



## UvA-DARE (Digital Academic Repository)

### Stepmothers are mean and academics are pretentious: What do pretrained language models learn about you?'

Choenni, R.; Shutova, E.; van Rooij, R.

**DOI**

[10.18653/v1/2021.emnlp-main.111](https://doi.org/10.18653/v1/2021.emnlp-main.111)

**Publication date**

2021

**Document Version**

Final published version

**Published in**

2021 Conference on Empirical Methods in Natural Language Processing

**License**

CC BY

[Link to publication](#)

**Citation for published version (APA):**

Choenni, R., Shutova, E., & van Rooij, R. (2021). Stepmothers are mean and academics are pretentious: What do pretrained language models learn about you?'. In M-C. Moens, X. Huang, L. Specia, & S. W. Yih (Eds.), *2021 Conference on Empirical Methods in Natural Language Processing: EMNLP 2021 : proceedings of the conference : November 7-11, 2021* (pp. 1477-1491). The Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.emnlp-main.111>

**General rights**

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

**Disclaimer/Complaints regulations**

If you believe that digital publication of certain material infringes any of your rights or (privacy) interests, please let the Library know, stating your reasons. In case of a legitimate complaint, the Library will make the material inaccessible and/or remove it from the website. Please Ask the Library: <https://uba.uva.nl/en/contact>, or a letter to: Library of the University of Amsterdam, Secretariat, Singel 425, 1012 WP Amsterdam, The Netherlands. You will be contacted as soon as possible.

*UvA-DARE is a service provided by the library of the University of Amsterdam (<https://dare.uva.nl>)*

# Stepmothers are mean and academics are pretentious: What do pretrained language models learn about you?

**Rochelle Choenni**

University of Amsterdam

r.m.v.k.choenni@uva.nl

**Ekaterina Shutova**

University of Amsterdam

e.shutova@uva.nl

**Robert van Rooij**

University of Amsterdam

r.a.m.vanrooij@uva.nl

## Abstract

*Warning: this paper contains content that may be offensive or upsetting.*

In this paper, we investigate what types of stereotypical information are captured by pretrained language models. We present the first dataset comprising stereotypical attributes of a range of social groups and propose a method to elicit stereotypes encoded by pretrained language models in an unsupervised fashion. Moreover, we link the emergent stereotypes to their manifestation as basic emotions as a means to study their emotional effects in a more generalized manner. To demonstrate how our methods can be used to analyze emotion and stereotype shifts due to linguistic experience, we use fine-tuning on news sources as a case study. Our experiments expose how attitudes towards different social groups vary across models and how quickly emotions and stereotypes can shift at the fine-tuning stage.

## 1 Introduction

Pretraining strategies for large-scale language models (LMs) require unsupervised training on large amounts of human generated text data. While highly successful, these methods come at the cost of interpretability as it has become increasingly unclear what relationships they capture. Yet, as their presence in society increases, so does the importance of recognising the role they play in perpetuating social biases. In this regard, [Bolukbasi et al. \(2016\)](#) first discovered that contextualized word representations reflect gender biases captured in the training data. What followed was a suite of studies that aimed to quantify and mitigate the effect of harmful social biases in word ([Caliskan et al., 2017](#)) and sentence encoders ([May et al., 2019](#)). Despite these studies, it has remained difficult to define what constitutes “bias”, with most work focusing on “gender bias” ([Manela et al., 2021](#); [Sun et al., 2019](#)) or “racial bias” ([Davidson et al., 2019](#);

[Sap et al., 2019](#)). More broadly, biases in the models can comprise a wide range of harmful behaviors that may affect different social groups for various reasons ([Blodgett et al., 2020](#)).

In this work, we take a different focus and study stereotypes that emerge within pretrained LMs instead. While bias is a personal preference that can be harmful when the tendency interferes with the ability to be impartial, stereotypes can be defined as a preconceived idea that (incorrectly) attributes general characteristics to all members of a group. While the two concepts are closely related i.e., stereotypes can evoke new biases or reinforce existing ones, stereotypical thinking appears to be a crucial part of human cognition that often emerges implicitly ([Hinton, 2017](#)). [Hinton \(2017\)](#) argued that implicit stereotypical associations are established through Bayesian principles, where the experience of their prevalence in the world of the perceiver causes the association. Thus, as stereotypical associations are not solely reflections of cognitive bias but also stem from real data, we suspect that our models, like human individuals, pick up on these associations. This is particularly true given that their knowledge is largely considered to be a reflection of the data they are trained on. Yet, while we consider stereotypical thinking to be a natural side-effect of learning, it is still important to be aware of the stereotypes that models encode. Psychology studies show that beliefs about social groups are transmitted and shaped through language ([Maass, 1999](#); [Beukeboom and Burgers, 2019](#)). Thus, specific lexical choices in downstream applications not only reflect the model’s attitude towards groups but may also influence the audience’s reaction to it, thereby inadvertently propagating the stereotypes they capture ([Park et al., 2020](#)).

Studies focused on measuring stereotypes in pretrained models have thus far taken supervised approaches, relying on human knowledge of common stereotypes about (a smaller set of) social groups

(Nadeem et al., 2020; Nangia et al., 2020). This, however, bears a few disadvantages: (1) due to the implicit nature of stereotypes, human defined examples can only expose a subset of popular stereotypes, but will omit those that human annotators are unaware of (e.g. models might encode stereotypes that are not as prevalent in the real world); (2) stereotypes vary considerably across cultures (Dong et al., 2019), meaning that the stereotypes tested for will heavily depend on the annotator’s cultural frame of reference; (3) stereotypes constantly evolve, making supervised methods difficult to maintain in practice. Therefore, similar to Field and Tsvetkov (2020), we advocate the need for implicit approaches to expose and quantify bias and stereotypes in pretrained models.

We present the first dataset of stereotypical attributes of a wide range of social groups, comprising  $\sim 2K$  attributes in total. Furthermore, we propose a stereotype elicitation method that enables the retrieval of salient attributes of social groups encoded by state-of-the-art LMs in an unsupervised manner. We use this method to test the extent to which models encode the human stereotypes captured in our dataset. Moreover, we are the first to demonstrate how training data at the fine-tuning stage can directly affect stereotypical associations within the models. In addition, we propose a complementary method to study stereotypes in a more generalized way through the use of emotion profiles, and systematically compare the emerging emotion profiles for different social groups across models. We find that all models vary considerably in the information they encode, with some models being overall more negatively biased while others are mostly positive instead. Yet, in contrast to previous work, this study is not meant to advocate the need for debiasing. Instead, it is meant to expose varying implicit stereotypes that different models incorporate and to bring awareness to how quickly attitudes towards groups change based on contextual differences in the training data used both at the pretraining and fine-tuning stage.

