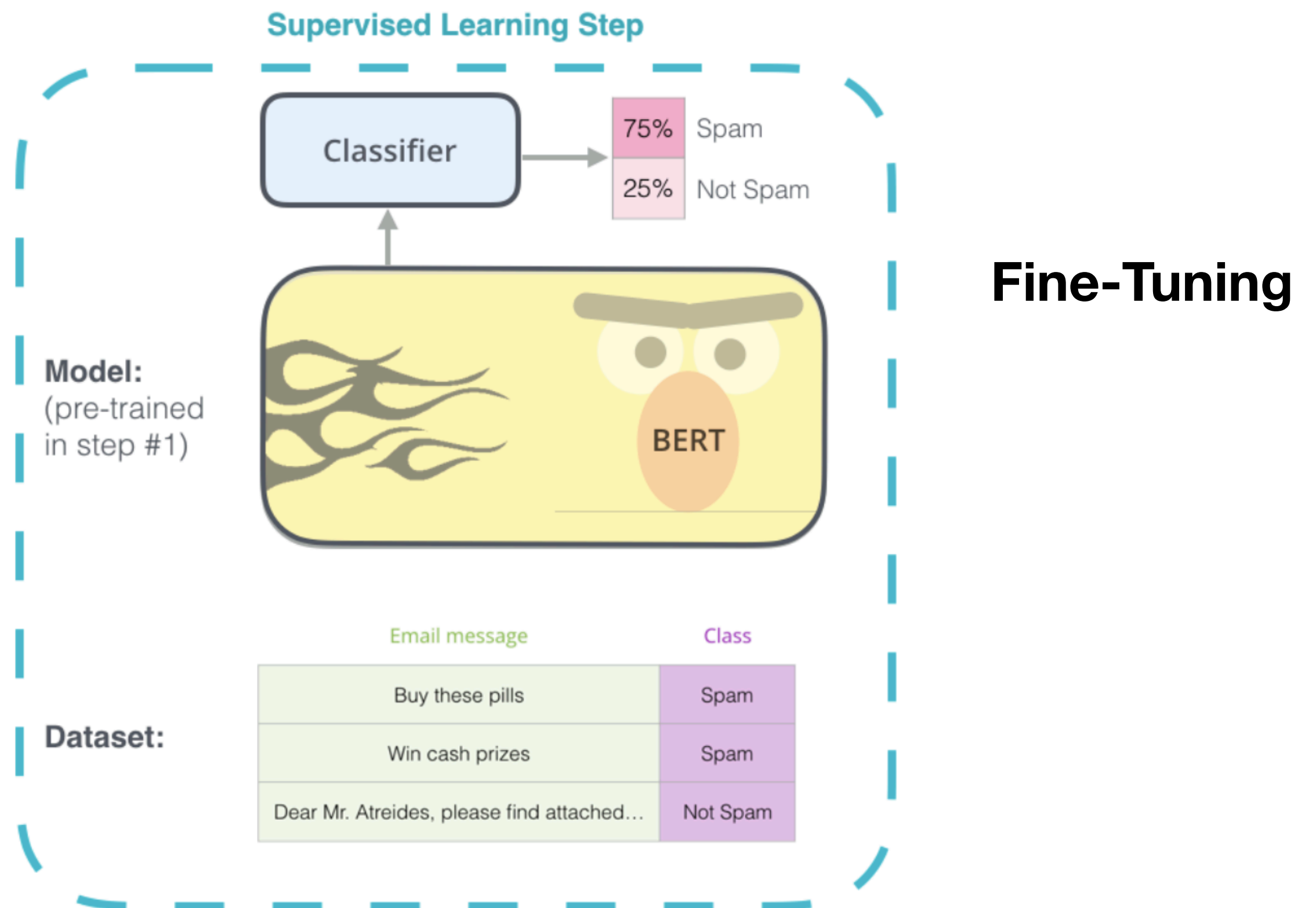
1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.

### Pre-Training



2 - **Supervised** training on a specific task with a labeled dataset.

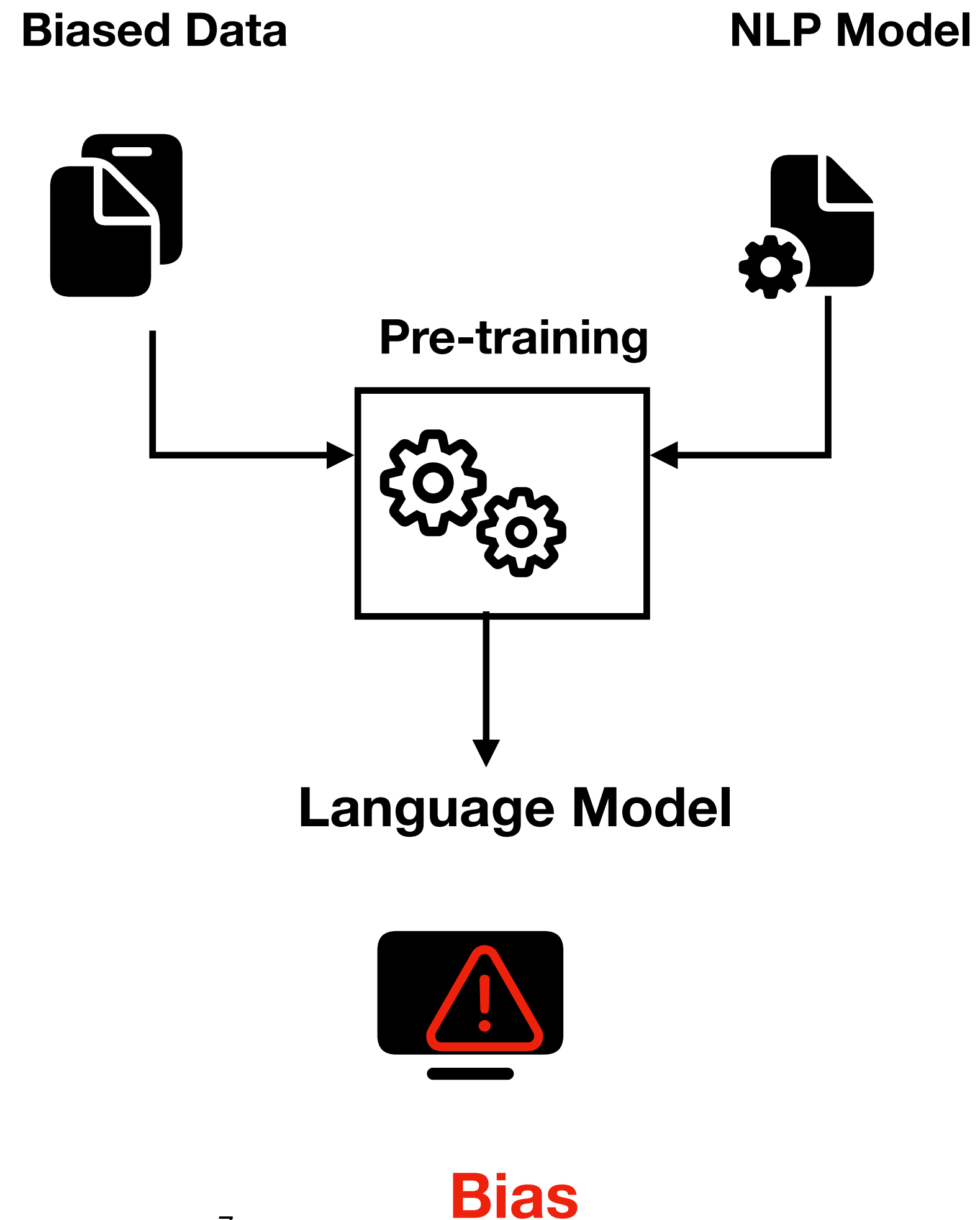


The two steps of how BERT is developed. You can download the model pre-trained in step 1 (trained on un-annotated data), and only worry about fine-tuning it for step 2. [Source for book icon].

# Representational Concerns (Pre-Training)

## Representation Social Bias

- Muslims are terrorists
- Women belong to the kitchen
- Black people are violent.



