

Causal Mediation Analysis for Interpreting Neural NLP: The Case of Gender Bias

Jesse Vig*

Salesforce Research

JVIG@SALESFORCE.COM

Sebastian Gehrmann*

Harvard University

GEHRMANN@SEAS.HARVARD.EDU

Yonatan Belinkov*

Technion – Israel Institute of Technology

BELINKOV@TECHNION.AC.IL

Sharon Qian

Harvard University

SHARONQIAN@SEAS.HARVARD.EDU

Daniel Nevo

Tel Aviv University

DANIELNEVO@TAU.EX.TAU.AC.IL

Simas Sakenis

SIMASSAKENIS@COLLEGE.HARVARD.EDU

Jason Huang

JASONHUANG@COLLEGE.HARVARD.EDU

Yaron Singer

YARON@SEAS.HARVARD.EDU

Stuart Shieber

Harvard University

SHIEBER@SEAS.HAVARD.EDU

Abstract

Common methods for interpreting neural models in natural language processing typically examine either their structure or their predictions, but not both. 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. It enables us to analyze the mechanisms by which information flows from input to output through various model components, known as mediators. We apply this methodology in a case study of gender bias in pre-trained Transformer language models. We analyze the role of individual neurons and attention heads in mediating gender bias across three datasets designed to gauge a model’s sensitivity to grammatical gender. Our mediation analysis reveals that gender bias effects are (i) sparse, concentrated in a small part of the network; (ii) synergistic, amplified or repressed by different components; and (iii) decomposable into effects flowing directly from the input and indirectly through the mediators.

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. Analysis methods may be categorized into structural and behavioral analyses (Tenney, Das, & Pavlick, 2019). Structural analyses aim to shed light on the internal structure of a neural model, for example through probing classifiers (Conneau, Kruszewski, Lample, Barrault, & Baroni, 2018; Hupkes, Veldhoen, & Zuidema, 2018; Adi, Kermay, Belinkov, Lavi, & Goldberg, 2017) that predict linguistic properties using representations from trained models. This methodology has been used for analyzing sentence

*. Equal contribution. Work conducted while J.V. was at Palo Alto Research Center.

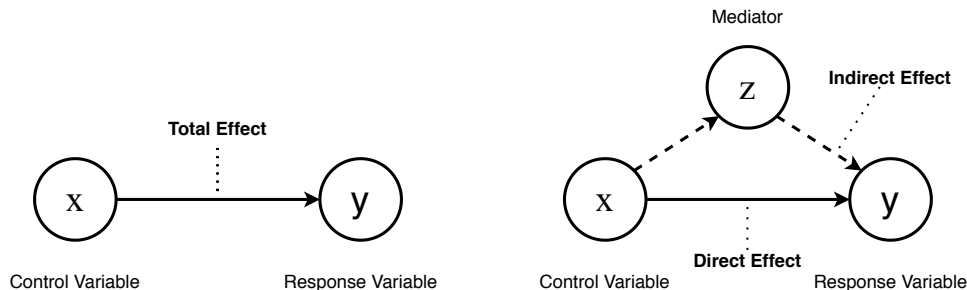


Figure 1: Causal mediation analysis introduces a mediator into the causal graph shown on the left. The mediator decouples the total effect into a direct and indirect effect.

embeddings, machine translation models, and contextual word representation models, among other models (Belinkov & Glass, 2019). *Behavioral* analyses, on the other hand, aim to assess a model’s behavior by its performance on constructed examples (e.g., Isabelle, Cherry, & Foster, 2017; Naik, Ravichander, Sadeh, Rose, & Neubig, 2018,), or by visualizing important input features via saliency methods (e.g., Li, Chen, Hovy, & Jurafsky, 2016; Murdoch, Liu, & Yu, 2018,).

Despite yielding interesting and useful insights, both types of analyses suffer from significant limitations. As pointed out by Belinkov and Glass (2019), probing classifiers only yield a correlational measure between a model’s representations and an external linguistic property, and are thus not *causally* connected to the model’s predictions. Barrett, Kementchedjhieva, Elazar, Elliott, and Søgaard (2019) further demonstrate that classifiers that aim to detect biases in learned representations focus on spurious correlations in their training data and fail to generalize to unseen data. Probing classifiers may thus fail to provide faithful interpretations. On the other hand, while behavioral analyses directly evaluate model predictions, they do not typically link them to the model’s internal structure.

To address these limitations, we introduce a methodology for interpreting neural NLP models based on *causal mediation analysis* (Pearl, 2001). Causal mediation analysis is a method from *causal inference*, which studies the change in a response variable following an intervention, or treatment, e.g., the health outcome of a drug treatment in a clinical trial. Causal mediation analysis (Figure 1) extends this approach by considering the indirect effect of intermediaries, or *mediators*, on the final outcome—e.g., a drug treatment causes headaches, which cause subjects to take aspirin (mediator), which in turn impacts the health outcome. We use mediation analysis to interpret neural networks by treating internal model components, e.g., specific neurons or attention heads, as mediators between model inputs and model outputs. We propose multiple controlled interventions on the model inputs and mediators, which reveal the causal role of specific components in a model’s behavior.

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, both static word embeddings (Caliskan, Bryson, & Narayanan, 2017; Bolukbasi, Chang, Zou, Saligrama, & Kalai, 2016) and contextualized word representations (Zhao, Wang, Yatskar, Cotterell, Ordonez, & Chang, 2019a; Basta, Costa-jussà, & Casas, 2019; Tan & Celis, 2019). We study how grammatical gender bias effects are mediated via different model components in Transformer-based language models, primarily several versions of GPT2 (Radford, Wu, Child, Luan, Amodei, & Sutskever, 2019), focusing on the role of individual neurons or attention heads in mediating these effects.

Our approach is a structural-behavioral analysis. It is structural in that our results highlight internal model components that are responsible for gender bias. It is behavioral in that said components are causally implicated in how gender bias manifests in the model outputs. In an experimental evaluation using several datasets designed to gauge a model’s gender bias, we find that larger models show larger gender bias effects, potentially absorbing more bias from the underlying training data. The causal mediation analysis further yields several insights regarding the role of different model components in mediating gender bias:

- Gender bias is *sparse*: Much of the effect is concentrated in relatively few model components.
- Gender bias is *synergistic*: Some model components interact to produce mutual effects that amplify their individual effects. Other components operate relatively independently, capturing complementary aspects of gender bias.
- Gender bias is *decomposable*: The total gender bias effect approximates the sum of the direct and indirect effect, a surprising result given the non-linear nature of the model.

We use GPT2 as a primary model for testing our framework, but also perform select analyses with two additional autoregressive models and three masked language models. Our experiments confirm that the insights outlined above apply to all three of the autoregressive models and, albeit to a lesser extent, to masked language models as well. This finding indicates that the causal mediation analysis framework is capable of capturing general patterns of causal structure in Transformer architectures, as opposed to only revealing model-specific characteristics.

In summary, this article 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>.¹

2. Related Work

2.1 Analysis Methods

Methods for interpreting neural network models in NLP can be broadly divided into two kinds. 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 (i.e., its predictions) on the task it was trained on (Belinkov & Glass, 2019; Tenney et al., 2019). The representation may thus have some information by coincidence, without it being used by the original model. In addition, it is challenging to differentiate the information learned by the probing classifier from that learned by the underlying model (Hewitt & Liang, 2019). Similarly, interactive methods can be useful to

1. This article expands upon our conference paper (Vig, Gehrmann, Belinkov, Qian, Nevo, Singer, & Shieber, 2020) in the following ways: (a) the conference paper only studied sparsity, while here we also study synergism and decomposition; (b) we extend the analysis to other models besides GPT2; (c) we consider various bias metrics; and (d) we draw broader connections to the causality literature.

identify network components that capture specific properties by relating them to similar training examples (Strobelt, Gehrmann, Pfister, & Rush, 2017).

An alternative approach is to assess how well a model captures different linguistic phenomena by evaluating its performance on curated examples (e.g., Sennrich, 2017; Isabelle et al., 2017; Naik et al., 2018,). This approach directly evaluates a model’s predictions but does not provide insight into the roles that the internal structure of the network played in arriving at the prediction. Another approach identifies important input features that contribute to a model’s prediction via saliency methods (Li et al., 2016; Arras, Horn, Montavon, Müller, & Samek, 2017; Murdoch et al., 2018), which typically ignore the model’s internal structure, although they may in principle be computed with respect to internal representations (Gehrmann, Dernoncourt, Li, Carlson, Wu, Welt, Foote Jr, Moseley, Grant, Tyler, et al., 2018; Montavon, Binder, Lapuschkin, Samek, & Müller, 2019).

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.

Causal approaches for interpreting models have very recently begun to be explored in NLP. Giulianelli, Harding, Mohnert, Hupkes, and Zuidema (2018) use gradients from probing classifiers to update recurrent language model hidden states, and study the effect of such an intervention on a subject-verb agreement task. Elazar, Ravfogel, Jacovi, and Goldberg (2020) remove linguistic information from Transformer hidden states and evaluate the effect of such removal on language modeling, effectively performing a sort of intervention at the layer level. Feder, Oved, Shalit, and Reichart (2020) add auxiliary adversarial tasks to language models in order to learn counterfactual representations with respect to a given concept. Our work is distinguished from these lines of research by our focus on mediation analysis and measurement of direct and indirect effect. Our approach is complementary to the counterfactual representations of Feder et al. (2020), which may be integrated in our causal mediation analysis.

This proposed strategy is rooted in theories of directed acyclic graphs (DAGs) from the causality literature, which consider DAGs as either the output of causal discovery (Pearl, 2009; Spirtes, Glymour, Scheines, & Heckerman, 2000) or as a means to encode prior knowledge to inform the design of analysis methods, e.g., variable selection methods (Pearl, 1995). Our approach views the neural network itself as a DAG, with a common *ancestor* to all nodes, the input, and a common *descendent* to all nodes, the output. Thus, for each input unit (e.g., a sentence), we can estimate the causal effect of an intervention (e.g., a text edit) by comparing the model output under the intervention to the output given the original input. Repeating this modification for a number of units and averaging over the obtained effects resembles the process of studying average causal effects in the population, say of a drug to treat a disease, with one major advantage. The so-called *fundamental problem of causal inference* (Holland, 1986) says that we cannot observe the counterfactual of two different interventions on the same unit, e.g., one cannot know what would have happened to a person had they been treated with drug B when in reality they were treated with drug A. Our adaption of the causal language and framework does not suffer from this problem, as we can manufacture outputs from the model under any conceivable interventions on the units.

2.2 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, May, & Van Durme, 2017), coreference resolution (Cao & Daumé III, 2020), and sentiment analysis (Kiritchenko & 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, Loftus, Russell, & Silva, 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, Mardziel, Wu, Amancharla, & Datta, 2020; Hall Maudslay, Gonen, Cotterell, & Teufel, 2019; Zhao, Wang, Yatskar, Ordonez, & Chang, 2018a; Kaushik, Hovy, & Lipton, 2019,), the training process (e.g., Huang, Zhang, Jiang, Stanforth, Welbl, Rae, Maini, Yogatama, & Kohli, 2020; Qian, Muaz, Zhang, & Hyun, 2019,), or the model itself (e.g., Madras, Creager, Pitassi, & Zemel, 2019; Romanov, De-Arteaga, Wallach, Chayes, Borgs, Chouldechova, Geyik, Kenthapadi, Rumshisky, & Kalai, 2019; Gehrmann, Strobel, Krüger, Pfister, & Rush, 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 (Bolkunov et al., 2016; Zhao, Zhou, Li, Wang, & Chang, 2018b; Yang & Feng, 2020), though recent work has questioned the efficacy of these debiasing techniques in removing both grammatical and societal biases (Elazar & Goldberg, 2018; Gonen & Goldberg, 2019). Biases may also be introduced in downstream tasks and representations in models where representations depend on additional context (Zhao, Wang, Yatskar, Cotterell, Ordonez, & Chang, 2019b; Kurita, Vyas, Pareek, Black, & Tsvetkov, 2019).