## 2 Related work

**Previous work on stereotypes** While studies that explicitly focus on stereotypes have remained limited in NLP, several works on bias touch upon this topic (Blodgett et al., 2020). This includes, for instance, studying specific phenomena such as the infamous ‘Angry Black Woman’ stereotype and the ‘double bind’ (Heilman et al., 2004) theory (Kir-

itchenko and Mohammad, 2018; May et al., 2019; Tan and Celis, 2019), or relating model predictions to gender stereotype lexicons (Field and Tsvetkov, 2020). To the best of our knowledge, Nadeem et al. (2020); Nangia et al. (2020) and Manela et al. (2021) are the first to explicitly study stereotypes in pretrained sentence encoders. While Manela et al. (2021) focus on gender stereotypes using the Wino-Bias dataset (Zhao et al., 2018), the other works propose new crowdsourced datasets (i.e. StereoSet and Crowspair) with stereotypes that cover a wide range of social groups. All datasets, however, have a similar set-up: they contain pairs of sentences of which one is more stereotypical than the other. Working in the language modeling framework, they evaluated whether the model “prefers” the stereotypical sentence over the anti-stereotypical one. In contrast, we propose a different experimental setup and introduce a new dataset that leverages search engines’ autocomplete suggestions for the acquisition of explicit stereotypical attributes. Instead of indirectly uncovering stereotypes through comparison, our elicitation method directly retrieves salient attributes encoded in the models. Our technique is inspired by Kurita et al. (2019), but while they measure the LM probability for completing sentences with the pronouns *she* and *he* specifically, we study the top  $k$  salient attributes without posing any restrictions on what these could be. Moreover, we are the first to include both monolingual and multilingual models in our analysis.

**Stereotype-driven emotions** Stereotypes are constantly changing and identifying negative ones in particular, is an inherently normative process. While some stereotypes clearly imply disrespect (e.g., women are incompetent), others emerge from excessive competence instead (e.g., Asians are good at math). Moreover, stereotypical content is heavily influenced by the social pressures of society at the time. Cuddy et al. (2009) argue that no stereotype remains stable and predictable from theoretical principles. Hence, many social psychologists have abandoned the study of stereotype content to focus on systematic principles that generalize across different specific instances of stereotypes instead, presumably making them more stable over time and place (Cuddy et al., 2009; Mackie et al., 2000; Weiner, 1993). Similarly, we explore a more robust approach to uncovering stereotypes in pretrained LMs by studying how stereotypes are more generally manifested as varying emotion profiles in

Search query	Category	Social group	Stereotypical attributes	Search engine
Why are black people so	race	black people	fast, athletic, hated, angry, loud	Yahoo
Why are british people so	race	british people	polite, pale, tall, reserved, cold	DuckDuckGo
Why are californians so	lifestyle	californians	entitled, rich, flaky, backstabby	multiple
Why are comedians so	profession	comedians	funny, sad, intelligent, depressed	Google
Why are millennials so	age	millennials	fragile, nostalgic, lonely, broke	Google
Why is Norway so	country	Norway	healthy, wealthy, happy, rich	multiple
Why are conservatives so	political	conservatives	angry, controlling, racist	Yahoo
Why are mormons so	religion	mormons	misunderstood, rich, succesful, nice	multiple

Table 1: A few (random) examples of stereotypical attributes sampled from the dataset. Note that the number of attributes vary as we do not retrieve the same number for each social group per search engine.

the models. Previous works show that groups evoke different emotional profiles (Cottrell and Neuberg, 2005; Tapias et al., 2007; Mackie et al., 2000), and a variety of theories link particular intergroup relations to distinct stereotype-driven emotions such as disgust and anger (Harris and Fiske, 2006, 2009).

### 3 Stereotypes from search engines

Retrieving human stereotypes in an implicit manner can be useful as people are likely to give more politically correct answers when asked for stereotypes explicitly. Questions we ask to search engines are often done in the comfort of our own homes, making them likely to reflect true stereotypes that are out there in the real world (Stephens-Davidowitz, 2018). When feeding search engines queries that already imply a stereotype about a specific social group, the suggested completions are often salient (i.e. stereotypical) attributes of that particular group e.g., search query: ‘Why are Asian parents so’ → completion: ‘strict’. As the autocomplete suggestions are based on popular human generated queries, we exploit this property for the acquisition of stereotypical attributes.

**Social groups** We use 382 social groups which we divide into 9 categories: age (11), gender (26), profession (115), race (99), country (73), religion (10), political (9), sexuality (12) and lifestyle (19). See Appendix C for the full lists of social groups per category. The social groups include, but are not limited to, the ones used by Nadeem et al. (2020) who collected Wikidata relation triples to automatically retrieve target terms.

**Templates** To increase coverage, we design 5 slightly varying templates to retrieve autocomplete suggestions for. For groups pertaining to people we use the templates:

- ‘Why are [TGT] so [ATTR]?’
- ‘What makes [TGT] so [ATTR]?’
- ‘How come [TGT] are so [ATTR]?’

- ‘Why are [TGT] always so [ATTR]?’
- ‘Why are all [TGT] so [ATTR]?’

For countries we use:

- ‘Why is [TGT] so [ATTR]?’
- ‘What makes [TGT] so [ATTR]?’
- ‘How come [TGT] is so [ATTR]?’
- ‘Why is [TGT] always so [ATTR]?’
- ‘Why are all people in [TGT] so [ATTR]?’

where [TGT] are social groups for which we search stereotypes and [ATTR] is the salient attribute with which the search engine completes the sequence. We tested other (longer and more elaborate) templates but we found that they did not produce many autocomplete suggestions. In fact, we believe that the above queries are so successful precisely because of their simplicity, given that people are likely to keep search queries concise.

**Search engines** Due to Google’s hate speech filtering system the autocomplete feature is disabled for frequently targeted groups e.g. black people, Jewish people and members of the LGBTQ+ community. Thus, we retrieve autocomplete suggestions from 3 search engines: Google, Yahoo and DuckDuckGo. In many cases, identical completions were given by multiple search engines. We sort these duplicate samples under the category ‘multiple engines’. We find that most negative (offensive) stereotypes are retrieved from Yahoo.

**Pre-processing** We clean up the dataset manually, using the following procedure:

1. Remove noisy completions that do not result in a grammatically correct sentence e.g. non adjectives.
2. Remove specific trend-sensitive references: e.g. to video games ‘why are asians so good at *league of legends*’.
3. Remove neutral statements not indicative of stereotypes e.g. ‘why are [TGT] so *called*’.
4. We filter out completions consisting of mul-

multiple words.<sup>1</sup> Yet, when possible, the input is altered such that only the key term has to be predicted by the model e.g., ‘Why are *russians* so  $x$ ’, where  $x$  = good at playing chess  $\rightarrow$  ‘Why are *russians* so good at  $x$ ’,  $x$  = chess.

The final dataset contains  $\sim$ 2K stereotypes about 274 social groups. The stereotypes are distributed across categories as follows – profession: 713, race: 412, country: 396, gender: 198, age: 171, lifestyle: 123, political: 50, religion: 36. None of the search engines produce stereotypical autocomplete suggestions for members of the LGBTQ+ community. In Table 1 we provide some examples from the dataset. See Appendix B for more details on the data acquisition and search engines. The full code and dataset are publicly available.<sup>2</sup>

#### 4 Correlating human stereotypes with salient attributes in pretrained models

To test for human stereotypes, we propose a stereotype elicitation method that is inspired by cloze testing, a technique that stems from psycholinguistics. Using our method we retrieve salient attributes from the model in an unsupervised manner and compute recall scores over the stereotypes captured in our search engine dataset.

