
Investigating Gender Bias in Language Models Using Causal Mediation Analysis

Jesse Vig^{*1} Sebastian Gehrmann^{*2} Yonatan Belinkov^{*2}
Sharon Qian² Daniel Nevo³ Yaron Singer² Stuart Shieber²
¹ Salesforce Research ² Harvard University ³ Tel Aviv University
jvig@salesforce.com danielnevo@tauex.tau.ac.il
{gehrmann, belinkov, sharonqian, yaron, shieber}@seas.harvard.edu

Abstract

Many interpretation methods for neural models in natural language processing investigate how information is encoded inside hidden representations. However, these methods can only measure whether the information exists, not whether it is actually used by the model. We propose a methodology grounded in the theory of causal mediation analysis for interpreting which parts of a model are causally implicated in its behavior. The approach enables us to analyze the mechanisms that facilitate the flow of information from input to output through various model components, known as mediators. As a case study, we apply this methodology to analyzing gender bias in pre-trained Transformer language models. We study the role of individual neurons and attention heads in mediating gender bias across three datasets designed to gauge a model’s sensitivity to gender bias. Our mediation analysis reveals that gender bias effects are concentrated in specific components of the model that may exhibit highly specialized behavior.

1 Introduction

The success of neural network models in various natural language processing tasks, coupled with their opaque nature, has led to much interest in interpreting and analyzing such models. One goal of these analyses is to identify whether a model utilizes latent information in its internal representations to arrive at a prediction. This is of particular importance when diagnosing the reasons for a *biased* prediction. A popular class of analysis methods, often called structural analysis, aims to extract this information using probing classifiers that predict linguistic properties from representations of trained models (e.g., Adi et al., 2017; Conneau et al., 2018; Hupkes et al., 2018; Tenney et al., 2019). However, since probing classifiers only yield a correlational measure between a model’s representations and an external property (Belinkov and Glass, 2019), they cannot show if the property is *causally* connected to the model’s predictions. Moreover, Barrett et al. (2019) showed that probing classifiers may generate unfaithful interpretations and fail to generalize to unseen data.

We introduce a methodology for interpreting neural models to address these limitations. We adapt *causal mediation analysis* (Pearl, 2001) for analyzing the mechanisms by which information flows from input to output through different model components. Mediation analysis relies on measuring the change in an output following a counterfactual intervention in an intermediate variable, or *mediator*. Through such interventions, one can measure the degree to which inputs influence outputs directly (*direct effect*), or indirectly through the mediator (*indirect effect*). In our case, the mediator can be any model components that we wish to study, such as neurons or attention heads. We propose multiple controlled interventions in these mediators, which reveal their causal role in a model’s behavior.

* Equal contribution. Y.B. is now at the Technion – Israel Institute of Technology. Work conducted while J.V. was at Palo Alto Research Center.

In a case study, we apply this framework to the analysis of gender bias in large pre-trained language models. Gender bias has surfaced as a major concern in word representations, with strong effects in both static word embeddings (Caliskan et al., 2017; Bolukbasi et al., 2016) and contextualized word representations (Zhao et al., 2019; Basta et al., 2019; Tan and Celis, 2019). Mediation analysis enables us to study how biased predictions arise from different model components. In our study, we focus on the role of individual neurons or attention heads in Transformer-based language models, in particular, several versions of GPT2 (Radford et al., 2019). In an experiment using several datasets designed to gauge a model’s gender bias, we find that gender bias effects increase with larger models, which potentially absorb more bias from the underlying training data. The causal mediation analysis further reveals that gender bias is sparse, with much of the effect concentrated in a relatively small proportion of neurons and attention heads.

In summary, this paper makes two broad contributions. First, we cast causal mediation analysis as an approach for analyzing neural NLP models, which may be applied to a variety of models and phenomena. Second, we demonstrate this methodology in the case of analyzing gender bias in pre-trained language models, revealing the internal mechanisms by which bias effects flow from input to output through various model components. The code for reproducing our results is available at <https://github.com/sebastianGehrmann/CausalMediationAnalysis>.

2 Methodology

2.1 Preliminaries

Consider a large pre-trained neural language model (LM), parameterized by θ , which predicts the probability of the next word given a prefix: $p_\theta(x_t | x_1, \dots, x_{t-1})$. We will focus on LMs based on Transformers (Vaswani et al., 2017), although much of the methodology applies to other architectures as well. Let $h_{l,i} \in \mathbb{R}^K$ denote the (contextual) representation of word i in layer l of the model, with neuron activations $h_{l,i,k}$ ($1 \leq k \leq K$). These representations are composed using so-called multi-headed attention. Let $\alpha_{l,h,i,j} \geq 0$ denote the attention directed from word i to word j by head h in layer l , such that $\sum_j \alpha_{l,h,i,j} = 1$.

2.2 Causal Mediation Analysis

Causal mediation analysis aims to measure how a treatment effect is mediated by intermediate variables (Robins and Greenland, 1992; Pearl, 2001; Robins, 2003). Pearl (2001) described an example where a side effect of a drug may cause patients to take aspirin, and the latter has a separate effect on the disease the drug was originally prescribed for. Thus, the drug has a direct effect through its standard mechanism and an indirect effect operating via aspirin taking.

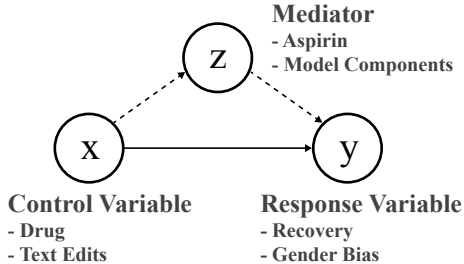


Figure 1: Mediation analysis illustration.

As illustrated in Figure 1, we may consider each neuron in a neural network to be analogous to aspirin in the example above – the neuron is influenced by the input and, in turn, affects the model output; however, there also exist direct pathways from the input to the output that do not pass through the neuron. We can thus decouple model components from the rest of the model by framing them as intermediaries in the causal path from inputs to outputs. Throughout this paper, we specifically focus on the use case of gender bias in language models, as past work suggests that gender is captured in specific model components, e.g., subspaces of contextual word representations (Zhao et al., 2019). By measuring the direct and indirect effects of targeted interventions, we can pinpoint how gender bias propagates through different parts of pre-trained LMs. While we use gender bias as a case study, the approach can be applied to other biases as well (race, ethnicity, etc.).

The example on the right illustrates a typical problem with biased LMs. Given a prompt u such as *The nurse said that*, a language model generates a continuation. A biased model may assign a higher likelihood to *she* than to *he*, such that $p_\theta(\textit{she} | u) > p_\theta(\textit{he} | u)$. In this case, *she* is the stereotypical candidate, while *he* is the anti-stereotypical candidate, which reflects a societal bias associating nurses with women more than men. Coming back to the binary setup, the relative

<p>Prompt u: The nurse said that __</p> <p>Stereotypical candidate: she</p> <p>Anti-stereotypical candidate: he</p>

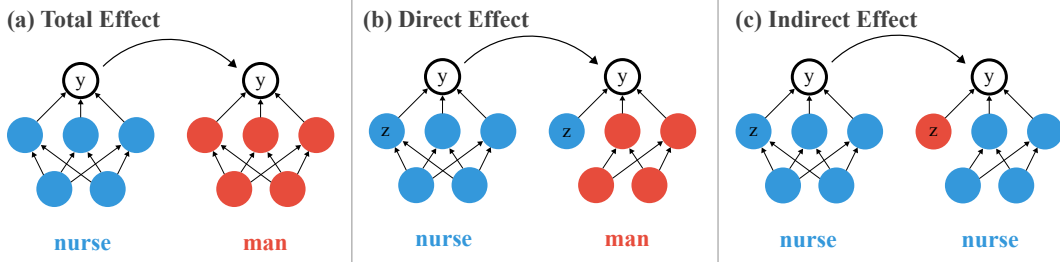


Figure 2: Mediation analysis illustration. Here the do -operation is $x = \text{set-gender}$, which changes u from *nurse* to *man* in this example. The **total effect** measures the change in y resulting from the intervention; the **direct effect** measures the change in y resulting from performing the intervention while holding a mediator z fixed; the **indirect effect** measures the change caused by setting z to its value under the intervention, while holding u fixed.

probabilities assigned to the candidates can be thought of as a measure of gender bias in the model:

$$\mathbf{y}(u) = \frac{p_\theta(\text{anti-stereotypical} \mid u)}{p_\theta(\text{stereotypical} \mid u)}. \quad (1)$$

In our example, $\mathbf{y}(u) = p_\theta(\textit{he} \mid \textit{The nurse said that})/p_\theta(\textit{she} \mid \textit{The nurse said that})$. If $\mathbf{y}(u) < 1$, the prediction is stereotypical; if $\mathbf{y}(u) > 1$, it is anti-stereotypical. A perfectly unbiased model would achieve $\mathbf{y}(u) = 1$ and exhibit bias toward neither the stereotypical nor the anti-stereotypical case. This binary simplification of grammatical gender does not capture the full spectrum, as argued by Cao and Daumé III (2019), and it is not defined, for example, what probability mass a gender-neutral reference should receive². While we leave extension of the framework to a continuous setup to future work, we report experimental results on the singular *they* compared to *he*.

We then apply causal mediation analysis by performing interventions on the input text, and measuring the effect on the gender bias measure defined above (Eq. 1), which we treat as the response variable. We define the following do -operations: (a) **set-gender**: replace the ambiguous profession with an anti-stereotypical gender-specific word (that is, replace *nurse* with *man*, *doctor* with *woman*, etc.); (b) **null**: leave the sentence as is. The population of units for this analysis is a set of example sentences such as the above prompt. We define $\mathbf{y}_x(u)$ as the value that \mathbf{y} attains in unit $u = u$ under the intervention $do(x = x)$.