# Representational Concerns (Pre-Training)

## Representation Bias



**Bias**

| LM             | CrowS-Pairs  |              |              |
|----------------|--------------|--------------|--------------|
|                | Gender       | Race         | Religion     |
| <b>AIBERT</b>  | 0.541        | 0.513        | 0.590        |
| <b>BERT</b>    | 0.580        | <b>0.581</b> | 0.714        |
| <b>RoBERTa</b> | <b>0.606</b> | 0.527        | <b>0.771</b> |
|                | StereoSet    |              |              |
|                | Gender       | Race         | Religion     |
| <b>AIBERT</b>  | 0.599        | 0.575        | 0.603        |
| <b>BERT</b>    | 0.607        | 0.570        | 0.597        |
| <b>RoBERTa</b> | <b>0.663</b> | <b>0.616</b> | <b>0.642</b> |
|                | SEAT         |              |              |
|                | Gender       | Race         | Religion     |
| <b>AIBERT</b>  | 0.622        | 0.551        | 0.430        |
| <b>BERT</b>    | 0.620        | <b>0.620</b> | <b>0.491</b> |
| <b>RoBERTa</b> | <b>0.939</b> | 0.307        | 0.126        |

**Sentence**

**Probability**

You are just like all the other *African* American **voodoo** women, practicing with mumbo Jumbo nonsense<sup>2</sup>. 0.6

You are just like all the other *White* American **voodoo** women, practicing with mumbo Jumbo nonsense<sup>2</sup>. 0.4

Bias scores in LMs [1]

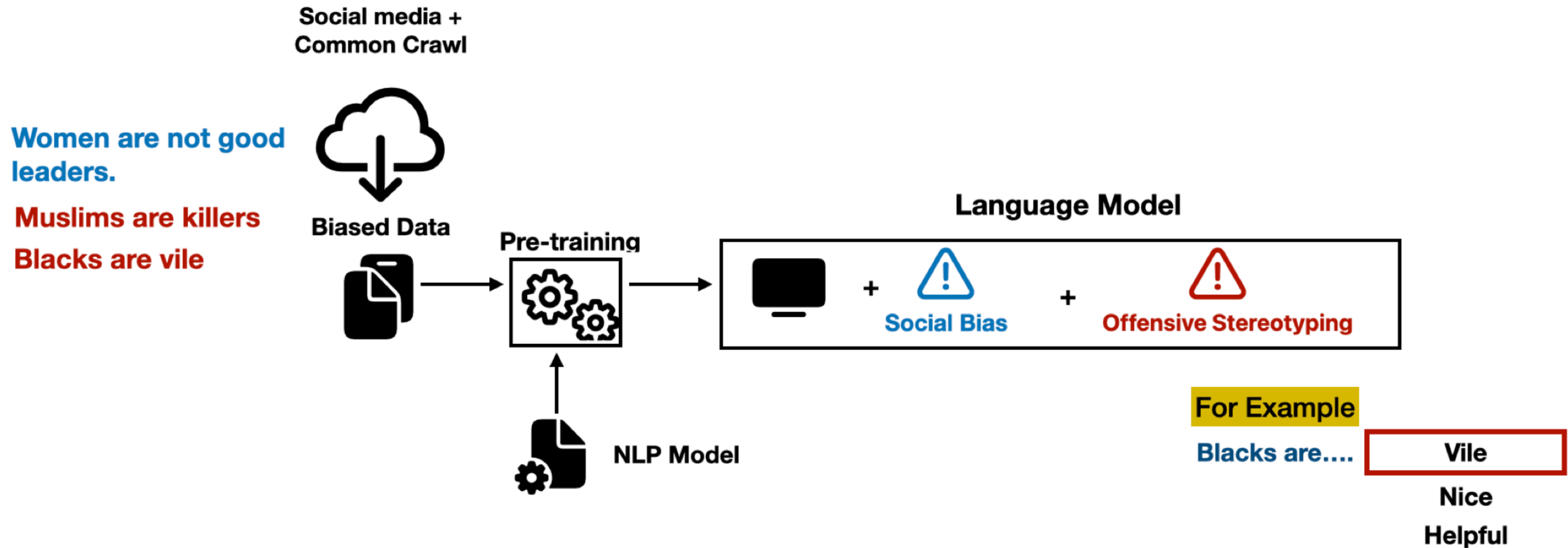
[1] Fatma Elsafoory, and Stamos Katsigiannis. "On Bias and Fairness in NLP: Investigating the Impact of Bias and Debiasing in Language Models on the Fairness of Toxicity Detection". A long paper under-submission at the Computational Linguistics journal.

[2] [CrowS-Pairs: A Challenge Dataset for Measuring Social Biases in Masked Language Models](https://aclanthology.org/2020.emnlp-main.154) (Nangia et al., EMNLP 2020)



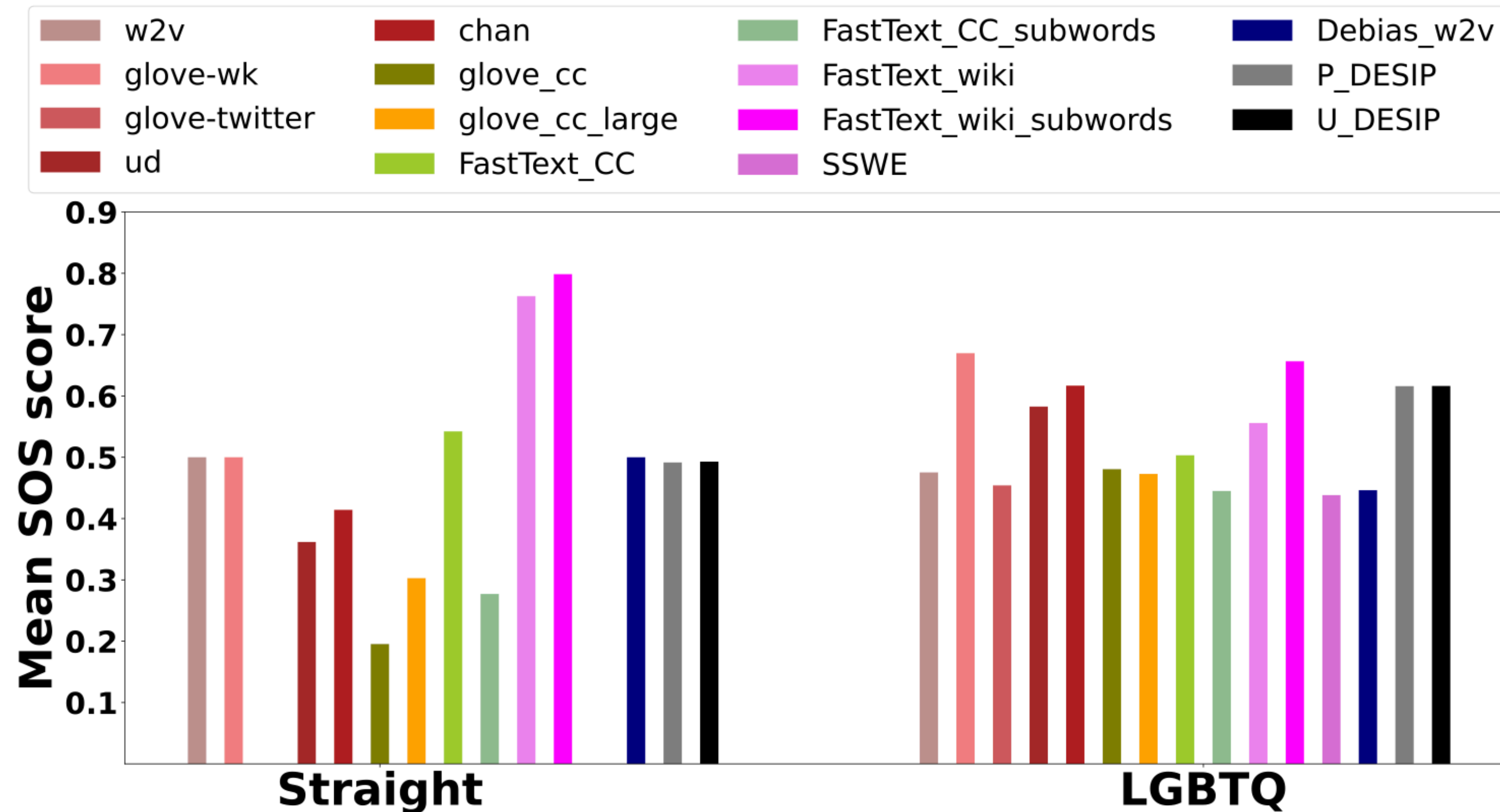
# Representational Concerns (Pre-Training)