3. Methodology

3.1 Preliminaries

Consider a large pre-trained neural language model (LM), parameterized by θ , which predicts the probability of the next word given a prefix: $p_\theta(x_t \mid x_1, \dots, x_{t-1})$. We will focus on LMs based on Transformers (Vaswani, Shazeer, Parmar, Uszkoreit, Jones, Gomez, Kaiser, & Polosukhin, 2017), although much of the methodology will apply to other architectures as well. Let $\mathbf{h}_{l,i} \in \mathbb{R}^K$ denote the (contextual) representation of word i in layer l of the model, with neuron activations $\mathbf{h}_{l,i,k}$ ($1 \leq k \leq K$). 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 $\sum_j \alpha_{l,h,i,j} = 1$.

3.2 Causal Mediation Analysis

Causal mediation analysis aims to measure how a treatment effect is mediated by intermediate variables (Robins & 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

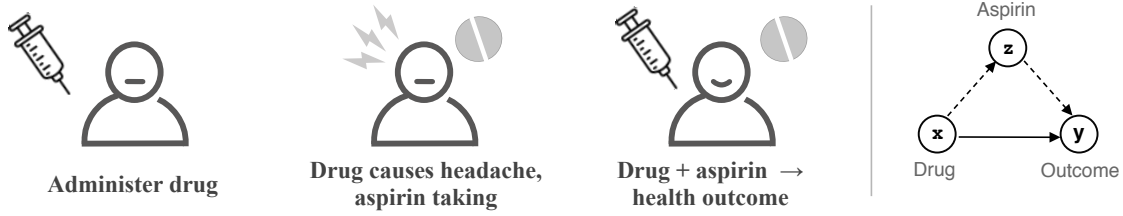


Figure 2: In this example, we want to estimate the causal effect of drug U on patient’s outcome (left panel). However, a side effect of the drug is a headache and patients with this side effect take aspirin, which also impacts the health outcome. In order to disentangle the effect of taking U with respect to aspirin status, we can carry out causal mediation analysis (right panel).

standard mechanism and an indirect effect operating via aspirin taking (the mediator) as illustrated in Figure 2.

We similarly frame internal model components, e.g., specific neurons, as mediators along the causal path between model inputs and outputs. We thus may consider a neuron to be analogous to aspirin in the example above: the neuron is influenced by the input and, in turn, affects the model output. By measuring the direct and indirect effects of targeted interventions on the model inputs, we can pinpoint the role of specific model components on model predictions. In this work, we 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., 2019a). While we use gender bias as a case study, the approach can be applied to controlled effect or bias.

The following example illustrates the problem:

Prompt u : The nurse said that __

Stereotypical candidate: she

Anti-stereotypical candidate: he

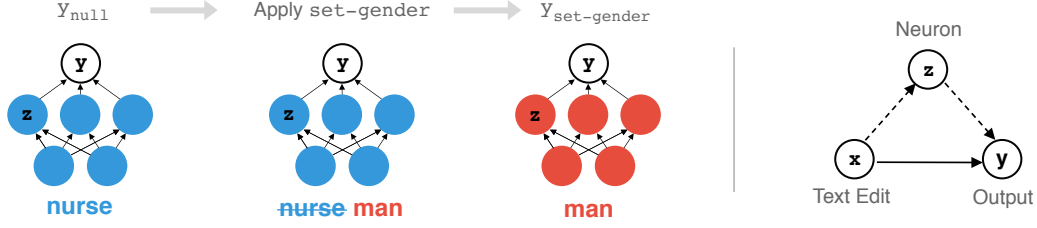
Given a prompt u such as *The nurse said that*, a language model is asked to generate a continuation. A biased model may assign a higher likelihood to *she* than to *he*, such that $p_\theta(\textit{she} \mid u) > p_\theta(\textit{he} \mid u)$. We say that *she* is the stereotypical candidate, while *he* is the anti-stereotypical candidate, reflecting a societal bias associating nurses with women more than men.

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

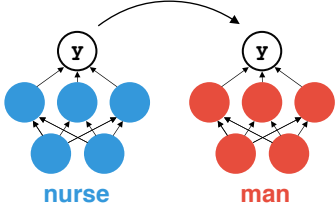
$$\mathbf{y}(u) = \frac{p_\theta(\text{anti-stereotypical} \mid u)}{p_\theta(\text{stereotypical} \mid u)}. \quad (1)$$

In our example, we have: $\mathbf{y}(u) = p_\theta(\textit{he} \mid \textit{The nurse said that}) / p_\theta(\textit{she} \mid \textit{The nurse said that})$. 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 thus exhibit bias toward neither the stereotypical nor the anti-stereotypical case. This experimental setup is based on a binary notion of a stereotypical and an anti-stereotypical candidate. Unfortunately, the datasets investigated in this work are designed for experiments with a binary grammatical gender instead of a gender-inclusive spectrum. While *they* should always be used until we know an individuals preferred pronouns, it is challenging to translate

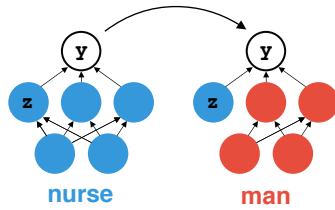
(a) Causal mechanism



(b) Total Effect



(c) Direct Effect



(d) Indirect Effect

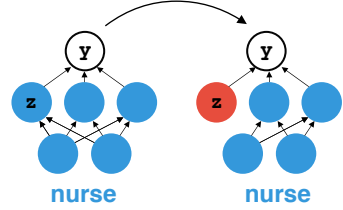


Figure 3: Mediation analysis illustration. Here the *do*-operation is $x = \text{set-gender}$, which changes u from *nurse* to *man* in this example. The **total effect** measures the change in y resulting from the intervention; the **direct effect** measures the change in y resulting from performing 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. Note that the outcomes y_{null} and $y_{\text{set-gender}}$ are observed for the same unit (input text), in contrast to the drug trial (Fig. 2) where each unit (patient) is associated with a single outcome.

this requirement into our experimental setup which aims to measure the extent to which a model is biased toward a societal stereotype. As an attempt to remedy these shortcomings, we will report experiments that treat the singular *they* as gender-neutral reference and *person* as the associated stereotypically neutral entity. These experiments measure the degree to which a model is biased against the more inclusive referential statement. For further discussion of this topic, we refer to Cao and Daumé III (2020) and Webster, Recasens, Axelrod, and Baldridge (2018).

In order to understand the role of individual model components on these biased predictions, we apply causal mediation analysis. We first perform targeted interventions on the input text and measure their effect on the gender bias measure defined above (Eq. 1), which serves as the response variable y . Specifically, we perform 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 $y_x(u)$ as the value that y attains in unit $u = u$ under the intervention $do(x = x)$.

Next, we define different kinds of effects of the intervention x on the response variable y .

Example

$u = \text{The nurse said that [blank]}$

1) Compute relative probabilities of the baseline.

$$p([he]|u) = p([he]|\text{the nurse said that}) \approx 3.1\%$$

$$p([she]|u) = p([she]|\text{the nurse said that}) \approx 22.4\%$$

$$y_{\text{null}}(u) = 3.1/22.4 \approx 0.14$$

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

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

$$p([he]|u, \text{set-gender}) = p([he]|\text{the man said that}) \approx 31.5\%$$

$$p([she]|u, \text{set-gender}) = p([she]|\text{the man said that}) \approx 2.4\%$$

$$y_{\text{set-gender}}(u) = 31.5/2.4 \approx 13.1$$

3) Compute the total effect

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

Figure 4: In this example, we present the setup to measure the **total effect** in an example with the prompt $u = \text{The nurse said that}$ with the control variable $x = \text{set-gender}$. As we compute the proportional probability prior to the intervention, we notice that 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.

3.2.1 TOTAL EFFECT

The unit-level **total effect** (TE) of $x = x$ on y in unit $u = u$ is the proportional difference² between the amount of bias under a gendered reading and under an ambiguous reading:

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

For our running example, this results in

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

An illustration of the total effect is provided in Figure 3a and an example computation is given in Figure 4.

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

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

3.2.2 DIRECT AND INDIRECT EFFECTS

We now analyze the causal role of specific mediators which lie between x and y . The mediator, denoted as z , might be a particular neuron, a full layer, an attention head, or a certain attention

2. We make the difference proportional to control for the high variance of y across examples. See Appendix A.1 for further evidence. While we limit the discussion to this metric for now, Section 3.5 expands to alternative metrics.

weight. Following Pearl’s definitions, we measure the direct and indirect effects of intervening in the model relative to z (Pearl, 2001).

The **natural direct effect** (NDE) measures how much an intervention x changes an outcome variable 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 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 3b):

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

The **natural indirect effect** (NIE) measures how much the intervention x changes 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 y through z (Figure 3c):

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

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, which are defined next.

3.3 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 4 provides more information on the dataset. As mentioned above, we pick *person* as target of the `set-gender` change for the gender-neutral reference and we measure the probability of the continuation *they*. 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* which we are not able to disambiguate from the singular case.

Throughout 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.

3.4 Attention Interventions

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

We align our intervention approach with two resources for assessing gender bias in pronoun resolution: Winobias (Zhao et al., 2018a) and Winogender (Rudinger, Naradowsky, Leonard, & Van Durme, 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 the following example from Winobias shows:

Prompt u : The nurse examined the farmer for injuries because she ____

Stereotypical candidate: was caring

Anti-stereotypical candidate: was screaming

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 $y(u) = p_\theta(\text{was screaming} \mid u) / p_\theta(\text{was caring} \mid u)$.³ In this case, we define the swap-gender operation, which changes *she* to *he*. The total effect is then

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

In the experiments, we study the effect of the attention from the last word (*she* or *he*) to the rest of the sentence.⁴ 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.

3.5 Alternate Metrics and Algorithmic Fairness

In this section, we consider alternate metrics to the ones discussed above and their connections to algorithmic fairness. A natural starting point would be to redefine bias as a simple difference between probabilities for the grammatical genders.

$$y^{\text{alt}}(u) = p_\theta(\text{anti-stereotypical} \mid u) - p_\theta(\text{stereotypical} \mid u)$$

To address the issue of high variance in the raw predicted probabilities for the stereotypical and anti-stereotypical candidates across different input sentences (Appendix A.1), they can be normalized to form a probability distribution such that $p_\theta(\text{anti-stereotypical} \mid u) + p_\theta(\text{stereotypical} \mid u) = 1$. Similarly, effects can then be computed as the difference between bias before and after intervention. This intuitive approach closely mirrors the original approach in Pearl’s work (Pearl, 2001).

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

4. One 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.

$$\text{Effect}(\text{set-gender}, \text{null}; \mathbf{y}^{alt}, u) = \mathbf{y}_{\text{intervention}}^{alt}(u) - \mathbf{y}_{\text{null}}^{alt}(u) \quad (8)$$

where $\text{Effect} \in \{\text{TE}, \text{NDE}, \text{NIE}\}$ and $\mathbf{y}_{\text{intervention}}$ would correspond to the intervention for computing the respective effect.

Since the causal effect is the primary result of interest, we can alternatively directly construct metrics that quantify the difference between prediction outcomes before and after the intervention. The disparate impact affected by gender bias can be considered a representational harm rather than an allocative harm (Barocas, Hardt, & Narayanan, 2019). Most existing work on fair classification focuses on allocative harms (Barocas et al., 2019; Dwork, Hardt, Pitassi, Reingold, & Zemel, 2011). On the other hand, bias in NLP is centered more on representational harms, particularly in word representations (Bolukbasi et al., 2016; Zhao et al., 2019b). We are nonetheless considering gender bias in a supervised prediction task, so the alternate measures of causal effect will be based on fairness in classifiers rather than fairness in word representations.