The unit-level **total effect** (TE) of $x = x$ on \mathbf{y} in unit $u = u$ is the proportional difference³ between the amount of bias under a gendered reading and under an ambiguous reading (Figure 2a):

$$\text{TE}(\text{set-gender}, \text{null}; \mathbf{y}, u) = \frac{\mathbf{y}_{\text{set-gender}}(u) - \mathbf{y}_{\text{null}}(u)}{\mathbf{y}_{\text{null}}(u)} = \frac{\mathbf{y}_{\text{set-gender}}(u)}{\mathbf{y}_{\text{null}}(u)} - 1. \quad (2)$$

For our running example, this results in

$$\frac{p_\theta(\textit{he} \mid \textit{The man said that})}{p_\theta(\textit{she} \mid \textit{The man said that})} \Big/ \frac{p_\theta(\textit{he} \mid \textit{The nurse said that})}{p_\theta(\textit{she} \mid \textit{The nurse said that})} - 1. \quad (3)$$

An illustrative example of the computation of the total effect is provided in Figure 3.

The average total effect of $x = x$ on \mathbf{y} is calculated by taking the expectation over the population u :

$$\text{TE}(\text{set-gender}, \text{null}; \mathbf{y}) = \mathbb{E}_u [\mathbf{y}_{\text{set-gender}}(u)/\mathbf{y}_{\text{null}}(u) - 1]. \quad (4)$$

We then analyze the causal role of specific mediators, or intermediary variables, which lie between x and \mathbf{y} . The mediator, denoted as z , might be a particular neuron, a full layer, an attention head, or a certain attention weight. Following Pearl’s definitions, we measure the direct and indirect effects of intervening in the model relative to z (Pearl, 2001).

²One could argue for the case where *he*, *she*, and *they* have same probability or where *they* has equal representation to the sum of *he* and *she*. Alternatively, one could argue that grammatical genders are inherently discriminatory and that we should change all of them to *they*, unless we know an individual’s preferred pronouns.

³We make the difference proportional to control for the high variance of \mathbf{y} across examples (appendix A.1).

Example

$u =$ The nurse said that [blank]

1) Compute relative probabilities of the baseline.

$$\begin{aligned} p([he]|u) &= p([he]|\text{the nurse said that}) \approx 0.03 \\ p([she]|u) &= p([she]|\text{the nurse said that}) \approx 0.22 \\ \mathbf{y}_{\text{null}}(u) &= 0.03/0.22 \approx 0.14 \end{aligned}$$

3) Compute the total effect

$$\begin{aligned} \text{TE}(\text{set-gender}, \text{null}; \mathbf{y}, u) \\ &= 13.1/0.14 - 1 \approx 92.6 \end{aligned}$$

2) Set u to an anti-stereotypical case and recompute.

$x =$ set-gender: change nurse \rightarrow man

$$\begin{aligned} p([he]|u, \text{set-gender}) &= \\ p([he]|\text{the man said that}) &\approx 0.32 \\ p([she]|u, \text{set-gender}) &= \\ p([she]|\text{the man said that}) &\approx 0.02 \\ \mathbf{y}_{\text{set-gender}}(u) &= 0.32/0.02 \approx 13.1 \end{aligned}$$

Figure 3: An example calculation of the **total effect** with the prompt $u =$ *The nurse said that* and the control variable $x =$ set-gender. Before the intervention, the model assigns a much higher probability to [she], the stereotypical example, than to [he]. By changing nurse to man, we compute the proportional probability of a definitionally gendered example. The total effect measures the effect of this intervention.

The **natural direct effect** (NDE) measures how much an intervention x changes an outcome variable \mathbf{y} directly, without passing through a hypothesized mediator z . It is computed by applying the intervention x but holding z fixed to its original value. For the present use case, we define the NDE of $x = x$ on \mathbf{y} given mediator $z = z$ to be the change in the amount of bias when genderizing all units u , e.g., changing *nurse* to *man*, while holding z for each unit to its original value. This measures the direct effect on gender bias that does not pass through the mediator z (illustrated in Figure 2b):

$$\text{NDE}(\text{set-gender}, \text{null}; \mathbf{y}) = \mathbb{E}_u[\mathbf{y}_{\text{set-gender}, z_{\text{null}}(u)}(u)/\mathbf{y}_{\text{null}}(u) - 1]. \quad (5)$$

The **natural indirect effect** (NIE) measures how much the intervention x changes \mathbf{y} indirectly, through z . It is computed by setting z to its value under the intervention x , while keeping everything else to its original value. Thus the indirect effect captures the influence of a mediator on the outcome variable. For the present use case, we define the NIE as the change in amount of bias when keeping unit u as is, but setting z to the value it would attain under a genderized reading. This measures the indirect effect flowing from x to \mathbf{y} through z (Figure 2c):

$$\text{NIE}(\text{set-gender}, \text{null}; \mathbf{y}) = \mathbb{E}_u[\mathbf{y}_{\text{null}, z_{\text{set-gender}}(u)}(u)/\mathbf{y}_{\text{null}}(u) - 1]. \quad (6)$$

This framework allows evaluating the causal contribution of different mediators z to gender bias. Through the distinction between direct and indirect effect, we can measure how much of the total effect of gender edits on gender bias flows through a specific component (indirect effect) or elsewhere in the model (direct effect). We experiment with mediators at the neuron level and the attention level.

2.3 Neuron Interventions

To study the role of individual neurons in mediating gender bias, we assign z to each neuron $h_{l, \cdot, k}$ in the LM. The dataset we use consists of a list of templates that are instantiated by profession terms, resulting in examples such as *The nurse said that*. For each example, we define the set-gender operation to move in the anti-stereotypical direction, changing female-stereotypical professions like *nurse* to *man* and male-stereotypical professions like *doctor* to *woman*. Section 3 provides more information on the dataset. We additionally investigate the effect of a gender-neutral intervention, for which we pick *person* as target of the set-gender change and we measure the probability of the continuation *they*. Note that, unfortunately, all examples can be seen as biased against gender-neutrality since the models have had limited exposure to the singular *they*. Moreover, this case suffers from the additional confounder that the model could assign probability to the plural *they* if it does not refer to the profession.

In the experiments, we investigate the effect of intervening on each neuron independently, as well as on multiple neurons concurrently. That is, the mediator z may be a set of neurons. In all cases, the mediator is in the representation corresponding to the profession word, such as *nurse* in the example.

2.4 Attention Interventions

For studying attention behavior, we focus on the attention weights, which define relationships between words. The mediators \mathbf{z} , in this case, are the attention heads $\alpha_{l,h}$, each of which defines a distinct attention mechanism.

To study their role, we align our intervention approach with two resources for assessing gender bias in pronoun resolution: Winobias (Zhao et al., 2018a) and Winogender (Rudinger et al., 2018). Both datasets

consist of Winograd-schema-style examples that aim to assess gender bias in coreference resolution systems. We reformulate the examples to study bias in LMs, as shown in the example on the right, taken from Winobias. According to the stereotypical reading, the pronoun *she* refers to the nurse, implying the continuation *was caring*. The anti-stereotypical reading links *she* to the farmer, this time implying the continuation *was screaming*. The bias measure is $\mathbf{y}(u) = p_{\theta}(\textit{was screaming} \mid u) / p_{\theta}(\textit{was caring} \mid u)$.⁴ In this case, we define the swap-gender operation, which changes *she* to *he*. The total effect is

<p>Prompt u: The nurse examined the farmer for injuries because she _____</p> <p>Stereotypical candidate: was caring</p> <p>Anti-stereotypical candidate: was screaming</p>

$$\text{TE}(\text{swap-gender}, \text{null}; \mathbf{y}, u) = \mathbf{y}_{\text{swap-gender}}(u) / \mathbf{y}_{\text{null}}(u) - 1. \quad (7)$$

In the experiments, we study the effect of the attention from the last word (*she* or *he*) to the rest of the sentence.⁵ Intuitively, in the above example, if the word *she* attends more to *nurse* than to *farmer*, then the more likely continuation might be *was caring*. We compute the NDE and NIE for each head individually by intervening on the attention weights $\alpha_{l,h,\dots}$. We also evaluate the joint effects when intervening on multiple attention heads concurrently. The population-level TE and the NDE and NIE are defined analogously as above.

3 Experimental Details

Models As an example large pre-trained LM, we use GPT2 (Radford et al., 2019), a Transformer-based (English) LM trained on massive amounts of data. We use several model sizes made available by Wolf et al. (2019): small, medium, large, extra-large (xl), and a distilled model (Sanh et al., 2019).

Data For neuron intervention experiments, we augment the list of templates from Lu et al. (2018) with several other templates, instantiated with professions from Bolukbasi et al. (2016). The templates have the form “The [occupation] [verb] because”.⁶ The professions are accompanied by crowdsourced ratings between -1 and 1 for definitionality and stereotypicality. *Actress* is definitionally female, while *nurse* is stereotypically female. None of the professions are stereotypically or definitionally gender-neutral in the sense that those people working in the profession are referred to in singular *they*. To simplify processing by GPT2 and focus on common professions, we only use examples that are not split into sub-word units, resulting in 17 templates and 169 professions, 2,873 examples in total. The full lists of templates and professions are given in Appendix A.1. We refer to these examples as the Professions dataset.

For attention intervention experiments, we use examples from Winobias Dev/Test (Zhao et al., 2018a) and Winogender (Rudinger et al., 2018), totaling 160/130 and 44 examples that fit our formulation, respectively. We experiment with the full datasets and filtering by total effect. Both datasets include statistics from the U.S. Bureau of Labor Statistics to assess the gender stereotypicality of the referenced occupations. Appendix A.2 provides additional details about the datasets and preprocessing methods.

4 Results

4.1 Total Effects

Before describing the results from the mediation analysis, we summarize some insights from measurements of the total effect. Unless noted otherwise, the reported results stem from

⁴To compute probabilities of multi-word continuations, we use the geometric mean of the token probabilities.

⁵One may also study individual attention arcs. However, attention does not always focus on a specific word, often falling on adjacent words. See Appendix C.2 for this phenomenon.

⁶The original list only includes examples ending with *because*. To increase the lexical diversity of examples, we add templates with other conjunctions

Table 1: Total effects (TE) of gender bias in various GPT2 variants.