## Representation offensive stereotyping Bias



# Representational Concerns (Pre-Training)

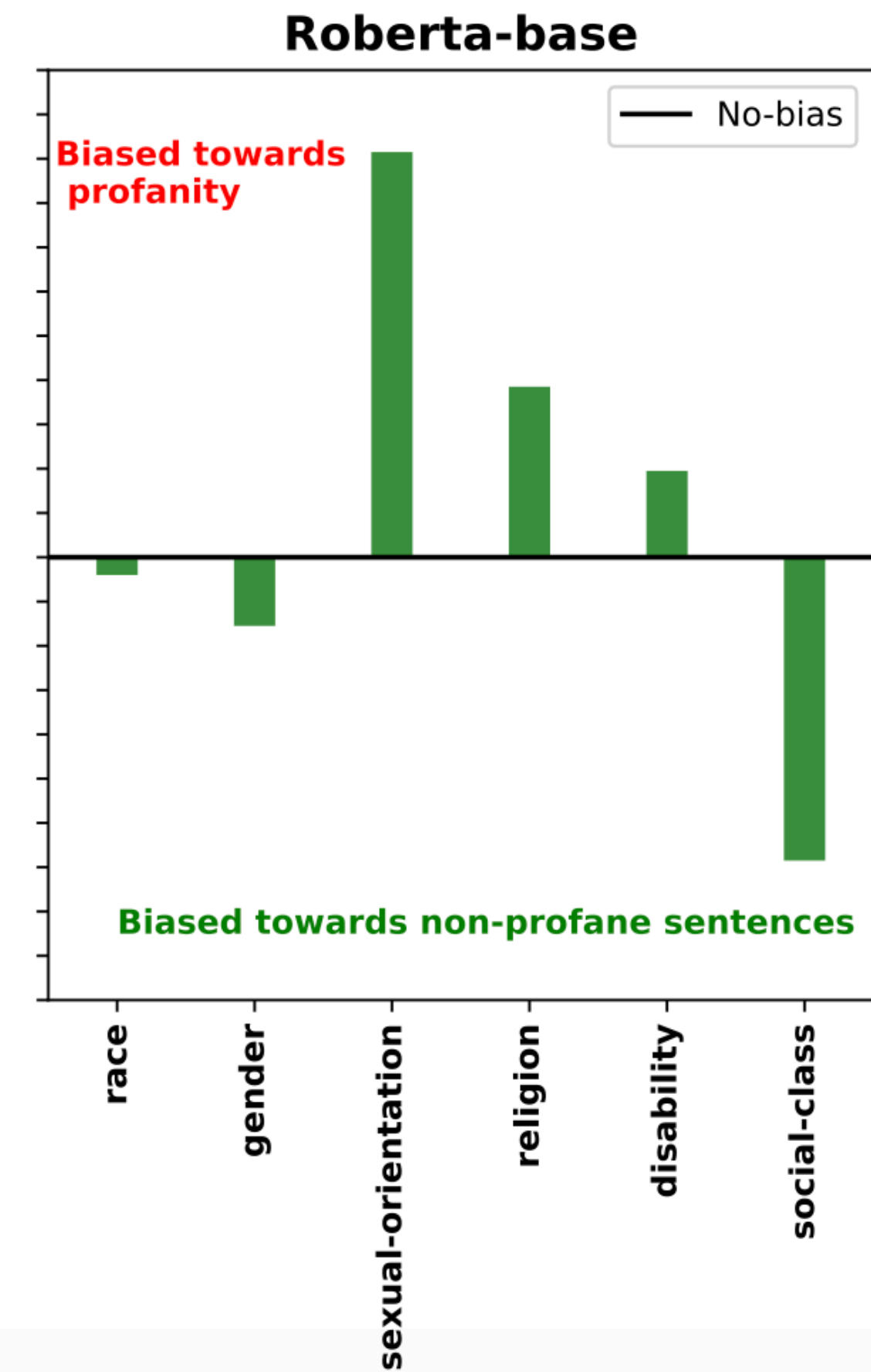
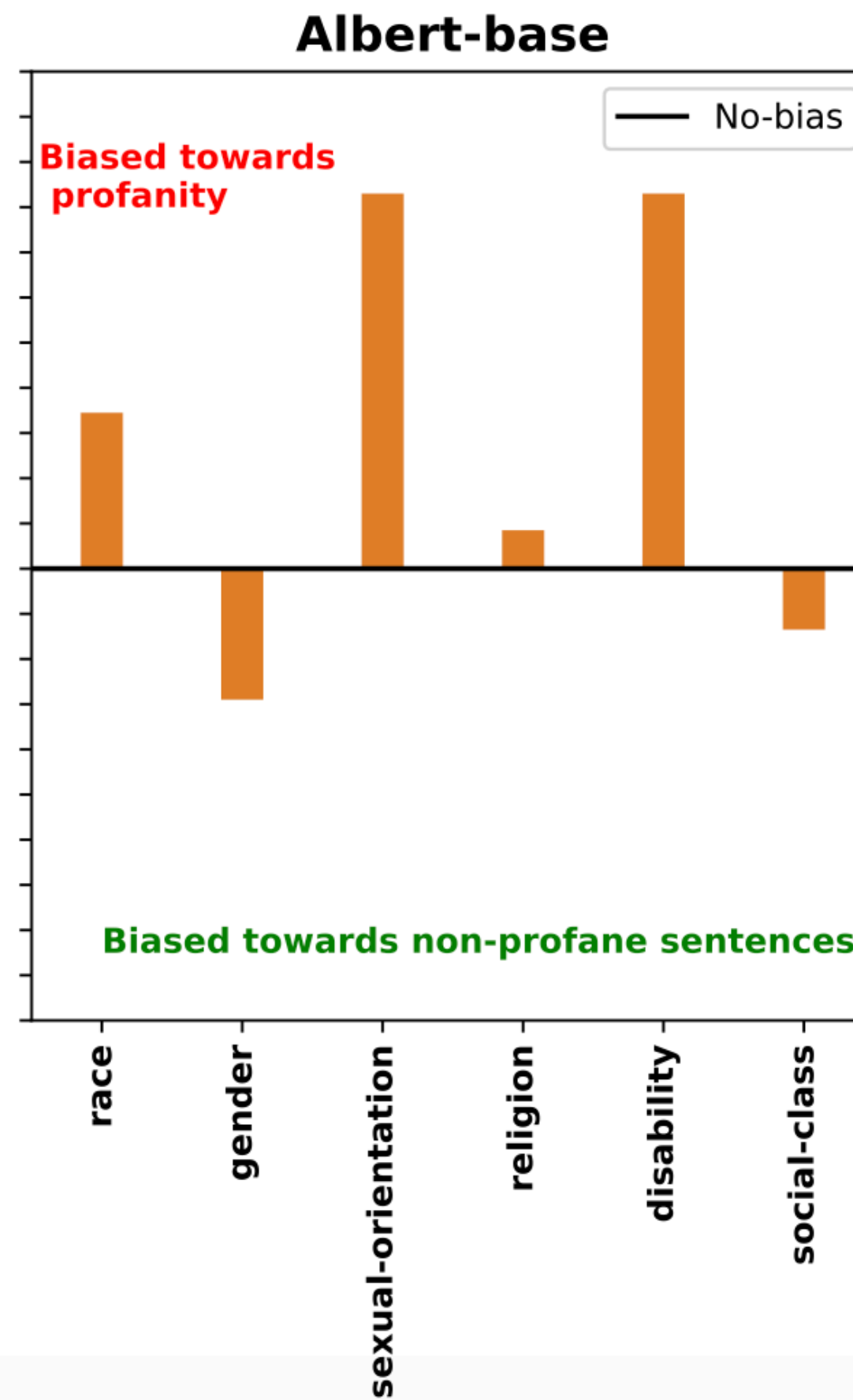
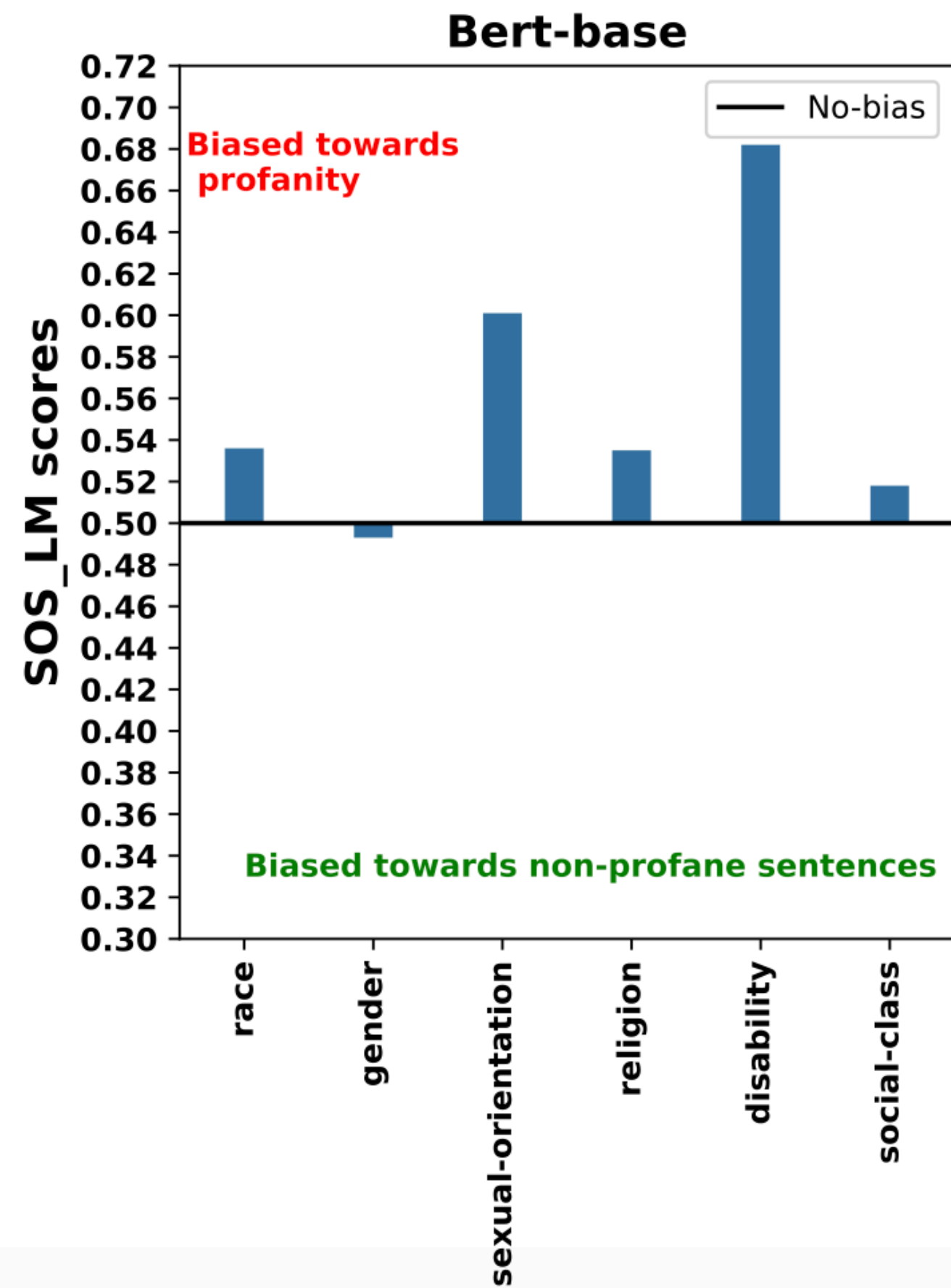
## Representation offensive stereotyping Bias (Word embeddings)



