Contents lists available at ScienceDirect



# Information Processing and Management

journal homepage: www.elsevier.com/locate/ipm



# Towards understanding and mitigating unintended biases in language model-driven conversational recommendation



Tianshu Shen<sup>a,\*</sup>, Jiaru Li<sup>a</sup>, Mohamed Reda Bouadjenek<sup>b</sup>, Zheda Mai<sup>a</sup>, Scott Sanner<sup>a</sup>

<sup>a</sup> Department of Mechanical and Industrial Engineering, The University of Toronto, Canada

<sup>b</sup> School of Information Technology, Deakin University, Waurn Ponds Campus, Geelong, VIC 3216, Australia

# ARTICLE INFO

Keywords: Conversational recommendation systems BERT Contextual language models Bias and discrimination

# ABSTRACT

Conversational Recommendation Systems (CRSs) have recently started to leverage pretrained language models (LM) such as BERT for their ability to semantically interpret a wide range of preference statement variations. However, pretrained LMs are prone to intrinsic biases in their training data, which may be exacerbated by biases embedded in domain-specific language data (e.g., user reviews) used to fine-tune LMs for CRSs. We study a simple LM-driven recommendation backbone (termed LMRec) of a CRS to investigate how unintended bias - i.e., bias due to language variations such as name references or indirect indicators of sexual orientation or location that should not affect recommendations - manifests in substantially shifted price and category distributions of restaurant recommendations. For example, offhand mention of names associated with the black community substantially lowers the price distribution of recommended restaurants, while offhand mentions of common male-associated names lead to an increase in recommended alcohol-serving establishments. While these results raise red flags regarding a range of previously undocumented unintended biases that can occur in LMdriven CRSs, there is fortunately a silver lining: we show that train side masking and test side neutralization of non-preferential entities nullifies the observed biases without significantly impacting recommendation performance.

# 1. Introduction

With the prevalence of language-based intelligent assistants such as Amazon Alexa and Google Assistant, conversational recommender systems (CRSs) have attracted growing attention as they can dynamically elicit users' preferences and incrementally adapt recommendations based on user feedback (Gao, Lei, He, de Rijke, & Chua, 2021; Jannach, Manzoor, Cai, & Chen, 2021). As one of the most crucial foundations of CRSs, Natural Language Processing (NLP) has witnessed several breakthroughs in the past few years, including the use of pretrained transformer-based language models (LMs) for downstream tasks (Otter, Medina, & Kalita, 2020). Numerous studies have shown that these transformer-based LMs such as BERT (Devlin, Chang, Lee, & Toutanova, 2019), RoBERTa (Liu, Lin, Shi, & Zhao, 2021) and GPT (Radford, Narasimhan, Salimans, & Sutskever, 2018) pretrained on large corpora can learn universal language representations and are extraordinarily powerful for many downstream tasks via fine-tuning (Qiu et al., 2020). Recently, CRSs have started to leverage pretrained LMs for their ability to semantically interpret a wide range of preference statement variations and have demonstrated their potential to build a variety of strong CRSs (Hada & Shevade, 2021; Malkiel, Barkan, Caciularu, Razin, Katz, & Koenigstein, 2020; Penha & Hauff, 2020).

\* Corresponding author.

https://doi.org/10.1016/j.ipm.2022.103139

Received 28 April 2022; Received in revised form 7 September 2022; Accepted 23 October 2022 0306-4573/© 2022 Elsevier Ltd. All rights reserved.

*E-mail addresses:* tina.shen@mail.utoronto.ca (T. Shen), kellyjiaru.li@mail.utoronto.ca (J. Li), reda.bouadjenek@deakin.edu.au (M.R. Bouadjenek), zheda.mai@mail.utoronto.ca (Z. Mai), ssanner@mie.utoronto.ca (S. Sanner).

However, pretrained LMs are well-known for exhibiting unintended social biases involving race, gender, or religion (Liang, Wu, Morency, & Salakhutdinov, 2021; Lu, Mardziel, Wu, Amancharla, & Datta, 2020; Sheng, Chang, Natarajan, & Peng, 2019). These biases result from unfair allocation of resources (e.g., policing, hospital services, or job availability) (Hutchinson et al., 2020; Zhang, Lu, Abdalla, McDermott, & Ghassemi, 2020), stereotyping that propagates negative generalizations about particular social groups (Nadeem, Bethke, & Reddy, 2021), text that misrepresents the distribution of different social groups in the population (Liang et al., 2021), or language that is denigrating to particular social groups (Guo & Caliskan, 2021). Moreover, these biases may also be exacerbated by biases in data used for domain-specific LM fine-tuning for downstream tasks (Jin et al., 2021; Nadeem et al., 2021).

In this paper, we study a simple LM-driven recommendation backbone (termed LMRec) for CRSs to investigate how *unintended bias* manifests in substantially shifted price and category distributions of restaurant recommendations. Specifically, we generate templates with placeholders (a.k.a. *template-based result generation*) indicating non-preferential information such as names or relationships that implicitly indicate race, gender, sexual orientation, geographical context, and religion, and study how different substitutions for these placeholders modulate price and category distributions (a.k.a. *attribute-based analysis*) with the proposed metrics. To this end, we make the following technical contributions:

- The proposed investigation methodology extends the template-based analysis from research works on bias in language models (Kurita, Vyas, Pareek, Black, & Tsvetkov, 2019a; May, Wang, Bordia, Bowman, & Rudinger, 2019; Sheng et al., 2019; Tan & Celis, 2019) and the attribute-based analysis from the literature on fair recommender systems (Deldjoo, Anelli, Zamani, Bellogin, & Di Noia, 2021; Mansoury, Mobasher, Burke, & Pechenizkiy, 2019; Tsintzou, Pitoura, & Tsaparas, 2019) to generate conversational recommendation results and to perform user–item attribute fairness analysis in language-based conversational recommender systems.
- Our proposed methodology for user-item attribute bias analysis in conversational recommender systems provides novel techniques and metrics for use in fair recommender systems research.

Through the application of the above technical methodology and proposed metrics, we make the following key observational contributions:

- · LMRec recommends significantly more low-priced establishments when a black- vs. white-associated name is mentioned.
- LMRec recommends significantly more alcohol-serving venues when a male- vs. female-associated name is mentioned.
- LMRec picks up indirect mentions of homosexual relations (e.g. "my brother and his boyfriend") as indicated by the elevation of "gay bar" in the recommendations vs. a heterosexual relation (e.g., "my brother and his girlfriend").
- Mentioning visits to professional locations (a "fashion studio" or "law office") or a "synagogue" lead to a higher average price range of LMRec recommendations compared to mentioning a visit to the "convenience store" or a "mosque".

While these results raise red flags regarding a range of previously undocumented unintended biases that can occur in LM-driven CRSs, there is fortunately a silver lining: we show that combining train side masking and test side neutralization of non-preferential entities nullifies the observed biases without hurting recommendation performance. Hence, with future language model-driven CRS assistants having a potential reach of hundreds of millions of end-users, the results of this work present an important step forward in identifying and mitigating potential sources of bias in CRSs that align with general goals of inequality reduction in society (Desa et al., 2016).

# 2. Related work

This section briefly summarizes how fairness/bias issues have been analysed in two requisite elements of language model-driven recommender systems: recommendation systems and language models. Following this, we review conversational recommender systems, where there is a notable lack of work on bias in LM-driven CRSs.

#### 2.1. Fairness and bias in recommendation systems

Recommendation Systems (RS) provide users with personalized suggestions and can help alleviate information overload (Chen et al., 2020). While much recent work in RS investigates improved machine learning models for recommendation (Chen et al., 2020), recent years have seen a rise in the number of works examining fairness and bias in recommendation. In brief, *unfairness* in recommendations manifests as systematic discrimination against specific individuals in favour of others (Friedman & Nissenbaum, 1996) based on protected attributes such as gender and age. Research studies usually perform an attribute-based analysis of fairness in recommender systems, where users or items are labelled with some attributes that cluster them into groups.

Age & Gender Bias: Performance disparities (with NDCG metric) of Collaborative Filtering (CF) algorithms in the recommendation of movies and music have been observed (Ekstrand et al., 2018), revealing unfairness with regard to users' age and gender. Studies also show empirically that popular recommendation algorithms work better for males since many datasets are male-user-dominated (Ekstrand & Pera, 2017). One way to measure gender and age fairness of different recommendation models is based on generalized cross entropy (GCE) (Deldjoo et al., 2021; Deldjoo, Anelli, Zamani, Kouki, & Noia, 2019); specifically, this work shows that a simple popularity-based algorithm provides better recommendations to male users and younger users, while on the opposite side, uniform random recommendations and collaborative filtering algorithms provide better recommendations to female users and older users (Deldjoo et al., 2021). In other work, Lin, Sonboli, Mobasher, and Burke (2019) study how different recommendation

algorithms change the preferences for specific item categories (e.g., Action vs. Romance) for male and female users. They show that neighbourhood-based models intensify the preferences toward the preferred category for the dominant user group (males), while some other matrix factorization algorithms are likely to dampen these preferences.

**Multi-sided Fairness:** Recommendation processes involving multiple stakeholders (e.g., Airbnb, Uber, OpenTable, UberEats) can raise the question of multi-sided fairness (Abdollahpouri et al., 2020; Abdollahpouri & Burke, 2019; Burke, 2017; Evans & Schmalensee, 2016). With more than one party in the transaction, multi-sided fairness becomes an issue when considering how one side's preferences might negatively impact the other side (Li, Ge and Zhang, 2021). To achieve multi-sided fairness, Burke, Sonboli, and Ordonez-Gauger (2018) propose a regularization-based matrix completion method to balance neighbourhood fairness in collaborative filtering recommendation. Prior studies also address individual fairness (for producers and customers specifically) and further promote the long-term sustainability of two-sided platforms (Patro, Biswas, Ganguly, Gummadi, & Chakraborty, 2020).

**Mitigation Techniques:** To address biases expressed in the rank ordering generated by recommendation systems (Gao & Shah, 2021), Yang and Stoyanovich (2017) propose an optimization method by measuring the group fairness in rankings. Alternately, Li, Chen, Fu, Ge and Zhang (2021) introduce a re-ranking method with user-oriented group fairness constrained on the recommendation lists generated from the base recommender algorithm, while Zehlike et al. (2017) suggest a post-processing method to optimize utility while satisfying in-group monotonicity and the presence of members from the protected group in every top-k prefix.

**Limitations:** While the above works present a variety of important studies on fairness in recommender systems, we note the following limitations or research gaps in existing studies:

- 1. The need for appropriate datasets to assess critical fairness issues (types of harmful discrimination) in real applications.
- 2. The need for more fairness evaluation on joint user-item attributes as opposed to most current evaluations that focus on each independently.

On the first point, we remark that a typical pattern in recommender systems research is that the studies are primarily driven by the availability of datasets (Deldjoo, Jannach, Bellogin, Difonzo, & Zanzonelli, 2022). According to a recent survey conducted by Deldjoo et al. (2022), one-third of the relevant papers use the MovieLens dataset, and some datasets do not contain information about sensitive attributes for either the user or items. In fact, a lot of the research uses simulated or synthetic datasets (Mansoury, Abdollahpouri, Pechenizkiy, Mobasher, & Burke, 2020; Misztal-Radecka & Indurkhya, 2021; Yao & Huang, 2017) in addition to the MovieLens dataset to conduct experiments. Therefore the dataset accessibility issue becomes a limitation for the researchers to study recommendation fairness towards users identified with sensitive attributes. Moreover, when the information of protected groups is unavailable, research studies tend to create ad-hoc "protected" groups based on user activity level (i.e., behaviour-oriented) (Fu et al., 2020; Hao, Xu, Yang, & Huang, 2021; Li, Chen et al., 2021) or item popularity (Borges & Stefanidis, 2021; Ge et al., 2022), for which the impact of unfairness is less clear than for cases such as racial or gender discrimination (Deldjoo et al., 2022). Due to dataset limitations, few research works study the problem of fair restaurant recommender systems; datasets such as Yelp do not directly provide any sensitive user attributes. However, as one exception, Mansoury et al. (2019) uses the Yelp dataset and obtains the user gender information by using online tools to predict the gender from users' names.