A recent line of work in algorithmic fairness is known as individual fairness, and is motivated by the concept that similar individuals should be treated similarly (Dwork et al., 2011). Unlike group fairness, which provides broad aggregate statements about fairness, individual fairness operates at the level of each individual input. In this case, each input sentence can be thought of as an individual, and the unit u before and after an intervention or gendered reading can be treated as similar individuals. In the neuron interventions, the examples *The nurse said that* and *The man said that* can be considered “similar” individuals, with the latter being the outcome of applying the `set-gender` operation on the former. The same concept applies in the attention intervention case, only with the `swap-gender` operation instead of the `set-gender` operation. The outcomes generated for each “individual” will be the predicted next word, and the difference between the outcome distributions for two similar input sentences can will be the unfairness. Statistical distance measures are a suggested choice for comparing these outcomes (Dwork et al., 2011). Two such examples are:

- Statistical distance, or total variation norm, between two probability measures P and Q on a domain A is defined as:

$$D_{tv}(P, Q) \triangleq \frac{1}{2} \sum_{a \in A} |P(a) - Q(a)| \quad (9)$$

- The relative l_∞ metric is similarly defined between two probability measures P and Q on a domain A by the following expression:

$$D_\infty(P, Q) \triangleq \sup_{a \in A} \log \left(\max \left(\frac{P(a)}{Q(a)}, \frac{Q(a)}{P(a)} \right) \right) \quad (10)$$

Here, we choose the domain A to be the set of predicted outcomes of interest, specifically the stereotypical and anti-stereotypical candidates. When working with neuron interventions, the outcomes of interest in response to a prompt like *The nurse said that* will be the set of possible genders predicted by the model, specifically $A = \{\text{he}, \text{she}\}$ in our case. Since these metrics quantify distributional differences, P and Q will be derived by normalizing the raw probabilities predicted by the neural network for each gender in A such that they sum to one and form a probability distribution.

P and Q would respectively refer to the outcome probability distribution under the original reading and under the alternate reading with different interventions depending on the causal effect (TE, NDE, or NIE) being computed.

The probability distributions used in the metrics discussed above can reasonably be extended to the full range of possible output words by the neural network. The inclusion of only pronouns in our construction, though, is due to the focus on the specific case of gender bias. Considering the larger subset, or even the full set, of possible output words and how their predicted probabilities change with various interventions would make it more difficult to concretely identify the effect on gender bias. However, it may also yield interesting insights on more subtle consequences of the interventions, such as how other non-pronoun words with associated stereotypes may be affected by these interventions.

We find that our results are robust to the metric used to quantify the effects and the work in this article focuses on measures of bias and effects originally described in the previous subsections. For reference, Appendix F contains the results of the same analyses but using the alternate metrics described here.

4. Experimental Details

4.1 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: small, medium, large, extra-large (xl), and a very small distilled model (Sanh, Debut, Chaumond, & Wolf, 2019). To test the universality of our findings across Transformer-based architectures, we perform a subset of experiments with two additional autoregressive models – Tranformer-XL (Dai, Yang, Yang, Carbonell, Le, & Salakhutdinov, 2019) and XLNet (Yang, Dai, Yang, Carbonell, Salakhutdinov, & Le, 2019) – as well as three masked LMs – BERT (Devlin, Chang, Lee, & Toutanova, 2019), DistilBERT (Sanh et al., 2019), and RoBERTa (Liu, Ott, Goyal, Du, Joshi, Chen, Levy, Lewis, Zettlemoyer, & Stoyanov, 2019). We use the Transformers library (Wolf, Debut, Sanh, Chaumond, Delangue, Moi, Cistac, Rault, Louf, Funtowicz, & Brew, 2019) for access to all of the listed models.

4.2 Data

For neuron intervention experiments, we augment the list of templates from Lu et al. (2020) with several other templates, instantiated with professions from Bolukbasi et al. (2016). 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 people working in the profession are referred to in singular *they*. To simplify processing by GPT2 and focus on common professions, we only take examples that are not split into sub-word units, resulting in 17 templates and 169 professions, or 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 with filtering by total effect. Both datasets include statistics from the U.S. Bureau of Labor Statistics to assess the gender stereotyp-

Model	Params	Layers	Heads	TE		
				WB	WG	Professions
GPT2-small rand.	117M	12	12	0.07	0.05	0.12
GPT2-distil	82M	6	12	0.12	0.08	130.86
GPT2-small	117M	12	12	0.25	0.10	112.28
GPT2-medium	345M	24	16	0.77	0.32	115.95
GPT2-large	774M	36	20	0.75	0.36	96.86
GPT2-xl	1558 M	48	25	1.05	0.34	225.22

Table 1: Model sizes and total effects (TE) of gender bias in various GPT2 variants evaluated on Winobias (WB), Winogender (WG), and the Professions dataset.

icality of the referenced occupations. Appendix A.2 provides additional details about the datasets and preprocessing methods.

5. Results

5.1 Total Effects

Before describing the results from the mediation analysis, we summarize some insights from measurements of the total effect. 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 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 the external gender bias⁵ and the log-total effect, ranging from 0.35 to 0.45 over the 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 either a smaller size compared to the professions dataset or the more complex relations in these datasets.

5. For this analysis, we add the stereotypicality and definitionality of each profession to capture the overall bias value.

The gender-neutral case leads to more consistent effects Throughout the templates in the professions experiments, 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 also much more consistent across models and templates. For the professions dataset, the GPT2 variants from distil to large have respective total effects of 8.3, 7.5, 9.6, and 12.5, all with standard deviations < 10 . 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”.

5.2 Sparsity

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 5a 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 did not find similar behavior in a randomly initialized model, indicating that these patterns do not occur by chance. We found this sparsity consistent in all model variants and datasets we examined. See Appendix C.1 for additional visualizations of indirect effects as well as direct effects.

To determine how many heads are required to achieve the full effect of intervening on all heads, we also intervene on groups of heads. We do so by selecting a subset of heads, using either a GREEDY approach, which iteratively selects the head with the maximal marginal contribution to the indirect effect, or a TOP-K approach, which selects the k elements with the strongest individual effects. 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 5b). The first 6 selected ones are from layers 4 and 5, further demonstrating the concentration of the effect in the middle layers.

Neurons Figure 6a 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. Figure 6b shows the effect from the top 5% neurons for the neutral intervention with GPT2-medium. We can observe that, while the variance of the embedding importance is much higher, the effect size is in line with the rest of the model. Additionally, the effect is much more evenly distributed across all layers. These two observations are further indications that the model has not learned a distinct and sparse representation of gender-neutral references.

Figure 7 shows the indirect effects when selecting neurons by the TOP-K algorithm.⁶ Similar to the attention result, a tiny fraction of neurons is sufficient for obtaining an effect equal to that of

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

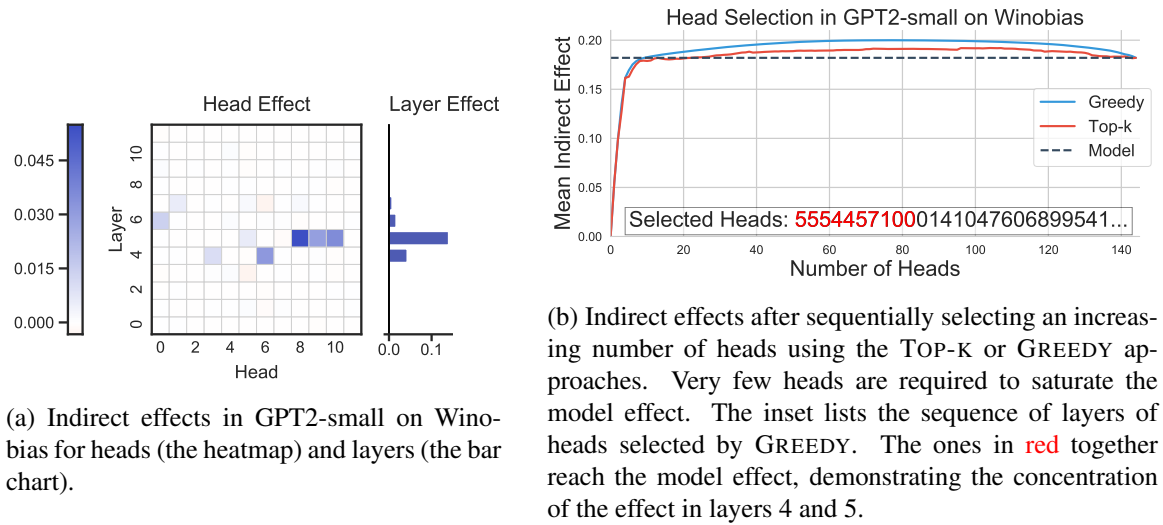


Figure 5: Sparsity effects in attention heads.

intervening on all neurons concurrently. Most of the top selected neurons are concentrated in the embedding layer and first hidden layer.

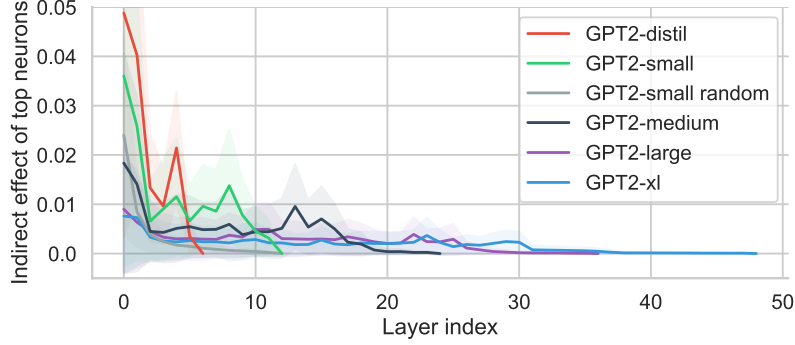
5.3 Synergism

How do different model components interact in capturing gender bias? Do different components work independently or jointly? Are gender bias effects amplified by different components or constrained?

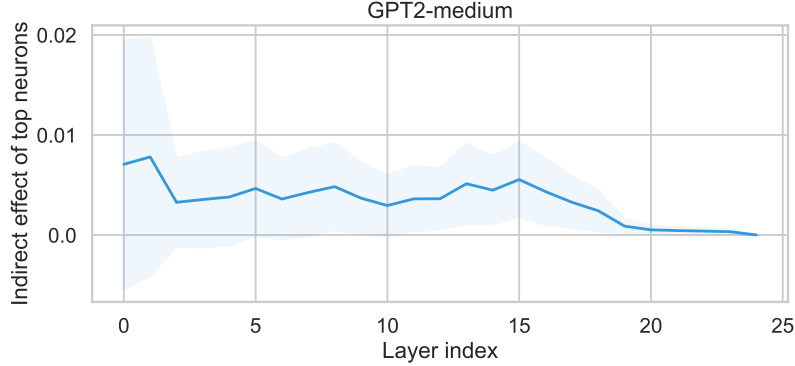
Attention Recent work found that attention heads in GPT2 and other Transformers play highly differentiated roles. For instance, some heads focus on adjacent tokens while others align with syntactic properties (Kovaleva, Romanov, Rogers, & Rumshisky, 2019; Hoover, Strobel, & Gehrmann, 2019; Clark, Khandelwal, Levy, & Manning, 2019; Vig & Belinkov, 2019). We use mediation analysis to study the interdependence of attention heads.

Figure 8 compares indirect effects of concurrent intervention on all heads (NIE-all) to summing the effects of independent interventions (NIE-sum). The differences are fairly small (maximum relative distance from NIE-all between 0.7% and 11.3%), indicating that heads operate primarily in an independent and complementary manner, capturing different aspects of gender bias. As Figure 5b shows, most heads do not contribute much to the indirect effect, and many reduce it. This trend is consistent across models and datasets (Appendix D).

Figure 9 shows the attention of the three heads with the highest indirect effects on Winobias. The figure demonstrates that they 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 (layer 5, head 10; noted as 5-10) as relating to coreference resolution based on visual inspection. 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.



(a) Indirect effects of top neurons in different models on the professions dataset. The embedding and first layer are much more significant than all others. We can additionally observe a “bump” in the middle layers where the effect increases after the initial decline.