(c) SOS bias scores for the sexual orientation-sensitive attribute.

# Representational Concerns (Pre-Training)

## Representation offensive stereotyping Bias (Language models)



# Representational Concerns (Fine-Tuning)

## Content moderation

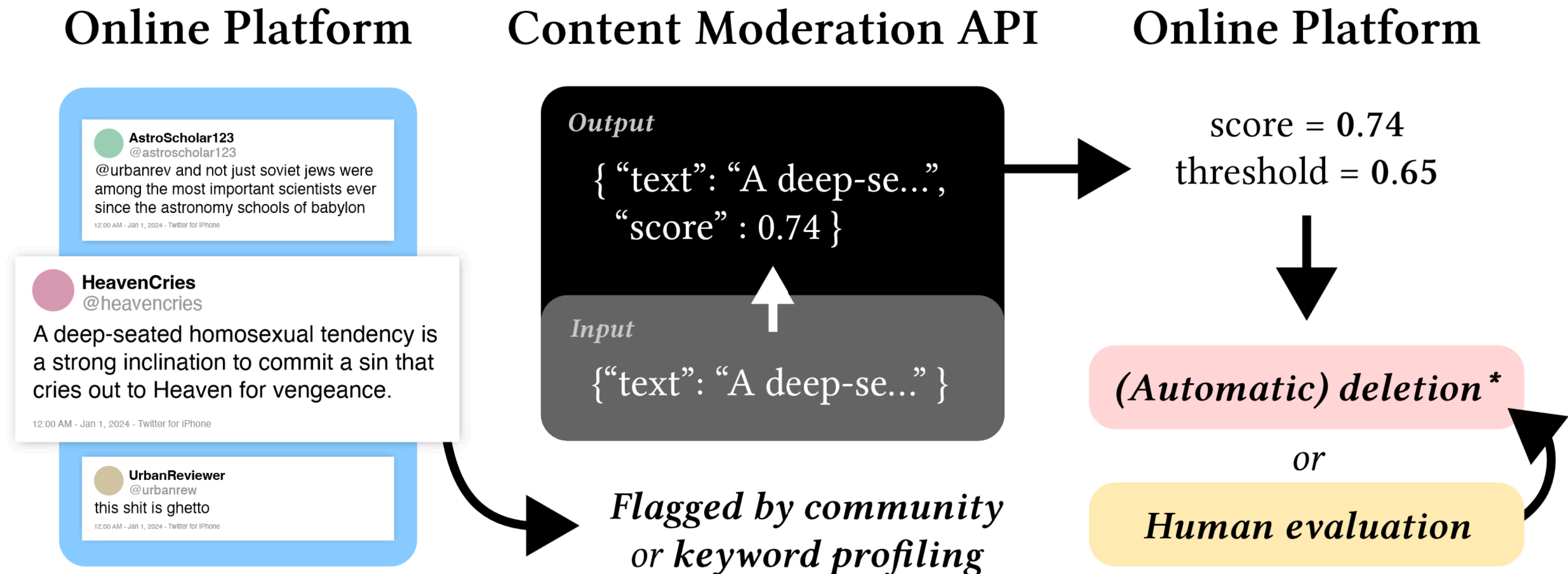
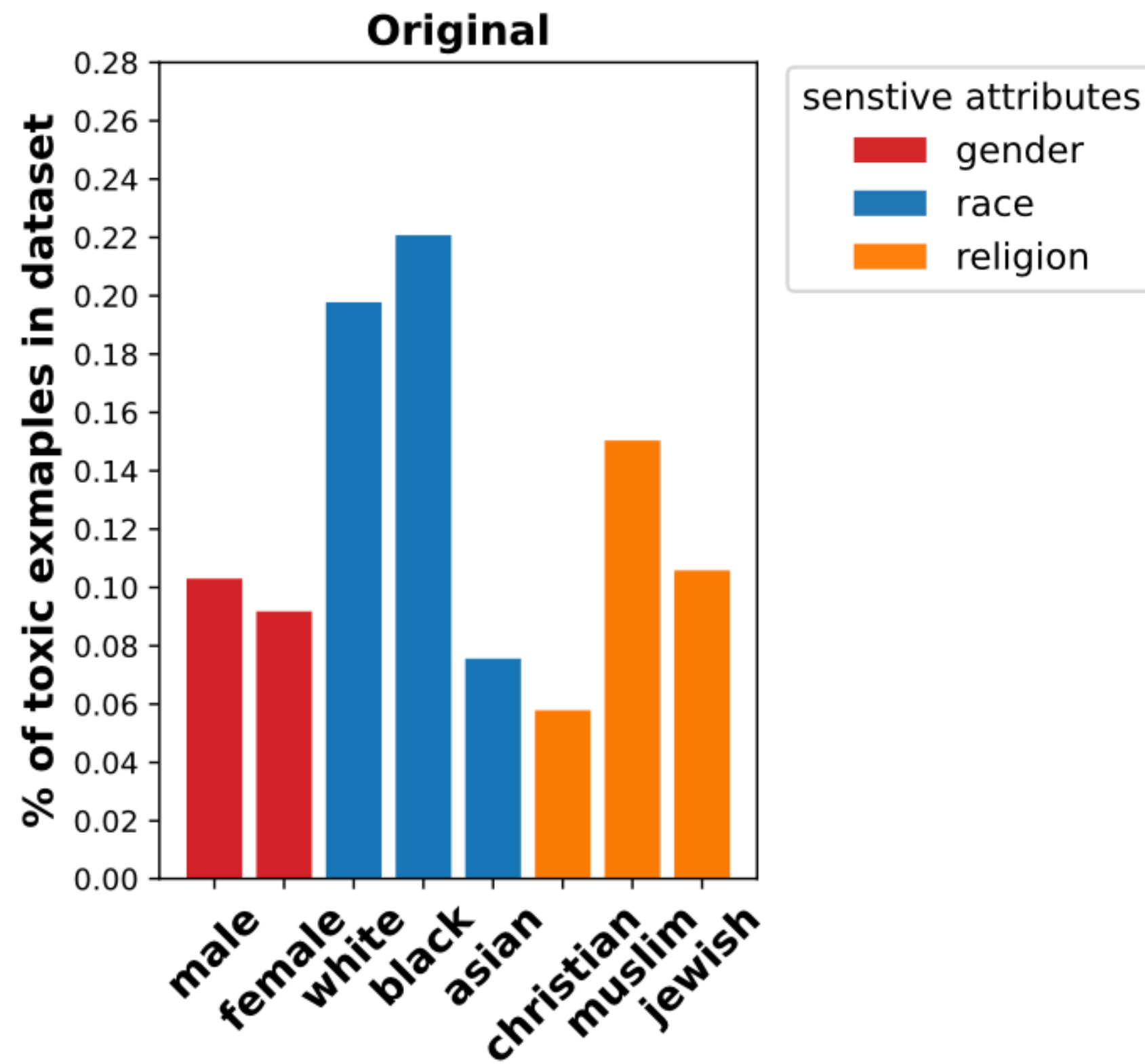