**Pretrained models** We study different types of pretrained LMs of which 3 are monolingual and 2 multilingual: **BERT** (Devlin et al., 2019) uncased trained on the BooksCorpus dataset (Zhu et al., 2015) and English Wikipedia; **RoBERTa** (Liu et al., 2019), the optimized version of BERT that is in addition trained on data from CommonCrawl News (Nagel, 2016), OpenWebTextCorpus (Gokaslan and Cohen, 2019) and STORIES (Trinh and Le, 2018); **BART**, a denoising autoencoder (Lewis et al., 2020) that while using a different architecture and pretraining strategy from RoBERTa, uses the same training data. Moreover, we use **mBERT**, that apart from being trained on Wikipedia in multiple languages, is identical to BERT. We use the uncased version that supports 102 languages. Similarly, **XLNet** is the multilingual variant of RoBERTa (Conneau et al., 2020) that is trained on cleaned CommonCrawl data (Wenzek et al., 2020) and

<sup>1</sup>Although incompatible with our set-up, we do not remove them from the dataset as they can be valuable in future studies.

<sup>2</sup>[https://github.com/RochelleChoenni/stereotypes\\_in\\_lms](https://github.com/RochelleChoenni/stereotypes_in_lms)

supports 100 languages. We include both versions of a model (i.e. **Base** and **Large**) if available. Appendix A provides more details on the models.

**Stereotype elicitation method** For each sample in our dataset we feed the model the template sentence and replace [ATTR] with the [MASK] token. We then retrieve the top  $k = 200$  model predictions for the MASK token, and test how many of the stereotypes found by the search engines are also encoded in the LMs. We adapt the method from Kurita et al. (2019) to rank the top  $k$  returned model outputs based on their typicality for the respective social group. We quantify typicality by computing the log probability of the model probability for the predicted completion corrected for by the prior probability of the completion e.g.:

$$P_{post}(y = \text{strict} | \text{Why are parents so } y ?) \quad (1)$$

$$P_{prior}(y = \text{strict} | \text{Why are [MASK] so } y ?) \quad (2)$$

$$p = \log(P_{post}/P_{prior}) \quad (3)$$

i.e., measuring association between the words by computing the chance of completing the template with ‘strict’ given ‘parents’ corrected by the prior chance of ‘strict’ given any other group. Note that Eq. 3 has been well-established as a measure for stereotypicality in research from both social psychology (McCauley et al., 1980) and economics (Bordalo et al., 2016). After re-ranking by typicality, we evaluate how many of the stereotypes are correctly retrieved by the model through recall@ $k$  for each of the 8 target categories.

**Results** Figure 1 shows the recall@ $k$  scores per model separated by category, showcasing the ability to directly retrieve stereotypical attributes of social groups using our elicitation method. While models capture the human stereotypes to similar extents, results vary when comparing across categories with most models obtaining the highest recall for country stereotypes. Multilingual models obtain relatively low scores when recalling stereotypical attributes pertaining to age, gender and political groups. Yet, XLNet-L is scoring relatively high on stereotypical profession and race attributes.

Table 2: Ranking: ‘why are old people so bad with’.

Prior	Post
1. memory	1. memory
2. math	2. alcohol
3. money	3. technology
4. children	4. dates

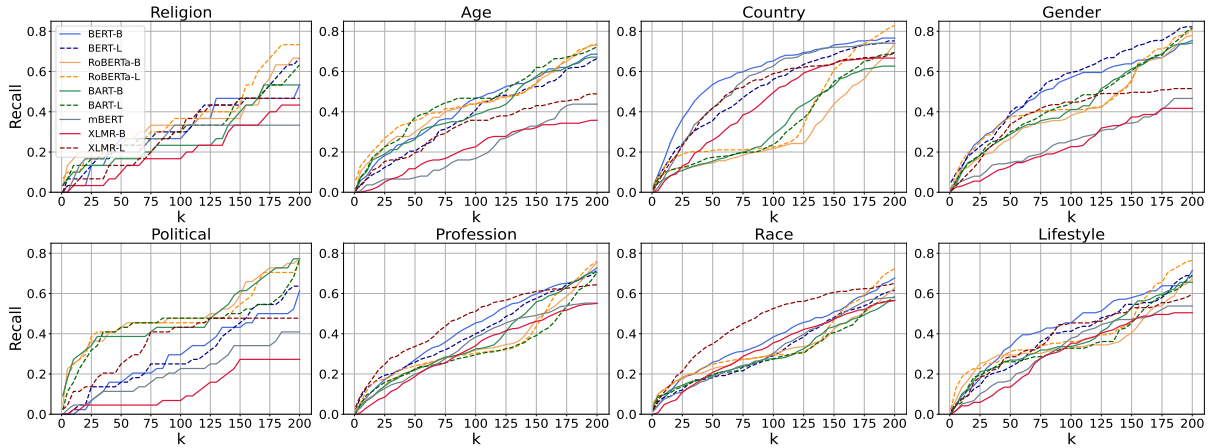


Figure 1: Recall@k scores for recalling the human-defined stereotypes captured in our dataset using our stereotype elicitation method on various pretrained LMs.

The suboptimal performance of multilingual models could be explained in different ways. For instance, as multilingual models are known to suffer from negative interference (Wang et al., 2020), their quality on individual languages is lower compared to monolingual models, due to limited model capacity. This could result in a loss of stereotypical information. Alternatively, multilingual models are trained on more culturally diverse data, thus conflicting information could counteract within the model with stereotypes from different languages dampening each other’s effect. Cultural differences might also be more pronounced when it comes to e.g. age and gender, whilst profession and race stereotypes might be established more universally.

## 5 Quantifying emotion towards different social groups

To study stereotypes through emotion, we draw inspiration from psychology studies showing that stereotypes evoke distinct emotions based on different types of perceived threats (Cottrell and Neuberg, 2005) or perceived social status and competitiveness of the targeted group (Fiske, 1998). For instance, Cottrell and Neuberg (2005) show that both feminists and African Americans elicit anger, but while the former group is perceived as a threat to social values, the latter is perceived as a threat to property instead. Thus, the stereotypes that underlie the emotion are likely different. Whilst strong emotions are not evidence of stereotypes per se, they do suggest the powerful effects of subtle biases captured in the model. Thus, the study into emotion profiles provides us with a good starting point to identify which stereotypes associated with

the social groups evoke those emotions. To this end, we (1) build emotion profiles for social groups in the models and (2) retrieve stereotypes about the groups that most strongly elicit emotions.

**Model predictions** To measure the emotions encoded by the model, we feed the model the 5 stereotype eliciting templates for each social group and retrieve the top 200 predictions for the [MASK] token (1000 in total). When taking the 1000 salient attributes retrieved from the 5 templates, we see that there are many overlapping predictions, hence we are left with only approx. between 300-350 unique attributes per social group. This indicates that the returned model predictions are robust with regard to the different templates.

**Emotion scoring** For each group, we score the predicted set of stereotypical attributes  $W_{TGT}$  using the NRC emotion lexicon (Mohammad and Turney, 2013) that contains  $\sim 14K$  English words that are manually annotated with Ekman’s eight basic emotions (fear, joy, anticipation, trust, surprise, sadness, anger, and disgust) (Ekman, 1999) and two sentiments (negative and positive). These emotions are considered basic as they are thought to be shaped by natural selection to address survival-related problems, which is often denoted as a driving factor for stereotyping (Cottrell and Neuberg, 2005). We use the annotations that consist of a binary value (i.e. 0 or 1) for each of the emotion categories; words can have multiple underlying emotions (e.g. *selfish* is annotated with ‘negative’, ‘anger’ and ‘disgust’) or none at all (e.g. *vocal* scores 0 on all categories). We find that the coverage for the salient attributes in the NRC lexicon is  $\approx 70-75\%$  per group.