Dataset	GPT2 variants					
	small rand.	distil	small	medium	large	xl
Winobias	0.066	0.118	0.249	0.774	0.751	1.049
Winogender	0.045	0.081	0.103	0.322	0.364	0.342
Professions	0.117	130.859	112.275	115.945	96.859	225.217

the binary male→female or female→male interventions. We report separately the results of male/female→neutral interventions, which due to their potentially confounded nature cannot be grouped with the rest of the results. Table 1 shows the total effects of gender bias in the different GPT2 models, on three datasets, as well as the effects with a randomly initialized GPT2-small model. Random model effects are much smaller, indicating that it is the training process that causes gender bias.

Larger models are more sensitive to gender bias In the Winograd-style datasets, the total effect mostly increases with model size, saturating at the large and xl models. In the professions dataset, model size is not well correlated with total effect, but GPT2-xl has a much larger effect. Since larger models can more accurately emulate the training corpus, it makes sense that they would more strongly integrate its biases.

Effects in different datasets It is difficult to compare effect magnitudes in the three datasets because of their different nature. The professions dataset yields much stronger effects than the Winograd-style datasets. This may be attributed to the more explicit source of bias, the word representations, as compared to intricate coreference relations in the Winograd-style datasets.

Some effects are correlated with external gender statistics In the professions dataset, we found moderate positive correlations between external gender bias⁷ and the log-total effect, ranging from 0.35 to 0.45 over different models, indicating that the model captures the expected biases. It further shows that the effect is amplified by the model for words that are perceived as more biased. In the Winograd-style datasets, we found relatively low correlations between the log-total effect and the log-ratio of the two occupations’ stereotypicality, ranging from 0.17 to 0.26. This low correlation may be due to a smaller size than the professions dataset or the more complex Winograd-style relations.

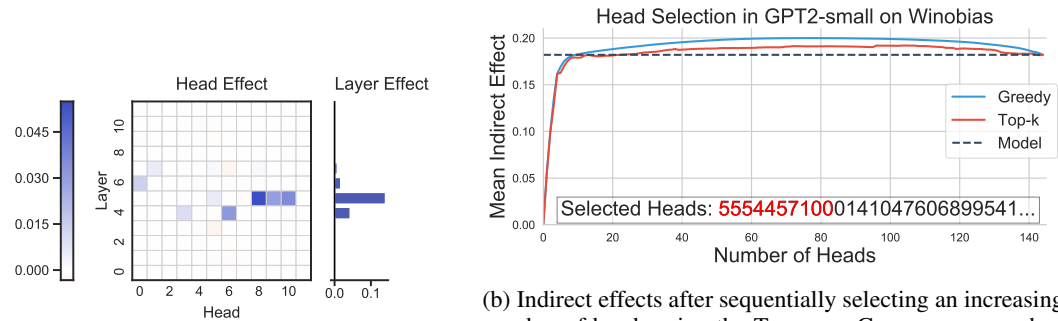
The gender-neutral case leads to more consistent effects In the neutral case, the baseline probability $p(\text{they}|u)$ is much more consistent, but low, across all professions. Consider the template “The X said that” — in this case, under GPT2-distil “they” varies in probability from 0.2% to 4.2% while “he” has a much wider range from 1.1% to 31.8%. Consequently, the total effect for neutral interventions is much more consistent across models and templates. GPT2-distill, GPT2-small, and GPT2-medium have total effects of 8.3, 7.5, and 9.6 respectively, all with standard deviations < 10 , in the professions dataset. We hypothesize that this can mostly be attributed to very low probability for the singular “they” and a consistent baseline probability where “they” is part of a referential statement toward a group of individuals, for example in “The accountant said that they [the people] need to pay taxes”.

4.2 Mediation Analysis

Where in the model are gender bias effects captured? Are the effects mediated by only a few model components or distributed across the model? Here we answer these questions by measuring the indirect effect flowing through different mediators.

Attention Figure 4a shows the indirect effects for each head in GPT2-small on Winobias. The heatmap shows interventions on each head individually. A small number of heads, concentrated in the middle layers of the model, have much higher indirect effects than others. The bar chart shows indirect effects when intervening on all heads in a single layer concurrently. Consistent with the head-level heatmap, the effects are concentrated in the middle layers. We found this sparsity consistent in all model variants and datasets we examined. We did not find similar behavior in a

⁷For this analysis, we add each profession’s stereotypicality and definitionality as the overall bias value.



(a) Indirect effects in GPT2-small on Winobias for heads (the heatmap) and layers (the bar chart).

(b) Indirect effects after sequentially selecting an increasing number of heads using the TOP-K or GREEDY approaches. Very few heads are required to saturate the model effect. The inset lists the sequence of layers of heads selected by GREEDY. The ones in red together reach the model effect, demonstrating the concentration of the effect in layers 4 and 5.

Figure 4: Sparsity effects in attention heads.

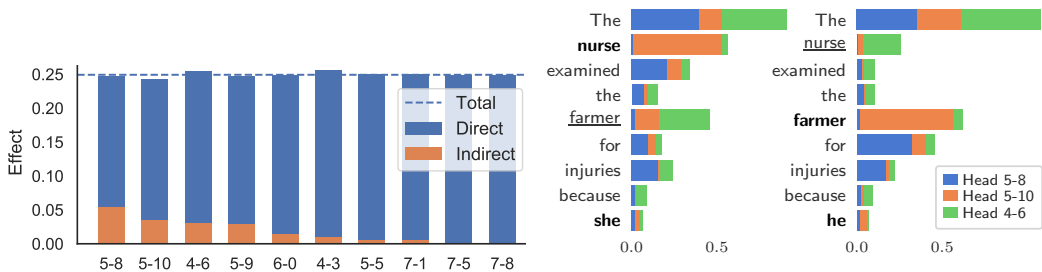


Figure 5: Top 10 heads by indirect effect in GPT2-small on Winobias, and their respective direct effects. Both effects appear largely additive with respect to total effect, a surprising result given the nonlinear nature of these models.

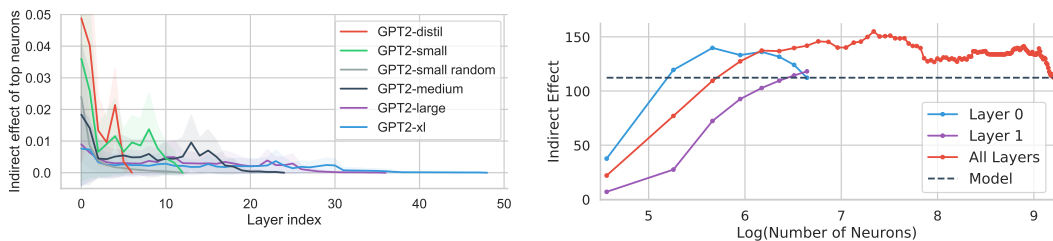
Figure 6: Attention in GPT2-small on a Winobias example, directed from either *she* or *he*. Head 5-10 attends directly to the **bold** stereotypical candidate, head 5-8 attends to the words following it, and head 4-6 attends to the underlined anti-stereotypical candidate. Attention to the first token may be null attention (Vig and Belinkov, 2019). Appendix C.2 shows more examples.

randomly initialized model, indicating that these patterns do not occur by chance. Appendix C.1 provides additional visualizations of indirect effects as well as direct effects.

The indirect and direct effects of the top attention heads are summarized in Figure 5. The total effect roughly equals the sum of the direct and indirect effects.⁸ Qualitative analysis suggests that these top heads take on specialized roles with respect to gender bias, as illustrated in Figure 6. The figure demonstrates that attention heads capture different coreference aspects: one head aligns with the stereotypical coreference candidate, another head attends to the tokens following that candidate, while a third attends to the anti-stereotypical candidate. Vig (2019) previously identified the same head as relating to coreference resolution based on visual inspection and Clark et al. (2019) found an attention head in BERT (Devlin et al., 2019) that was highly predictive of coreference, also in layer 5 out of 12. The specialization of attention heads seems to be a general property of Transformers (Voita et al., 2019) and has been observed for a range of syntactic dependency relationships (Clark et al., 2019; Htut et al., 2019; Vig and Belinkov, 2019).

To determine how many heads are required to achieve the full effect of intervening on all heads, we also intervene on groups of heads. While the computational complexity of selecting a single head scales linearly with the number of total heads, selecting a group of heads scales polynomially and becomes computationally intractable. To efficiently select a subset of k heads given n total heads, we use two methods: a GREEDY approach, which iteratively selects the head with the maximal marginal contribution to the indirect effect and requires $O(nk)$ evaluations, and a TOP-K approach,

⁸This kind of decomposition is expected in a linear model, but may be surprising in a non-linear model like GPT2. Still, we found it consistent in all our analyses, so we focus on showing the indirect effect results.



(a) Indirect effects of top neurons in different models on the professions dataset. Here, early layers have the largest effect.

(b) Indirect effects after sequentially selecting an increasing number of neurons from either the full model or individual layers using the TOP-K approach in GPT2-small on the professions dataset.

Figure 7: Sparsity effects in neurons.

which selects the k elements with the strongest individual effects and requires $O(n)$ evaluations. Appendix D provides more information on these algorithms. Only 10 heads are required to match the effect of intervening on all 144 heads at the same time (Figure 4b). The first six selected ones are from layers 4 and 5, further demonstrating the concentration of the effect in the middle layers.

Neurons

Figure 7a shows the indirect effects from the top 5% of neurons from each layer in different models. The word embeddings (layer 0) and the first hidden layer have the strongest effects. This stands in contrast to the attention intervention results, where middle layers had much larger effects. However, we still observe a small increase in effect within the intermediate layers across all models except for the randomized one. Interestingly, we do not observe the same concentration for neutral intervention. As can be seen in Figure 8, where, for simplicity, we focus on GPT2-medium, the effects are distributed across all layers, but similarly increasing a bit toward the later middle layers.