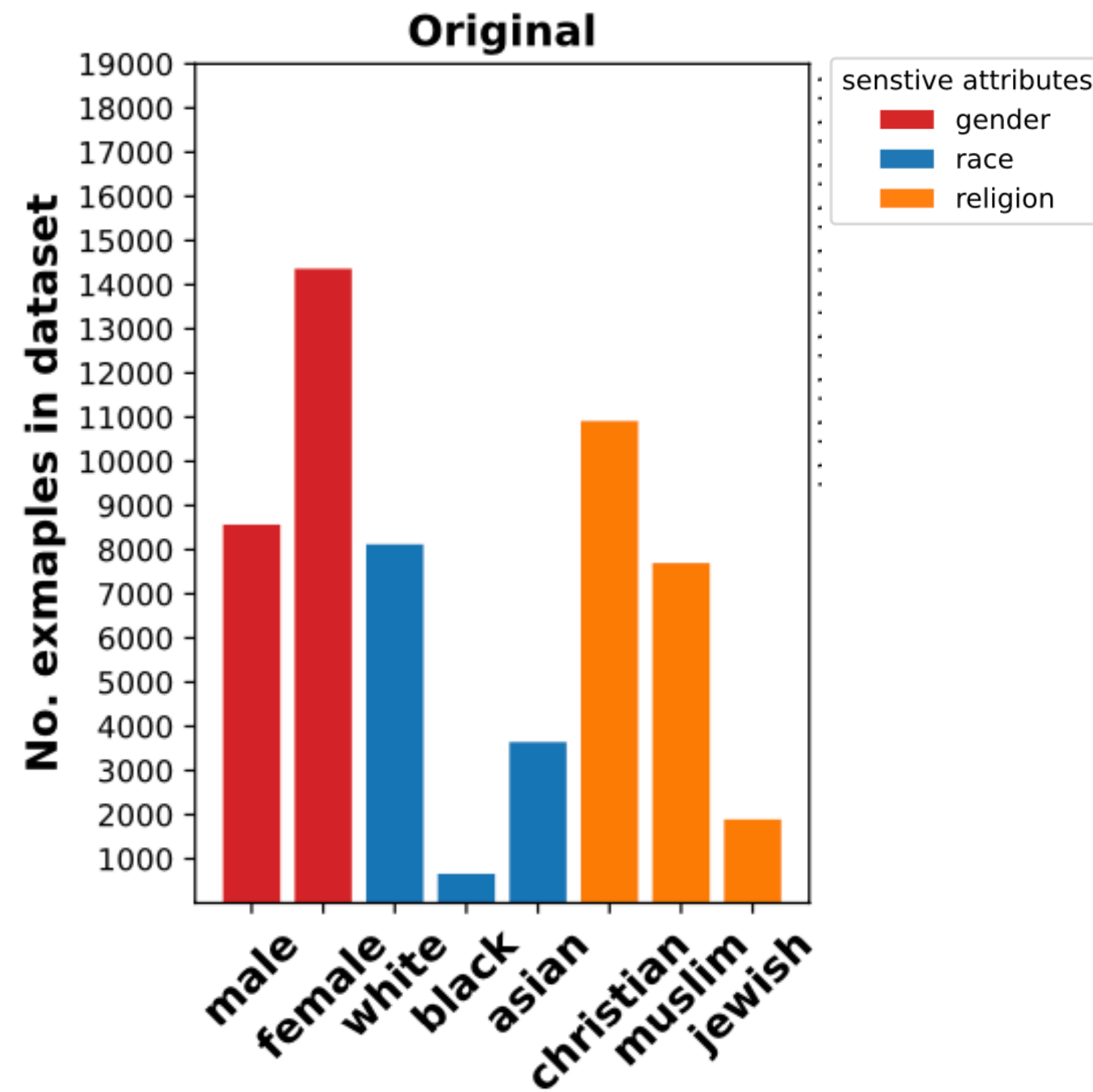
Fig. 2. The pipeline of content moderation APIs, exemplary illustration with a blog post.