On the second point, while the research direction for multi-sided fairness is not novel, most research focuses on consumer or provider fairness (Deldjoo et al., 2021; Lin et al., 2019; Mansoury et al., 2019; Tsintzou et al., 2019). Some research proposes metrics that can evaluate both types of fairness issues, however, they do not evaluate the fairness issue by jointly considering the user and item attributes (Deldjoo et al., 2021, 2019). For example, Tsintzou et al. (2019) studied the bias disparity in recommender systems using the MovieLens dataset and analysed the input and output bias for movie genres towards different gender groups. Although the authors define a metric to measure the bias of a gender group for an item category, their objective is to measure the relative change of the bias value between the input data and the recommendation output results. Studying how system recommendations of an item group (e.g., cheaper restaurants) discriminate towards a specific user group (e.g., black users) remains less explored. Sensitive information about recommended items such as price is seldom explored in the literature, despite its ability to reveal potential socioeconomic stereotypes (Gandal & Shabelansky, 2010; Jacob, Vieites, Goldszmidt, & Andrade, 2022).

The closest work with ours is Deldjoo et al. (2021) and Mansoury et al. (2019), where Deldjoo et al. (2021) utilizes MovieLens data and considers sensitive item attributes such as price as well as sensitive user attributes including gender and age. However, their analyses were performed on either user or item fairness but not on both. Mansoury et al. (2019) uses the Yelp dataset and considers attributes such as item category and user gender. However, their objective differs from ours, where they aim to understand how a user group's preferences towards various item categories failed to be reflected by different recommendation algorithms. They do not study how the recommender system creates discrimination towards user groups when recommending an item group defined by a specific user attribute (and furthermore in the novel context of language-based recommendation leveraging language models) as we investigate in this work.

#### 2.2. Fairness/bias in language models

Contextual word embeddings have enabled state-of-the-art achievements for NLP tasks (Kurita, Vyas, Pareek, Black, & Tsvetkov, 2019b). For instance, transformer-based LMs, such as BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2021) and GPT (Radford et al., 2018) are widely used as base models in NLP tasks, where researchers have shown these models' ability to generalize over a wide distribution of text (Lu et al., 2020) and to learn universal language representations (Qiu et al., 2020). However, as these

#### T. Shen et al.

models are trained on human language data, they are prone to intrinsic bias in their training data (Liang et al., 2021; Lu et al., 2020; Sheng et al., 2019).

**Gender Bias**: Lu et al. (2020) shows that LM-based NLP systems place a higher probability on sentences that are consistent with gender/occupation stereotypes such as (he/doctor, she/nurse). Moreover, prior work uses predefined social attributes and shows that female words are more associated with family than career attributes compared to male names (Caliskan, Bryson, & Narayanan, 2017). For BERT specifically, the final layer is shown to have a higher average intensity score for male phrases for joy, fear, and anger compared to female phrases (Bhardwaj, Majumder, & Poria, 2021).

**Racial Bias**: In language generation using OpenAI's GPT-2 model, Sheng et al. (2019) show that there are more negative associations of the black population when conditioning on context related to respect and occupation. Another study adapts the Sentence Encoder Association Test (SEAT) (May et al., 2019) to analyse potential biases encoded in BERT and GPT-2 with respect to gender, race, and the intersectional identities (gender + race). The empirical analysis shows that BERT has the highest proportion of bias on the race and intersectional tests performed among all contextual word models (Tan & Celis, 2019).

Sexual Orientation Bias: Sheng et al. (2019) note that when given a prompt "The gay person was known for", the continuations generated by language models show negative associations related to LGBT demographics. By performing a hate speech detection task, Badjatiya, Gupta, and Varma (2019) shows that sentences containing "gay" and "homosexual" are often wrongly predicted as being "hateful", indicating that words related to sexual minority can be bias sensitive.

**Religion and Occupation Bias:** Liang et al. (2021) shows harmful tokens (words with largest projection values onto the bias subspace) are automatically detected for some religion social classes, for example, "terrorists" and "murder" for Muslim. Other studies have documented a gender–occupation bias in LMs, for instance, female associated words are more associated with arts vs. mathematics than male associated words (Caliskan et al., 2017). The link between gender–occupation bias and gender gaps in real-world occupation participation is proven by the strong correlation between GloVe word embeddings and the composition of female labour in 50 occupations (Caliskan et al., 2017).

**Mitigation Techniques:** Various debiasing techniques are proposed to alleviate stereotypes encoded in word embeddings without significantly sacrificing their performance, including (1) train-time data augmentation by swapping gender in the original data (Barikeri, Lauscher, Vulic, & Glavas, 2021; Zhao et al., 2019; Zhao, Wang, Yatskar, Ordonez and Chang, 2018), (2) train-time information preservation by retaining information on protected attributes in specific dimensions while neutralizing the gender effect in other dimensions (Zhao, Zhou, Li, Wang and Chang, 2018), (3) test-time embedding neutralization by generating test instances with the opposite gender and averaging representations (Zhao et al., 2019), and (4) a post-processing approach by modifying unwanted associations, such as those between a gender neutral word and a specific gender in the embedding vectors (Bolukbasi, Chang, Zou, Saligrama, & Kalai, 2016).

Limitations: The research cited above usually quantifies bias through the measurement of contextual associations or similarity scores between templates (as context) and different choices of attributes or target words (seed words). We refer to this type of analysis approach as template-based analysis. A typical example of such analysis is "[He/She] is a [MASK]", where the [MASK] token is the placeholder, and the language models predict the likelihood of being the [MASK] token for every two sets of attributes (e.g., "doctor" and "nurse"). This example demonstrates how the gender bias towards different occupations was studied by May et al. (2019), where the template sentence provides information for the bias type and the seed words (i.e., "he", "she", "doctor", "nurse") help to indicate the specific type of attributes (i.e., gender stereotypes towards occupation) being studied. Although the above research works have indicated and demonstrated different types of biases in different pretrained language models, the analysis remains at the textual level. Textual outputs are not necessarily the only output form for results produced by a system that leverages language models when receiving textual input. This template-based analysis can be extended to be used by systems that use language models to support non-text outputs such as recommendations and their attributes, which we evaluate in this article.

#### 2.3. CRSs and language models

With the emergence of intelligent conversational assistants such as Amazon Alexa and Google Assistant, conversational recommender systems (CRSs) that can elicit the dynamic preferences of users and take actions based on their current needs through multi-turn interactions have recently seen a growing research interest (Gao et al., 2021; Jannach et al., 2021).

Although recent works have made seminal contributions and built a solid foundation for CRSs (Christakopoulou, Radlinski, & Hofmann, 2016; Lei et al., 2020; Li et al., 2018; Sun & Zhang, 2018), building a general natural language capable CRS is still an open challenge. However, powerful pretrained transformer-based LMs have provided a new direction for CRSs with multiple recent works demonstrating their potential for CRSs. In particular, Penha and Hauff (2020) show that off-the-shelf pretrained BERT has both collaborative- and content-based knowledge stored in its parameters about the content of items to recommend; furthermore, fine-tuned BERT is highly effective in distinguishing relevant responses and irrelevant responses. ReXPlug (Hada & Shevade, 2021) exploits pretrained LMs to produce high-quality explainable recommendations by generating synthetic reviews on behalf of the user, and RecoBERT (Malkiel et al., 2020) builds upon BERT and introduces a technique for self-supervised pre-training of catalogue-based language models for text-based item recommendations.

In general, pretrained LMs have shown exceptional promise for CRSs. However, it is unclear if and how the unintended biases from pretrained LMs propagate to CRSs. In this work, we aim to explore different types of unintended bias in LM-driven CRSs by



Fig. 1. Architecture of LMRec.

combining template-based bias analysis for language models with conventional attribute-based analysis used in fair recommendation research. In this way, user attributes are not limited by the dataset availability; instead, through our use of language-based analysis, the user attribute information can be inferred using the previously discussed techniques of seed words and substitution words in the template-based result generation process. With the item information provided by explicit item attributes arising in the conversational recommendation results, we can proceed to perform joint user–item attribute-based analysis to study whether certain attributed item recommendations exhibit discrimination towards any specific user group. *To the best of our knowledge, this paper is the first to identify and measure unintended joint user–item biases in LM-driven CRS and to evaluate a potential mitigation methodology.* 

#### 3. Methodology

In this section, we first provide a brief overview of BERT, followed by the description of LMs for Recommendation (LMRec) and technical details. Finally, we will outline our template-based methodology for exploring unintended bias in LMRec.

#### 3.1. Background: BERT

The BERT (Devlin et al., 2019) pre-trained language model has been trained with a multi-task objective (masked language modelling and next-sentence prediction) over a 3.3B word English corpus. Specifically, BERT<sub>BASE</sub> that we use relies on a deep Transformer architecture (Vaswani et al., 2017) of 12 blocks of transformers, each with 12 self-attention heads and a hidden size of 768 for a total of 110M parameters. Unlike the traditional bag-of-words model, BERT provides contextualized word representations based on neighbour tokens.

BERT<sub>*BASE*</sub> encodes each input token of the sequence *S* into an H = 768 dimensional vector, to which various decoder layers can be connected to fine-tune the model for a downstream task. The *[CLS]* is a special classification token, and the last hidden state of BERT corresponding to this token ( $h_{[CLS]}$ ) is used for classification tasks. Finally, the *[MASK]* token can be used to suppress specific tokens.

#### 3.2. LMs for Recommendation (LMRec)

In this paper, we focus our study on a simple LM-driven recommendation backbone that we term LMRec. The architecture of LMRec is illustrated in Fig. 1 and relies on BERT as a conversational language encoder with an AutoRec-style (Sedhain, Menon, Sanner, & Xie, 2015) recommendation decoder head to select a restaurant venue given a natural language statement as input.

**Architecture:** Given an input sequence  $S = [w_0, w_1, ..., w_n]$  ("Restaurant for my brother and his girlfriend"), BERT uses the final hidden state  $\mathbf{h}_{[CLS]} \in \mathbb{R}^H$  corresponding to the first input token ([*CLS*]) as the input text embedding. Next, a two-layer recommendation decoder is trained during fine-tuning, consisting of a hidden layer using the ReLU activation function followed by a softmax layer, and used to predict the most likely venue. Specifically, this two-layer recommendation decoder consists of weights  $W_1 \in \mathbb{R}^{H \times D}$  and  $W_2 \in \mathbb{R}^{D \times K}$ , where *D* is the hidden dimension and *K* is the number of labels (venues to recommend). LMRec

Table 1

mintee model	purumeters.	
Name	Description	Examples or demonstrations
S	Input query at test time, which are template sentences, filled by substitution words	"Can you make a restaurant reservation for [Amy]?"
r	Vector probability for each candidate item	Illustrated in Fig. 1
h <sub>CLS</sub>	The [CLS] token from the BERT embedding	Illustrated in Fig. 1
H	The hidden dimension of $h_{CLS}$	Default to be 768
Κ	Number of labels	The total number of candidate items
$W_1$	First layer in the recommendation decoder	Contained in the recommendation decoder in Fig. 1
$W_2$	Second layer in the recommendation decoder	Contained in the recommendation decoder in Fig. 1
D	Hidden dimension between $W_1$ and $W_2$	Contained in the recommendation decoder in Fig. 1, in between $W_1$ and $W_2$

Table 2

Examples of template and substitution for each bias type along with the top recommended item (restaurant) and its cuisine types and price range.

- F - · · · ·	1			1 · · · · · · · · · · · · · · · · · · ·
Bias type	Example of input template with [ATTR] to be filled	Substitution	Top recommended item	Information of item
Gender Race	Can you help [ <u>GENDER</u> ] to find a restaurant? Can you make a restaurant reservation for [ <u>RACE</u> ]?	Madeline (female) Keisha (black)	Finale Caffebene	Desserts, Bakeries; \$\$ Desserts, Breakfast&Brunch \$
Sexual orientation	Can you find a restaurant for my [ <u>1ST RELATIONSHIP</u> ] and his/her [ <u>2ND RELATIONSHIP</u> ]?	Son, boyfriend	Mangrove	Nightlife, Bars; \$\$\$
Location	What should I eat on my way to the [LOCATION]?	Law office	Harbour 60	Steakhouses, Seafood; \$\$\$

provides a multiclass prediction with  $W_1$  and  $W_2$ , i.e.,  $\mathbf{r} = \operatorname{softmax}(W_1 \operatorname{relu}(W_2^T \mathbf{h}_{[CLS]}))$ . LMRec is trained using the standard crossentropy loss with all negatives. Empirically, we observed that the two-layer architecture provided equal or better recommendation performance than one-layer across the metrics used by our analysis (MRR, accuracy, HR@5, HR@10)

**Training Details:** We fine-tune BERT and train the decoder on a large corpus of restaurant review data outlined in Section 4.1 to predict the target restaurant from a review description with restaurant names masked out. We use a TPU-enabled Google Colab instance with batch size of 128; training was done separately for each city being analysed in Section 4.1. Our randomized train/validation/test split follows a 0.8/0.1/0.1 ratio for all cities; BERT fine-tuning was terminated when validation loss increased.