(b) Indirect effects of top neurons in GPT2-medium for gender-neutral interventions on the professions dataset. The effect is lower and distributed across all layers.

Figure 6: Sparsity effects in neurons.

Neurons Similar to the case of attention, Figure 7 shows that after a few neurons (4%) match the model-wise concurrent effect, most neurons do not contribute much, and many even diminish the effect. This result suggests that neurons may be as specialized as the individual attention heads. However, an analogous qualitative analysis is challenging due to the large number of neurons.

By definition, concurrent intervention on all neurons entails $TE = NIE\text{-all}$, since then $\mathbf{y}_{\text{set-gender}}(u) = \mathbf{y}_{\text{null}, \mathbf{z}_{\text{set-gender}}(u)}(u)$. Notably, the sum of independent indirect effects (NIE-sum) is much smaller than the concurrent intervention (NIE-all), as shown in the following table. Thus, neurons combine synergistically to compound independent effects.

Residual Connections Visualizing the indirect effect of each neuron in a heatmap (Figure 10 top) reveals vertical stripes when a neuron at the same index, but different layers, has a similar effect. While sparse, this effect sometimes continues over multiple layers. Two possible explanations for

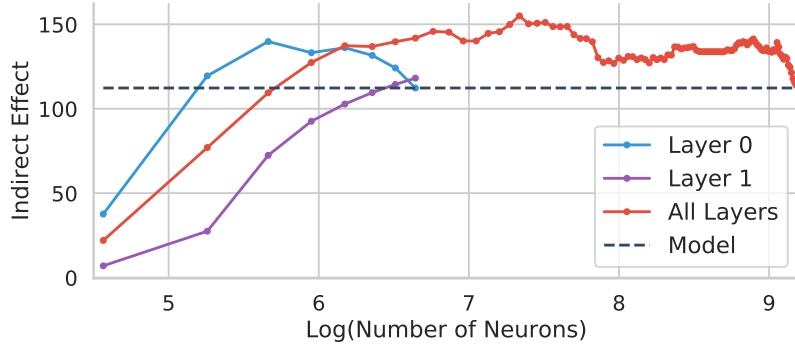


Figure 7: Indirect effects after sequentially selecting an increasing number of neurons from either the full model or individual layers using the TOP-K approach in GPT2-small on the professions dataset. Very few neurons (4%) are required to saturate the model effect. Of those, 57% are from layers 0 and 1.

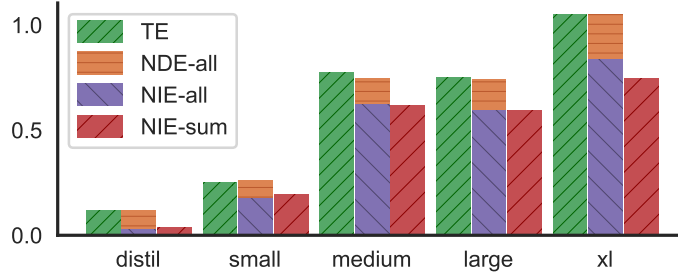


Figure 8: Effects of intervening on all heads concurrently (all) vs. independently and summing (sum) in various GPT2 variants evaluated on Winobias.

this are random alignments of two effective cells or the residual connections between the layers. To analyze this, we computed the number of stripes between layer pairs across the professions dataset, with and without randomizing neuron indices. As Figure 10 (bottom) shows, the stripes are less random in higher layers. This implies that, as the information gets transformed, the model converges on a representation. This is akin to gated recurrent networks, except that those transform across time steps instead of layers. This result may partially explain the higher neuron importance in earlier layers since those neurons have not yet converged to a representation and thus have a higher variance and contribution to the representation in other neurons.

5.4 Decomposition of the Total Effect

Attention heads mediate most of the effect Figure 8 also shows the concurrent direct and indirect effects, when intervening on all heads. In all but the smallest model (distil), the concurrent indirect effect is larger than the direct effect, indicating that most of the effect is mediated through the attention heads. Other model components (e.g., word representations) are nonetheless responsible for a portion of the total effect. This might be due to biased word embeddings predisposing the model towards certain continuations. For instance, the representation of *he* might lead the model to

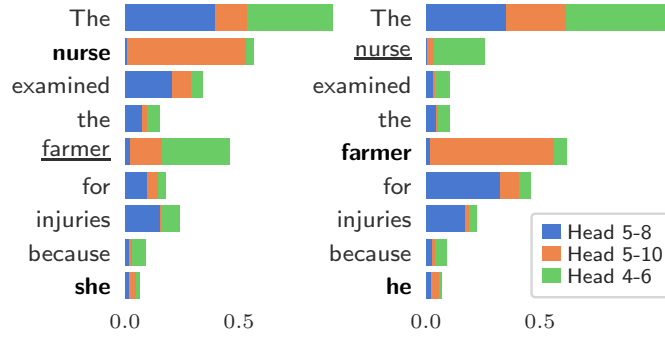


Figure 9: Attention of different heads in GPT2-small on a Winobias example, directed from either *she* or *he*. Colors correspond to different heads. Head 5-10 attends directly to the **bold** stereotypical candidate, head 5-8 attends to the words following it, and head 4-6 attends to the underlined anti-stereotypical candidate. Attention to the first token may be null attention (Vig & Belinkov, 2019). Appendix C.2 shows more examples.

	distil	small	medium	large	xl
NIE-sum	6.8	4.0	3.5	2.1	2.9
NIE-all	130.9	112.3	116.0	96.9	225.2

predict a lower probability for *was caring* compared to *she*, irrespective of any previous occupation mention.

TE \approx NDE + NIE In linear models, it is known that the linear total effect decomposes to direct and indirect effects (Pearl, 2001). Intuitively, intervention effects either flow through a mediator or directly. In our case, we have a highly non-linear model and this decomposition is not guaranteed.⁷ Nevertheless, Figure 11 shows such approximate decomposition for the top heads in GPT2-small.⁸ The same holds for concurrent interventions (Figure 8), where $TE \approx NDE\text{-all} + NIE\text{-all}$. To understand this phenomenon, observe that under our formulation of the effects using a proportional difference, a decomposition of the form $TE = NDE + NIE$ is expected if the following equality holds for all u :

$$y_{\text{set-gender}}(u) - y_{\text{set-gender}, z_{\text{null}}(u)}(u) = y_{\text{null}, z_{\text{set-gender}}(u)}(u) - y_{\text{null}}(u). \quad (11)$$

See Appendix E for intuition, a proof, and evidence that Eq. 11 approximately holds in our results.

7. VanderWeele and Vansteelandt (2009) discuss which decompositions can be guaranteed. While causal effect definitions are model-free, some decompositions are possible even for non-linear models and effects, in the presence of, for example, interaction between the intervention and the mediator. However, the $TE = NDE + NIE$ decomposition is not guaranteed without further assumptions (e.g., under linear models). In our case, an additional no-interaction condition was needed for the decomposition to hold, as shown and discussed in Appendix E.

8. In the neuron intervention case, by definition $TE = NIE\text{-all}$ and $NDE\text{-all} = 0$, so the decomposition trivially holds.

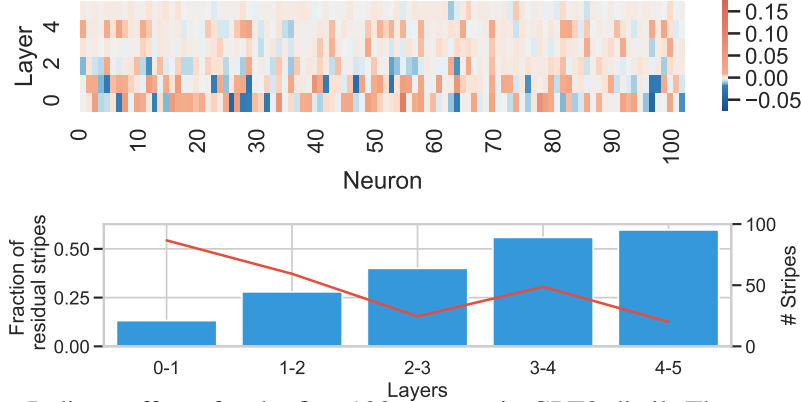


Figure 10: **Top**: Indirect effects for the first 100 neurons in GPT2-distil. There are distinct vertical stripes where the effect of a neuron at an index continues to the next layer. **Bottom**: The fraction of the continued effect over layer pairs in GPT2-distil that can be explained by residual connections. In higher layers, the model more strongly relies on the connections to refine representations across its layers.

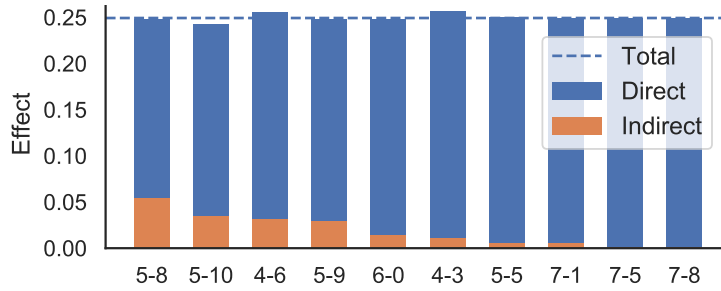


Figure 11: Top 10 heads by indirect effect in GPT2-small on Winobias, and their respective direct effects.

5.5 Experiments with other models

How specific are the findings outlined above to GPT2? In this section, we examine whether causal mediation analysis has yielded any general insights about gender bias effects in Transformer-based LMs.

Attention interventions Figure 12 shows the indirect effects for each head and layer in Transformer-XL and XLNet-base on Winobias (top) along with bar charts of top heads by indirect effect and their respective direct effects for both models (bottom). These results support our main findings with GPT2: the heat maps show the sparsity of indirect effects in Transformer-XL and XLNet while the bar charts demonstrate decomposability. The agreement between all autoregressive models we have tested indicates that our framework has revealed general patterns in gender bias effects manifesting across different architectures rather than features of a particular model.

The results of attention intervention experiments with masked language models – BERT, DistilBERT, and RoBERTa – show mixed levels of agreement with GPT2 results. For example, Figure 13 illustrates the lack (left) and apparent presence (right) of sparsity of indirect effects in BERT-large-uncased when taking two different approaches for scoring candidates. In contrast to the autoregres-

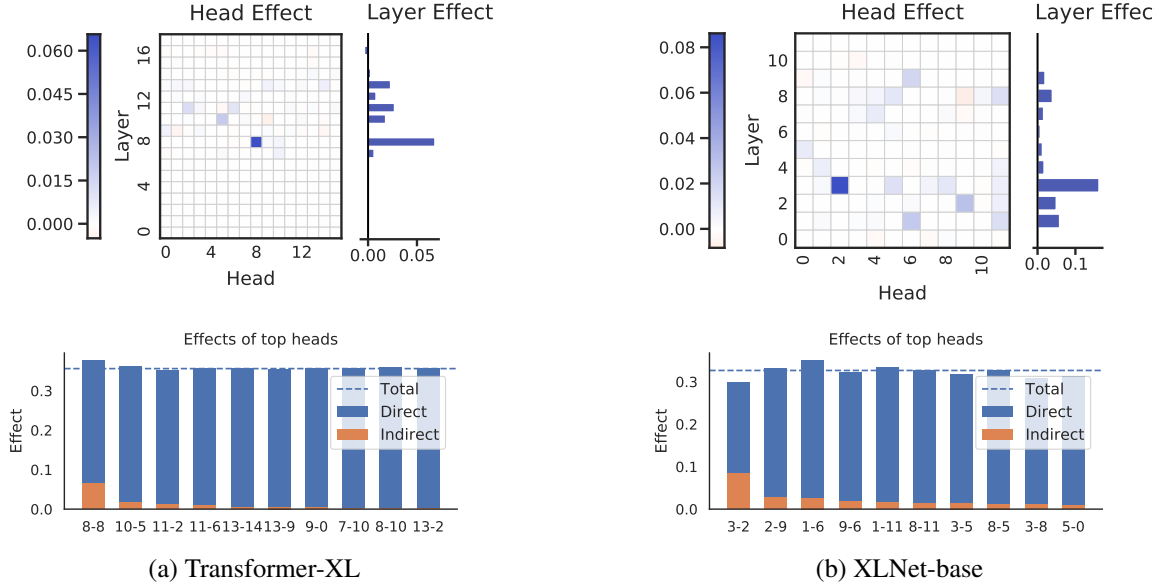


Figure 12: Sparsity and decomposability of gender bias effects in Transformer-XL and XLNet-base for Winobias.