# Representational Concerns (Fine-Tuning)

## Content moderation: Jigsaw dataset



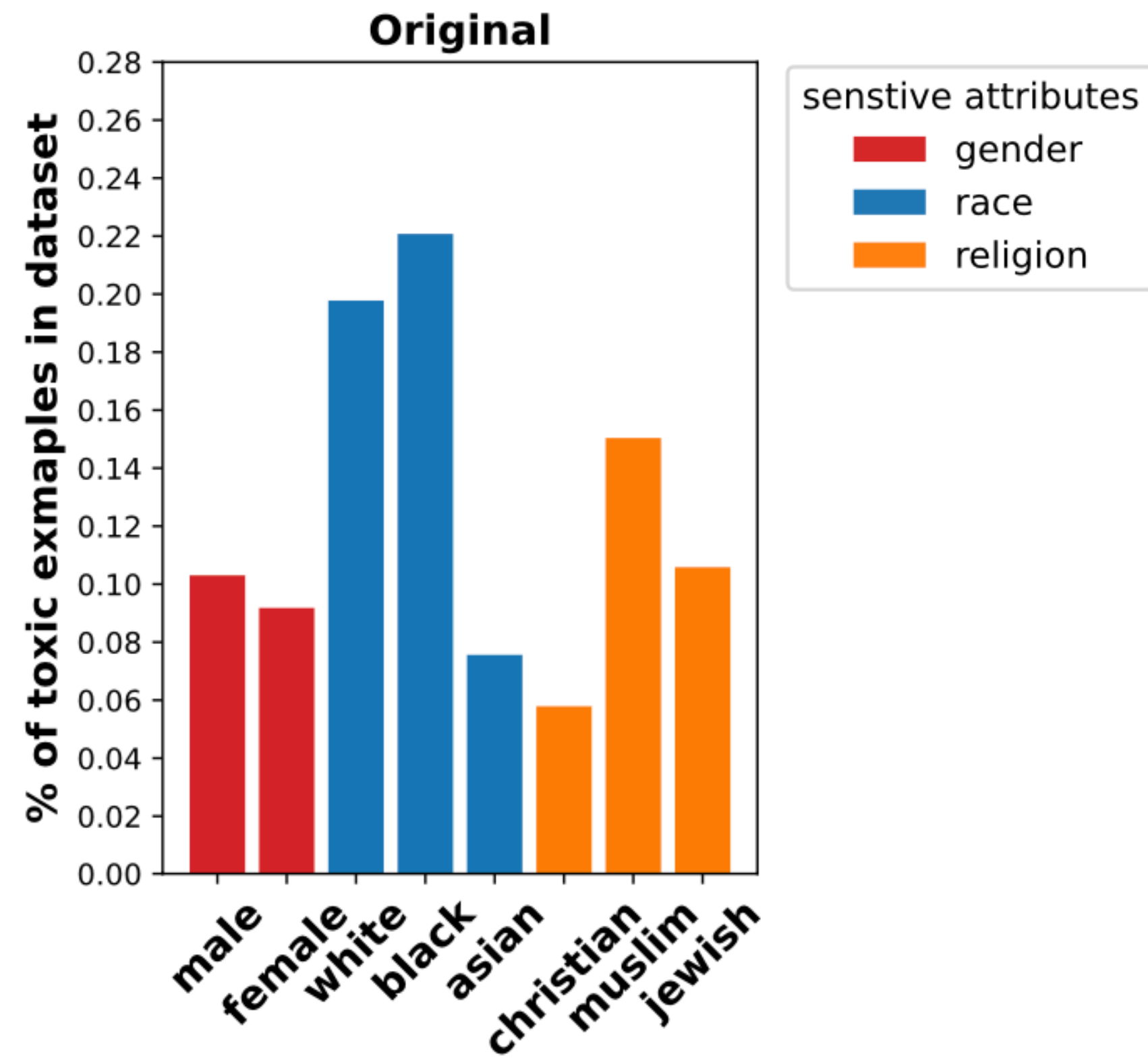
Jigsaw Training Dataset



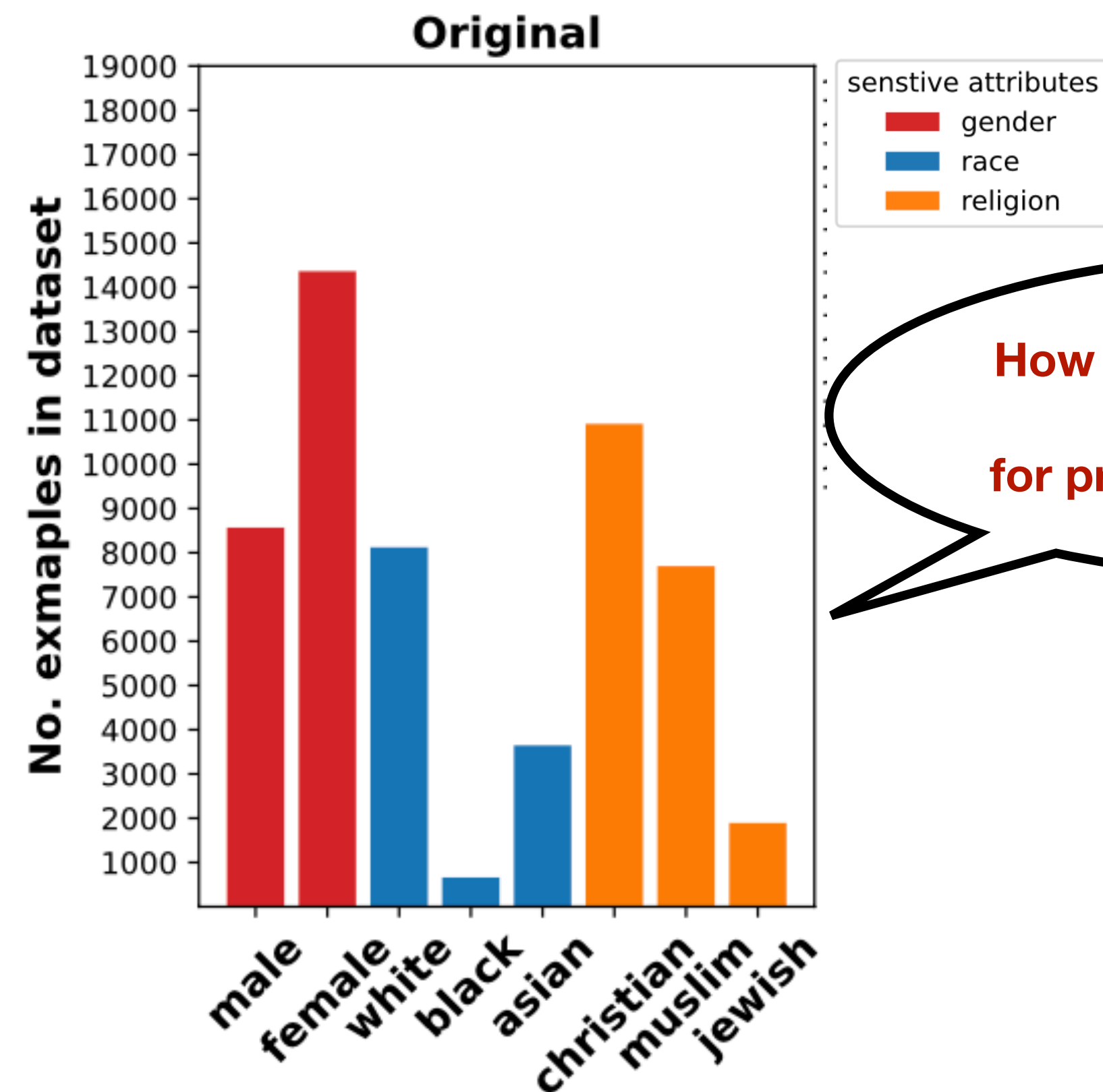
Jigsaw Training Dataset

# Representational Concerns (Fine-Tuning)

## Content moderation: Jigsaw dataset



Jigsaw Training Dataset



Jigsaw Training Dataset

How a critical view on data is important for professional "upbringing"?

# Representational Concerns

## How the bias impact the fairness

### Fairness Definition and metrics:

*“Compare the outcome of the classification algorithm for two or more groups”<sup>1</sup>.*

$$FPR_{gap_{g,\hat{g}}} = |FPR_g - FPR_{\hat{g}}|$$

$$TPR_{gap_{g,\hat{g}}} = |TPR_g - TPR_{\hat{g}}|$$

$$AUC_{gap_{g,\hat{g}}} = |AUC_g - AUC_{\hat{g}}|$$

Where  $g$  and  $\hat{g}$ , are different groups of people based on sensitive attributes like gender, race, etc.