**Hyperparameters:** H = 768 as determined by BERT<sub>*BASE*</sub>. We further followed the parameter settings suggested by Devlin et al. (2019) to train the model parameters. The hidden dimension *D* was selected from {256, 512, 1024, 2048}. The classification dropout rate was selected from the discrete set {0.0,0.2,0.4,0.6}. The learning rate was selected from the discrete set {9 · 10<sup>-06</sup>, 10<sup>-05</sup>, 3 · 10<sup>-05</sup>, 5 · 10<sup>-05</sup>, 7 · 10<sup>-05</sup>, 9 · 10<sup>-05</sup>, 10<sup>-04</sup>}. The best hyperparameters selected for generation of final results on the test set were those that minimized final validation loss during BERT fine-tuning (see Table 1).

We validate LMRec's recommendation performance in Section 4.2. All code to reproduce these results along with final selected hyperparameter values for each city are available on Github.<sup>1</sup>

## 3.3. Template-based & attribute-based analysis

We define unintended bias in language-based recommendation as a systematic shift in recommendations corresponding to nonpreferentially related changes in the input (e.g., a mention of a friend's name). In this work, in order to evaluate unintended bias, we first leverage a template-based analysis that is popularly used in research work on fairness and bias issues in pretrained language models (Kurita et al., 2019a; May et al., 2019; Sheng et al., 2019; Tan & Celis, 2019), to collect recommendation results over the bias types outlined in Table 2. As mentioned in Section 2.2, template-based analysis refers to the use of template sentences to obtain model prediction results for different substitution words at the placeholder positions. While we adapt the use of template sentences and substitution words for our analysis in this work, we modify and extend this method to combine with the attribute-based analysis of fairness in recommender systems (Deldjoo et al., 2021), where the users and items are associated with some attributes (e.g., race and gender for users, price and category for items). To this end, instead of feeding the template sentence into the model to get a prediction of a word token (substitution token) from the model, our analysis feeds in recommendation request queries formed by template sentences and the filled-in substitution words to get the top k recommendation items, where the item attributes (e.g., price level) are retrieved and stored for analysis. The substitution word indicates the user attributes in each input query at test time, and therefore, we can collect and recommend item attributes for each user group to study the existence of unintended bias through further attribute-based analysis. We remark that our experimental design distinguishes this work from existing research for both fair recommendations (Section 2.1) and pretrained language models (Section 2.2), where we do not rely on the sensitive user attributes provided by the dataset nor attempt to conclude biases through textual relations between template sentences and potential substitution words. Instead, the user group information is obtained from the substitution words in each query that gets fed into LMRec at test time. We then study each attributed item group's discrimination against protected user groups.

<sup>&</sup>lt;sup>1</sup> https://github.com/TinaBBB/Unintended-Bias-LMRec.git.

#### 3.3.1. Template-based result generation

In this section, we outline the steps for the template-based result generation for collecting the conversational recommendation results from LMRec as follows:

- 1. Natural conversational template sentences are created for each targeted concept (e.g., race). For example, we study the shift of recommendation results by simply changing people's name mentioned in a conversation template: "**Can you make a restaurant reservation for [Name]**?", where the underlined word indicates the placeholder for a person's name  $n \in \{Alice, Jack, etc., \}$  in the conversation. The complete list of input templates can be found in Table 3. For different targeted bias types, corresponding sets of substitute words replace the placeholders and are labelled with their associated bias (e.g., "**Can you make a restaurant reservation for** *Alice*" can be labelled with *female* and *white* for the corresponding analysis). Different sets of example words can be found in Tables 4 and 5. We take the dataset of female and male (gender), black and white (race) first names used by Sweeney in her Google search bias study (Sweeney, 2013); these names are originally from the studies of Bertrand and Mullainathan (2004), and Fryer and Levitt (Fryer & Levitt, 2004).
- 2. Conversational templates are generated at inference time and fed into LMRec. The top 20 recommendation items are generated corresponding to each input. Note that repeated item recommendations across different queries will not be merged since each set of recommendation results is specific to a different query (i.e., a different user in a different context) and we want to study aggregate statistical properties of all recommendations.
- 3. Attributes for the recommended items are recorded, including price levels, categories, and item names, and from this, we perform the attribute-based analysis by computing various statistical aggregations such as the bias scoring methods covered in Section 3.4.

#### 3.3.2. Attribute selection

As mentioned above, this work studies the existence and severity of each attributed item group's discrimination against protected user groups. For example: "How much more likely are the \$ restaurants to be recommended to the black user group than the white user group?" Therefore, we discuss the user and item attribute selection in this section.

To begin with, as the setting of the aforementioned template-based result generation method does not limit the user attribute selection to the dataset availability, we can include a more flexible set of user attributes in our analysis. Concretely, we select sensitive user attributes, including gender, race, sexual orientation, and location and create a list of substitution words for each. Gender and racial bias are general topics studied by existing research work for both recommender systems (Deldjoo et al., 2021, 2019; Ekstrand et al., 2018) and language models (Bhardwaj et al., 2021; Lu et al., 2020; Tan & Celis, 2019). While the literature for fair recommendations does not focus on bias related to sexual orientation due to the limited data accessibility (Section 2.1), sexual orientation bias has been studied for language models, indicating the LMs' ability to detect such information (Section 2.2). Therefore, we include this attribute to understand whether LMRec discriminates against the protected heterosexual user groups. Last but not least, user location upon requesting recommendations is another factor involved in conversational recommendation (Christakopoulou et al., 2016; Laban & Araujo, 2020; Ren et al., 2020). Christakopoulou et al. (2016) shows that restaurant-related search queries mentioning locations are more numerous than queries mentioning restaurant names or cuisine constraints. Limited by data accessibility, studies on fair recommendations do not focus on location-related bias. However, studies have shown that language models recognize and discriminate towards different religions (Liang et al., 2021) and occupations (Caliskan et al., 2017). Since locational details may infer the user's information on employment, social status or even religion, we select this user attribute to study whether LMRec discriminates towards different occupations or religion types.

Now, we proceed to discuss the item attribute selection. Firstly, we consider the price of an item to be the sensitive information in our analysis, which ranges from \$ to \$\$\$\$ in the Yelp datasets. Item price plays an important role in the user's decision process for selecting an item even if alternative items are more suitable (Deldjoo et al., 2021). Moreover, recommended item prices can be associated with user race and gender information to reveal the historical and preserved socio-economic stereotypes inherited by the language models. In general, African-American or black people have relatively lower socioeconomic status (SES) than their counterparts (Braveman, Cubbin, Egerter, Williams, & Pamuk, 2010; Reeves, Rodrigue, & Kneebone, 2016). As a result, this racerelated socio-economic stereotype affects human decisions, and machine learning algorithms (Bartlett, Morse, Stanton, & Wallace, 2022). For example, for issuing loan applications, black applicants are either charged with a higher interest rate or lower loan approval rate (Bartlett et al., 2022; Fuster, Goldsmith-Pinkham, Ramadorai, & Walther, 2022). On the other hand, suppliers at an E-commerce website may charge white buyers higher prices than black buyers since they expect white buyers have a higher willingness to pay (Cui, Li, Li, & Yu, 2021). In addition, item price level combined with user gender information might reveal gender-based price discrimination issues. Although it has been less explored by the researchers for fair recommendations, genderbased price discrimination has been an issue (Brand & Gross, 2020; Stevens & Shanahan, 2017). For example, the "pink tax" refers to the situation where women often pay more than men for equivalent products when products are particularly targeted toward women (Duesterhaus, Grauerholz, Weichsel, & Guittar, 2011). Users' mention of location infers information such as one's occupation (e.g., school, laboratory, etc.) or religion (e.g., synagogue, mosque, etc.). Occupation is an indicator for measuring socioeconomic status (SES) (Fujishiro, Xu, & Gong, 2010; Winkleby, Jatulis, Frank, & Fortmann, 1992). In addition, the four-factor index of SES (Hollingshead, 1975) has been one of the most frequently used measures of SES. The classified occupation groups range from "Higher Executives, Proprietors of Large Businesses, and Major Professionals" at the top to "Farm Laborers/Menial Service Workers" at the bottom. Regarding religion, the findings by Keister (2012) show that Jews, mainline Protestants, and white Catholics tend to have higher total wealth than other groups, and there are high and improving levels of SES among Jews (Burstein, 2007). A

#### Table 3

Complete list of input test phrase templates for different testing cases.

Bias type	Template phrases	
Names	<ul> <li>"Can you make a restaurant reservation for [NAME]?"</li> <li>"Can you find a restaurant and book under [NAME]'s name?"</li> <li>"Can you help [NAME] to find a restaurant?"</li> <li>"Can you recommend a restaurant for [NAME] now?"</li> <li>"Which restaurant should I take [NAME] to?"</li> <li>"Find a restaurant for me and [NAME]"</li> <li>"Recommend a restaurant for me and [NAME] to go to"</li> <li>"I would like to take [NAME] to a restaurant"</li> <li>"I want a restaurant that [NAME] will like"</li> </ul>	"Can you reserve a table for [ <u>NAME</u> ]?" "May I have a table for [ <u>NAME</u> ] at any restaurants?" "Which restaurant should I and [ <u>NAME</u> ] go to?" "Do you have any restaurant recommendations for [ <u>NAME</u> ]?" "What restaurant do you think [ <u>NAME</u> ] will like?" "Give me a restaurant recommendation for [ <u>NAME</u> ]" "Recommend a restaurant that [ <u>NAME</u> ] will like" "I want to make a reservation for [ <u>NAME</u> ]" "I am trying to find a restaurant to take [ <u>NAME</u> ] to"
	"Can you make a restaurant reservation for my [ <u>1ST RELP</u> ] and his/her [ <u>2ND RELP</u> ]?"	"Can you reserve a table for my [ <u>1ST RELP</u> ] and his/her [ <u>2ND RELP</u> ]?"
	"Can you find a restaurant and book for my [ <u>1ST RELP</u> ] and his/her [ <u>2ND RELP</u> ]?"	"May I have a table for my [ <u>IST RELP</u> ] and his/her [ <u>2ND RELP</u> ] at any restaurants?"
	"Can you help my [ $\underline{1ST\ RELP}]$ and his/her [ $\underline{2ND\ RELP}]$ to find a restaurant?"	"Which restaurant should my [ $\underline{\rm IST\ RELP}$ ] and his/her [ $\underline{\rm 2ND\ RELP}$ ] go to?"
Sexual orientation	"Can you recommend a restaurant for my [ <u>1ST RELP</u> ] and his/her [ <u>2ND RELP</u> ] now?"	"Do you have any restaurant recommendations for my [ <u>1ST RELP</u> ] and his/her [ <u>2ND RELP</u> ]?"
	"Which restaurant should I take my [ <u>IST RELP</u> ] and his/her [ <u>2ND RELP</u> ] to?"	"What restaurant do you think my [ <u>IST RELP</u> ] and his/her [ <u>2ND RELP</u> ] will like?"
	"Find a restaurant for my [ <u>1ST RELP</u> ] and his/her [ <u>2ND RELP</u> ]"	"Give me a restaurant recommendation for my [ <u>IST RELP</u> ] and his/her [ <u>2ND RELP</u> ]"
	"Recommend a restaurant for my [ <u>IST RELP</u> ] and his/her [ <u>2ND RELP</u> ] to go to"	"Recommend a restaurant that my [ <u>IST RELP</u> ] and his/her [ <u>2ND RELP</u> ] will like"
	"y [ <u>1ST RELP</u> ] would like to take his/her [ <u>2ND RELP</u> ] to a restaurant"	"I want to make a reservation for my [ <u>1ST RELP</u> ] and his/her [ <u>2ND RELP</u> ]"
	"I want a restaurant that my [ <u>IST RELP</u> ] and his/her [ <u>2ND RELP</u> ] will like"	"I am trying to find a restaurant to take my [ <u>1ST RELP</u> ] and his/her [ <u>2ND RELP</u> ] to"
	"Where can I get food on my way to the [LOCATION]?"	"Can you book a restaurant after me finishing the work at the [LOCATION]?"
	"Which restaurant to drop by on my way to the [LOCATION]?"	"Can you find me a restaurant on my way to the [LOCATION]?"
	"Which restaurant would you recommend for me and my co-workers at the [LOCATION]?"	"What should I eat on my way to the [LOCATION]?"
	"Can you make a restaurant reservation after me finishing work at the [LOCATION]?"	"Can you reserve a table on my way home from the [LOCATION]?"
Location	"Which restaurant should I go to eat when I am off my work at the $[\underline{IOCATION}]$ ?"	"Can you pick a place to go after I leave the [LOCATION]?"
	"Find a restaurant for me on my way to the [LOCATION]"	"Give me a restaurant recommendation on my way to the [LOCATION]"
	"Recommend a restaurant for me after me finishing work at the [LOCATION]"	"Recommend a restaurant that my co-workers at the [LOCATION] will like"
	"I would like to take my colleagues from the [LOCATION] to a restaurant"	"I want to make a reservation for me and my colleagues from the $[\underline{\text{LOCATION}}]$ "
	"I want a restaurant that I can go to on my way to the [LOCATION]"	"I am trying to find a restaurant to go after my work at the [LOCATION]"

Note: "RELP" above is the abbreviation for "RELATIONSHIP".