We score groups by counting the frequencies with which the predicted attributes  $W_{TGT}$  are associated with the emotions and sentiments. For each group, we remove attributes from  $W_{TGT}$  that are not covered in the lexicon. Thus, we do not extract emotion scores for the exact same number of attributes per group (number of unique attributes and coverage in the lexicon vary). Thus, we normalize scores per group by the number of words for which we are able to retrieve emotion scores ( $\approx 210$ -250 per group). The score of an emotion-group pair is computed as follows:

$$s_{emo}(TGT) = \sum_{i=w}^{|W_{TGT}|} \text{NRC}_{emo}(i) / (|W_{TGT}|) \quad (4)$$

We then define emotion vectors  $\hat{v} \in \mathcal{R}^{10}$  for each group  $TGT$ :  $\hat{v}_{TGT} = [s_{fear}, s_{joy}, s_{sadness}, s_{trust}, s_{surprise}, s_{anticipation}, s_{disgust}, s_{anger}, s_{negative}, s_{positive}]$ , which we use as a representation for the emotion profiles within the model.

**Analysis** Figure 2, provides examples of the emotion profiles encoded for a diverse set of social groups to demonstrate how these profiles allow us to expose stereotypes. For instance, we see that in RoBERTa-B religious people and liberals are primarily associated with attributes that underlie anger. Towards homosexuals, the same amount of anger is accompanied by disgust and fear as well. As a result, we can detect distinct salient attributes that contribute to these emotions e.g.: Christians are *intense*, *misguided* and *perverse*, liberals are *phony*, *mad* and *rabid*, whilst homosexuals are *dirty*, *bad*, *filthy*, *appalling*, *gross* and *indecent*. The finding that homosexuals elicit relatively much disgust can be confirmed by studies on humans as well (Cottrell and Neuberg, 2005). Similarly, we find that Greece and Puerto Rico elicit relatively much fear and sadness in RoBERTa-B. Whereas Puerto Rico is *turbulent*, *battered*, *armed*, *precarious* and *haunted*, for Greece we find attributes such as *failing*, *crumbling*, *inefficient*, *stagnant* and *paralyzed*.

Emotion profiles elicited in BART-B differ considerably, showcasing how vastly sentiments vary across models. In particular, we see that overall the evoked emotion responses are weaker. Moreover, we detect relative differences such as liberals being more negatively associated than homosexuals, encoding attributes such as *cowardly*, *greedy* and *hypocritical*. We also find that BART-B encodes more positive associations e.g., *committed*,

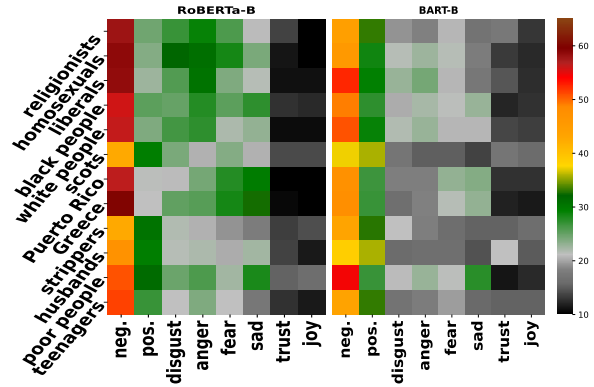


Figure 2: Examples of emotion profiles for a diverse set of social groups from RoBERTa-B and BART-B.

*reliable*, *noble* and *responsible* contributing to trust for husbands. Interestingly, all multilingual models encode vastly more positive attributes for all social groups (see Appendix D). We expect that this might be an artefact of the training data, but leave further investigation of this for future work.

**Comparison across models** We systematically compare the emotion profiles elicited by the social groups across different models by adapting the Representational Similarity Analysis (RSA) from Kriegeskorte et al. (2008). We opted for this method as it takes the relative relations between groups within the same model into account. This is particularly important as we have seen that some models are overall more negatively or positively biased. Yet, when it comes to bias and stereotypicality, we are less interested in absolute differences across models, but rather in how emotions differ towards groups in relation to the other groups. First, the representational similarity within each model is defined using a similarity measure to construct a representational similarity matrix (RSM). We define a similarity vector  $\hat{w}_{TGT}$  for a social group such that every element  $\hat{w}_{ij}$  of the vector is determined by the cosine similarity between  $\hat{v}_i$ , where  $i = TGT$ , and the vector  $\hat{v}_j$  for the  $j$ -th group in the list. The RSM is then defined as the symmetric matrix consisting of all similarity vectors. The resulting matrices are then compared across models by computing the Spearman correlation ( $\rho$ ) between the similarity vectors corresponding to the emotion profiles for a group in a model  $a$  and  $b$ . To express the similarity between the two models we take the mean correlation over all social groups in our list.

**Results** Computing RSA over all categories combined, shows us that RoBERTa-B and BART-B ob-

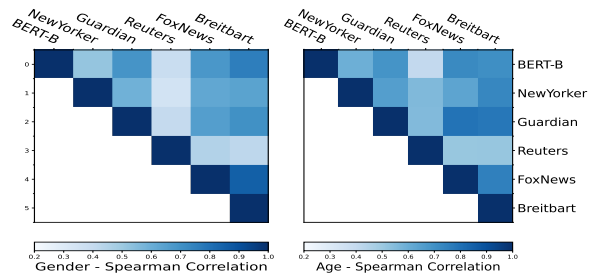


Figure 3: Correlations in emotion profiles for gender and age groups across news sources (BERT-B).

tain the highest correlation ( $\rho = 0.44$ ). While using different architectures and pretraining strategies, the models rely on the same training data. Yet, we included base and large versions of models in our study and find that these models show little to no correlation (see Appendix E, Fig.10). This is surprising, as they are pretrained on the same data and tasks as their base versions (but contain more model parameters e.g. through additional layers). This shows how complex the process is in which associations are established and provides strong evidence that other modelling decisions, apart from training data, contribute to what models learn about groups. Thus, carefully controlling training content can not fully eliminate the need to analyze models w.r.t. the stereotypes that they might propagate.

## 6 Stereotype shifts during fine-tuning

Many debiasing studies intervene at the data level e.g., by augmenting imbalanced datasets (Manela et al., 2021; Webster et al., 2018; Dixon et al., 2018; Zhao et al., 2018) or reducing annotator bias (Sap et al., 2019). These methods are, however, dependent on the dataset, domain, or task, making new mitigation needed when transferring to a new set-up (Jin et al., 2020). This raises the question of how emotion profiles and stereotypes are established through language use, and how they might shift due to new linguistic experience at the fine-tuning stage. We take U.S. news sources from across the political spectrum as a case study, as media outlets are known to be biased (Baron, 2006). By revealing stereotypes learned as an effect of fine-tuning on a specific source, we can trace the newly learned stereotypes back to the respective source.