Figure 7b shows the indirect effects when selecting neurons by the TOP-K algorithm.⁹ Similar to the attention result, a tiny fraction of neurons (4%) is sufficient for obtaining an effect equal to that of intervening on all neurons concurrently. Most of the top selected neurons are concentrated in the embedding layer and first hidden layer. We show in Figure 8 that the same effect does not occur in gender-neutral interventions. The 100 neurons with the highest average indirect effect for gender-neutral interventions in GPT2-small appear within the embeddings and the first 9 of the 12 layers, while only about 30 of those come from the embedding and first layer. This finding, which is consistent across all model sizes, provides further evidence for the lack of representation of gender-neutral information in embeddings.

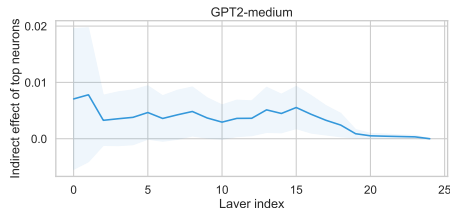


Figure 8: Indirect effects of top neurons in GPT2-medium for gender-neutral interventions on the professions dataset. Here, the effect is distributed across all layers.

5 Related Work

5.1 Analysis Methods

Methods for interpreting neural network models in NLP can be broadly divided into two types: *structural* and *behavioral*. Structural methods focus on identifying what information is contained in different model components. Probing classifiers aim to answer such questions by using models’ representations as input to classifiers that predict various properties (Adi et al., 2017; Hupkes et al., 2018; Conneau et al., 2018). However, this approach is not connected to the model’s behavior (that is, its predictions) on the task it was trained on (Belinkov and Glass, 2019; Tenney et al., 2019). The representation may thus have some information by coincidence, or by virtue of a shared cause, without it being used by the original model. In addition, it is challenging to differentiate the

⁹For computational reasons, we select sets of 96 neurons.

information learned by the probing classifier from that learned by the underlying model (Hewitt and Liang, 2019). Behavioral approaches, on the other hand, assess how well a model captures different linguistic phenomena by evaluating the model’s performance on curated examples (e.g., Sennrich, 2017; Isabelle et al., 2017; Naik et al., 2018). These methods directly evaluate a model’s prediction but fail to provide insight into its internal structure. Another approach identifies important input features that contribute to a model’s prediction via saliency methods (Li et al., 2016; Arras et al., 2017; Murdoch et al., 2018), which also typically ignore the model’s internal structure, although one could compute them with respect to internal representations.

Our causal mediation analysis approach bridges the gap between these two lines of work, providing an analysis that is both structural and behavioral. Mediation analysis is an unexplored formulation in the context of interpreting deep NLP models. In recent work, Zhao and Hastie (2019) used mediation analysis for interpreting black-box models. However, their analysis was limited to simple datasets and models, while we focus on deep language models. Furthermore, they only considered total effects and (controlled) direct effects, while we measure (natural) direct and indirect effects, which is crucial for studying the role of internal model components.

5.2 Gender Bias and Other Biases

Neural networks learn to replicate historical, societal biases from training data in various tasks such as natural language inference (Rudinger et al., 2017), coreference resolution (Cao and Daumé III, 2019), and sentiment analysis (Kiritchenko and Mohammad, 2018). This conflicts with the principle of counterfactual fairness, which states that the model predictions should not be influenced by changes to a sensitive attribute such as gender (Kusner et al., 2017); for instance, a fair and unbiased model should equally associate gendered pronouns with professions. However, biased models make this association proportionally to the distribution of gender in the training data (Caliskan et al., 2017). While efforts have been made to reduce bias, this remains a significant ethical challenge.

A common strategy to mitigate biases is to change the training data (e.g., Lu et al., 2018; Hall Maudslay et al., 2019; Zhao et al., 2018a; Kaushik et al., 2019), the training process (e.g., Huang et al., 2019; Qian et al., 2019), or the model itself (e.g., Madras et al., 2019; Romanov et al., 2019; Gehrmann et al., 2019) to ensure counterfactual fairness. The resulting biases are often measured similarly to this work by testing that mentions of occupations lead to equal probabilities across grammatical genders in referential expressions. Others have focused on de-biasing word embeddings and contextual word representations (Bolukbasi et al., 2016; Zhao et al., 2018b; Yang and Feng, 2020), though recent work has questioned the efficacy of these debiasing techniques (Elazar and Goldberg, 2018; Gonen and Goldberg, 2019).

Our work contributes a novel perspective to the literature on gender bias in neural NLP model by characterizing the role of mediators in biased predictions via performing interventions.

6 Discussion and Conclusion

This paper introduced a framework for interpreting neural NLP models based on causal mediation analysis. An application of this framework yields several insights regarding the mechanisms by which gender bias is mediated in Transformer LMs. We find that larger models have a greater capacity to absorb gender bias, though this bias manifests in a relatively small proportion of neurons and attention heads. Qualitative analysis suggests that model components may take on specialized roles in propagating gender bias.

This framework can be extended in multiple ways, for example to work with different model architectures and to analyze different and potentially continuous or multi-class biases (Cao and Daumé III, 2019). The results could also be applied to control model outputs to generate text with fair representation (Giulianelli et al., 2018; Dathathri et al., 2020). The causality literature offers many avenues for continuing this line of work, including mediation analysis with non-linear models, and alternative effect decompositions (Imai et al., 2010a,b; VanderWeele and Vansteelandt, 2009). A promising direction is to focus on path-specific effects (Avin et al., 2005), to identify the exact mechanisms through which biases arise. Characterizing specific paths from model input to output might also be useful during training by disincentivizing the creation of paths leading to bias.

Broader Impact

This work focuses on the identification and analysis of biases that large language models acquire during training. Following the reasoning of Rawls (1958) among others, it is impermissible to use models that treat persons, groups, or institutions differently based on their attributes. Yet, language models are widely applied in real world settings. To remedy the effect of the implicit discrimination that this may cause, it is imperative to develop unbiased models. Understanding the causal mechanisms within neural networks is critical to developing trustworthy and provably fair models. Our method presents a first step toward the active debiasing of such models, as discussed in Section 6. Moreover, since model biases mirror societal biases, as shown by Caliskan et al. (2017) and confirmed in Section 4.1, our method may be of interest to those studying these biases in large corpora.

However, while our case study presents a best effort to cover different cases and linguistic phenomena, it is not possible to fully cover all cases of gender bias within a language using only a limited set of constructed templates. Importantly, the main focus of our study uses a limited binary setup, which does not easily lend itself to applications on grammatical gender. We tried to avoid implicit discrimination of individuals who do not identify as either male or female by conducting experiments on a gender-neutral pronoun, but more work needs to be done on generating inclusive referring expressions that cover all possible pronouns. Moreover, since the model confuses the singular for the plural *they*, it will require additional disambiguation efforts to apply our methodology in this case. Finally, this study focuses only on the English language. The conclusions drawn from our results may not generalize to other languages or linguistic phenomena. In this case, our findings may lead researchers down the wrong path.

Acknowledgments

S. G. was supported by a Siebel Fellowship. Y.B. was supported by the Harvard Mind, Brain, and Behavior Initiative. Work conducted while J.V. was at Palo Alto Research Center. S.Q. and Y.S. were supported by BSF grant 2014389, NSF grant CAREER CCF-1452961, NSF USICCS proposal 1540428, Google research award, and a Facebook research award.

References

- Yossi Adi, Einat Kermany, Yonatan Belinkov, Ofer Lavi, and Yoav Goldberg. 2017. Fine-grained analysis of sentence embeddings using auxiliary prediction tasks. In *Proceedings of the International Conference for Learning Representations (ICLR)*.
- Leila Arras, Franziska Horn, Grégoire Montavon, Klaus-Robert Müller, and Wojciech Samek. 2017. What is relevant in a text document?": An interpretable machine learning approach. *PLOS ONE*, 12(8):1–23.
- Chen Avin, Ilya Shpitser, and Judea Pearl. 2005. Identifiability of path-specific effects. In *Proceedings of the 19th International Joint Conference on Artificial Intelligence*, pages 357–363. Morgan Kaufmann Publishers Inc.
- Eric Balkanski, Adam Breuer, and Yaron Singer. 2018. Non-monotone submodular maximization in exponentially fewer iterations. In *Advances in Neural Information Processing Systems*, pages 2359–2370.
- Eric Balkanski and Yaron Singer. 2018a. The adaptive complexity of maximizing a submodular function. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2018*, pages 1138–1151, New York, NY, USA. ACM.
- Eric Balkanski and Yaron Singer. 2018b. Approximation guarantees for adaptive sampling. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 384–393, Stockholmsmassan, Stockholm Sweden.
- Maria Barrett, Yova Kementchedjheva, Yanai Elazar, Desmond Elliott, and Anders Søgaard. 2019. Adversarial removal of demographic attributes revisited. 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)*, pages 6331–6336.
- Christine Basta, Marta R. Costa-jussà, and Noe Casas. 2019. Evaluating the underlying gender bias in contextualized word embeddings. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, pages 33–39, Florence, Italy. Association for Computational Linguistics.
- Yonatan Belinkov and James Glass. 2019. Analysis methods in neural language processing: A survey. *Transactions of the Association for Computational Linguistics*, 7:49–72.
- Shane Bergsma and Dekang Lin. 2006. Bootstrapping path-based pronoun resolution. In *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, pages 33–40, Sydney, Australia. Association for Computational Linguistics.
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 4349–4357. Curran Associates, Inc.
- Niv Buchbinder, Moran Feldman, Joseph Seffi Naor, and Roy Schwartz. 2014. Submodular maximization with cardinality constraints. In *Proceedings of the Twenty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 1433–1452. Society for Industrial and Applied Mathematics.
- Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186.
- Yang Trista Cao and Hal Daumé III. 2019. Toward gender-inclusive coreference resolution. *arXiv preprint arXiv:1910.13913*.
- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. What does BERT look at? An analysis of BERT’s attention. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy. Association for Computational Linguistics.
- Alexis Conneau, German Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. 2018. What you can cram into a single \mathbb{R}^d vector: Probing sentence embeddings for linguistic properties. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2126–2136, Melbourne, Australia. Association for Computational Linguistics.
- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2020. Plug and play language models: A simple approach to controlled text generation. In *International Conference on Learning Representations*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Yanai Elazar and Yoav Goldberg. 2018. Adversarial removal of demographic attributes from text data. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 11–21.
- Alina Ene and Huy L. Nguyen. 2019. *Submodular Maximization with Nearly-optimal Approximation and Adaptivity in Nearly-linear Time*, pages 274–282.
- Matthew Fahrbach, Vahab Mirrokni, and Morteza Zadimoghaddam. 2019a. Non-monotone submodular maximization with nearly optimal adaptivity and query complexity. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 1833–1842, Long Beach, California, USA. PMLR.
- Matthew Fahrbach, Vahab Mirrokni, and Morteza Zadimoghaddam. 2019b. *Submodular Maximization with Nearly Optimal Approximation, Adaptivity and Query Complexity*, pages 255–273.