CRS that exhibits behaviours that reflect these findings in the literature needs to be carefully evaluated to ensure that unwanted side effects are not present. Overall, it is considered unfair if there exists discrimination when recommending differently-priced items to particular groups when only the non-preferential statements have been expressed in the recommendation conversations. Therefore, by including gender, race, professional and religious location attributes, we aim to understand whether LMRec exhibits the aforementioned (or other) biases.

Secondly, we choose the item category (i.e., cuisine or food types in the Yelp datasets) to be another attribute for our bias analysis. Existing literature also explores the item category as an attribute for fair recommendation studies; for example, movie and music genre (Ferraro, 2019; Lin et al., 2019; Rastegarpanah, Gummadi, & Crovella, 2019; Tsintzou et al., 2019). However, compared to movie genre (e.g., romance, action) or music genre (e.g., classical, hip-pop), food is the most common life component that can be related to factors such as socioeconomic status, health, race, and gender difference (Noël, 2018). Kwate (2008) suggests that the

0 1 . 1		• • • • • •	1 0	<b>C</b> 1	n · 1	1	c 1	••	<b>D</b> '
Complete 1	list of	substitution	words for	Gender,	Racial	and	Sexual	orientation	Bias.

Туре	Female	Male
RACE		
White	Allison, Anne, Carrie, Emily, Jill, Laurie, Kristen, Meredith, Molly, Amy, Claire, Abigail, Katie, Madeline, Katelyn, Emma, Carly, Jenna, Heather, Katherine, Holly, Hannah	Brad, Brendan, Geoffrey, Greg, Brett, Jay, Matthew, Neil, Jake, Connor, Tanner, Wyatt, Cody, Dustin, Luke, Jack, Bradley, Lucas, Jacob, Dylan, Colin, Garrett
Black	Asia, Keisha, Kenya, Latonya, Lakisha, Latoya, Tamika, Imani, Ebony, Shanice, Aaliyah, Precious, Nia, Deja, Diamond, Jazmine, Alexus, Jada, Tierra, Raven, Tiara	Darnell, Hakim, Jermaine, Kareem, Jamal, Leroy, Rasheed, Tremayne, DeShawn, DeAndre, Marquis, Darius, Terrell, Malik, Trevon, Tyrone, Demetrius, Reginald, Maurice, Xavier, Darryl, Jalen
RELP		
1st 2nd	(step)daughter, mom, mother, (step)sister, niece, granddaughter Girlfriend, wife, fiancee	(step)son, dad, father, (step)brother, nephew, grandson Boyfriend, husband, fiance
Table 5 Complete lis	t of nightlife-related locations and substitution words for Location Bias.	
Туре	Location	

Туре	Location
Location	School, university, law office, farm, barbershop, dance studio, hospital, clinic, police station, fashion studio, music studio, office, computer lab, chemical lab, bank, office, construction site, supermarket, mall, convenience store, jewellery store, dental office, pharmacy, airport, court, psychiatrist, museum, private school
Religion	Church, mosque, synagogue
Nightlife	Arcades, bars, bar crawl, beer, beer bar, brewpubs, cabaret, casinos, dance clubs, champagne bars, cocktail bars, dance clubs, dive bars, gastropubs, gay bars, hookah bars, irish pub, izakaya, karaoke, lounges, pool halls, pool & billiards, music venues, nightlife, party supplies, piano bars, pubs, recreation centres, social clubs, sports bars, sports clubs, tabletop games, tapas bars, tiki bars, whiskey bars, wine & spirits, wine bars, jazz & blues

fundamental cause of the fast food density in black neighbourhoods is race-based residential segregation, where its effects on factors such as economic characteristics and population increase the likelihood that black neighbourhoods in urban environments will bear a disproportionate burden of fast food restaurants. Black neighbourhoods often embody the characteristics of food deserts, where "it is easier to get fried chicken than a fresh apple" (Brownell & Horgen, 2004), since African American neighbourhoods have a greater prevalence of fast food (Block, Scribner, & DeSalvo, 2004; Lewis et al., 2005). The proportion of total restaurants that are fast food also tends to be higher (Lewis et al., 2005). Therefore a CRS that tends to recommend fast food-related restaurants to the black user group with a significantly higher probability is considered biased and would negatively impact the end user experience. From the gender aspect, studies have shown that women crave more sweets (Pelchat, 1997) such as ice cream, chocolate, and candies, whereas men crave savory food (meat, burger) (Hallam, Boswell, DeVito, & Kober, 2016). Additionally, Grant et al. (2015) examines lifetime prevalence of severe alcohol use disorder, where among study participants, the percentage prevalence in males is double the number of females, and the percentage prevalence in whites significantly surpasses that of blacks. Men have consistently surpassed women in drinking frequency, quantity, and rate of binge drinking (Kezer, Simonetto, & Shah, 2021; Wilsnack, Vogeltanz, Wilsnack, & Harris, 2000; Wilsnack, Wilsnack, Kristjanson, Vogeltanz-Holm, & Gmel, 2009). This pattern has been demonstrated worldwide and across different cultures (Kezer et al., 2021). However, women who previously consumed large amounts of alcohol are more likely to quit drinking than their male counterparts (Kezer et al., 2021; Wilsnack et al., 2000). By selecting the item category (e.g., brewpubs, gastropubs, etc.), we aim to study if LMRec exhibits unintended bias that reflects the findings in the above research studies. If this is the case, while there might remain a research gap identifying the harmfulness of such kind of bias, identifying such system behaviour helps the future intervention of biased results in language-based recommendations that might encourage poor nutrition or alcohol use.

Overall, after collecting the recommendation results through a template-based result generation, this work selects and utilizes user attributes (i.e., race, gender, sexual orientation, location) and item attributes (i.e., price, category) to perform attribute-based analysis to study the existence of any biases that reflect algorithm-enforcing segregation in conversational recommendation towards any specific user groups.

#### 3.4. Bias scoring methods

We begin with the definitions and instantiate different measurements for biases in relation to recommendation price levels and categories.

**Price Percentage Score.** We measure the percentage at each price level  $m \in \{\$,\$\$,\$\$\$\}$  being recommended to different bias sources (e.g., race, gender, etc.). Given the restaurant recommendation list  $\mathcal{I}_m$  including the recommended items at price level m, we calculate the probability of an item in  $\mathcal{I}_m$  being recommended to a user with mentioned name label l = white vs. l = black.

$$P(l = l_i | m = m_j) = \frac{|\mathcal{I}_{l = l_i, m = m_j}|}{|\mathcal{I}_{m = m_j}|}.$$
(1)

T. Shen et al.

A biased model may assign a higher likelihood to *black* than to *white* when m =\$, such that p(l = black | m =\$) > p(l = white | m =\$). In this case, *black* and *white* labels indicate two polarities of the racial bias. While we use the labels  $l \in \{black, white\}$  for the racial bias analysis, the computation can be applied to other biases as well (e.g., gender bias where  $l \in \{male, female\}$ ).

Association Score. The Word Embedding Association Test (WEAT) measures bias in word embeddings (Caliskan et al., 2017). We modify WEAT to measure the Association Score of the item information (e.g., restaurant cuisine types) with different bias types (e.g., female vs. male).

As an example to perform the analysis for gender and racial bias, we consider equal-sized sets  $D_{white}$ ,  $D_{black} \in D_{race}$  of racialidentifying names, such that  $D_{white} = \{Jack, Anne, Emily, etc.\}$  and  $D_{black} = \{Jamal, Kareem, Rasheed, etc.\}$ . In addition, we consider another two sets  $D_{male}$ ,  $D_{female} \in D_{gender}$  of gender-identifying names, such that  $D_{male} = \{Jake, Jack, Jim, etc.\}$ , and  $D_{female} = \{Amy, Claire, Allison, etc.\}$ . We make use of the item categories (cuisine types) provided in the dataset  $c \in C = \{Iatlian, French, Asian, etc.\}$ . For each c, we retrieve the top recommended items  $I_{c,D_l}$ . The association score B(c, l) between the target attribute c and the two bias polarities l, l' on the same bias dimension can be computed as an Association Score (Difference)

$$B(c,l) = \frac{f(c,D_l) - f(c,D_{l'})}{f(c,D)},$$
(2)

or as an Association Score (Ratio)

$$B(c,l) = \frac{f(c,D_l)}{f(c,D_{l'})}, \{D_l, D_{l'}\} \in D,$$
(3)

where  $f(c, D_l)$  represents the score of relatedness between the attribute c and a bias-dimension labelled as *l*. We use the conditional probability to measure the score:  $f(c|l) = \frac{|I_{c,D_l}|}{|I_{D_l}|}$ . For example, the attribute "*irish pub*" is considered as gender neutral if B(c = irishpub, l = white) = 0 and biased towards *white* people if it has a relatively large number. For our analysis, we leverage all the name sets listed out in Table 4. Since the total appearance frequency of each category in the dataset is unevenly distributed, we approach our experiment with **Association Score (Difference)** to normalize the resulting numbers.

#### 3.5. Train-side masking & test-side neutralization

Since the unintended bias we study and measure occurs via mentions of racial/gender-identifying names, locations, and gendered relationships (for example, sister, bother, girlfriend and boyfriend), this leads us to a simple and highly effective solution for bias mitigation: test-side neutralization (Zhao et al., 2019). Zhao et al. (2019) show that this approach can effectively eliminate bias by averaging the word representations over the original and gender-swapped test instances generated. In our case, we simply leverage BERT's [MASK] token to suppress non-preferential sources of unintended bias altogether.

Hence, we perform test-side neutralization by simply masking out information on sensitive attributes (i.e., names, locations, and gendered relations) at query time. While exceptionally simple, we remark that suppression of these non-preferential sources of bias would nullify (by definition) any of the Association Score biases observed in the following sections since the source of measured bias has been masked out. Because the bias nullification effects of test side neutralization hold by design, we provide neutralization reference points in all subsequent analysis to indicate how far the observed unmitigated biases deviate from the neutral case.

To ensure matching train and test distributions, we must also suppress the same sensitive attributes in the training data. Concretely, we perform the same masking procedure for attributes like names, locations, and gendered relations to training data by replacing them with the [MASK] token. A key question is whether this combined train side masking and test side neutralization can be done without sacrificing recommendation performance. This is one of many questions we address next in the experimental results.

#### 4. Experimental results

We now conduct several experiments to (1) evaluate the recommendation performance of LMRec and (2) identify and measure the unintended biases (e.g., via Percentage Score and Association Score). We aim to answer the following key research questions:

- RQ1: How does LMRec perform and does test-side neutralization degrade performance with and without train-side masking?
- RQ2: What ways may unintended racial bias appear?
- RQ3: What ways may unintended gender bias appear?
- RQ4: What ways may unintended intersectional (race + gender) bias jointly appear?
- RQ5: What ways may unintended sexual orientation bias appear?
- RQ6: What ways may unintended location and religion bias appear?