We rely on the political bias categorisation of news sources from the *AllSides*<sup>3</sup> media bias rating website. These ratings are retrieved using multiple

<sup>3</sup><https://www.allsides.com/media-bias/media-bias-ratings>

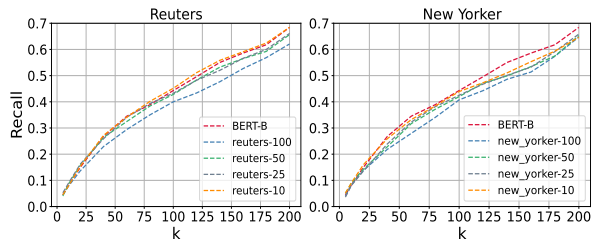


Figure 4: Effect on recall@k when fine-tuning BERT-B on 10, 25, 50 and 100 % of the data

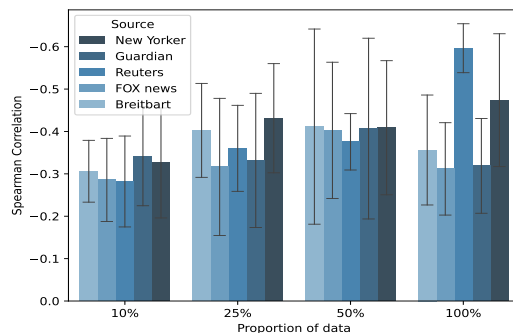


Figure 5: Decrease in Spearman correlation ( $\Delta\rho$ ) after fine-tuning the pretrained models compared to no fine-tuning ( $\Delta\rho = 1$ ) (no correlation left:  $\Delta\rho = -1$ ). We show results for models trained on varying proportions of the data. Results are averaged over categories and standard deviations are indicated by error bars.

methods, including editorial reviews, blind bias surveys, and third party research. Based on these ratings we select the following sources: New Yorker (*far left*), The Guardian (*left*), Reuters (*center*), FOX News (*right*) and Breitbart (*far right*). From each news source we take 4354 articles from the *All-The-News*<sup>4</sup> dataset that contains articles from 27 American Publications collected between 2013 and early 2020. We fine-tune the 5 base models<sup>5</sup> on these news sources using the MLM objective for only 1 training epoch with a learning rate of  $5e-5$  and a batch size of 8 using the HuggingFace library (Wolf et al., 2020). We then quantify the emotion shift after fine-tuning using RSA.

**Results** We find that fine-tuning on news sources can directly alter the encoded stereotypes. For instance, for  $k = 25$ , fine-tuning BERT-B on Reuters informs the model that Croatia is good at *sports* and Russia is good at *hacking*, at the same time, associations such as Pakistan is bad at *football*, Romania is good at *gymnastics* and South Africa at

<sup>4</sup>Available at: <https://tinyurl.com/bx3r3de8>

<sup>5</sup>Training the large models was computationally infeasible.





illustrate the salient attributes that are removed, added and remained constant after fine-tuning. For instance, the role of news media in shaping public opinion about police has received much attention in the wake of the growing polarization over high-profile incidents (Intravia et al., 2018; Graziano, 2019). We find clear evidence of this polarization as fine-tuning on New Yorker results in attributes such as *cold*, *unreliable*, *deadly* and *inept*, yet, fine-tuning on FOX news yields positive associations such as *polite*, *loyal*, *cautious* and *exceptional*. In addition, we find evidence for other stark contrasts such as the model picking up on sexist (e.g. women are not *interesting* and *equal* but *late*, *insecure* and *entitled*) and racist stereotypes (e.g. black people are not *misunderstood* and *powerful*, but *bitter*, *rude* and *stubborn*) after fine-tuning on FOX news.

## 7 Conclusion

We present the first dataset containing stereotypical attributes of a range of social groups. Importantly, our data acquisition technique enables the inexpensive retrieval of similar datasets in the future, enabling comparative analysis on stereotype shifts over time. Additionally, our proposed methods could inspire future work on analyzing the effect of training data content, and simultaneously contribute to the field of social psychology by providing a testbed for studies on how stereotypes emerge from linguistic experience. To this end, we have shown that our methods can be used to identify stereotypes evoked during fine-tuning by taking news sources as a case study. Moreover, we have exposed how quickly stereotypes and emotions shift based on training data content, and linked stereotypes to their manifestations as emotions to quantify and compare attitudes towards groups within LMs. We plan to extend our approach to more languages in future work to collect different, more culturally dependent, stereotypes as well.

## 8 Ethical consideration

The examples given in the paper can be considered offensive but are in no way a reflection of the authors' own values and beliefs and should not be taken as such. Moreover, it is important to note that for the fine-tuning experiments only a few interesting examples were studied and showcased. Hence, more thorough research should be conducted before drawing any hard conclusions about the news papers and the stereotypes they propagate. In ad-

dition, our data acquisition process is completely automated and did not require the help from human subjects. While the stereotypes we retrieve stem from real humans, the data we collect is publicly available and completely anonymous as the specific stereotypical attributes and/or search queries can not be traced back to individual users.

## References

- David P Baron. 2006. Persistent media bias. *Journal of Public Economics*, 90(1-2):1–36.
- Camiel J Beukeboom and Christian Burgers. 2019. How stereotypes are shared through language: a review and introduction of the social categories and stereotypes communication (sesc) framework. *Review of Communication Research*, 7:1–37.
- Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. Language (technology) is power: A critical survey of “bias” in nlp. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5454–5476.
- Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pages 4356–4364.
- Pedro Bordalo, Katherine Coffman, Nicola Gennaioli, and Andrei Shleifer. 2016. Stereotypes. *The Quarterly Journal of Economics*, 131(4):1753–1794.
- Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.
- Catherine A Cottrell and Steven L Neuberg. 2005. Different emotional reactions to different groups: a sociofunctional threat-based approach to “prejudice”. *Journal of personality and social psychology*, 88(5):770.
- Amy JC Cuddy, Susan T Fiske, Virginia SY Kwan, Peter Glick, Stéphanie Demoulin, Jacques-Philippe Leyens, Michael Harris Bond, Jean-Claude Croizet, Naomi Ellemers, Ed Sleebos, et al. 2009. Stereotype content model across cultures: Towards universal similarities and some differences. *British Journal of Social Psychology*, 48(1):1–33.