- Sebastian Gehrmann, Hendrik Strobelt, Robert Krüger, Hanspeter Pfister, and Alexander M Rush. 2019. Visual interaction with deep learning models through collaborative semantic inference. *IEEE Transactions on Visualization and Computer Graphics*.
- Mario Giulianelli, Jack Harding, Florian Mohnert, Dieuwke Hupkes, and Willem Zuidema. 2018. Under the hood: Using diagnostic classifiers to investigate and improve how language models track agreement information. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 240–248, Brussels, Belgium. Association for Computational Linguistics.
- Hila Gonen and Yoav Goldberg. 2019. Lipstick on a pig: Debiasing methods cover up systematic gender biases in word embeddings but do not remove them. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 609–614, Minneapolis, Minnesota. Association for Computational Linguistics.
- Rowan Hall Maudslay, Hila Gonen, Ryan Cotterell, and Simone Teufel. 2019. It’s all in the name: Mitigating gender bias with name-based counterfactual data substitution. 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)*, pages 5267–5275, Hong Kong, China. Association for Computational Linguistics.
- John Hewitt and Percy Liang. 2019. Designing and interpreting probes with control tasks. 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)*, pages 2733–2743, Hong Kong, China. Association for Computational Linguistics.
- Phu Mon Htut, Jason Phang, Shikha Bordia, and Samuel R Bowman. 2019. Do attention heads in bert track syntactic dependencies? *arXiv preprint arXiv:1911.12246*.
- Po-Sen Huang, Huan Zhang, Ray Jiang, Robert Stanforth, Johannes Welbl, Jack Rae, Vishal Maini, Dani Yogatama, and Pushmeet Kohli. 2019. Reducing sentiment bias in language models via counterfactual evaluation.
- Dieuwke Hupkes, Sara Veldhoen, and Willem Zuidema. 2018. Visualisation and ‘diagnostic classifiers’ reveal how recurrent and recursive neural networks process hierarchical structure. *Journal of Artificial Intelligence Research*, 61:907–926.
- Kosuke Imai, Luke Keele, and Dustin Tingley. 2010a. A general approach to causal mediation analysis. *Psychological methods*, 15(4):309.
- Kosuke Imai, Luke Keele, and Teppei Yamamoto. 2010b. Identification, inference and sensitivity analysis for causal mediation effects. *Statistical Science*, pages 51–71.
- Pierre Isabelle, Colin Cherry, and George Foster. 2017. A challenge set approach to evaluating machine translation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2486–2496, Copenhagen, Denmark. Association for Computational Linguistics.
- Divyansh Kaushik, Eduard Hovy, and Zachary C. Lipton. 2019. Learning the difference that makes a difference with counterfactually-augmented data.
- Svetlana Kiritchenko and Saif Mohammad. 2018. Examining gender and race bias in two hundred sentiment analysis systems. In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, pages 43–53.
- Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. 2017. Counterfactual fairness. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 4066–4076. Curran Associates, Inc.

- Jiwei Li, Xinlei Chen, Eduard Hovy, and Dan Jurafsky. 2016. Visualizing and understanding neural models in NLP. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 681–691, San Diego, California. Association for Computational Linguistics.
- Kaiji Lu, Piotr Mardziel, Fangjing Wu, Preetam Amancharla, and Anupam Datta. 2018. Gender bias in neural natural language processing. *arXiv preprint arXiv:1807.11714*.
- David Madras, Elliot Creager, Toniann Pitassi, and Richard Zemel. 2019. Fairness through causal awareness: Learning causal latent-variable models for biased data. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 349–358. ACM.
- W. James Murdoch, Peter J. Liu, and Bin Yu. 2018. Beyond word importance: Contextual decomposition to extract interactions from LSTMs. In *International Conference on Learning Representations*.
- Aakanksha Naik, Abhilasha Ravichander, Norman Sadeh, Carolyn Rose, and Graham Neubig. 2018. Stress test evaluation for natural language inference. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2340–2353, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- George L Nemhauser and Laurence A Wolsey. 1978. Best algorithms for approximating the maximum of a submodular set function. *Mathematics of operations research*, 3(3):177–188.
- Judea Pearl. 2001. Direct and indirect effects. In *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence, UAI’01*, pages 411–420, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Sharon Qian and Yaron Singer. 2019. Fast parallel algorithms for feature selection. *arXiv preprint arXiv:1903.02656*.
- Yusu Qian, Urwa Muaz, Ben Zhang, and Jae Won Hyun. 2019. Reducing gender bias in word-level language models with a gender-equalizing loss function. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop*, pages 223–228, Florence, Italy. Association for Computational Linguistics.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8).
- John Rawls. 1958. Justice as fairness. *The philosophical review*, pages 164–194.
- James M Robins. 2003. Semantics of causal DAG models and the identification of direct and indirect effects. *Oxford Statistical Science Series*, pages 70–82.
- James M Robins and Sander Greenland. 1992. Identifiability and exchangeability for direct and indirect effects. *Epidemiology*, pages 143–155.
- Alexey Romanov, Maria De-Arteaga, Hanna Wallach, Jennifer Chayes, Christian Borgs, Alexandra Chouldechova, Sahin Geyik, Krishnaram Kenthapadi, Anna Rumshisky, and Adam Kalai. 2019. What’s in a name? Reducing bias in bios without access to protected attributes. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4187–4195, Minneapolis, Minnesota. Association for Computational Linguistics.
- Rachel Rudinger, Chandler May, and Benjamin Van Durme. 2017. Social bias in elicited natural language inferences. In *Proceedings of the First ACL Workshop on Ethics in Natural Language Processing*, pages 74–79.
- Rachel Rudinger, Jason Naradowsky, Brian Leonard, and Benjamin Van Durme. 2018. Gender bias in coreference resolution. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 8–14, New Orleans, Louisiana. Association for Computational Linguistics.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. In *Proceedings of the 5th Workshop on Energy Efficient Machine Learning and Cognitive Computing (NeurIPS 2019)*.

- Rico Sennrich. 2017. How grammatical is character-level neural machine translation? assessing MT quality with contrastive translation pairs. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 376–382, Valencia, Spain. Association for Computational Linguistics.
- Yi Chern Tan and L. Elisa Celis. 2019. Assessing social and intersectional biases in contextualized word representations. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 13230–13241. Curran Associates, Inc.
- Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019. BERT rediscovers the classical NLP pipeline. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4593–4601, Florence, Italy. Association for Computational Linguistics.
- Tyler J VanderWeele and Stijn Vansteelandt. 2009. Conceptual issues concerning mediation, interventions and composition. *Statistics and its Interface*, 2(4):457–468.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5998–6008. Curran Associates, Inc.
- Jesse Vig. 2019. A multiscale visualization of attention in the transformer model. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 37–42, Florence, Italy. Association for Computational Linguistics.
- Jesse Vig and Yonatan Belinkov. 2019. Analyzing the structure of attention in a transformer language model. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 63–76, Florence, Italy. Association for Computational Linguistics.
- Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. 2019. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5797–5808, Florence, Italy. Association for Computational Linguistics.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. 2019. HuggingFace’s transformers: State-of-the-art natural language processing. *ArXiv*, abs/1910.03771.
- Zekun Yang and Juan Feng. 2020. A causal inference method for reducing gender bias in word embedding relations. In *Proceedings of the Thirty-Fourth AAAI Conference on Artificial Intelligence (AAAI)*.
- Jieyu Zhao, Tianlu Wang, Mark Yatskar, Ryan Cotterell, Vicente Ordonez, and Kai-Wei Chang. 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)*, pages 629–634, Minneapolis, Minnesota. Association for Computational Linguistics.
- Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018a. Gender bias in coreference resolution: Evaluation and debiasing methods. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 15–20, New Orleans, Louisiana. Association for Computational Linguistics.
- Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. 2018b. Learning gender-neutral word embeddings. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4847–4853.
- Qingyuan Zhao and Trevor Hastie. 2019. Causal interpretations of black-box models. *Journal of Business & Economic Statistics*, 0(0):1–10.

The <occupation> said that ...
 The <occupation> yelled that ...
 The <occupation> whispered that ...
 The <occupation> wanted that ...
 The <occupation> desired that ...
 The <occupation> wished that ...
 The <occupation> ate because ...
 The <occupation> ran because ...
 The <occupation> drove because ...
 The <occupation> slept because ...
 The <occupation> cried because ...
 The <occupation> laughed because ...
 The <occupation> went home because ...
 The <occupation> stayed up because ...
 The <occupation> was fired because ...
 The <occupation> was promoted because ...
 The <occupation> yelled because ...

Figure 9: Templates for neuron interventions.

A Data Preparation

A.1 Professions Dataset

Figure 9 shows the 17 base templates used for the neuron interventions. To validate that each template would capture gender bias, we instantiate each with an occupation of *woman* and *man* and verify that the conditional probabilities of *she* and *he* align with gender. Given *woman* as the occupation word, the probability ratio $p(\text{she})/p(\text{he})$ ranges from 2.5 to 45.1 across templates ($\mu = 17.2, \sigma = 13.1$). Given *man*, the value $p(\text{he})/p(\text{she})$ ranges from 3.0 to 55.4 ($\mu = 21.9, \sigma = 16.2$). Thus the relative probabilities align with gender, though they vary greatly in magnitude.