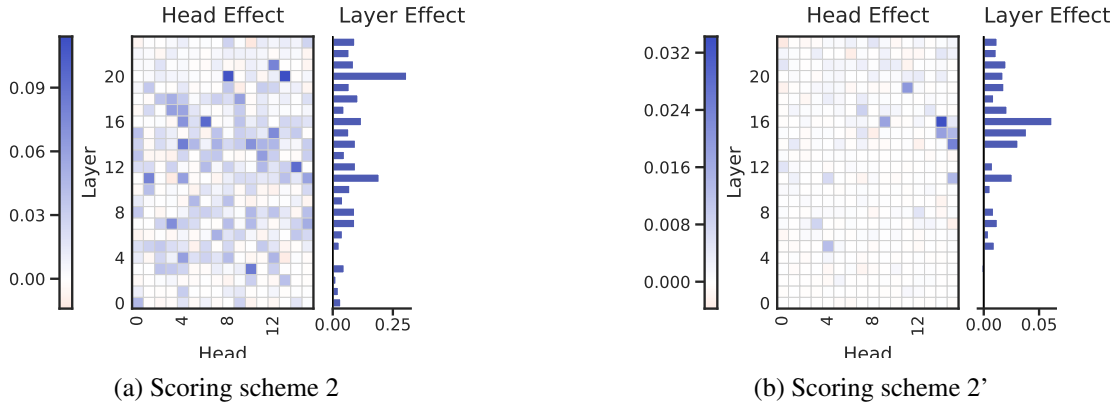


Figure 13: Indirect effects in BERT-large-uncased on Winobias for different scoring schemes.

sive case, there does not seem to be a single obvious way to score multi-token continuations with masked LMs since they do not directly output probability distributions for the next word given a prefix. This led us to try out several different scoring schemes in our experiments (see Appendix G for details). We suspect that this difference in how candidates are scored in autoregressive models on the one hand and masked LMs on the other may be the underlying cause for the discrepancy in results for the two kinds of models.

One general trend that seems to manifest across all of the different models we have tested is that of larger variants of the same models having larger total effects. Appendix G contains a table illustrating this, as well a more detailed exposition of the results with masked LMs.

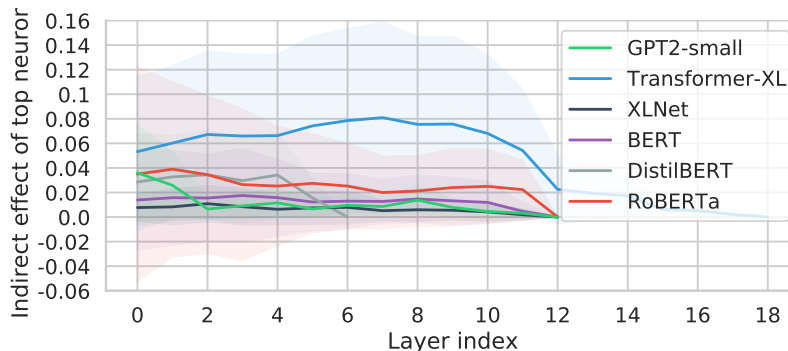


Figure 14: Indirect effects of top neurons in different models on the professions dataset.

Neuron interventions Figure 14 shows the indirect effects from the top 5% of neurons from each layer in all of our five additional models along with the same plot for GPT2-small. None of the non-GPT2 models seem to share the exact same pattern that we have observed with GPT2 variants, where the strongest effects in the word embedding layer used to be followed by a rapid drop in the successive layers. Here, a general trend seems to be one of roughly gradual decrease to 0 in the strength of effects, albeit with more and less significant fluctuations, the only exception being Transformer-XL whose effects are steadily increasing in the first third of the layers. As can be seen from Figure 14, we also observe considerably higher variances with some of the new models than in the case of GPT2. We do not have a compelling theory to explain these differences in results.

6. Discussion and Conclusion

This article introduced a structural-behavioral framework for interpreting neural models based on causal mediation analysis. An application of this framework to gender bias in NLP yielded several insights regarding the mechanisms by which gender bias is mediated in large Transformer LMs, revealing that gender bias effects are sparse, synergistic, and decomposable to direct and indirect effects.

This work can be extended in multiple ways. While the framework may apply to any property expressed as a function of model predictions, our experimental design focuses on gender bias in a binary setup. Preliminary experiments with gender-neutral references showed that models have not learned the concept of these references, at least for our chosen templates. Applying the methodology to more gender-inclusive setups and other kinds of biases is thus especially important since it can uncover these shortcomings in models.

Our results can also guide model selection for limiting the amount of bias. In related work, Giulianelli et al. (2018) demonstrated that changing hidden representations in an LSTM based on the output of diagnostic classifiers can decrease the error on a subject-verb agreement classification task. Inspired by these results, mediation analysis could be used to determine where and how to intervene in similar ways and thus not only assess the success of debiasing techniques, but also motivate new debiasing methods that establish counterfactual fairness for protected groups. Another extension of this work for training a fair model would be to use the individual fairness framework (Dwork et al., 2011) that inspired some of the alternate metrics discussed in this article. This approach would

involve maximizing a traditional objective function through the typical training process subject to an additional Lipschitz constraint.

This work is a first attempt to adopt mediation analysis for interpreting NLP models. The causality literature often focuses on assumptions needed for identification of mediation effects from observed data (Pearl, 2001; Avin, Shpitser, & Pearl, 2005; Imai, Keele, & Yamamoto, 2010b). The challenge with inferring causality from observational data is that for each unit, the outcome is observed under a single intervention. However, in this work we use the language of causal mediation to study the structure of NLP models, utilizing the fact that the outcome of the same unit (e.g., a sentence) can be observed under any intervention given a trained model. As a result, causal effects can be computed in a relatively simple manner. While we observed consistent results under multiple metrics, our definitions of causal effects to quantify bias could be refined, and alternative definitions might be advantageous for NLP research.

The causality literature offers many avenues for continuing this line of work, including mediation analysis with non-linear models, and alternative effect decompositions (Imai, Keele, & Tingley, 2010a; Imai et al., 2010b; VanderWeele & 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. We believe the present work sets the ground for employing this literature in the ongoing effort to analyze neural NLP models.

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.

Appendix A. Data Preparation

A.1 Professions Dataset

Figure 15 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,

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 15: Templates for neuron interventions.

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 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.

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.⁹ 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, 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 & 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.

Appendix 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.

9. An 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.*”

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

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

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

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

Appendix C. Additional Attention Results

C.1 Indirect and Direct Effects

Figure 16 complements Figure 5a by visualizing the indirect effects for additional GPT2 models. As with Figure 5a, the attention heads with the largest indirect effects lie in the middle layers of each model. Figure 17 shows the indirect effects for a model with randomized weights. Figures 18 and 19 visualize the indirect effects for other dataset variations for the GPT2-small model from Figure 5a. The attention heads with largest indirect effect have significant overlap across the dataset variations.

Figure 20 visualizes *direct* effects on Winobias for GPT2-small and GPT2-large. As discussed in Section 5.4, the sum of direct and indirect effects approximate the total effect.

C.2 Examples

Figure 21 visualizes attention for the Winobias examples with the greatest total effect in GPT2-small, complementing the example shown in Figure 9. Figure 22 visualizes attention for additional models for the same example shown in Figure 9.

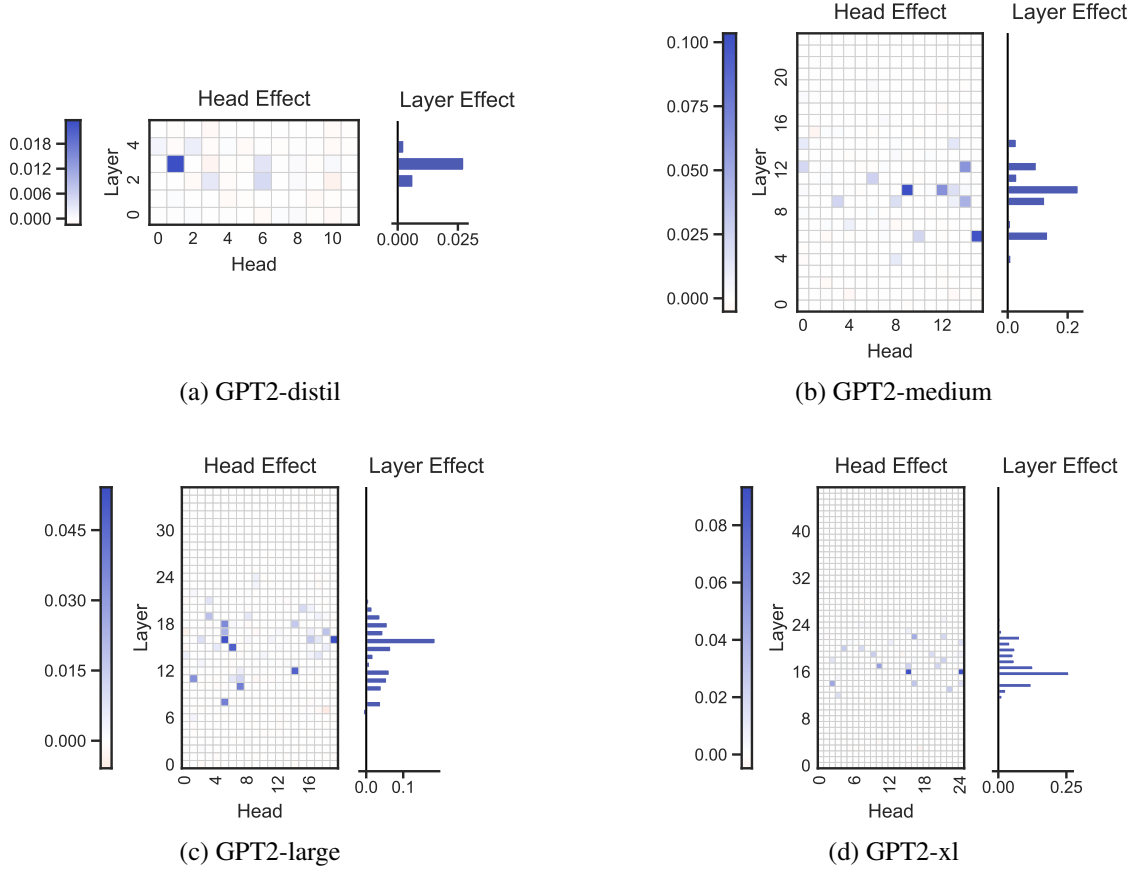


Figure 16: Mean indirect effect on Winobias for heads (the heatmap) and layers (the bar chart) over additional GPT2 variants.

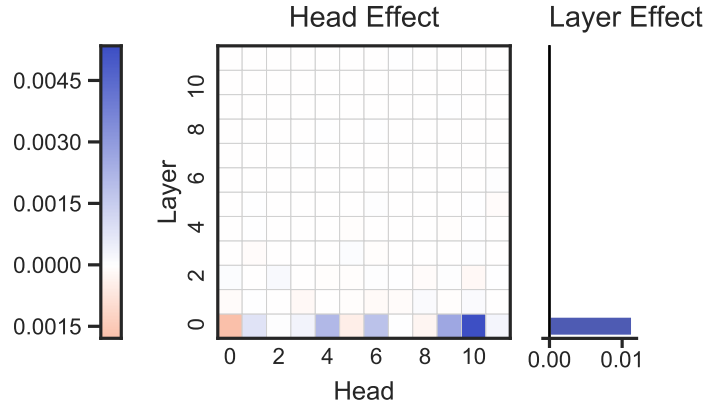


Figure 17: Indirect effect when using a randomly initialized GPT2-small model on Winobias.