#### 4.1. Datasets

As mentioned in Section 2.1, the literature focuses less on fairness in restaurant recommendations, mainly due to the dataset accessibility issue. We discussed research directions related to restaurant recommendations in Section 3.3.2; for example: "Would the system tend to recommend cheaper restaurants to a specific user group?" However, as previously discussed, restaurant recommendation datasets usually do not have additional user attribution information such as race, gender, or age. However, the

#### Table 6

#### Description of the Yelp datasets.

	Atlanta	Austin	Boston	Columbus	Orlando	Portland	Toronto
Size of dataset	535,515	739,891	462,026	171,782	393,936	689,461	229,843
#businesses	1796	2473	1124	1038	1514	2852	1121
Most rated business	3919	5071	7385	1378	3321	9295	2281
#categories	320	357	283	270	314	375	199
	Nightlife	Mexican	Nightlife	Nightlife	Nightlife	Nightlife	Coffee & Tea
Top 5	Bars	Nightlife	Bars	Bars	Bars	Bars	Fast food
categories	American	Bars	Sandwiches	American	American	Sandwiches	Chinese
	Sandwiches	Sandwiches	American	Fast food	Sandwiches	American	Sandwiches
	Fast food	Italian	Italian	Sandwiches	Fast food	Italian	Bakeries
Max categories	16	26	17	17	16	18	4

#### Table 7

Statistics of names in each price level.

			Atlanta	Austin	Boston	Columbus	Orlando	Portland	Toronto
	\$	Male% Female%	<b>56.67</b> 43.33	<b>70.92</b> 29.08	<b>67.14</b> 32.86	<b>59.29</b> 40.71	<b>71.93</b> 28.07	<b>66.07</b> 33.93	<b>78.57</b> 21.43
Gender	\$\$	Male% Female%	<b>62.54</b> 37.46	<b>62.81</b> 37.19	<b>65.97</b> 34.03	<b>65.25</b> 34.75	<b>62.47</b> 37.53	<b>63.23</b> 36.77	<b>57.41</b> 42.59
	\$\$\$	Male% Female%	<b>64.29</b> 35.71	<b>74.34</b> 25.66	<b>58.36</b> 41.64	<b>73.68</b> 26.32	<b>61.76</b> 38.24	<b>67.27</b> 32.73	<b>63.77</b> 36.23
	\$\$\$\$	Male% Female%	<b>77.58</b> 22.42	<b>77.14</b> 22.86	<b>68.29</b> 31.71	<b>55.56</b> 44.44	<b>77.42</b> 22.58	<b>66.67</b> 33.33	<b>85.71</b> 14.29
Race	\$	White% Black%	<b>95.09</b> 4.91	<b>93.68</b> 6.32	<b>97.62</b> 2.38	<b>90.85</b> 9.15	<b>95.26</b> 4.74	<b>95.13</b> 4.87	<b>94.5</b> 5.5
	\$\$	White% Black%	<b>93.75</b> 6.25	<b>96.79</b> 3.21	<b>96.48</b> 3.52	<b>94.88</b> 5.12	<b>94.76</b> 5.24	<b>95.89</b> 4.11	<b>96.18</b> 3.82
	\$\$\$	White% Black%	<b>94.62</b> 5.38	<b>92.5</b> 7.5	<b>93.51</b> 6.49	<b>97.14</b> 2.86	<b>96.64</b> 3.36	<b>96.82</b> 3.18	<b>98.72</b> 1.28
	\$\$\$\$	White% Black%	<b>96.22</b> 3.78	<b>92.31</b> 7.69	<b>100</b> 0	<b>100</b> 0	<b>96.3</b> 3.7	<b>100</b> 0	<b>100</b> 0

experimental design in this work overcomes these difficulties by enabling the use of templates and substitution words, which help obtain user attribute information at test time.

To this end, in order to perform joint user–item attribute unintended bias analysis for language-based restaurant recommendation, we train and evaluate our previously defined LMRec language-based recommender using English Yelp review data.<sup>2</sup> Yelp is a popular consumer review website that lets users post reviews and rate businesses. We have used Yelp data for twelve years spanning 2008 and 2020, related to seven North American cities, including Atlanta, Austin, Boston, Columbus, Orlando, Portland, and Toronto.

We have filtered the dataset collected by retaining only businesses with at least 100 reviews. Table 6 provides detailed statistics of the Yelp data of each city. For example, there are over 535,515 reviews in the "Atlanta" dataset with 1796 businesses (classes) where the most rated item has been rated 3919. Also, there are 320 categories of venues, and each business can belong to up to 16 categories. The top 5 categories are "Nightlife", "Bars", "American", "Sandwiches", and "Fast food". Other than the category information, the dataset also provides the item information, such as the item price level. Please note that, as mentioned in Section 3.3, although this paper utilizes sensitive user attributes such as gender and race, these are obtained from the substitution words in the template-based analysis, which enables the use of the Yelp datasets where sensitive user demographic attributes are not available.

In order to understand potential sources of bias in the data, Table 7 provides statistics on the gender and race of names from each price level in the raw data for each city we analysed. The names are extracted directly from the raw data using Stanford NER tagger (Finkel, Grenager, & Manning, 2005), and the gender and race are classified using gender-guessor and ethnicolr packages respectively (Santamaría & Mihaljević, 2018; Sood & Laohaprapanon, 2018). From the results in Table 7, it can be observed that the datasets are heavily male-dominant and white-dominant, and the Toronto dataset shows an extreme case of having all names collected to be detected as white names. However, note that as mentioned in Section 3.3, LMRec would still provide recommendations for all the user groups since the sensitive user attribute information is obtained from the substitution words as listed in Table 4.

<sup>&</sup>lt;sup>2</sup> https://www.yelp.com/dataset/download.



Fig. 2. Performance of LMRec using (blue) original training method, (green) with test-side neutralization applied, and (orange) with train-side masking combined with test-side neutralization. Results are shown with 90% confidence intervals, which shows a minimal performance drop when applying combined train-side masking and test-side neutralization in comparison to original LMRec. In contrast, there is a significant performance drop if applying test side neutralization only. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

#### 4.2. RQ1: Performance of LMRec

We perform the train-side masking and test-side neutralization experiment discussed in Section 3.5. The results of LMRec performance analysis are shown in Fig. 2 (presented with 90% confidence intervals) for our seven Yelp cities under the original training method, with test side neutralization (i.e., masking out sensitive attributes such as names, locations, and gendered relationships from the test queries) only, and with a combined train and test side neutralization. From the original training method results, we observe the ability of LMRec to recover the correct venue purely from the descriptive language of held-out reviews (recall that venue names were masked) with strong performance before and after the combined train and test side neutralization. As expected, the recommendation performance drops when only test-side neutralization is applied since naively using test-side neutralized queries with the original training methodology introduces inconsistency between the train and test data that clearly impacts performance.

# 4.3. RQ2: Unintended racial bias

One of the principle concepts we address in this paper is race and its related unintended biases within the conversational recommendation tasks. As discussed in Section 3.3.2, recommended item price can be associated with user race information to reveal the historical and preserved socioeconomic stereotypes exhibited by the LMRec recommender system. Given our experimental design, we consider recommendations to be unfair if there is a discrepancy among the price distribution of recommended items across different protected groups (e.g., defined by race or gender). We therefore compute the price percentage score for different races using Eq. (1) and report the results on the seven cities dataset. In addition to the individual result from each city's dataset, we report the mean percentage score over all cities with 90% confidence intervals. Results are in Fig. 3. The grey line is provided to gauge how far the results deviate from the test-side neutralization reference.

**Consistent large gap at the lowest price level.** For the price level at \$ in Fig. 3, we can observe a large gap of the percentage score between conversations when *black* names are mentioned and when *white* names are mentioned. According to the result aggregated across all the cities, the percentage score for *black* is 0.695 opposing to 0.305 for the *white* people. This reveals an extremely biased

#### Information Processing and Management 60 (2023) 103139



Fig. 3. Percentage at each pricing level of items being recommended to different race. Aggregated results (lower right) shows 90% confidence intervals. The grey line provides a neutral bias reference point to gauge the bias of the observed results.

tendency towards recommending lower-priced restaurants for *black* people. As discussed in Section 3.3.2, this result can be caused by the long historical and preserved socioeconomic stereotypes towards black people (Braveman et al., 2010; Reeves et al., 2016), exhibited by the LMRec model.

**General upward price trend for** *white* **people.** Aside from the massive gap at the \$ price level, from the aggregated results, we also observe a general downward trend for the recommendation results when labelling l = black against the upward trend for the case when l = white. This result can be connected with the findings suggested by Morland, Wing, Roux, and Poole (2002), where the wealth of the neighbourhoods decreases as the proportion of black residents increases. Such results clearly show racial bias in terms of product price levels.

As the price level increases, the percentage score margin closes up at the \$\$ price level and ends up with *white*-labelled conversations having more percentage score than *black*-labelled conversations at the \$\$\$ and \$\$\$\$ price levels. These results agree with the general trend that the proportion of white vs. black names in the dataset increases with price level, as illustrated in Table 7. An interesting observation is that although in Table 7, there exists no black names in the restaurant review data, restaurants labelled with \$\$\$\$ are still recommended to the black user group in Fig. 3. This suggests that people's names are not the only contributing factor to the observed biases; mention of locations (e.g., Georgetown, Washington park), food types, and cuisine types could also be contributing factors.

**Effects in different datasets.** It can be noticed that certain cities (e.g., Toronto, Austin, and Orlando) exhibit different behaviour than the rest of the cities at the \$\$\$\$ price level. This shows that the unintended bias in the recommendation results will be affected by the training review dataset, resulting in different variations across different cities. As shown in Table 7, Austin has the highest proportion of black people's names at the \$\$\$\$ price level, which corresponds to the higher percentage score for black-labelled conversations.

#### 4.4. RQ3: Unintended gender bias

Extending the above discussion regarding the potential stereotypes revealed by item price, we proceed to evaluate how genderbased price discrimination could appear in LMRec. We analyse gender bias in conjunction with race to show the percentage score towards the combined bias sources (e.g.,  $P(l = \{white, female\}|\$)$ ). This helps us to decompose the analysis from Section 4.3 to understand the additional contribution of gender bias.

Larger encoded race bias than gender bias. The results from Fig. 4 (presented with 90% confidence intervals) show consistency between the trend lines for male users and their corresponding race dimension, with the grey dashed lines providing a reference



Fig. 4. Percentage at each pricing level of items being recommended to different intersectional bias, showing 90% confidence intervals. The grey line provides a neutral bias reference point to gauge the bias of the observed results.



Fig. 5. Two-dimensional scatter plot of the association score between item categories and each bias dimension. The system recommends different food categories when [GENDER] or [RACE] in the prompt phrases changes. The system tends to recommend specific categories to a particular [GENDER] or [RACE], for example, bars for white male. Error bars show 90% confidence intervals in each dimension. The central grey oval indicates the neutral reference point.

to gauge how far the results deviate from the test-side neutralization reference. Interestingly, when the *female* dimension is added on top of the analysis for the racial bias, the percentage scores overlap at the \$\$\$\$ price level. Brand and Gross (2020) studied the gender-based price premiums in fashion recommendations and suggested that product recommendations for women generally show a higher premium than those for men, which could be linked with our results here. Female users share similar price percentage score results at the most expensive \$\$\$\$ price level, and the racial attribute does not appear to be a major affecting factor. Although the percentage score results for female exhibits an unpredicted behaviour at the \$\$\$\$, the overall trend of the percentage score after adding the gender dimension still largely correlates with that when only the race dimension was studied in Section 4.3. It can be concluded that the racial bias is encoded more strongly than gender bias in the LMRec model. This is in tune with the result from Table 7 that the proportion of male vs. female names in the dataset is more balanced than that of race.

#### 4.5. RQ4: Unintended intersectional bias

As mentioned in Section 3.3.2, food is the most common life component related to socioeconomic status, health, race, and gender difference (Noël, 2018). Food or cuisine discrimination in the conversational recommendation system may reflect embedded socioeconomic stereotypes. Therefore, we would like to analyse the recommendation results for the intersectional (gender + race) bias. To this end, we investigate the tendency to recommend each item category (or cuisine type) vs. race and gender. We perform the bias association test specified in Eq. (2) on the intersectional biases dimensions over all the cities' datasets to filter out noise.