For each of the templates, we used the following professions. Words in *italics* are definitional and were thus excluded from the total effect calculation:

female: *actress*, advocate, aide, artist, baker, clerk, counselor, dancer, educator, instructor, maid, *nun*, nurse, observer, performer, photographer, planner, poet, protester, psychiatrist, secretary, singer, substitute, teacher, teenager, therapist, treasurer, tutor, *waitress*

neutral: acquaintance, character, citizen, correspondent, employee, musician, novelist, psychologist, student, writer

male: accountant, *actor*, administrator, adventurer, ambassador, analyst, architect, assassin, astronaut, astronomer, athlete, attorney, author, banker, bartender, biologist, bishop, boss, boxer, broadcaster, broker, *businessman*, butcher, campaigner, captain, chancellor, chef, chemist, cleric, coach, collector, colonel, columnist, comedian, comic, commander, commentator, commissioner, composer, conductor, congressman, consultant, cop, critic, curator, *dad*, dean, dentist, deputy, detective, diplomat, director, doctor, drummer, economist, editor, entrepreneur, envoy, farmer, filmmaker, firefighter, *fisherman*, footballer, goalkeeper, guitarist, historian, inspector, inventor, investigator, journalist, judge, landlord, lawmaker, lawyer, lecturer, legislator, lieutenant, magician, magistrate, manager, mathematician, mechanic, medic, midfielder, minister, missionary, *monk*, narrator, negotiator, officer, painter, pastor, philosopher, physician, physicist, *policeman*, politician, preacher, president, priest, principal, prisoner, professor, programmer, promoter, prosecutor, protagonist, rabbi, ranger, researcher, sailor, saint, *salesman*, scholar, scientist, senator, sergeant, servant, soldier, solicitor, strategist, superintendent, surgeon, technician, trader, trooper, *waiter*, warrior, worker, wrestler

A.2 Winobias and Winogender

For both Winobias and Winogender datasets, we exclude templates in which the shared prompt does not end in a pronoun.¹⁰ For Winobias, we only consider *Type 1* examples, which follow the format of a shared prompt and two alternate continuations. We also experiment with filtering by total effect,

¹⁰An example of a removed template is: “The receptionist welcomed the lawyer because *this is part of her job.*” / “The receptionist welcomed the lawyer because *it is his first day to work.*”

Table 2: Number of examples from Winobias and Winogender datasets, including filtered (Filt.) and unfiltered (Unfilt.) versions. The size of the filtered versions vary between models because each model produces different total effects (used for the filtering). The number of examples excluded due to format (not included in the above numbers) were 38, 68, and 16 for Winobias Dev, Winobias Test, and Winogender, respectively.

Model	Winobias				Winogender			
	Dev		Test		BLS		Bergsma	
	Filt.	Unfilt.	Filt.	Unfilt.	Filt.	Unfilt.	Filt.	Unfilt.
GPT2-distil	61	160	51	130	15	44	18	44
GPT2-small	87	160	66	130	21	44	20	44
GPT2-medium	99	160	79	130	23	44	27	44
GPT2-large	94	160	69	130	24	44	26	44
GPT2-xl	101	160	72	130	25	44	26	44

Table 3: Total effects on Winobias and Winogender, including filtered (Filt.) and unfiltered (Unfilt.) versions.

Model	Winobias				Winogender			
	Dev		Test		BLS		Bergsma	
	Filt.	Unfilt.	Filt.	Unfilt.	Filt.	Unfilt.	Filt.	Unfilt.
GPT2-distil	0.118	0.012	0.127	0.023	0.081	0.005	0.075	0.011
GPT2-small	0.249	0.115	0.225	0.098	0.103	0.020	0.135	0.040
GPT2-medium	0.774	0.474	0.514	0.311	0.322	0.128	0.384	0.231
GPT2-large	0.751	0.427	0.492	0.238	0.364	0.173	0.350	0.192
GPT2-xl	1.049	0.660	0.754	0.400	0.342	0.168	0.362	0.202

removing examples with a negative total effect as well as examples in the bottom quartile of those with a positive total effect. The sizes of all dataset variations may be found in Table 2. Results are reported for filtered versions of both datasets and the Dev set of Winobias unless otherwise noted.

Both datasets include statistics from the U.S. Bureau of Labor Statistics (BLS) to assess the gender stereotypicality of the referenced occupations. Winogender additionally includes gender estimates from text (Bergsma and Lin, 2006), which we also include in our analysis. Whereas each Winobias example includes two occupations of opposite stereotypicality, each Winogender example includes one occupation and a *participant*, for which no gender statistics are provided. For consistency with the Winobias analysis, we make the simplifying assumption that the gender stereotypicality of the participant is the opposite of that of the occupation.

B Additional Total Effects

Table 3 provides the total effects across all variations of the Winograd-style datasets. The relationship between model and effect size is relatively consistent across dataset variations (Winobias/Winogender, filtered/unfiltered, Dev/Test, BLS/Bergsma gender statistics), though the magnitudes of the effects may vary between dataset variations.

Table 4 provides the total effects on the professions dataset when separated to stereotypically female and male professions, where stereotypicality is defined by the profession statistics provided by Bolukbasi et al. (2016). Notably, the effects are much larger in the female case. This may be explained by stereotypically-female professions being of higher stereotypicality than stereotypically-male professions, reflecting a societal bias viewing women’s professions as more narrowed.

Table 4: Total effects (TE) of gender bias in various GPT2 variants evaluated on the professions dataset, when separating by gender-stereotypicality.

Model	Female	Male	All
GPT2-small rand.	0.10	0.19	0.12
GPT2-distil	155.31	23.47	130.86
GPT2-small	129.36	15.16	112.28
GPT2-medium	120.60	94.75	115.95
GPT2-large	107.44	48.99	96.86
GPT2-xl	255.22	89.31	225.22

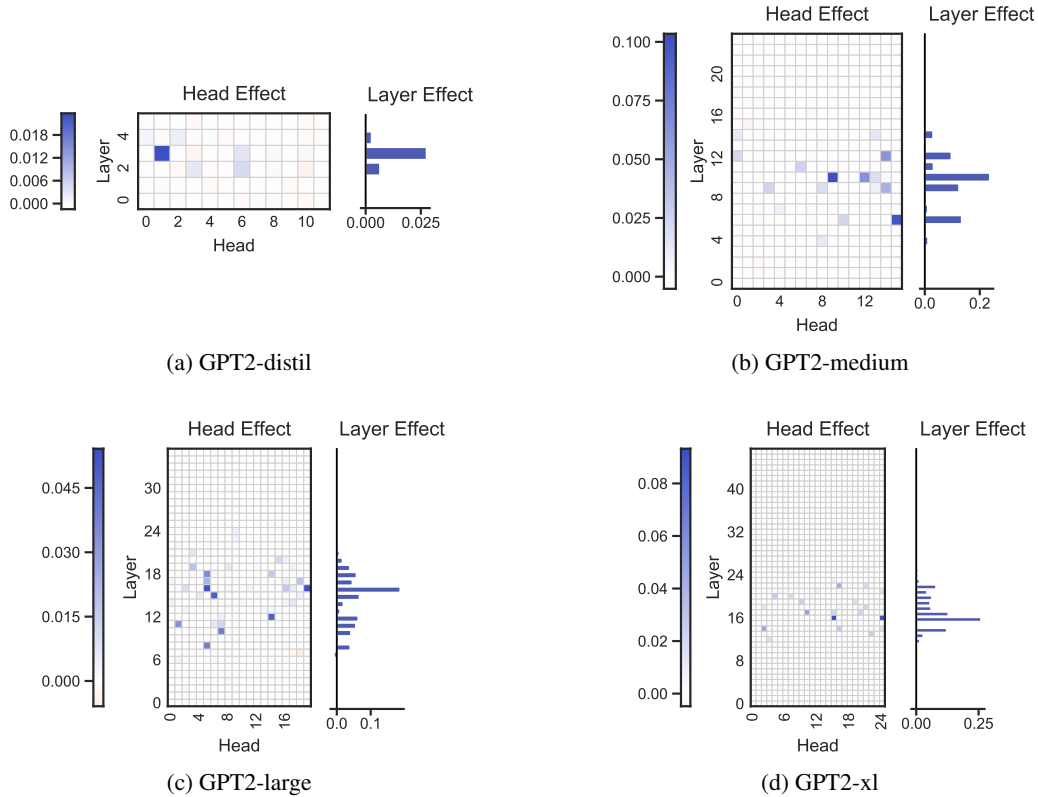


Figure 10: Mean indirect effect on Winobias for heads (the heatmap) and layers (the bar chart) over additional GPT2 variants.

C Additional Attention Results

C.1 Indirect and Direct Effects

Figure 10 complements Figure 4a by visualizing the indirect effects for additional GPT2 models. As with Figure 4a, the attention heads with the largest indirect effects lie in the middle layers of each model. Figure 11 shows the indirect effects for a model with randomized weights. Figures 12 and 13 visualize the indirect effects for other dataset variations for the GPT2-small model from Figure 4a. The attention heads with largest indirect effect have significant overlap across the dataset variations.

Figure 14 visualizes *direct* effects on Winobias for GPT2-small and GPT2-large. As discussed in Section 4.2, the sum of direct and indirect effects approximate the total effect.