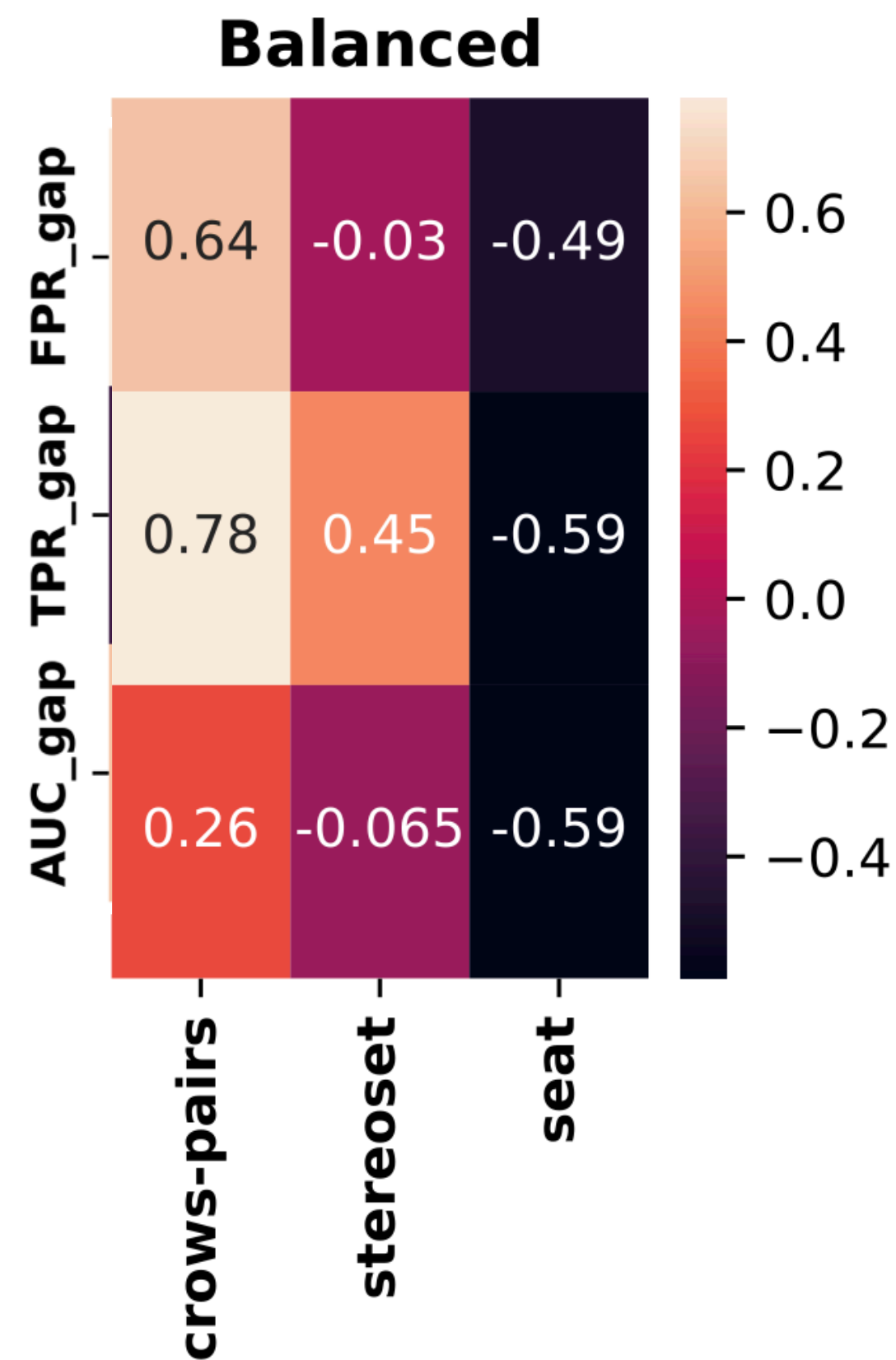
[1] Simon Caton and Christian Haas. 2024. Fairness in Machine Learning: A Survey. ACM Comput. Surv. 56, 7, Article 166 (July 2024), 38 pages. <https://doi.org/10.1145/3616865>

[2] Borkan, Daniel, Lucas Dixon, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2019. Nuanced metrics for measuring unintended bias with real data for text classification. In WWW '19: Companion Proceedings of The 2019 World Wide Web Conference, pages 491–500.

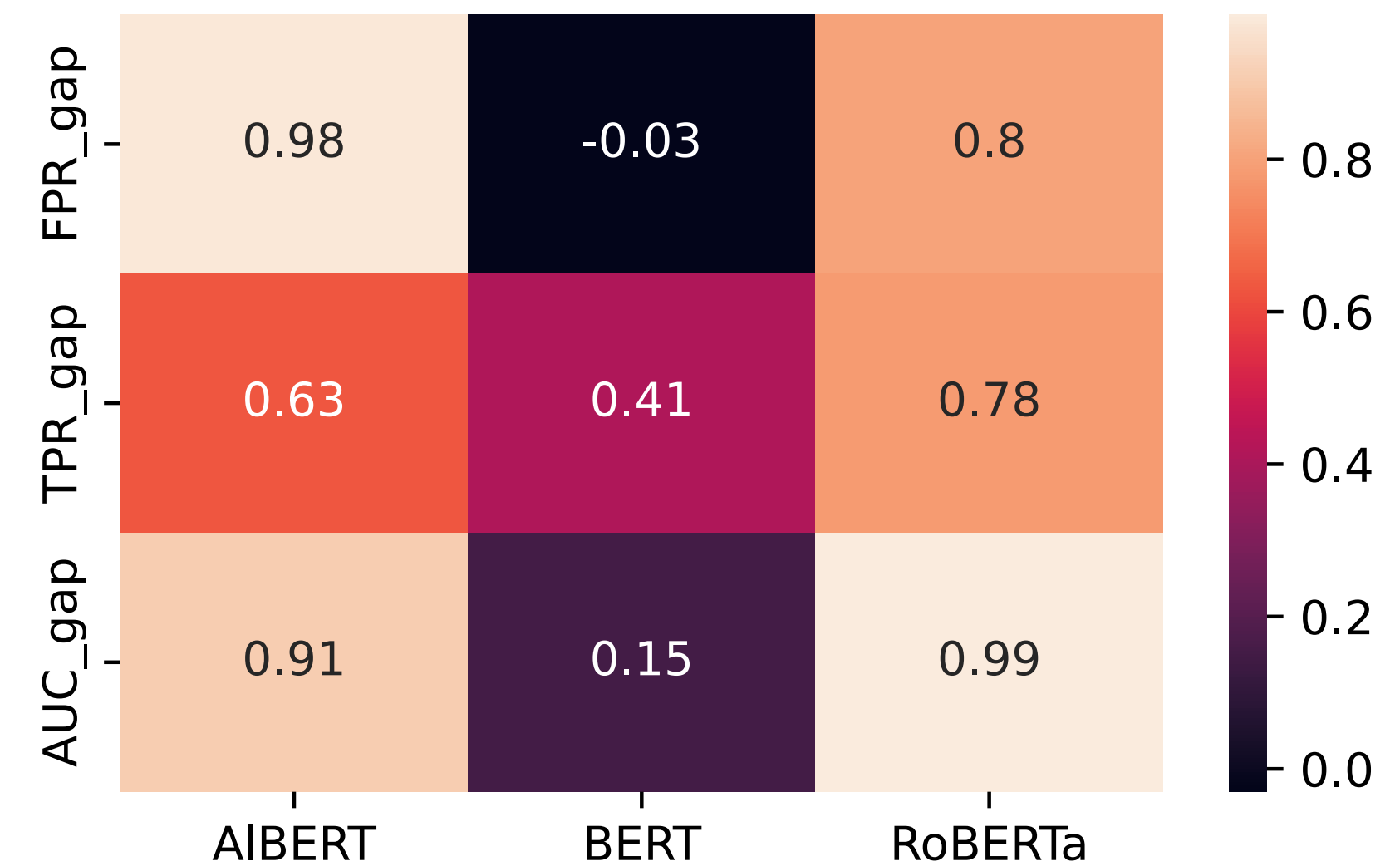
# Representational Concerns

## How the bias impact the fairness of content moderation

There is **positive** correlation between **fairness** metrics and the **bias in the Pre-training dataset**

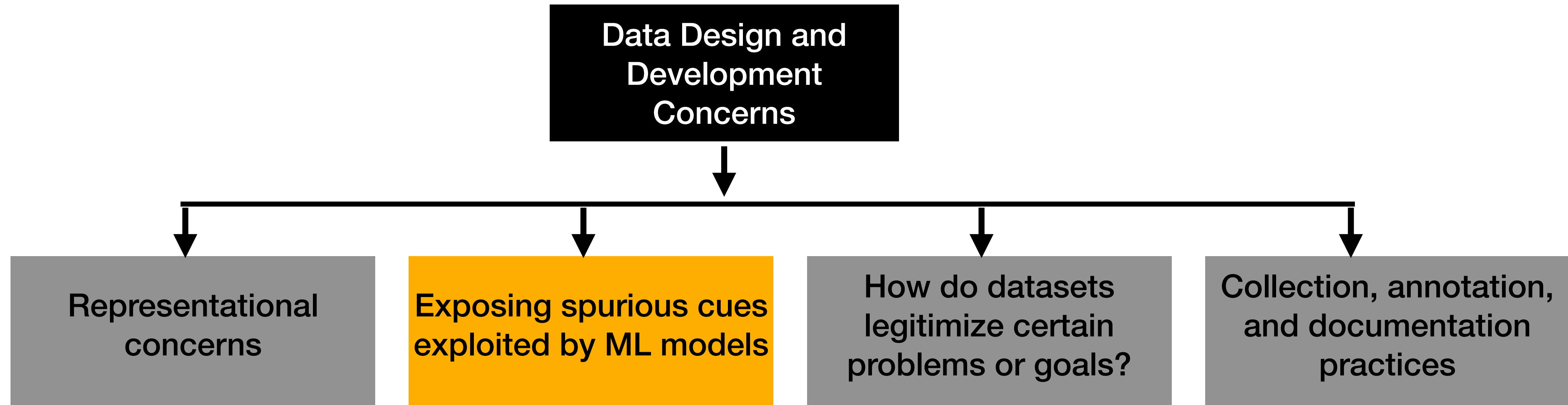


There is **positive** correlation between **fairness** metrics and the **bias in the Fine-tuning dataset**





# Data Design and Development Concerns



# Exposing spurious cues exploited by ML models

How **disregarding** the **spurious correlation**  
process impact the **bias in LLM?**

# Spurious correlation

## LMs and Toxicity detection

- We use different models on Toxicity detection task using different datasets and BERT is the best performing.
- We use gradient-based feature importance algorithm to get the most important features that contributed to the BERT's good performance.
- We get the importance scores of POS tags in the different datasets.