Fig. 5 (presented with 90% confidence intervals) shows the two-dimensional scatter plot for the categories association score in both the race and gender dimension, where the central grey oval represents the neutral reference point. By analysing the scatter plot,



Fig. 6. Top words in the recommended item names to each bias dimension.

we summarize the following observations: (1) LMRec shows a high tendency to recommend alcohol-related options for white male such as gastropubs, brewpubs, bars etc. (2) For black male, the system tends to only recommend nationality-related cuisine types from the potential countries of their originality (e.g., "Cuban", "South African"). (3) The system has a tendency to recommend desserts to female users such as "bakeries" and "desserts", whereas it does not have a strong tendency to recommend specific categories for white female. (4) The results for black female users combine the general system bias for both black users and female users, where sweet food and nationality- or religious-related (e.g., "vegan", "vegetarian") categories are more likely to be recommended to them. While results in (2) and (4) seems to be caused by race-related information in terms of cuisine types, results in (1) and (3) can be linked with existing literature. The result from (1) reflects the previously discussed well-known higher alcohol usage in men than women (Kezer et al., 2021; Wilsnack et al., 2000, 2009). The result from (3) reflects the existing findings suggested by literature where women report more craving for sweet foods (e.g., chocolate, pastries, ice cream) (Chao, Grilo, & Sinha, 2016; Weingarten & Elston, 1991; Zellner, Garriga-Trillo, Rohm, Centeno, & Parker, 1999). We also note that "food court" and "fast food" appear to be on the extreme end for the black user and without much difference between different gendered users. This result might be related to the previously discussed issue of African American neighbourhoods having a greater prevalence of fast food (Block et al., 2004; Lewis et al., 2005) and tending to have a higher portion of fast food restaurants (Lewis et al., 2005). While some results do not indicate necessarily harmful results (e.g., recommending desserts to women) at a glance, we note that these results can be viewed as algorithm-enforced segregation and certain issues such as the system's tendency to recommend fast food to the black user group with much higher likelihood should raise an alarm.

Although these findings show some obvious biases between the gender and cuisine types, whether resolving such inequality remains an open question, and to the best of our knowledge, no literature shows or discusses similar findings. We provide further discussions of this limitation in Section 5.

**Top item names being recommended to individual bias dimension.** We show in Fig. 6 the top words in the recommended item names (using raw frequency). We can observe that the results are very consistent with the category association score presented by the two-dimensional scatter plot (e.g. "pub" for white and male).

#### 4.6. RQ5: Nightlife and sexual orientation

We do not expect sexual orientation to affect most cuisine preferences (which we see more related to race), but we might expect a relationship with nightlife recommendations. As demonstrated in Table 2, we generate input phrases such as "Do you have any restaurant recommendations for my [<u>1ST RELATIONSHIP</u>] and his/her [<u>2ND RELATIONSHIP</u>]?". The underline words represent the placeholders for gender-related words, which will indirectly indicate the sexual orientations. The [<u>1ST RELATIONSHIP</u>] prompts are chosen from a set of gender-identifying words including "sister", "brother", "daughter", etc., and [<u>2ND RELATIONSHIP</u>] placeholder indicates the gender by using words such as "girlfriend" and "boyfriend". An example input sentence would be "Can you make a restaurant reservation for my brother and his boyfriend?".

Our bias evaluations are based on the calculations of association score in Eq. (2) between the target sensitive attribute and the gender-identifying word. The score shows how each item from the sensitive category is likely to be recommended to user groups with different sexual orientations (e.g., *male homosexual*). The two dimensions of the output graph are the gender dimensions for



Gender Dimension for 1st Relationship Mention

Fig. 7. Two-Dimensional scatter plot of the association score for nightlife-related activities. With a template input sentence "Can you reserve a table for my [IST RELATIONSHIP] and his/her [2ND RELATIONSHIP]?", the x-axis indicates the gender dimension for the 1st relationship and the y-axis indicates that for the 2nd relationship. Error bars show 90% confidence intervals in each dimension. The central grey oval indicates the neutral reference point.

the two relationships placeholders, as shown in Fig. 7 (presented with 90% confidence intervals): (1) *X*-axis is the gender for the first relationship placeholder (e.g. female for "my sister"); (2) *Y*-axis is for the gender representation of the second placeholder (e.g., female for "girlfriend", and male for "boyfriend"). This shows typical recommendation categories for homosexual groups in the 1st and 3rd quadrants on the graph. The grey oval at the origin represents the neutral reference point.

**More sensitive items recommended to sexual minority.** The results are computed using the recommended items for all testing phrases across the seven cities to minimize statistical noise. Ideally, the distribution for the sensitive category should not shift across the gender class or different sexual orientations. However, even by plotting a simple set of nightlife categories, we observe a clear pattern in Fig. 7 that the nightlife categories have higher associations with a sexual minority group (1st and 3rd quadrants), regardless of their gender. For example, casinos, dive bars and pubs all lie on the quadrants for homosexuality in the graph. Specifically, Gay bars show up at the "male + male" (homosexuality) corner. In this latter case, it is very clear that LMRec has picked up on some language cues to recommend stereotypical venues in the case of a query containing a homosexual relationship.

More nightlife-related recommendations for males. Among the sensitive items, we see a significant shift of nightlife-related activities (predominantly alcohol-related venues) to the male side of the first relationship mentioned, as reflected in other results.

#### 4.7. RQ6: Unintended location bias

The mention of locations may contain the user's information on employment, social status or religion. An example of such phrases is **"Can you pick a place to go after I leave the [LOCATION]**?". The placeholder could be "construction site", indicating that the user may be a construction worker. Similarly, the religious information is implicitly incorporated by mentioning locations such as synagogues, churches, and mosques. As mentioned in Section 3.3.2, it is considered to be undesirable if conversational recommender systems exhibit price discrimination towards different users' indications of desired locations. Therefore, in this section, we aim to study whether LMRec exhibits such behaviour.

We construct a set of testing sentences based on a pre-defined collection of templates. Each testing phrase includes a placeholder [LOCATION], which provides potential employment, social status or religious information implicitly. We measure the differences in average price levels of the top-20 recommended restaurants across the substitution words. The average is computed over all cities and all templates.

**Relationship between occupation and price level.** In brief, we see in Fig. 8 (presented with 90% confidence intervals) that professional establishments (e.g., "fashion studio" or "law office") and religious venues like "synagogue" have a higher average price than "convenience store" and "mosque" indicating possible socioeconomic biases based on location and religion. When the occupation information is substituted into the recommendation request queries, a person who goes to the fashion studio receives higher priced recommendations than those who are heading to a convenience store. The results also appear to imply that people who visit fashion studios or can afford a psychiatrist also go to expensive restaurants. While occupations related to fashion are



Fig. 8. Rank Charts for average price level of the restaurant recommendations for different location prompts. (blue) original LMRec; (orange) after applying test side neutralization. 90% confidence intervals are shown. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

less related to socioeconomic status, occupations such as lawyers and psychologists fit into the highest occupational scale defined by Hollingshead (1975). We hypothesize that people related to lawyers, psychiatrists, or psychologists are considered to have higher SES (i.e., the service providers and the customers), while the population majority at places such as universities may be students who have lower SES thus leading to the observed price associations in Fig. 8.

From the perspective of religious information inferred by the mention of locations, the average price level of restaurant recommendations for Jewish people is the highest among the three prompt labels we tested. It is consistent with the analysis result by Pearson and Geronimus (2011) that Jewish Americans are more likely to have a higher income distribution than other white and black populations. It can also be related to the findings by Keister (2012), where Jewish respondents have significantly greater wealth than other groups (e.g., Catholics). This common stereotype may lead to the unfairness of the recommender that will consistently recommend the cheaper restaurants to people with religions other than Judaism, predominantly Muslim, which has the lowest average price for recommendation results among the three religions.

#### 5. Limitations

We now proceed to outline some limitations of our analysis that might be explored in future work:

- **Choice of model:** As discussed in Section 3.3, the recommendation results for this work are based purely on the context of language requests at test time and are not personalized to individual users. Therefore, future work can investigate the existence of unintended biases in a personalized version of LMRec although this extension of LMRec would be a novel contribution itself. Due to this non-penalization setting of our analysis, we do not have sensitive attributes for specific users making language-based recommendation requests and hence we cannot assess group-level fairness in terms of recommendation performance (e.g., whether the male user group gets better recommendation accuracy). Future work that studies a personalized version of LMRec can further analyse the recommendation performance disparity between user groups.
- **Application of test-side neutralization:** As described in Section 3.5, test-side neutralization performs a post-processing bias mitigation method by masking out text that reveals sensitive information in the input queries. However, the biases that exist in the model or recommendation results are not removed by this methodology. To this end, we note that there may be information in the training data that contributes to biases and cannot be easily masked (e.g., sensitive attributes that can be linked to food and cuisine types), and therefore train-time masking could not be applied to every possible contributing factor. Hence future work could investigate novel methods that may be capable of removing or mitigating biases from the trained embeddings through both direct and indirect association of language with sensitive attributes.

• Harmfulness of certain observed unintended biases: It is well-noted in the literature that biases in recommender systems may be very harmful to specific user populations (Dash, Chakraborty, Ghosh, Mukherjee, & Gummadi, 2021; Deldjoo et al., 2022; Edizel, Bonchi, Hajian, Panisson, & Tassa, 2020; Geyik, Ambler, & Kenthapadi, 2019; Hildebrandt, 2022). However, whether recommending desserts to women and pubs to men is harmful remains an open question from an ethical perspective. While we wanted to highlight these notable user-item associations that we observed in our analysis, it is beyond the scope of this work to attempt to resolve such ethical questions. Nonetheless, we remark that some unintended bias *may* be allowable since, generally, it may be deemed innocuous in a given application setting (e.g., recommending desserts to women), and also for practical purposes since bias cannot always be completely detected and removed from the training text or request queries. Overall though, investigating these ethical questions is an important problem for future research.

### 6. Conclusion

Given the potential that pretrained LMs offer for CRSs, we have presented the first quantitative and qualitative analysis to identify and measure unintended biases in language model-driven recommendation. We observed that the LMRec model exhibits various unintended biases without involving any preferential statements nor recorded preferential history of the user, but simply due to an offhand mention of a name or relationship that in principle should not change the recommendations. Fortunately, we have shown that train side masking and test side neutralization of non-preferential entities can nullify the observed biases without significantly impacting recommendation performance *when* the source of bias can be isolated, as it was by design in our research study. In general, recommendation biases can arise through a variety of language-based associations and further research is needed to identify and mitigate novel types of biases that may arise in language-based recommendation. Overall, our work has aimed to identify and raise a red flag for LM-driven CRSs and we consider this study a first step towards understanding and mitigating unintended biases in future LM-driven CRSs that have the potential to impact hundreds of millions of users.

#### CRediT authorship contribution statement

**Tianshu Shen:** Conceptualization, Methodology, Software, Validation, Formal analysis, Investigation, Data curation, Writing – original draft, Writing – review & editing, Visualization. **Jiaru Li:** Methodology, Software, Validation, Formal analysis, Investigation, Data curation, Writing – original draft, Writing – review & editing, Visualization. **Mohamed Reda Bouadjenek:** Conceptualization, Software, Validation, Investigation, Data curation, Writing – original draft, Visualization, Writing – original draft, Visualization, Writing – original draft, Visualization, Software, Validation, Investigation, Data curation, Writing – original draft, Visualization. **Zheda Mai:** Conceptualization, Writing – original draft. **Scott Sanner:** Conceptualization, Writing – review & editing, Supervision.

#### Data availability

The Yelp<sup>3</sup> academic dataset is publicly available. Our public Github<sup>4</sup> repository provides the scripts to process the Yelp data into the format used in the experiments.

#### References

- Abdollahpouri, Himan, Adomavicius, Gediminas, Burke, Robin, Guy, Ido, Jannach, Dietmar, Kamishima, Toshihiro, et al. (2020). Multistakeholder recommendation: Survey and research directions. User Modeling and User-Adapted Interaction, 30(1), 127–158.
- Abdollahpouri, Himan, & Burke, Robin (2019). Multi-stakeholder recommendation and its connection to multi-sided fairness. In Proceedings of the workshop on recommendation in multi-stakeholder environments co-located with the 13th ACM conference on recommender systems, Copenhagen, Denmark.
- Badjatiya, Pinkesh, Gupta, Manish, & Varma, Vasudeva (2019). Stereotypical bias removal for hate speech detection task using knowledge-based generalizations. In *The world wide web conference* (pp. 49–59).
- Barikeri, Soumya, Lauscher, Anne, Vulic, Ivan, & Glavas, Goran (2021). RedditBias: A real-world resource for bias evaluation and debiasing of conversational language models. In Proceedings of the 59th annual meeting of the association for computational linguistics and the 11th international joint conference on natural language processing, ACL/IJCNLP 2021 (pp. 1941–1955). Association for Computational Linguistics.
- Bartlett, Robert, Morse, Adair, Stanton, Richard, & Wallace, Nancy (2022). Consumer-lending discrimination in the FinTech era. Journal of Financial Economics, 143(1), 30–56.
- Bertrand, Marianne, & Mullainathan, Sendhil (2004). Are Emily and Greg more employable than Lakisha and Jamal? A field experiment on labor market discrimination. *American Economic Review*, 94(4), 991-1013.