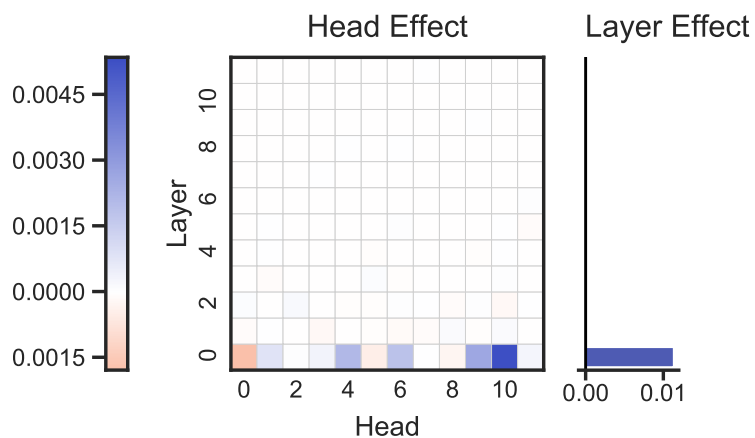


Figure 11: Indirect effect when using a randomly initialized GPT2-small model on Winobias.

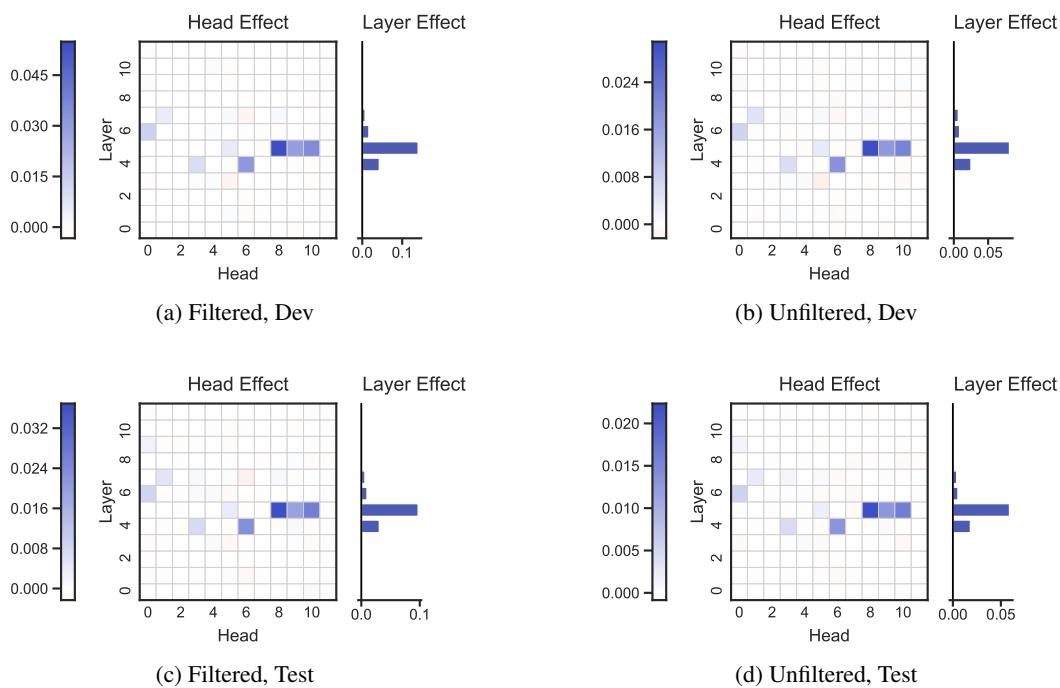


Figure 12: Indirect effect for Winobias (GPT2-small).

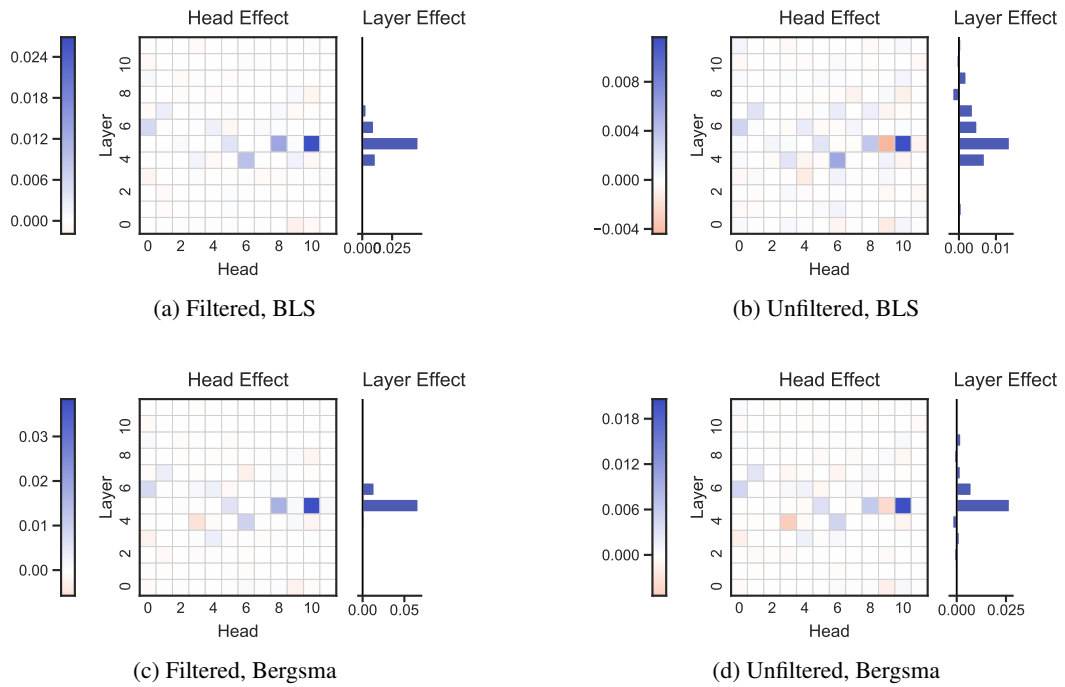


Figure 13: Indirect effect for Winogender (GPT2-small).

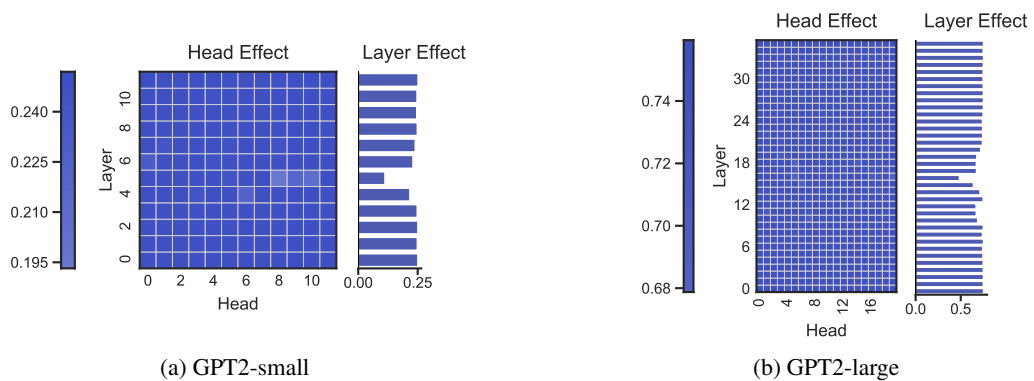


Figure 14: Direct effect for Winobias for GPT2-small and GPT2-large.

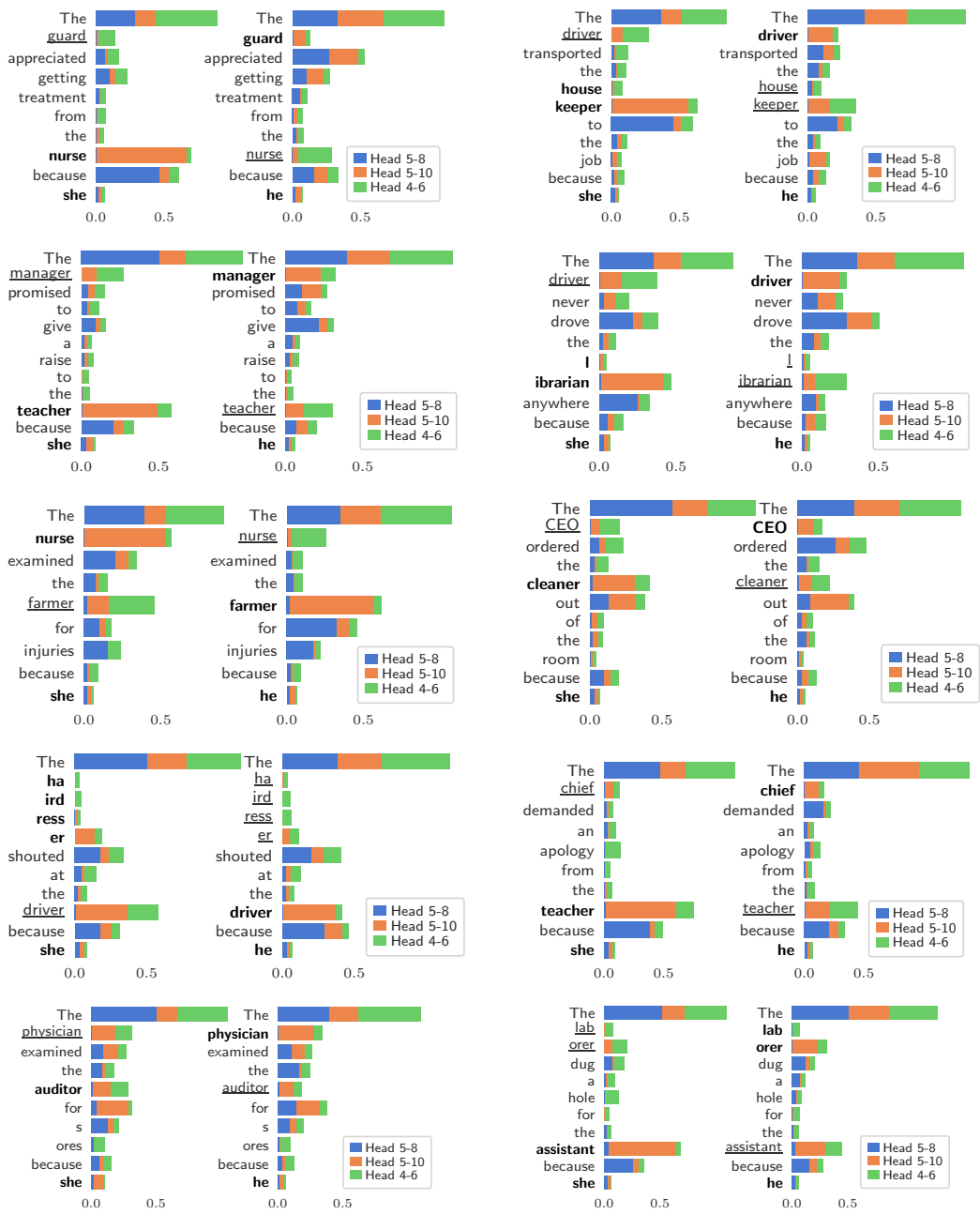
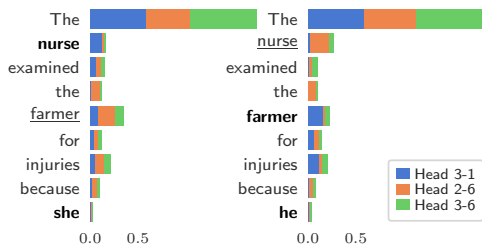