- Thomas Davidson, Debasmita Bhattacharya, and Ingmar Weber. 2019. Racial bias in hate speech and abusive language detection datasets. In *Proceedings of the Third Workshop on Abusive Language Online*, pages 25–35.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2018. Measuring and mitigating unintended bias in text classification. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, pages 67–73.
- MeiXing Dong, David Jurgens, Carmen Banea, and Rada Mihalcea. 2019. Perceptions of social roles across cultures. In *International Conference on Social Informatics*, pages 157–172. Springer.
- Paul Ekman. 1999. Basic emotions. *Handbook of Cognition and Emotion*, pages 45–60.
- Anjalie Field and Yulia Tsvetkov. 2020. Unsupervised discovery of implicit gender bias. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 596–608.
- Susan T Fiske. 1998. Stereotyping, prejudice, and discrimination. *The handbook of social psychology*, 2(4):357–411.
- Aaron Gokaslan and Vanya Cohen. 2019. Openwebtext corpus. <https://skylion007.github.io/OpenWebTextCorpus/>.
- Lisa M Graziano. 2019. News media and perceptions of police: a state-of-the-art-review. *Policing: An International Journal*.
- Lasana T Harris and Susan T Fiske. 2006. Dehumanizing the lowest of the low: Neuroimaging responses to extreme out-groups. *Psychological science*, 17(10):847–853.
- Lasana T Harris and Susan T Fiske. 2009. Social neuroscience evidence for dehumanised perception. *European review of social psychology*, 20(1):192–231.
- Madeline E Heilman, Aaron S Wallen, Daniella Fuchs, and Melinda M Tamkins. 2004. Penalties for success: reactions to women who succeed at male gender-typed tasks. *Journal of applied psychology*, 89(3):416.
- Perry Hinton. 2017. Implicit stereotypes and the predictive brain: cognition and culture in “biased” person perception. *Palgrave Communications*, 3(1):1–9.
- Jonathan Intravia, Kevin T Wolff, and Alex R Piquero. 2018. Investigating the effects of media consumption on attitudes toward police legitimacy. *Deviant Behavior*, 39(8):963–980.
- Xisen Jin, Francesco Barbieri, Brendan Kennedy, Aida Mostafazadeh Davani, Leonardo Neves, and Xiang Ren. 2020. On transferability of bias mitigation effects in language model fine-tuning. *arXiv preprint arXiv:2010.12864*.
- Svetlana Kiritchenko and Saif Mohammad. 2018. Examining gender and race bias in two hundred sentiment analysis systems. In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, pages 43–53.
- Nikolaus Kriegeskorte, Marieke Mur, and Peter A Bandettini. 2008. Representational similarity analysis-connecting the branches of systems neuroscience. *Frontiers in systems neuroscience*, 2:4.
- Keita Kurita, Nidhi Vyas, Ayush Pareek, Alan W Black, and Yulia Tsvetkov. 2019. Measuring bias in contextualized word representations. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, pages 166–172.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Anne Maass. 1999. Linguistic intergroup bias: Stereotype perpetuation through language. In *Advances in experimental social psychology*, volume 31, pages 79–121. Elsevier.
- Diane M Mackie, Thierry Devos, and Eliot R Smith. 2000. Intergroup emotions: Explaining offensive action tendencies in an intergroup context. *Journal of personality and social psychology*, 79(4):602.
- Daniel de Vassimon Manela, David Errington, Thomas Fisher, Boris van Breugel, and Pasquale Minervini. 2021. Stereotype and skew: Quantifying gender bias in pre-trained and fine-tuned language models. *arXiv preprint arXiv:2101.09688*.
- Chandler May, Alex Wang, Shikha Bordia, Samuel Bowman, and Rachel Rudinger. 2019. On measuring social biases in sentence encoders. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 622–628.

- Clark McCauley, Christopher L Stitt, and Mary Segal. 1980. Stereotyping: From prejudice to prediction. *Psychological Bulletin*, 87(1):195.
- Saif M Mohammad and Peter D Turney. 2013. Nrc emotion lexicon. *National Research Council, Canada*, 2.
- Moin Nadeem, Anna Bethke, and Siva Reddy. 2020. Stereoset: Measuring stereotypical bias in pretrained language models. *arXiv preprint arXiv:2004.09456*.
- Sebastian Nagel. 2016. Cc-news dataset. <https://commoncrawl.org/2016/10/news-dataset-available/>.
- Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel Bowman. 2020. Crows-pairs: A challenge dataset for measuring social biases in masked language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1953–1967.
- Chan Young Park, Xinru Yan, Anjalie Field, and Yulia Tsvetkov. 2020. Multilingual contextual affective analysis of lgbt people portrayals in wikipedia. *arXiv preprint arXiv:2010.10820*.
- Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A Smith. 2019. The risk of racial bias in hate speech detection. In *Proceedings of the 57th annual meeting of the association for computational linguistics*, pages 1668–1678.
- Seth Stephens-Davidowitz. 2018. Everybody Lies: What the internet can tell us about who we really are. In *Bloomsbury Publishing Plc*.
- Tony Sun, Andrew Gaut, Shirlyn Tang, Yuxin Huang, Mai ElSherief, Jieyu Zhao, Diba Mirza, Elizabeth Belding, Kai-Wei Chang, and William Yang Wang. 2019. Mitigating gender bias in natural language processing: Literature review. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1630–1640.
- Yi Chern Tan and L Elisa Celis. 2019. Assessing social and intersectional biases in contextualized word representations. *arXiv preprint arXiv:1911.01485*.
- Molly Parker Tapias, Jack Glaser, Dacher Keltner, Kristen Vasquez, and Thomas Wickens. 2007. Emotion and prejudice: Specific emotions toward outgroups. *Group Processes & Intergroup Relations*, 10(1):27–39.
- Trieu H Trinh and Quoc V Le. 2018. A simple method for commonsense reasoning. *arXiv preprint arXiv:1806.02847*.
- Zirui Wang, Zachary C Lipton, and Yulia Tsvetkov. 2020. On negative interference in multilingual language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4438–4450.
- Kellie Webster, Marta Recasens, Vera Axelrod, and Jason Baldridge. 2018. Mind the gap: A balanced corpus of gendered ambiguous pronouns. *Transactions of the Association for Computational Linguistics*, 6:605–617.
- Bernard Weiner. 1993. On sin versus sickness: A theory of perceived responsibility and social motivation. *American psychologist*, 48(9):957.
- Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzmán, Armand Joulin, and Édouard Grave. 2020. Ccnet: Extracting high quality monolingual datasets from web crawl data. In *Proceedings of The 12th Language Resources and Evaluation Conference*, pages 4003–4012.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. **Transformers: State-of-the-art natural language processing**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.
- Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. Gender bias in coreference resolution: Evaluation and debiasing methods. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, volume 2.
- Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *Proceedings of the IEEE international conference on computer vision*, pages 19–27.

## A Pretrained model details

Model	tokenization	L	dim	H	params	V	D	task	#lgs
BERT-B	WordPiece	12	768	12	110M	30K	16GB	MLM+NSP	1
BERT-L	WordPiece	24	1024	16	336M	30K	16GB	MLM+NSP	1
RoBERTa-B	BPE	12	768	12	125M	50K	160GB	MLM	1
RoBERTa-L	BPE	24	1024	16	335M	50K	160GB	MLM	1
BART-B	BPE	12	768	16	139M	50K	160GB	Denosing	1
BART-L	BPE	24	1024	16	406M	50K	160GB	Denosing	1
mBERT	WordPiece	12	768	12	168M	110K	-	MLM+NSP	102
XLMR-B	SentencePiece	12	768	8	270M	250K	2.5TB	MLM	100
XLMR-L	SentencePiece	24	1024	16	550M	250K	2.5TB	MLM	100

Table 3: Summary statistics of the model architectures: tokenization method, number of layers  $L$ , hidden state dimensionality  $dim$ , number of attention heads  $H$ , number of model parameters  $params$ , vocabulary size  $V$ , training data size  $D$ , pretraining tasks, and number of languages used  $\#lgs$ .