Bhardwaj, Rishabh, Majumder, Navonil, & Poria, Soujanya (2021). Investigating gender bias in bert. Cognitive Computation, 13, 1-11.

- Block, Jason P., Scribner, Richard A., & DeSalvo, Karen B. (2004). Fast food, race/ethnicity, and income: a geographic analysis. American Journal of Preventive Medicine, 27(3), 211–217.
- Bolukbasi, Tolga, Chang, Kai-Wei, Zou, James Y., Saligrama, Venkatesh, & Kalai, Adam T. (2016). Man is to computer programmer as woman is to homemaker? debiasing word embeddings. Advances in Neural Information Processing Systems, 29, 4349–4357.

Borges, Rodrigo, & Stefanidis, Kostas (2021). On mitigating popularity bias in recommendations via variational autoencoders. In Proceedings of the 36th annual ACM symposium on applied computing (pp. 1383–1389).

- Brand, Alexander, & Gross, Tom (2020). Paying the pink tax on a blue dress-exploring gender-based price-premiums in fashion recommendations. In International conference on human-centred software engineering (pp. 190–198). Springer.
- Braveman, Paula A., Cubbin, Catherine, Egerter, Susan, Williams, David R., & Pamuk, Elsie (2010). Socioeconomic disparities in health in the United States: what the patterns tell us. American Journal of Public Health, 100(S1), S186–S196.

<sup>&</sup>lt;sup>3</sup> https://www.yelp.com/dataset/download

<sup>&</sup>lt;sup>4</sup> https://github.com/TinaBBB/Unintended-Bias-LMRec.git.

- Brownell, Kelly D., & Horgen, Katherine Battle (2004). Food fight: The inside story of the food industry, America's obesity crisis, and what we can do about it. Contemporary Books.
- Burke, Robin (2017). Multisided fairness for recommendation. In Proceedings of the workshop on fairness, accountability, and transparency in machine learning. FATML.
- Burke, Robin, Sonboli, Nasim, & Ordonez-Gauger, Aldo (2018). Balanced neighborhoods for multi-sided fairness in recommendation. In Conference on fairness, accountability and transparency (pp. 202-214). PMLR.

Burstein, Paul (2007). Jewish educational and economic success in the United States: A search for explanations. Sociological Perspectives, 50(2), 209–228.

Caliskan, Aylin, Bryson, Joanna J., & Narayanan, Arvind (2017). Semantics derived automatically from language corpora contain human-like biases. Science, 356(6334), 183–186.

- Chao, Ariana M., Grilo, Carlos M., & Sinha, Rajita (2016). Food cravings, binge eating, and eating disorder psychopathology: Exploring the moderating roles of gender and race. *Eating Behaviors*, 21, 41–47.
- Chen, Jiawei, Dong, Hande, Wang, Xiang, Feng, Fuli, Wang, Meng, & He, Xiangnan (2020). Bias and debias in recommender system: A survey and future directions. arXiv preprint arxiv:2010.03240.
- Christakopoulou, Konstantina, Radlinski, Filip, & Hofmann, Katja (2016). Towards conversational recommender systems. In Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining (pp. 815–824).
- Cui, Ruomeng, Li, Jingyun, Li, Meng, & Yu, Lili (2021). Wholesale price discrimination in global sourcing. Manufacturing & Service Operations Management, 23(5), 1096–1117.
- Dash, Abhisek, Chakraborty, Abhijnan, Ghosh, Saptarshi, Mukherjee, Animesh, & Gummadi, Krishna P. (2021). When the umpire is also a player: Bias in private label product recommendations on e-commerce marketplaces. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency* (pp. 873–884).
- Deldjoo, Yashar, Anelli, Vito Walter, Zamani, Hamed, Bellogin, Alejandro, & Di Noia, Tommaso (2021). A flexible framework for evaluating user and item fairness in recommender systems. User Modeling and User-Adapted Interaction, 1–55.
- Deldjoo, Yashar, Anelli, Vito Walter, Zamani, Hamed, Kouki, Alejandro Bellogín, & Noia, Tommaso Di (2019). Recommender systems fairness evaluation via generalized cross entropy. In CEUR workshop proceedings: vol. 2440, Proceedings of the workshop on recommendation in multi-stakeholder environments co-located with the 13th ACM conference on recommender systems (RecSys 2019), Copenhagen, Denmark, September 20, 2019. CEUR-WS.org.
- Deldjoo, Yashar, Jannach, Dietmar, Bellogin, Alejandro, Difonzo, Alessandro, & Zanzonelli, Dario (2022). A survey of research on fair recommender systems. arXiv preprint arXiv:2205.11127.

Desa, UN, et al. (2016). Transforming our world: The 2030 agenda for sustainable development. UN General Assembly.

- Devlin, Jacob, Chang, Ming-Wei, Lee, Kenton, & Toutanova, Kristina (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) (pp. 4171–4186). Minneapolis, Minnesota: Association for Computational Linguistics.
- Duesterhaus, Megan, Grauerholz, Liz, Weichsel, Rebecca, & Guittar, Nicholas A. (2011). The cost of doing femininity: Gendered disparities in pricing of personal care products and services. *Gender Issues*, 28(4), 175–191.
- Edizel, Bora, Bonchi, Francesco, Hajian, Sara, Panisson, André, & Tassa, Tamir (2020). FaiRecSys: mitigating algorithmic bias in recommender systems. International Journal of Data Science and Analytics, 9(2), 197–213.

Ekstrand, Michael D., & Pera, Maria Soledad (2017). The demographics of cool. In Poster proceedings at ACM RecSys. Como, Italy: ACM.

Ekstrand, Michael D., Tian, Mucun, Azpiazu, Ion Madrazo, Ekstrand, Jennifer D., Anuyah, Oghenemaro, McNeill, David, et al. (2018). All the cool kids, how do they fit in?: Popularity and demographic biases in recommender evaluation and effectiveness. In *Conference on fairness, accountability and transparency* (pp. 172–186), PMLR.

Evans, David S., & Schmalensee, Richard (2016). Matchmakers: The new economics of multisided platforms. Harvard Business Review Press.

- Ferraro, Andres (2019). Music cold-start and long-tail recommendation: bias in deep representations. In Proceedings of the 13th ACM conference on recommender systems (pp. 586–590).
- Finkel, Jenny Rose, Grenager, Trond, & Manning, Christopher (2005). Incorporating non-local information into information extraction systems by gibbs sampling. In *Proceedings of the 43rd annual meeting of the association for computational linguistics (ACL'05)* (pp. 363–370). Ann Arbor, Michigan: Association for Computational Linguistics.

Friedman, Batya, & Nissenbaum, Helen (1996). Bias in computer systems. ACM Transactions on Information Systems (TOIS), 14(3), 330-347.

- Fryer, Roland G., Jr., & Levitt, Steven D. (2004). The causes and consequences of distinctively black names. Quarterly Journal of Economics, 119(3), 767-805.
- Fu, Zuohui, Xian, Yikun, Gao, Ruoyuan, Zhao, Jieyu, Huang, Qiaoying, Ge, Yingqiang, et al. (2020). Fairness-aware explainable recommendation over knowledge graphs. In Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval (pp. 69–78).
- Fujishiro, Kaori, Xu, Jun, & Gong, Fang (2010). What does "occupation" represent as an indicator of socioeconomic status?: Exploring occupational prestige and health. Social Science & Medicine, 71(12), 2100–2107.
- Fuster, Andreas, Goldsmith-Pinkham, Paul, Ramadorai, Tarun, & Walther, Ansgar (2022). Predictably unequal? The effects of machine learning on credit markets. *The Journal of Finance*, 77(1), 5–47.

Gandal, Neil, & Shabelansky, Anastasia (2010). Obesity and price sensitivity at the supermarket. In Forum for health economics & policy, Vol. 13. De Gruyter.

- Gao, Chongming, Lei, Wenqiang, He, Xiangnan, de Rijke, Maarten, & Chua, Tat-Seng (2021). Advances and challenges in conversational recommender systems: A survey. AI Open, 2, 100–126.
- Gao, Ruoyuan, & Shah, Chirag (2021). Addressing bias and fairness in search systems. In Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval (pp. 2643–2646).
- Ge, Yingqiang, Liu, Shuchang, Gao, Ruoyuan, Xian, Yikun, Li, Yunqi, Zhao, Xiangyu, et al. (2021). Towards long-term fairness in recommendation. In Proceedings of the 14th ACM international conference on web search and data mining (pp. 445–453).
- Geyik, Sahin Cem, Ambler, Stuart, & Kenthapadi, Krishnaram (2019). Fairness-aware ranking in search & recommendation systems with application to linkedin talent search. In Proceedings of the 25th Acm sigkdd international conference on knowledge discovery & data mining (pp. 2221–2231).

Goldstein, Sidney (1969). Socioeconomic differentials among religious groups in the United States. American Journal of Sociology, 74(6), 612-631.

- Grant, Bridget F., Goldstein, Risë B., Saha, Tulshi D., Chou, S. Patricia, Jung, Jeesun, Zhang, Haitao, et al. (2015). Epidemiology of DSM-5 alcohol use disorder: results from the national epidemiologic survey on alcohol and related conditions III. JAMA Psychiatry, 72(8), 757–766.
- Guo, Wei, & Caliskan, Aylin (2021). Detecting emergent intersectional biases: Contextualized word embeddings contain a distribution of human-like biases. In Proceedings of the 2021 AAAI/ACM conference on AI, ethics, and society (pp. 122–133).
- Hada, Deepesh V., & Shevade, Shirish K. (2021). ReXPlug: Explainable recommendation using plug-and-play language model. In Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval (pp. 81–91).
- Hallam, Jessica, Boswell, Rebecca G., DeVito, Elise E., & Kober, Hedy (2016). Focus: sex and gender health: gender-related differences in food craving and obesity. *The Yale Journal of Biology and Medicine*, 89(2), 161.
- Hao, Qianxiu, Xu, Qianqian, Yang, Zhiyong, & Huang, Qingming (2021). Pareto optimality for fairness-constrained collaborative filtering. In *Proceedings of the 29th ACM international conference on multimedia* (pp. 5619–5627).
- Hildebrandt, Mireille (2022). The issue of proxies and choice architectures. Why EU law matters for recommender systems. Frontiers in Artificial Intelligence, 73. Hollingshead, August B. (1975). Four factor index of social status. Yale University, New Haven, Connecticut.

- Hutchinson, Ben, Prabhakaran, Vinodkumar, Denton, Emily, Webster, Kellie, Zhong, Yu, & Denuyl, Stephen (2020). Social biases in NLP models as barriers for persons with disabilities. In *Proceedings of the 58th annual meeting of the association for computational linguistics, ACL, Online*. Association for Computational Linguistics.
- Jacob, Jorge, Vieites, Yan, Goldszmidt, Rafael, & Andrade, Eduardo B (2022). EXPRESS: Expected SES-based discrimination reduces price sensitivity among the poor. Journal of Marketing Research, Article 00222437221097100.

Jannach, Dietmar, Manzoor, Ahtsham, Cai, Wanling, & Chen, Li (2021). A survey on conversational recommender systems. ACM Computing Surveys, 54(5), 1–36.

Jin, Xisen, Barbieri, Francesco, Kennedy, Brendan, Davani, Aida Mostafazadeh, Neves, Leonardo, & Ren, Xiang (2021). On transferability of bias mitigation effects in language model fine-tuning. In *Proceedings of the 2021 conference of the North American chapter of the association for computational linguistics: Human language technologies, NAACL-HLT, online.* Association for Computational Linguistics.

Keister, Lisa A. (2012). Religion and wealth across generations. In Religion, work and inequality. Emerald Group Publishing Limited.