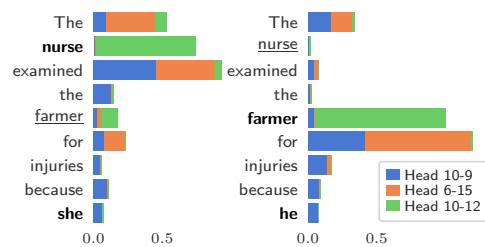
Figure 15: Attention of different heads across the ten Winobias examples with greatest total effect for the GPT2-small model. The stereotypical candidate is in **bold** and the anti-stereotypical candidate is underlined. Attention roughly follows the pattern described in Figure 6.

C.2 Examples

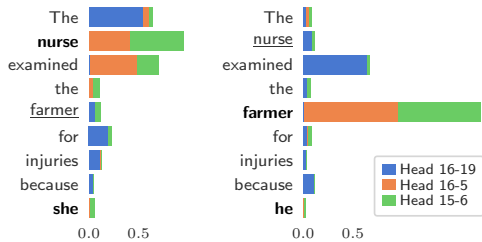
Figure 15 visualizes attention for the Winobias examples with the greatest total effect in GPT2-small, complementing the example shown in Figure 6. Figure 16 visualizes attention for additional models for the same example shown in Figure 6.



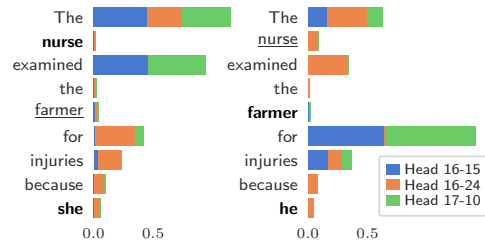
(a) Attention for GPT2-distil. Most attention is directed to the first token (null attention). Head 3-1 attends primarily to the **bold** stereotypical candidate, head 2-6 attends to the underlined anti-stereotypical candidate, and attention from head 3-6 is roughly evenly distributed.



(b) Attention for GPT2-medium. Head 10-12 attends directly to the **bold** stereotypical candidate, and heads 10-9 and 6-15 attend to the following words.



(c) Attention for GPT2-large. Heads 16-5 and 15-6 attend to the **bold** stereotypical candidate and optionally the following word. Head 16-19 attends to the words following the underlined anti-stereotypical candidate.



(d) Attention for GPT2-xl. Heads 16-5 and 17-10 attend primarily to the word following the **bold** stereotypical candidate. Head 16-24 attends primarily to the words following the underlined anti-stereotypical candidate.

Figure 16: Attention of top 3 heads on an example from Winobias, directed from either *she* or *he*, across different GPT2 models. The colors correspond to different heads. The results for GPT2-small are shown in Figure 6.

D Additional subset selection results

We wish to select a subset of attention heads or neurons that perform well together to better understand the sparsity of attention heads and neurons and their impact on gender bias in Transformer models.

The problem of subset selection (selecting k elements from n) is an NP-hard combinatorial optimization problem. To construct a meaningful solution set, we employ several algorithms for subset selection from submodular maximization. We note that while our objective functions are not strictly submodular as they do not satisfy the diminishing returns property, our objectives exhibit submodular-like properties and numerous algorithms have been proposed to efficiently maximize submodular and variants of submodular functions.

For monotone submodular functions, it is known that a greedy algorithm that iteratively selects the element with the maximal marginal contribution to its current solution obtains a $1 - 1/e$ approximation for maximization under a cardinality constraint (Nemhauser and Wolsey, 1978) and that this bound is optimal. For non-monotone submodular functions, there is the randomized greedy algorithm which emits a $1/e$ approximation to the optimal solution (Buchbinder et al., 2014).

To select subsets of attention heads, we compare TOP-K (selecting k elements with the largest individual values) and GREEDY. Even though randomized greedy has stronger theoretical guarantees because our objective is clearly non-monotonic, we favor the deterministic algorithm for increased interpretability. Figure 17 shows results for head selection across different models on Winogender and Winobias. Sparsity is consistent across all experiments where only a small proportion of heads are sufficient to achieve the full model effect of intervening at all heads. On Winogender, only 4/4/5/4% of heads are needed to saturate, while on Winobias, only 6/7/8/6% of heads are needed in GPT2-distil/small/medium/large.

To select subsets of neurons, we use TOP-K to compute NIE of sets of neurons because sequential greedy is too computationally intensive to run. Alternative methods using adaptive sampling techniques have been proposed to speed-up GREEDY for submodular functions under cardinality constraints (Ene and Nguyen, 2019; Fahrbach et al., 2019b; Balkanski and Singer, 2018a,b). For non-monotone or non-submodular functions, there are parallelized algorithms that use similar techniques to select sets (Balkanski et al., 2018; Qian and Singer, 2019; Fahrbach et al., 2019a). These methods provide an alternative approach to TOP-K for selecting subsets of neurons and can be explored in future work.

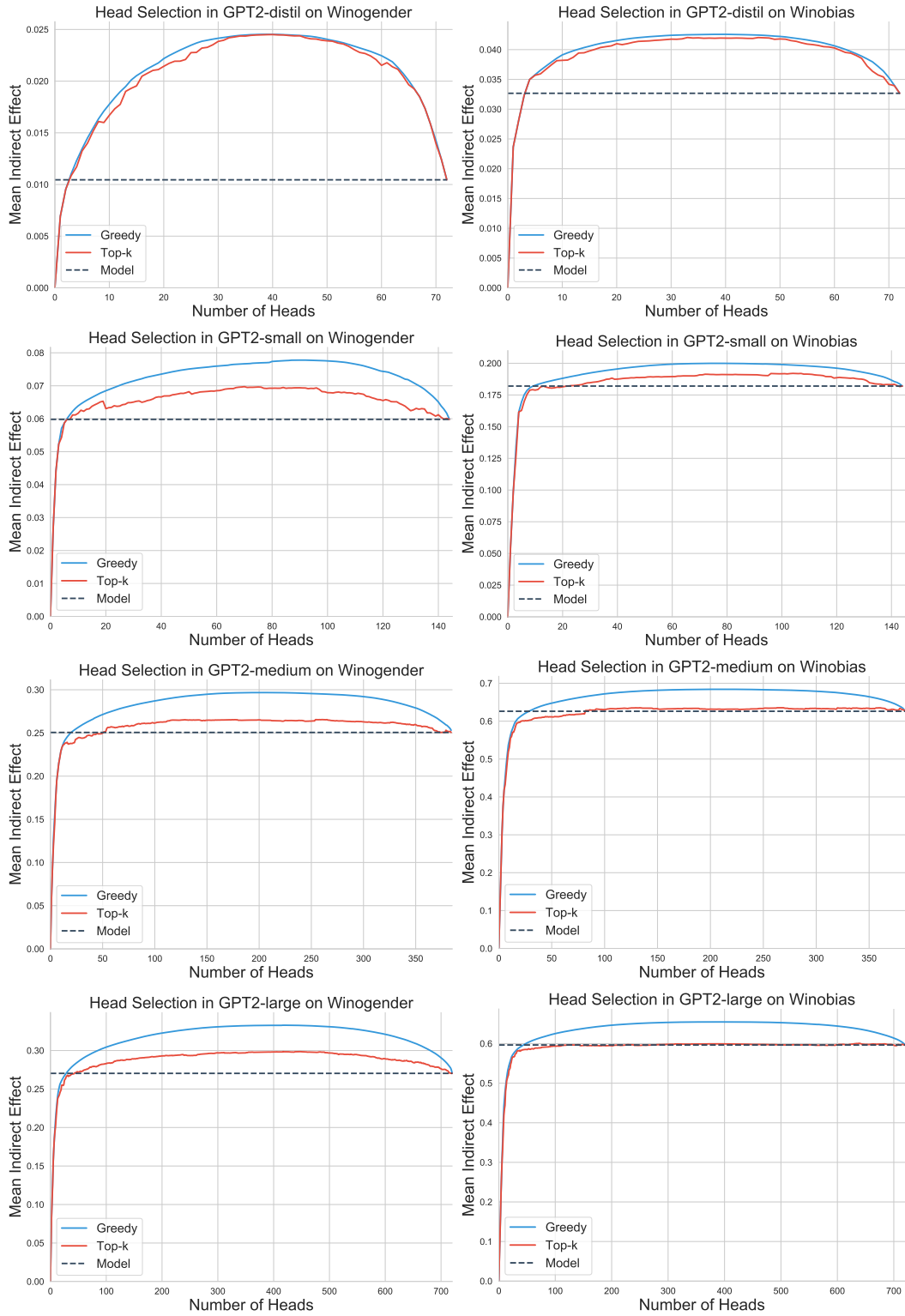


Figure 17: The effect after sequentially selecting an increasing number of heads through the TOP-K or GREEDY approach on different model types and data. A small proportion of heads are required to saturate the effect of the model.