## B Data acquisition

For the collection of autocomplete suggestions we rely on the free publicly available API’s from the respective engines using the following base url’s:

- Google: <http://suggestqueries.google.com/complete/search>
- Yahoo: <http://sugg.search.yahoo.net/sg>
- DuckDuckGo: <https://duckduckgo.com/ac>

All search engine suggestions are automatically generated by an algorithm without human involvement. These suggestions are supposed to be based on factors like popularity and similarity. We enter the search queries anonymously such that the resulting suggestions are mainly based on common queries from other people’s search histories. Unfortunately, however, exact details about the workings of the algorithms are not publicly available, but an extensive explanation of Google’s search predictions can be found here: [Google’s documentation on autocomplete suggestions](#). Moreover, Figure B illustrates the contribution of each search engine to the datasets. We see that while each search engine relies on a different algorithm, in many cases the engines predict similar stereotypical attributes regardless. Moreover, the dataset was constructed during the period January-May 2021. However, given that the algorithms behind these engines are constantly evolving, it is not guaranteed that the same approach will yield identical results in the future. We will make the dataset and corresponding code available upon publication.

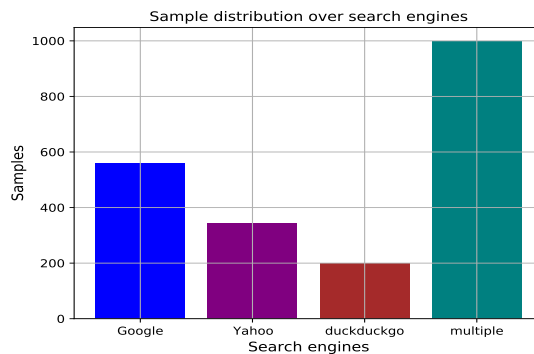


Figure 8: Distribution of the contribution of each search engine to the final dataset.

## C Lists of social groups per category

### Age (11)

"elders", "kids", "millenials", "teens", "boomers", "old people"  
 "old women", "children", "teenagers", "daughters", "sons"

### Gender (25)

"girls", "women", "men", "females", "males", "boys"  
 "boyfriends", "girlfriends", "stepmothers", "stepfathers", "ladies", "gentlemen"  
 "brothers", "sisters", "mothers", "fathers", "grandfathers", "grandmothers"  
 "wives", "husbands", "brides", "grooms", "schoolgirls", "schoolboys"  
 "mommies", "daddies"

### Race (99)

"Asians", "Americans", "Europeans", "Jews", "Indians", "Russians", "Africans", "Black people"  
 "Mexicans", "Whites", "Blacks", "White people", "Germans", "blondes", "blonde girls", "Arabs"  
 "White Americans", "Black Americans", "Hispanics", "Native Americans", "Black men", "White men"  
 "Asian women", "Asian men", "Black women", "the Dutch", "Irish people", "Irish men", "White women"  
 "Indian men", "Indian women", "Chinese men", "Chinese women", "Japanese women", "Japanese men"  
 "Indian parents", "Asian parents", "White parents", "Black parents", "Black fathers", "Latinas", "Latinos"  
 "Latin people", "Brazilian women", "Asian kids", "Black kids", "White kids", "African Americans"  
 "Nigerians", "Ethiopians", "Ukrainians", "Sudanese people", "Afghans", "Iraqis", "Hispanic men"  
 "Hispanic women", "Italians", "Italian men", "Italian women", "Somalis", "Iranian people", "Iranians"  
 "Australians", "Australian men", "Australian women", "Aussies", "Ghanaians", "Swedes", "Finns"  
 "Venezuelans", "Moroccans", "Syrians", "Pakistanis", "British people", "French people", "Greeks"  
 "Indonesians", "Vietnamese people", "Romanians", "Ecuadorians", "Norwegians", "Nepalis", "Scots"  
 "Bengalis", "Polish people", "Taiwanese people", "Albanians", "Colombians", "Egyptians", "Koreans"  
 "Persian people", "Portuguese men", "Portuguese women", "Turkish people", "Austrians"  
 "South Africans", "Dutch people", "Chileans", "Lebanese people"

### Countries (73)

"Ethiopia", "Europe", "Ukraine", "Sudan", "Afghanistan", "Iraq", "Yemen", "Russia", "Italy", "Somalia"  
 "Iran", "Venezuela", "Oman", "Syria", "China", "India", "Britain", "France", "Greece", "Scotland"  
 "Paraguay", "Brazil", "Africa", "Jordan", "Indonesia", "Vietnam", "Pakistan", "Romania", "Ecuador"  
 "Mexico", "Puerto Rico", "Cameroon", "Kenya", "Liberia", "Sierra Leon", "Qatar", "Syria", "Kenya"  
 "Saudi Arabia", "Lebanon", "South Africa", "Korea", "Singapore", "Germany", "Ireland", "Morocco"  
 "Turkey", "Laos", "Bangladesh", "Guatemala", "Ghana", "Nepal", "Albania", "Spain", "Sweden"  
 "Argentina", "Chile", "Taiwan", "Finland", "Australia", "Egypt", "Peru", "Poland", "Columbia", "Bolivia"  
 "Japan", "Norway", "Cape Verde", "Portugal", "Austria", "the Netherlands"  
 "Croatia", "Malta", "Belgium"

### Profession (115)

"students", "politicians", "doctors", "business men", "librarians", "artists", "professors", "priests"  
 "bosses", "police", "police officers", "soldiers", "scientists", "physicians", "cashiers", "housekeepers"  
 "teachers", "janitors", "models", "actresses", "pilots", "strippers", "brokers", "hairdressers", "bartenders"  
 "diplomats", "receptionists", "realtors", "mathematicians", "barbers", "coaches", "business people"  
 "construction workers", "managers", "accountants", "commanders", "firefighters", "movers", "attorneys"  
 "bakers", "athletes", "dancers", "carpenters", "mechanics", "handymen", "musicians", "detectives"  
 "entrepreneurs", "opera singers", "chiefs", "lawyers", "software developers", "farmers", "writers"  
 "real-estate developers", "butchers", "electricians", "prosecutors", "bankers", "cooks", "plumbers"  
 "football players", "boxers", "chess players", "swimmers", "tennis players", "supervisors", "attendants"  
 "producers", "researchers", "economists", "physicists", "psychologists", "sales people", "assistants"  
 "engineers", "comedians", "painters", "civil servants", "guitarists", "linguists", "laborers", "historians"  
 "chemists", "pensioners", "performing artists", "singers", "secretaries", "auditors", "counselors"  
 "dentists", "analysts", "nurses", "waiters", "authors", "architects", "academics", "directors", "illustrators"  
 "clerks", "photographers", "cleaners", "composers", "pharmacists", "sociologists", "journalists"  
 "guards", "actors", "midwives", "sheriffs", "editors", "designers", "judges", "poets", "maids"

### Religion (10)

"Religious people", "Muslims", "Christians", "Hindus", "atheists", "Buddhists"  
 "Catholics", "Protestants", "Sikhs", "Mormons"

<b>Political (9)</b>
"immigrants", "conservatives", "liberals", "trump supporters", "voters", "communists"
"capitalists", "populists", "socialists"

<b>Sexuality (12)</b>
"gay people", "lesbians", "queer people", "transgenders", "homosexuals", "pansexual people"
"queers", "faggots", "bisexual people", "asexual people", "crossdressers", "fags"

<b>Lifestyle (19)</b>
"hipsters", "nerds", "rednecks", "homeless people", "feminists", "rich people", "poor people", "criminals"
"frats", "frat boys", "sorority girls", "hippies", "geeks", "goths", "punks", "Californians"
"celebrities", "redheads", "gingers"