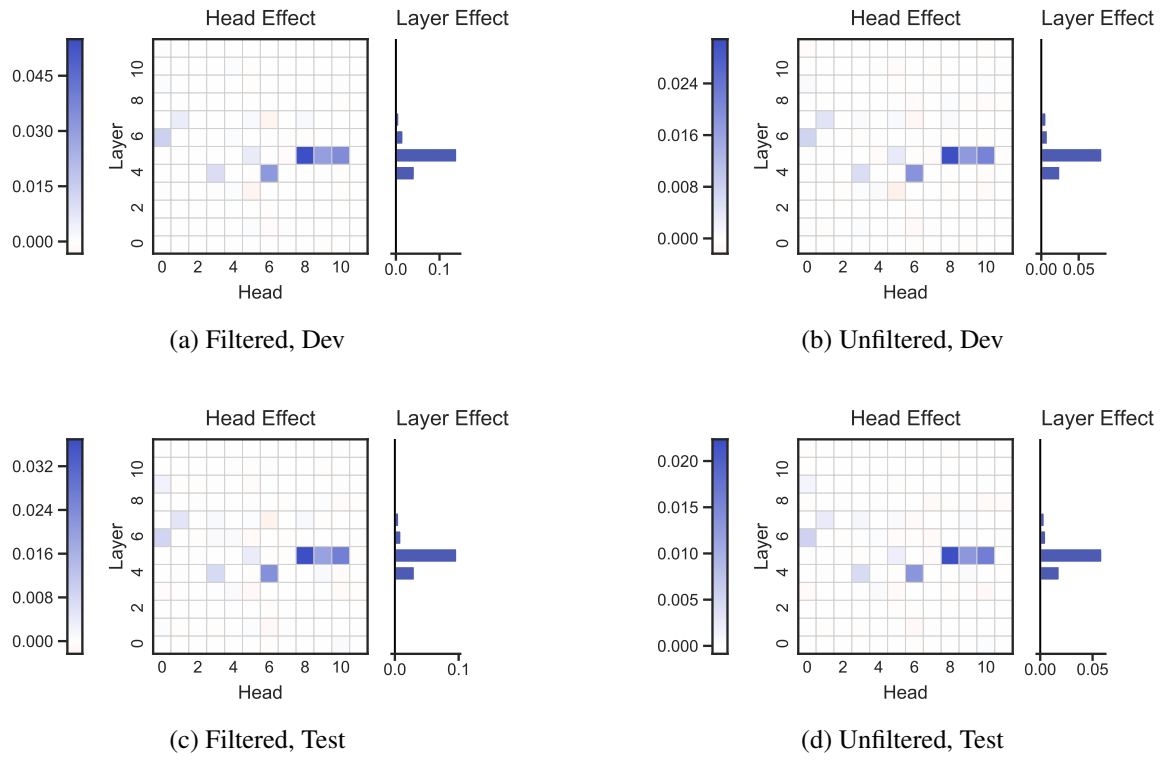


Figure 18: Indirect effect for Winobias (GPT2-small).

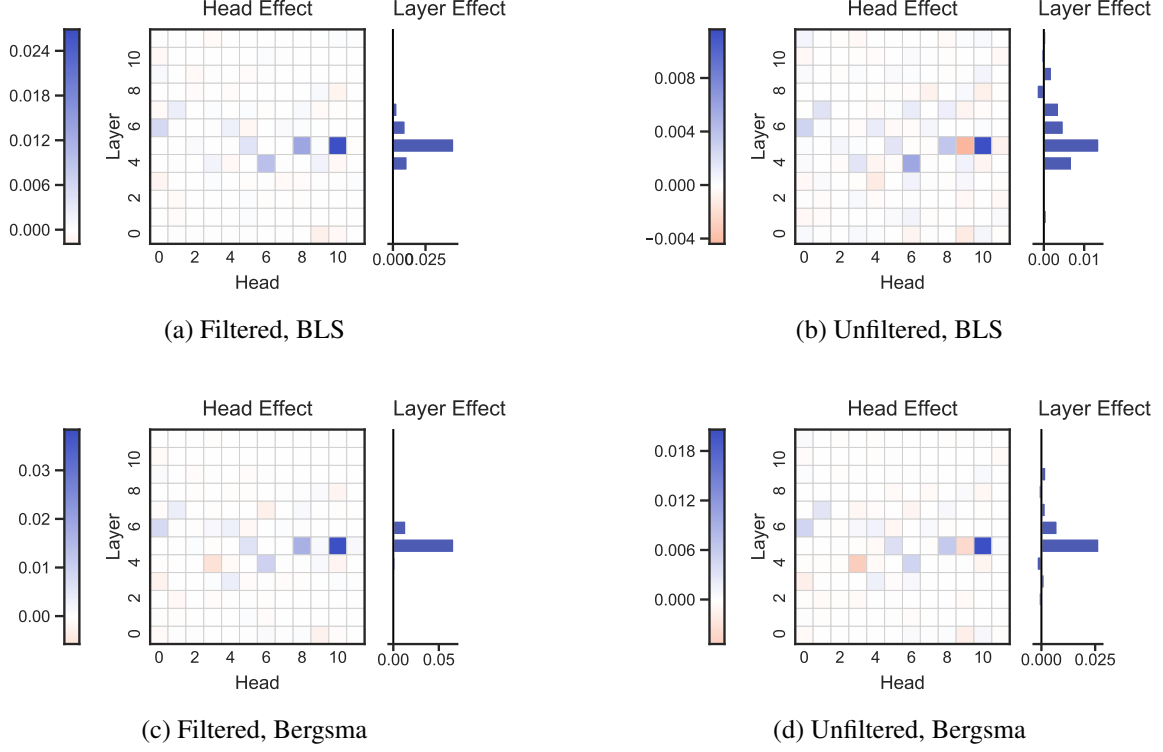


Figure 19: Indirect effect for Winogender (GPT2-small).

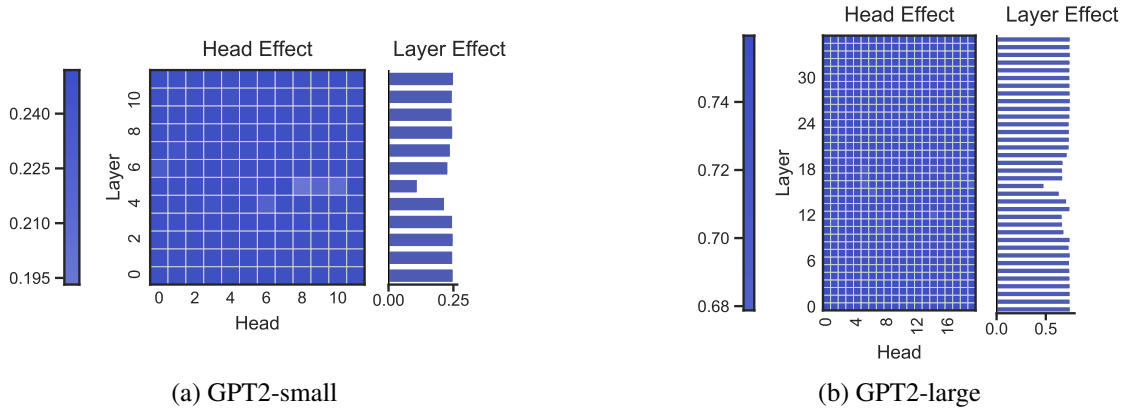


Figure 20: Direct effect for Winobias for GPT2-small and GPT2-large.

MEDIATION ANALYSIS FOR INTERPRETING NLP: GENDER BIAS

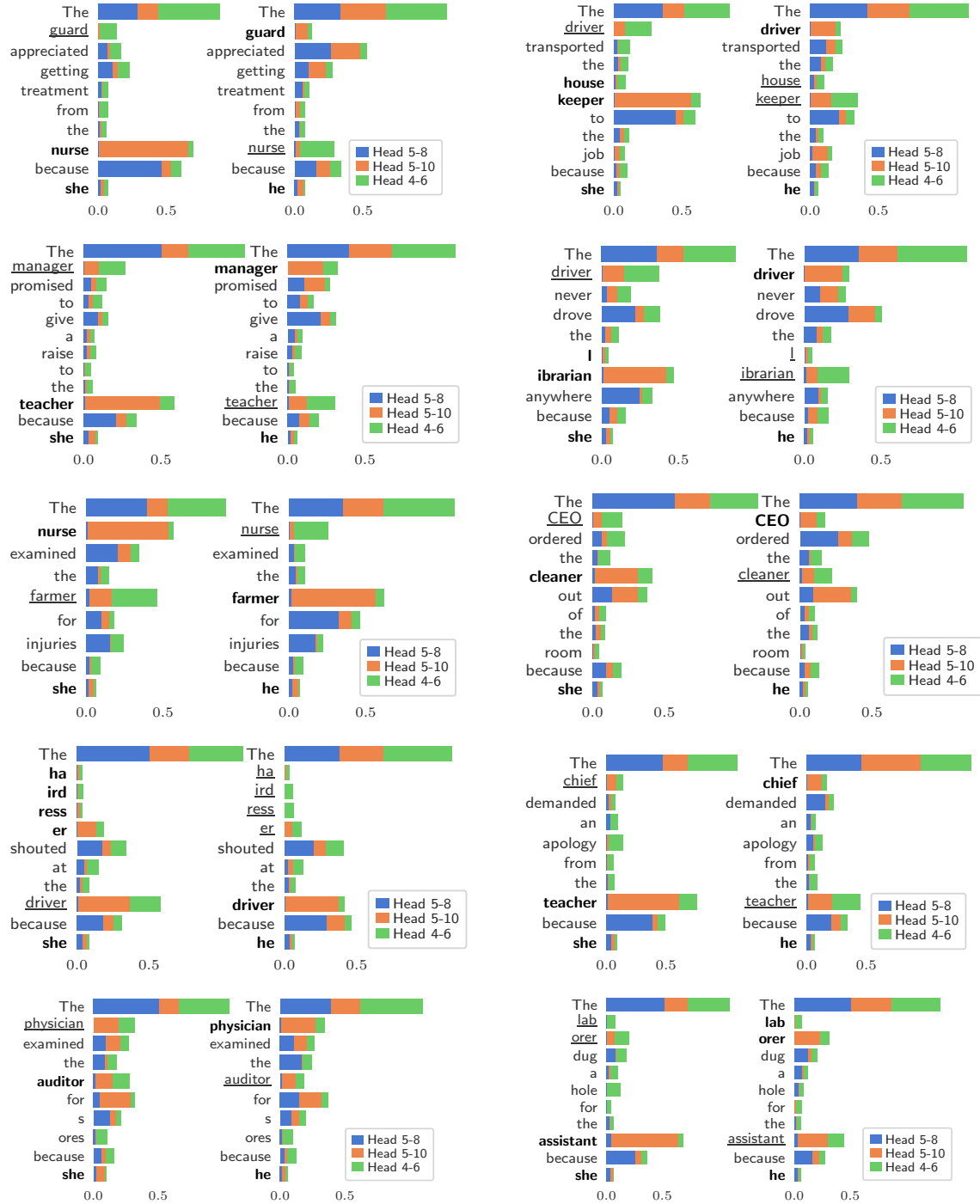
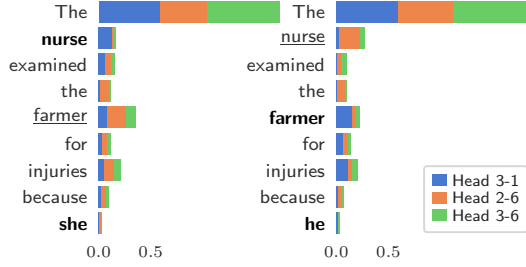
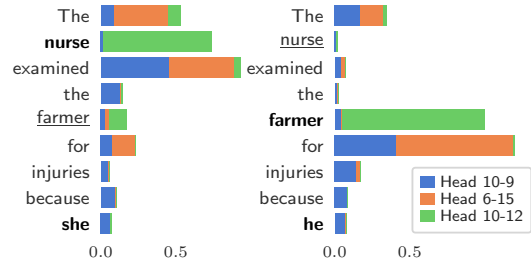


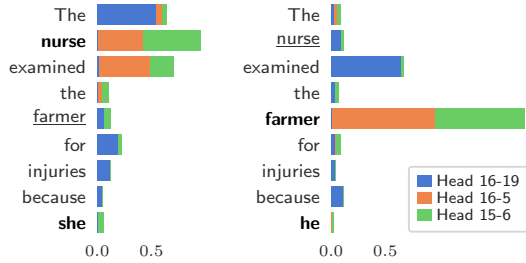
Figure 21: Attention of different heads across the 10 Winobias examples with greatest total effect for the GPT2-small model. The stereotypical candidate is in **bold** and the anti-stereotypical candidate is underlined. Attention roughly follows the pattern described in Figure 9.