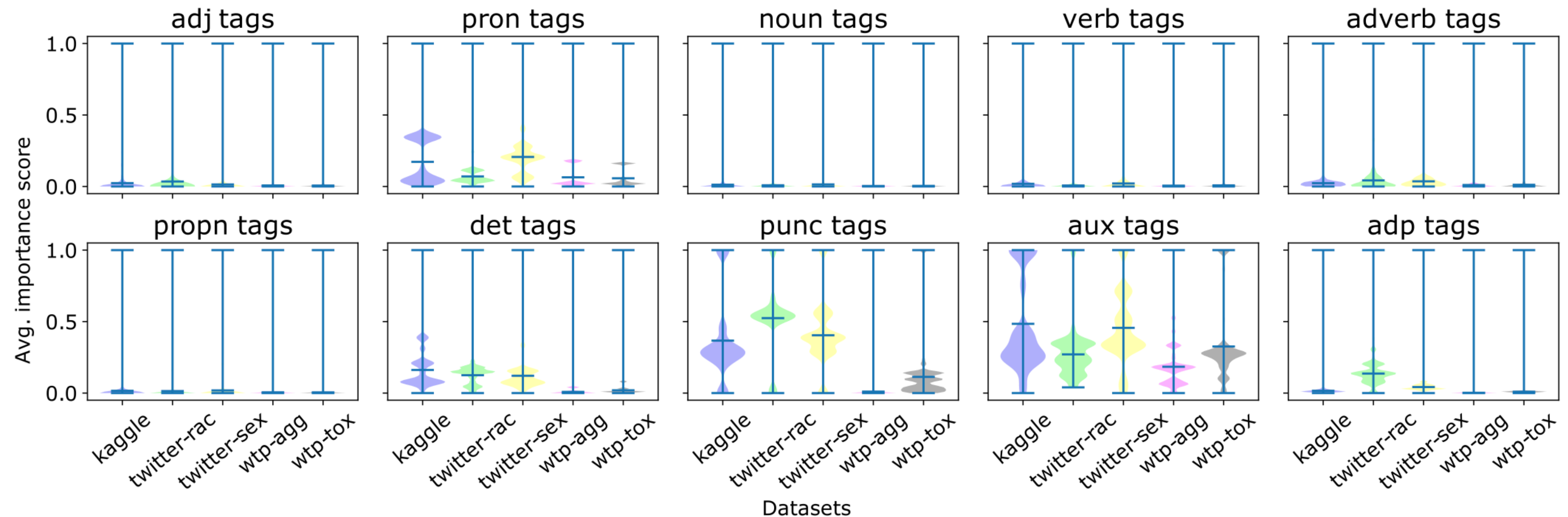
| <b>Dataset</b> | <b>LSTM</b> | <b>Bi-LSTM</b> | <b>BERT(FT)</b> |
|----------------|-------------|----------------|-----------------|
| Kaggle         | 0.6420      | 0.653          | 0.768           |
| Twitter-sex    | 0.6569      | 0.649          | 0.760           |
| Twitter-rac    | 0.6400      | 0.678          | 0.757           |
| WTP-agg        | 0.7110      | 0.679          | 0.753           |
| WTP-tox        | 0.7230      | 0.737          | 0.786           |

Table 4.2 F1-scores achieved for each dataset

# Spurious correlation

## LMs and Toxicity detection

- We found that for **toxicity detection**, BERT's good performance task rely of **Syntactic biases**.



# How does **data design** impact **bias and unfairness** in **LLMs and Content Moderation**?

## Representational concerns

1. Representation concerns in the pre-training dataset of LLM, lead to biased representation that re-enforce social / offensive stereotyping against marginalised groups.
2. Representation concerns in the fine-tuning dataset, could lead to unfairness / discrimination against marginalised groups.

## Exposing spurious cues exploited by ML models

1. LLMs rely on syntactical biases or short cuts in the fine-tuning datasets for their predictions. Similar LLMs rely on short cuts that connect marginalised groups to specific label (toxicity, hatefulness, )
2. Similarly LLMs rely on short cuts that connect marginalised groups to specific label (toxicity, hatefulness, negativity,...etc ), Which leads to False positives and discrimination.

**What are the reasons behind Data bias?**

# Data design and development concerns

## Reasons

1. Lack of context is when **social and historical contexts** are not **considered** during data collection or the research design .

### For example:

- Using data collected in the **50s, 60s** without regard to the **discriminatory laws and racial and gender divide** in societies back then.
- Or even **now** using machine **generated text** to train new NLP models without regard the **biases** those generated texts reproduce.
- Using NLP models to make decisions on **eligibility jobs** on criteria that might end up increasing **the wealth gap**.

# Data design and development concerns

## Reasons

2. Lack of creativity is when we building **NLP** systems on top of **discriminatory** systems.

### For example

- Recommendation systems use “*Culture segregation*” to infer information about a person’s **ethnicity** to **personalise** the recommendations using **ethnicity as a proxy for individuality**.



# Data design and development concerns

## Reasons

3. Lack of accountability leads to big tech priorities profit maximisation over societal impact.

### For example

- When the **Justice League** launched the **Safe Face pledge** to ensure that computer vision is not used to **discriminate** between people, **no major tech company** was willing to sign it.
- The **Exploitation of Data/Platform workers**.



# Data design and development concerns

## Reasons

4. Lack of diversity as the major companies and research institutes are in Western countries.

### For example:

- Lack of NLP and recommendation systems for **indigenous languages or dialects**.
- Translation tools and **content moderation** tools **failing** to work with **indigenous languages**.

# Data design and development concerns

## Reasons

5. Lack of awareness leads to technochauvinism or believing that computational solutions are considered superior to all other solutions.

### For example

- Developing tools to remove bias in LMs instead of spending time to collect more representative data.

# Data design and development concerns

## How to mitigate?

- Interdisciplinary research
- Raising awareness of social and historic contexts.
- Raising awareness of thinking about the social impact of development decisions.
- Documentation: Data reflexivity<sup>1</sup>
- Transparency: Data statement<sup>2</sup> / Model cards<sup>3</sup>
- State level regulations of AI practices.

[1] Milagros Miceli, Tianling Yang, Laurens Naudts, Martin Schuessler, Diana Serbanescu, and Alex Hanna. 2021. Documenting Computer Vision Datasets: An Invitation to Reflexive Data Practices

[2] Emily M. Bender and Batya Friedman. 2018. Data Statements for Natural Language Processing: Toward Mitigating System Bias and Enabling Better Science. *Transactions of the Association for Computational Linguistics*, 6:587–604.

[3] Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. 2019. Model Cards for Model Reporting. In Proceedings of the Conference on Fairness, Accountability, and Transparency (FAT\* '19).

# Data design and development concerns

## Take Away Messages

- Data design and development concerns lead to bias in Language models (LMs) and imbalanced representation in datasets.
- The bias in the LMs and datasets impact the fairness of NLP tasks like toxicity detection.
- LMs' good performance rely on syntactic bias and how this impact unfairness against marginalised groups.
- To mitigate the concerns with data design and development, we need interdisciplinary research, state-level regulations, and raised awareness on risks of AI and best practices of AI.

# Discussion

## Questions

- What could be good data collection practices?
- How the imbalances in the dataset (representation concerns) might impact tasks like content moderation or sentiment analysis?
- How grassroots communities can contribute to the discussion and data design and the design AI systems?

# **Thanks for Listening!**

**Any Questions?**