## D Emotion profiles from multilingual models

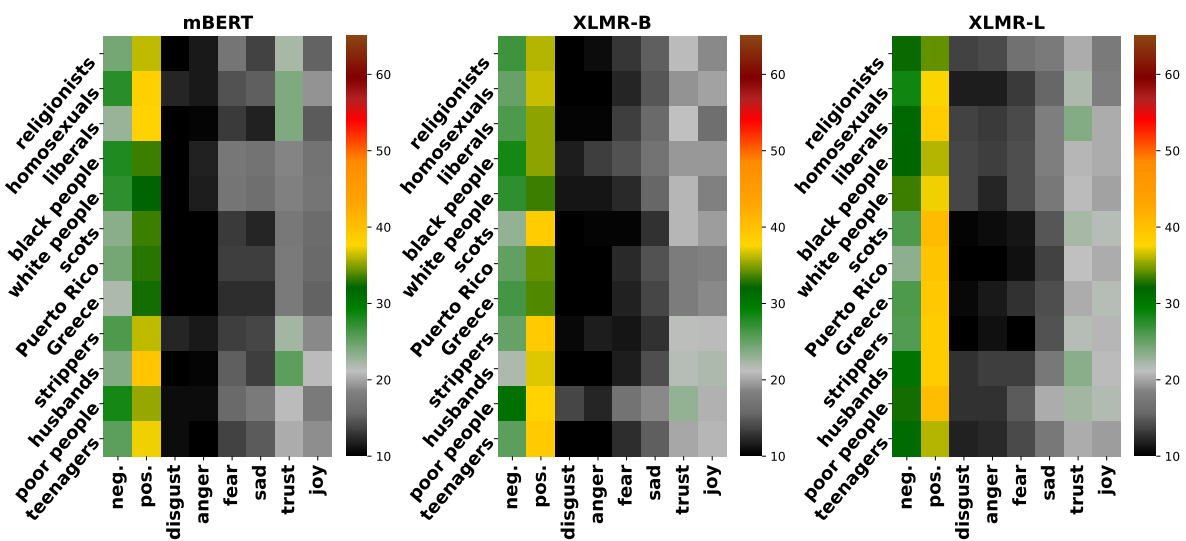


Figure 9: Examples of emotion profiles for the multilingual models. It showcases that these models are much more positive about all social groups in comparison to the monolingual models. Whereas we observed that monolingual models primarily encode negative associations for most groups, associations encoded within the multilingual models are more balanced between positive and negative sentiments.

## E Additional quantitative results of systematic shifts in emotion profiles across models

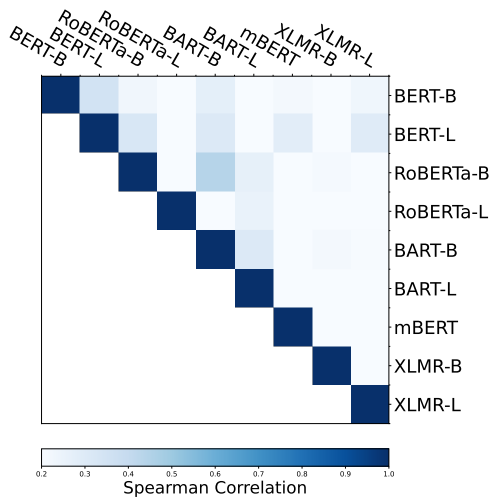


Figure 10: Spearman correlation between each pair of models computed over all social groups. This figure illustrates that there is fairly little correlation between any of the models when it comes to the emotion profiles that they capture.

$\Delta\rho$	Source	Religion	Profession	Lifestyle	Sexuality	Race	Gender	Country	Age	Political
BERT-B	NewYorker	-.56	-.34	-.25	-.23	-.39	-.47	-.47	-.43	<b>-.72</b>
	Guardian	<b>-.49</b>	-.34	-.08	-.23	-.37	-.31	-.43	-.31	<b>-.49</b>
	Reuters	<b>-.71</b>	-.53	-.43	-.65	-.53	-.63	-.69	-.60	-.54
	FOX news	-.46	-.30	-.16	-.22	-.35	-.30	-.44	-.33	<b>-.51</b>
	BreitBart	-.39	-.25	-.11	-.21	-.33	-.23	-.40	-.34	<b>-.66</b>
RoBERTa-B	NewYorker	-.20	-.22	-.20	<b>-.29</b>	-.21	-.24	-.16	-.08	-.38
	Guardian	-.19	-.20	-.19	-.20	-.22	-.18	-.16	-.13	<b>-.24</b>
	Reuters	-.25	-.32	-.33	-.21	-.33	<b>-.49</b>	-.37	-.24	-.40
	FOX news	-.10	-.18	-.14	<b>-.37</b>	-.16	-.12	-.16	-.25	-.25
	BreitBart	-.15	-.23	-.21	-.41	-.18	-.27	-.22	-.18	<b>-.43</b>
BART-B	NewYorker	-.56	-.48	-.40	<b>-.60</b>	-.44	-.55	-.43	-.48	-.49
	Guardian	-.49	-.48	-.32	-.41	-.37	-.50	-.47	<b>-.67</b>	-.33
	Reuters	-.43	-.51	-.45	-.51	-.53	-.54	-.54	<b>-.70</b>	-.29
	FOX news	-.27	-.50	-.32	-.44	-.37	-.44	-.42	<b>-.65</b>	-.50
	BreitBart	-.37	-.48	-.42	-.35	-.37	-.51	-.44	<b>-.56</b>	-.50
mBERT	NewYorker	-.58	-.64	-.33	-.44	-.64	-.63	<b>-.80</b>	-.59	-.38
	Guardian	-.58	-.49	-.30	-.50	-.63	-.72	<b>-.77</b>	-.53	-.37
	Reuters	-.50	-.56	-.29	-.46	-.37	-.59	<b>-.85</b>	-.33	-.42
	FOX news	-.35	-.64	-.36	-.54	-.68	<b>-.71</b>	<b>-.71</b>	-.49	-.60
	BreitBart	-.39	-.66	-.36	-.43	-.51	-.61	<b>-.75</b>	-.40	-.55
XLMR-B	NewYorker	-.44	-.76	-.45	-.66	-.61	<b>-.86</b>	-.66	-.72	-.58
	Guardian	-.52	-.72	-.49	-.46	-.68	<b>-.83</b>	-.53	-.63	-.38
	Reuters	-.53	<b>-.74</b>	-.69	-.55	-.67	-.73	-.53	-.69	-.57
	FOX news	-.40	<b>-.71</b>	-.47	-.57	-.58	-.69	-.51	-.69	-.30
	BreitBart	-.60	-.76	-.47	-.56	-.75	<b>-.79</b>	-.60	-.65	-.51

Table 4: Emotion shifts after fine-tuning for 1 training epoch on  $\pm 4.5K$  articles from the respective news sources. We quantify shift as the decrease in similarity after fine-tuning, i.e. change in averaged Spearman correlation ( $\Delta\rho$ ), between the pretrained and fine-tuned model respectively. If the emotion profiles do no change  $\rho = 1$  and thus  $\Delta\rho = 0$ , on the other hand, if no correlation remains after fine-tuning  $\Delta\rho = -1$ . Biggest changes are indicated by bold letters.