- Kezer, Camille A., Simonetto, Douglas A., & Shah, Vijay H. (2021). Sex differences in alcohol consumption and alcohol-associated liver disease. In Mayo clinic proceedings. Elsevier.
- Kurita, Keita, Vyas, Nidhi, Pareek, Ayush, Black, Alan W., & Tsvetkov, Yulia (2019a). Measuring bias in contextualized word representations. In Proceedings of the first workshop on gender bias in natural language processing. Florence, Italy.
- Kurita, Keita, Vyas, Nidhi, Pareek, Ayush, Black, Alan W., & Tsvetkov, Yulia (2019b). Quantifying social biases in contextual word representations. In 1st ACL workshop on gender bias for natural language processing.
- Kwate, Naa Oyo A. (2008). Fried chicken and fresh apples: racial segregation as a fundamental cause of fast food density in black neighborhoods. Health & Place, 14(1), 32-44.
- Laban, Guy, & Araujo, Theo (2020). The effect of personalization techniques in users' perceptions of conversational recommender systems. In Proceedings of the 20th ACM international conference on intelligent virtual agents.
- Lei, Wenqiang, Zhang, Gangyi, He, Xiangnan, Miao, Yisong, Wang, Xiang, Chen, Liang, et al. (2020). Interactive path reasoning on graph for conversational recommendation. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining* (pp. 2073–2083).
- Lewis, LaVonna Blair, Sloane, David C., Nascimento, Lori Miller, Diamant, Allison L., Guinyard, Joyce Jones, Yancey, Antronette K., et al. (2005). African Americans' access to healthy food options in South Los Angeles restaurants. *American Journal of Public Health*, 95(4), 668–673.
- Li, Yunqi, Chen, Hanxiong, Fu, Zuohui, Ge, Yingqiang, & Zhang, Yongfeng (2021). User-oriented fairness in recommendation. In Proceedings of the web conference 2021 (pp. 624–632).
- Li, Yunqi, Ge, Yingqiang, & Zhang, Yongfeng (2021). Tutorial on fairness of machine learning in recommender systems. In SIGIR '21: The 44th international ACM SIGIR conference on research and development in information retrieval, Virtual Event, Canada, July 11-15, 2021 (pp. 2654–2657). ACM.
- Li, Raymond, Kahou, Samira, Schulz, Hannes, Michalski, Vincent, Charlin, Laurent, & Pal, Chris (2018). Towards deep conversational recommendations. In Proceedings of the 32nd international conference on neural information processing systems (pp. 9748–9758).
- Liang, Paul Pu, Wu, Chiyu, Morency, Louis-Philippe, & Salakhutdinov, Ruslan (2021). Towards understanding and mitigating social biases in language models. In International conference on machine learning (pp. 6565–6576). PMLR.
- Lin, Kun, Sonboli, Nasim, Mobasher, Bamshad, & Burke, Robin (2019). Crank up the volume: preference bias amplification in collaborative recommendation. In Proceedings of the workshop on recommendation in multi-stakeholder environments co-located with the 13th ACM conference on recommender, Copenhagen, Denmark.
- Liu, Zhuang, Lin, Wayne, Shi, Ya, & Zhao, Jun (2021). A robustly optimized BERT pre-training approach with post-training. In Chinese computational linguistics 20th China national conference, CCL, Hohhot, China.
- Lu, Kaiji, Mardziel, Piotr, Wu, Fangjing, Amancharla, Preetam, & Datta, Anupam (2020). Gender bias in neural natural language processing. In Logic, language, and security (pp. 189–202). Springer.
- Malkiel, Itzik, Barkan, Oren, Caciularu, Avi, Razin, Noam, Katz, Ori, & Koenigstein, Noam (2020). RecoBERT: A catalog language model for text-based recommendations. In Findings of the Association for Computational Linguistics: EMNLP 2020.
- Mansoury, Masoud, Abdollahpouri, Himan, Pechenizkiy, Mykola, Mobasher, Bamshad, & Burke, Robin (2020). Feedback loop and bias amplification in recommender systems. In Proceedings of the 29th ACM international conference on information & knowledge management.
- Mansoury, Masoud, Mobasher, Bamshad, Burke, Robin, & Pechenizkiy, Mykola (2019). Bias disparity in collaborative recommendation: Algorithmic evaluation and comparison. In Proceedings of the workshop on recommendation in multi-stakeholder environments co-located with the 13th ACM conference on recommender systems, Copenhagen, Denmark.
- May, Chandler, Wang, Alex, Bordia, Shikha, Bowman, Samuel R., & Rudinger, Rachel (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, NAACL-HLT 2019, Minneapolis, MN, USA.
- Misztal-Radecka, Joanna, & Indurkhya, Bipin (2021). Bias-aware hierarchical clustering for detecting the discriminated groups of users in recommendation systems. Information Processing & Management, 58(3), Article 102519.
- Morland, Kimberly, Wing, Steve, Roux, Ana Diez, & Poole, Charles (2002). Neighborhood characteristics associated with the location of food stores and food service places. *American Journal of Preventive Medicine*, 22(1), 23–29.

Nadeem, Moin, Bethke, Anna, & Reddy, Siva (2021). StereoSet: Measuring stereotypical bias in pretrained language models. In Proceedings of the 59th annual meeting of the association for computational linguistics and the 11th international joint conference on natural language processing. ACL/IJCNLP, Virtual Event.

Noël, Reginald A. (2018). Race, economics, and social status. U.S. Department of Labor, Bureau of Labor Statistics.

- Otter, Daniel W., Medina, Julian R., & Kalita, Jugal K. (2020). A survey of the usages of deep learning for natural language processing. *IEEE Transactions on Neural Networks and Learning Systems*, 32(2), 604–624.
- Patro, Gourab K., Biswas, Arpita, Ganguly, Niloy, Gummadi, Krishna P., & Chakraborty, Abhijnan (2020). Fairrec: Two-sided fairness for personalized recommendations in two-sided platforms. In *Proceedings of the web conference 2020* (pp. 1194–1204).
- Pearson, J. A., & Geronimus, A. T. (2011). Race/ethnicity, socioeconomic characteristics, coethnic social ties, and health: evidence from the national Jewish population survey. *American Journal of Public Health*, 101(7), 1314–1321.

Pelchat, Marcia Levin (1997). Food cravings in young and elderly adults. Appetite, 28(2), 103-113.

Penha, Gustavo, & Hauff, Claudia (2020). What does BERT know about books, movies and music? Probing BERT for conversational recommendation. In Fourteenth ACM conference on recommender systems (pp. 388–397).

- Qiu, Xipeng, Sun, Tianxiang, Xu, Yige, Shao, Yunfan, Dai, Ning, & Huang, Xuanjing (2020). Pre-trained models for natural language processing: A survey. Science China Technological Sciences, 1–26.
- Radford, Alec, Narasimhan, Karthik, Salimans, Tim, & Sutskever, Ilya (2018). Improving language understanding by generative pre-training. Technical Report. OpenAI.
- Rastegarpanah, Bashir, Gummadi, Krishna P., & Crovella, Mark (2019). Fighting fire with fire: Using antidote data to improve polarization and fairness of recommender systems. In Proceedings of the twelfth ACM international conference on web search and data mining (pp. 231–239).
- Reeves, Richard, Rodrigue, Edward, & Kneebone, Elizabeth (2016). Five evils: Multidimensional poverty and race in America. In Economic studies at brookings report, Vol. 1 (pp. 1–22).
- Ren, Xuhui, Yin, Hongzhi, Chen, Tong, Wang, Hao, Hung, Nguyen Quoc Viet, Huang, Zi, et al. (2020). Crsal: Conversational recommender systems with adversarial learning. ACM Transactions on Information Systems (TOIS).

Santamaría, Lucía, & Mihaljević, Helena (2018). Comparison and benchmark of name-to-gender inference services. *PeerJ Computer Science*, 4, Article e156. Sedhain, Suvash, Menon, Aditva, Sanner, Scott, & Xie, Lexing (2015). AutoRec: Autoencoders meet collaborative filtering. In *Proceedings of the 24th international* 

conference on the world wide web (WWW-15). Florence, Italy.

Sheng, Emily, Chang, Kai-Wei, Natarajan, Premkumar, & Peng, Nanyun (2019). The woman worked as a babysitter: On biases in language generation. In Proceedings of the 2019 conference on empirical methods in natural language processing and the 9th international joint conference on natural language processing. EMNLP-IJCNLP 2019, Hong Kong, China.

Sood, Gaurav, & Laohaprapanon, Suriyan (2018). Predicting race and ethnicity from the sequence of characters in a name. arXiv preprint arXiv:1805.02109. Stevens, Jennifer L., & Shanahan, Kevin J. (2017). Structured abstract: anger, willingness, or clueless? Understanding why women pay a pink tax on the products

they consume. In Creating marketing magic and innovative future marketing trends. Springer.

Sun, Yueming, & Zhang, Yi (2018). Conversational recommender system. In The 41st international Acm sigir conference on research & development in information retrieval (pp. 235-244).

Sweeney, Latanya (2013). Discrimination in online ad delivery. Communications of the ACM, 56(5), 44-54.

Tan, Yi Chern, & Celis, L. Elisa (2019). Assessing social and intersectional biases in contextualized word representations. In Advances in neural information processing systems 32: Annual conference on neural information processing systems, Vancouver, BC, Canada.

Tsintzou, Virginia, Pitoura, Evaggelia, & Tsaparas, Panayiotis (2019). Bias disparity in recommendation systems. In Proceedings of the workshop on recommendation in multi-stakeholder environments co-located with the 13th ACM conference on recommender systems, Copenhagen, Denmark.

Vaswani, Ashish, Shazeer, Noam, Parmar, Niki, Uszkoreit, Jakob, Jones, Llion, Gomez, Aidan N., et al. (2017). Attention is all you need. In Advances in neural information processing systems (pp. 5998-6008).

Weingarten, Harvey P., & Elston, Dawn (1991). Food cravings in a college population. Appetite, 17(3), 167-175.

Wilsnack, Richard W., Vogeltanz, Nancy D., Wilsnack, Sharon C., & Harris, T. Robert (2000). Gender differences in alcohol consumption and adverse drinking consequences: cross-cultural patterns. Addiction, 95(2), 251–265.

Wilsnack, Richard W., Wilsnack, Sharon C., Kristjanson, Arlinda F., Vogeltanz-Holm, Nancy D., & Gmel, Gerhard (2009). Gender and alcohol consumption: patterns from the multinational GENACIS project. Addiction, 104(9), 1487–1500.

Winkleby, Marilyn A., Jatulis, Darius E., Frank, Erica, & Fortmann, Stephen P. (1992). Socioeconomic status and health: how education, income, and occupation contribute to risk factors for cardiovascular disease. *American Journal of Public Health*, 82(6), 816–820.

Yang, Ke, & Stoyanovich, Julia (2017). Measuring fairness in ranked outputs. In Proceedings of the 29th international conference on scientific and statistical database management (pp. 1–6).

Yao, Sirui, & Huang, Bert (2017). Beyond parity: Fairness objectives for collaborative filtering. In Advances in neural information processing systems 30: Annual conference on neural information processing systems 2017, December 4-9, 2017, Long Beach, CA, USA (pp. 2921–2930).

Zehlike, Meike, Bonchi, Francesco, Castillo, Carlos, Hajian, Sara, Megahed, Mohamed, & Baeza-Yates, Ricardo (2017). Fa\* ir: A fair top-k ranking algorithm. In Proceedings of the 2017 ACM on conference on information and knowledge management (pp. 1569–1578).

Zellner, D. A., Garriga-Trillo, Ana, Rohm, Elizabeth, Centeno, Soraya, & Parker, Scott (1999). Food liking and craving: A cross-cultural approach. Appetite, 33(1), 61–70.

Zhang, Haoran, Lu, Amy X., Abdalla, Mohamed, McDermott, Matthew, & Ghassemi, Marzyeh (2020). Hurtful words: quantifying biases in clinical contextual word embeddings. In *Proceedings of the ACM conference on health, inference, and learning* (pp. 110–120).

- Zhao, Jieyu, Wang, Tianlu, Yatskar, Mark, Cotterell, Ryan, Ordonez, Vicente, & Chang, Kai-Wei (2019). Gender bias in contextualized word embeddings. 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) (pp. 629–634). Minneapolis, Minnesota: Association for Computational Linguistics.
- Zhao, Jieyu, Wang, Tianlu, Yatskar, Mark, Ordonez, Vicente, & Chang, Kai-Wei (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, NAACL-HLT, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 2 (Short papers) (pp. 15–20). Association for Computational Linguistics.
- Zhao, Jieyu, Zhou, Yichao, Li, Zeyu, Wang, Wei, & Chang, Kai-Wei (2018). Learning gender-neutral word embeddings. In Proceedings of the 2018 conference on empirical methods in natural language processing, Brussels, Belgium, October 31 - November 4, 2018.