(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 22: 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 9.

Appendix 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 & 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, Feldman, Naor, & Schwartz, 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 23 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 & Nguyen, 2019; Fahrback, Mirrokni, & Zadimoghaddam, 2019b; Balkanski & Singer, 2018a, 2018b). For non-monotone or non-submodular functions, there are parallelized algorithms that use similar techniques to select sets (Balkanski, Breuer, & Singer, 2018; Qian & Singer, 2019; Fahrback, Mirrokni, & Zadimoghaddam, 2019a). These methods provide an alternative approach to TOP-K for selecting subsets of neurons and can be explored in future work.

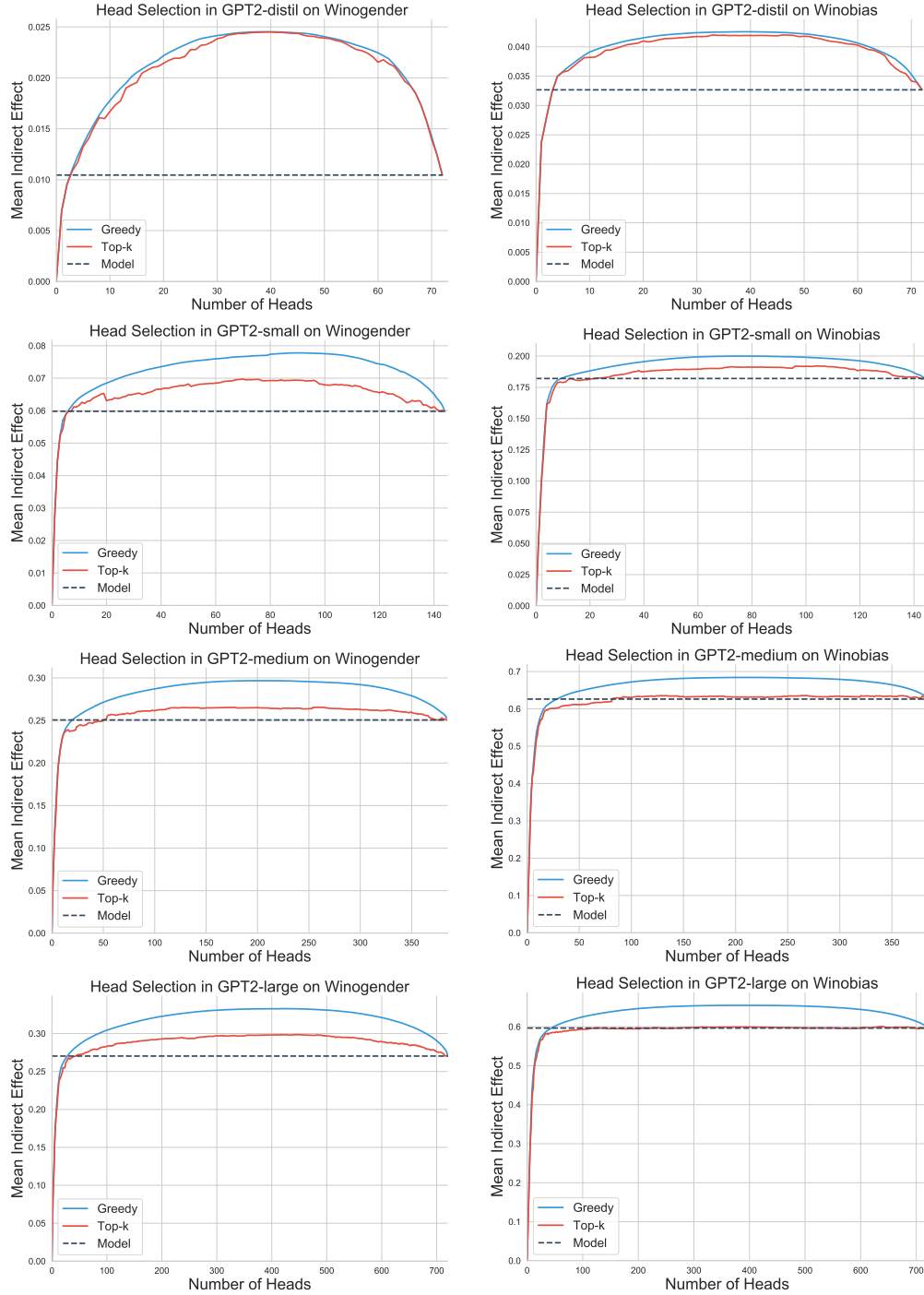


Figure 23: The indirect effect after sequentially selecting an increasing number of heads through the TOP-K or GREEDY approach on different model types and data. A small proportion of heads are required to saturate the effect of the model.

Appendix E. Proof that no-interaction in the difference NIE implies decomposition of the TE

Since by definition $y_{\text{set-gender}}(u) = y_{\text{set-gender}, z_{\text{set-gender}}(u)}(u)$ and $y_{\text{null}}(u) = y_{\text{null}, z_{\text{null}}(u)}(u)$, it can be seen that both sides of Eq. 11 describe a form of NIE (defined on the difference scale), where each contrasts y under two different interventions on z while keeping the sentence u the same (left side, under `set-gender`; right side, under `null`). This equation parallels a previously-described assumption in the causal mediation analysis literature that ascertains that the NIE is the same regardless of the fixed value at which the intervention (the analogue of `set-gender/null`) is held, known as a no-interaction assumption (Imai et al., 2010b). We show that no-interaction in the indirect effect on the difference scale implies that the $\text{TE} = \text{NDE} + \text{NIE}$ under our scale. Eq. 11 can be rewritten as

$$\begin{aligned} y_{\text{set-gender}}(u) - y_{\text{null}}(u) = & \quad (12) \\ y_{\text{set-gender}, z_{\text{null}}(u)}(u) - y_{\text{null}}(u) \\ + \quad y_{\text{null}, z_{\text{set-gender}}(u)}(u) - y_{\text{null}}(u). \end{aligned}$$

Now, dividing both sides of the equation by $y_{\text{null}}(u)$ and taking expectations over u yields

$$\begin{aligned} \mathbb{E}_u [y_{\text{set-gender}}(u)/y_{\text{null}}(u) - 1] = & \quad (13) \\ \mathbb{E}_u [y_{\text{set-gender}, z_{\text{null}}(u)}(u)/y_{\text{null}}(u) - 1] + \\ \mathbb{E}_u [y_{\text{null}, z_{\text{set-gender}}(u)}(u)/y_{\text{null}}(u) - 1], \end{aligned}$$

which is exactly

$$\begin{aligned} \text{TE}(\text{set-gender}, \text{null}; y) = & \quad (14) \\ \text{NDE}(\text{set-gender}, \text{null}; y) + \\ \text{NIE}(\text{set-gender}, \text{null}; y). \end{aligned}$$

It should be noted that even if Eq. 11 does not hold, but the equation approximately holds upon dividing both sides by y_{null} , we would expect the decomposition $\text{TE} \approx \text{NDE} + \text{NIE}$ to hold. Indeed, further inspection of the trained model revealed that the left side and right side of Eq. 11 were very close in the case of the attention intervention. Figure 24 shows a plot of the values attained by the two sides of the equation (normalized by y_{null} to make consistent with the earlier analyses) for all attention heads across all examples in the Winobias dataset. Fitting the data to a linear model yields a coefficient of 1.04 and an intercept of 0.00 ($R^2 = 0.78$).

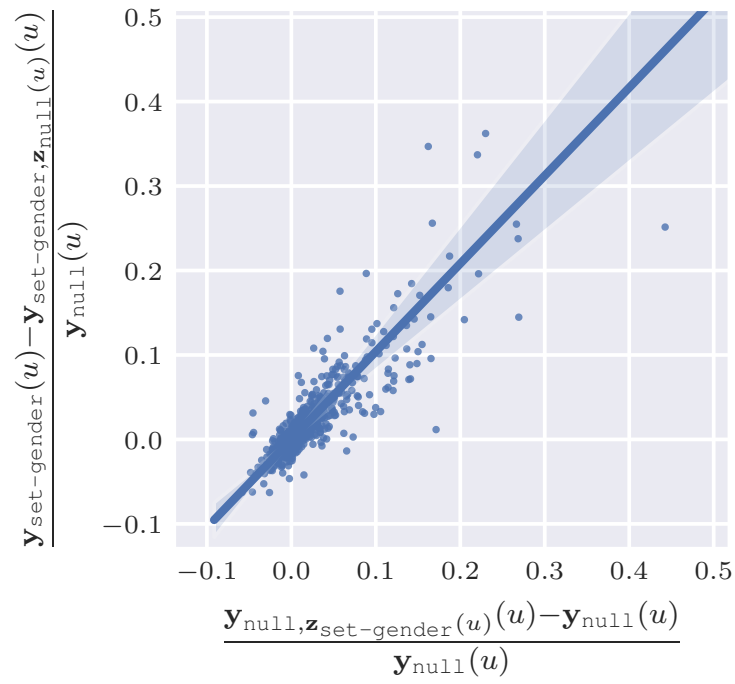


Figure 24: Plot of right side of Eq. 11 (x axis) against left side (y axis), normalized by y_{null} . For visualization purposes, we exclude a single outlier at (0.60, 1.07).

Appendix F. Alternate Metrics

Here are the results of select analyses from throughout the article, replicated using the different alternate metrics described in Section 3.5. Specifically, we investigate a total of four metrics, consisting of the primary metric and the three alternate metrics, referred to and defined by the following:

- Original metric: The primary measurement of bias and effects as described in the setup in Sections 3.2, 3.3, and 3.4 and as used throughout the article.
- Relative ℓ_∞ metric: Effects based on the distance measure described in Equation 10.
- Normalized difference: Effects based on the distance measure described in Equation 8.
- Total Variation distance (or TV-distance): Effects based on the distance measure described in Equation 9.

Figures 25, 26, and 27 evaluate all the metrics on the GPT2-small model using the filtered Winobias Dev, filtered Winogender Bergsma, and Professions data sets respectively.

Figures 28, 29, and 30 evaluate the three alternate metrics – the relative ℓ_∞ metric, normalized difference, and TV-distance, respectively – on the filtered Winobias Dev data set for different GPT2 model variants. These can be contrasted with the comparable results of the original metric in Figure 16.

Tables 5 and 6 reporting total effects are also included. Overall, the main takeaways are that the results of our analyses are quite robust to the metric that is being used to compute the effects. There are minor visible differences, but the relative behavior overall is consistent, which is not necessarily an intuitive result.

For instance, one notable difference is that there appears to be less absolute sparsity, but relative sparsity holds, which suggests that these alternate metrics are more sensitive to the effects of causal mediation analysis while demonstrating that the conclusions drawn by our article continue to hold in all of these situations. Similarly, larger models appear to exhibit a more pronounced effect than smaller models even under alternate metrics, which is another consistent result.

It is worth noting that because of the way that TV-distance is constructed in terms of difference between magnitudes rather than ratios, the effect values for each model and each data set are on different levels. However, when normalizing relative to a constant (i.e., the minimum effect across layers for the neuron line plots or across heads for the attention heatmaps), then the emerging behavior is typically consistent with the broader patterns that have been identified.

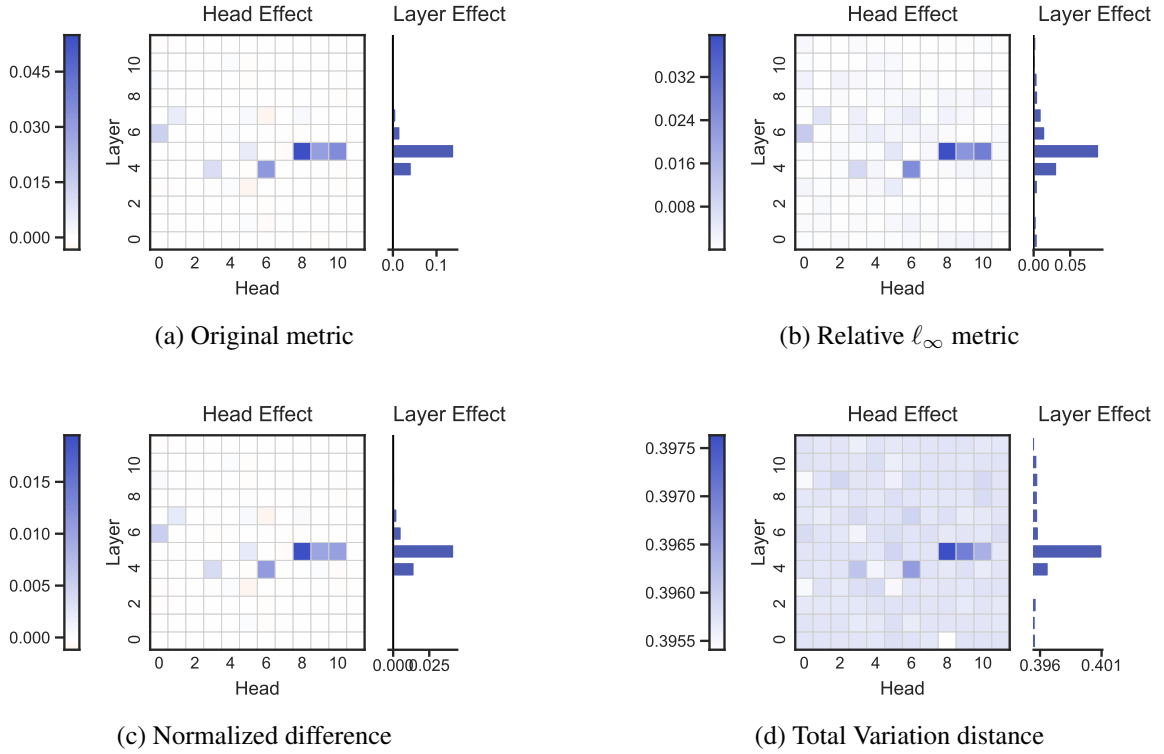


Figure 25: Indirect effect for the filtered Winobias Dev data set on GPT2-small using various alternate metrics.

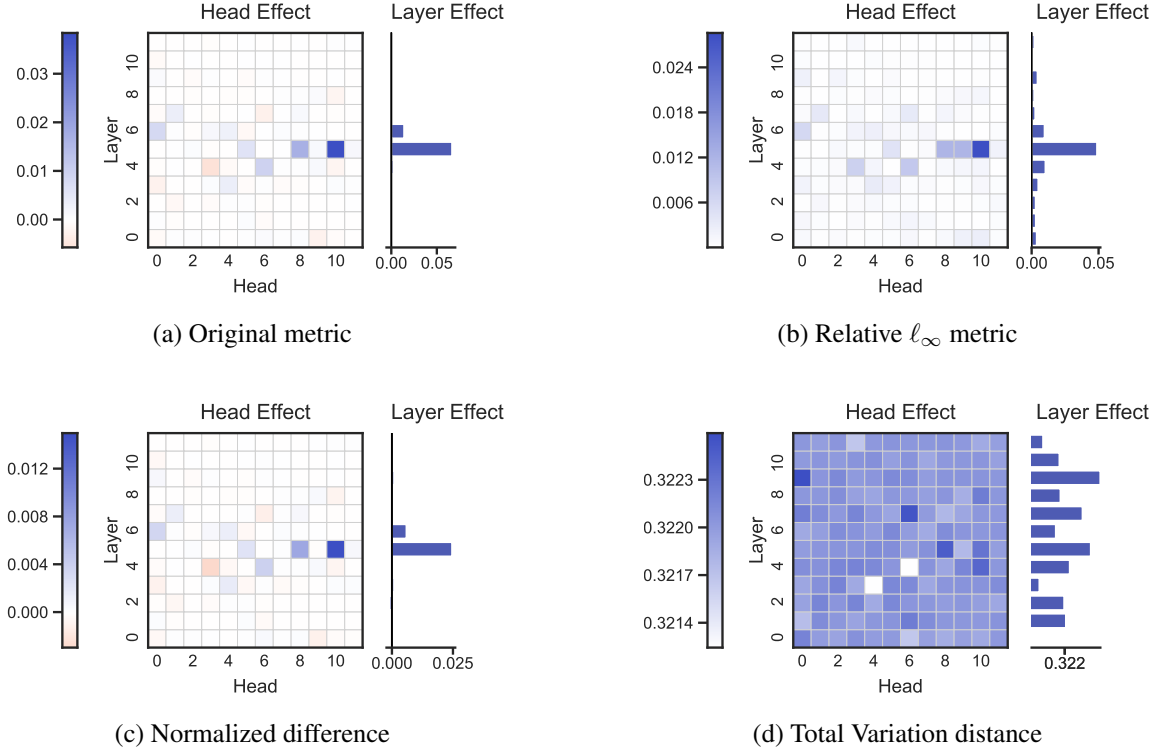


Figure 26: Indirect effect for the filtered Winogender Bergsma data set on GPT2-small using various alternate metrics.

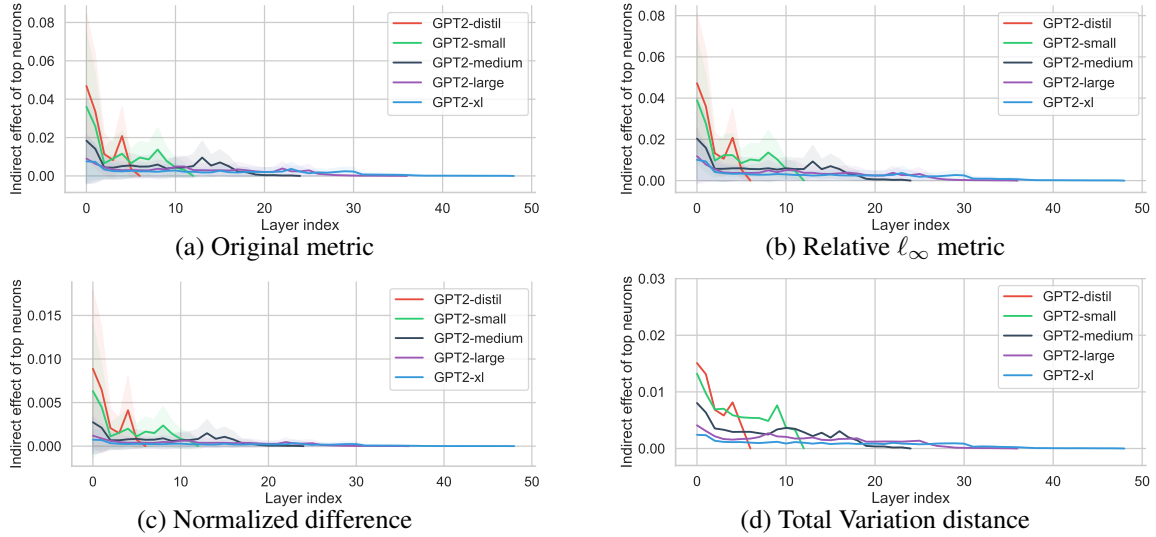


Figure 27: Indirect effect for the Professions data set on GPT2-small using various alternate metrics.

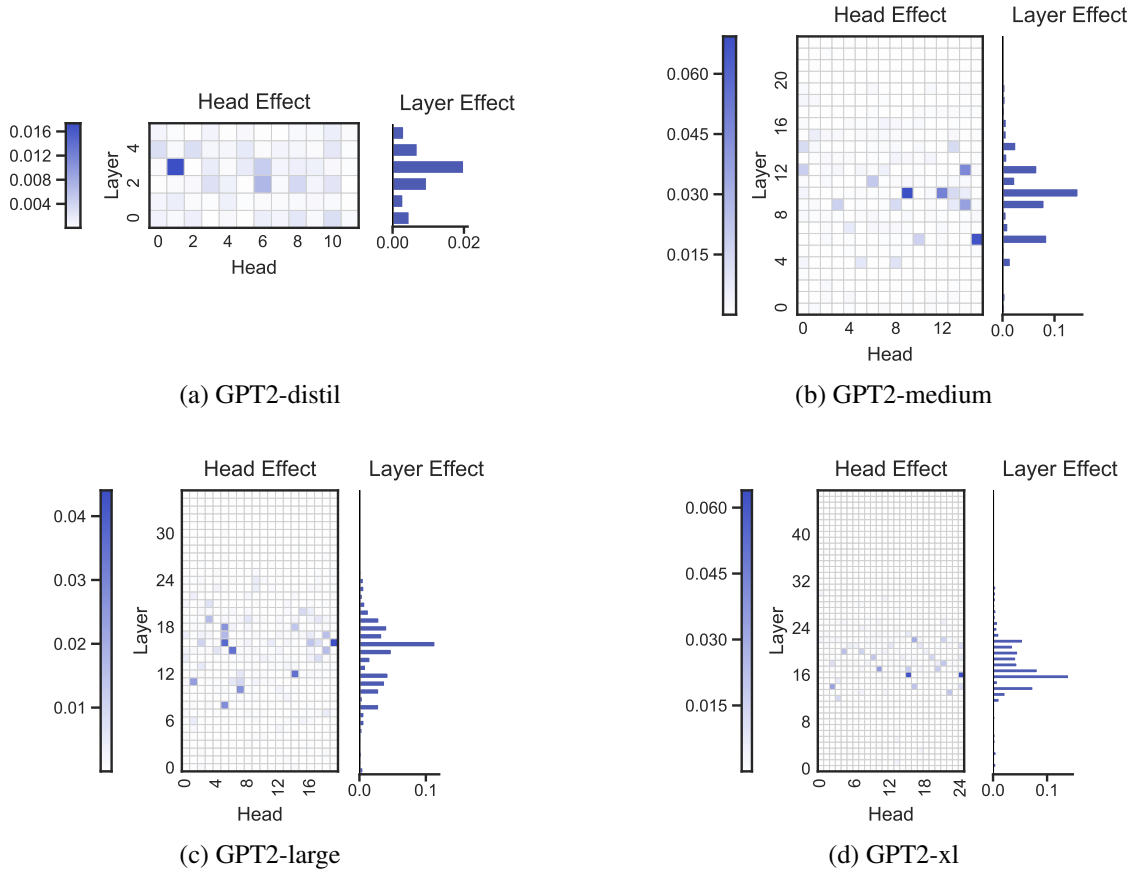


Figure 28: Mean indirect effect on Winobias for heads (the heatmap) and layers (the bar chart) over additional GPT2 variants as measured using the alternate metric, relative ℓ_∞ .

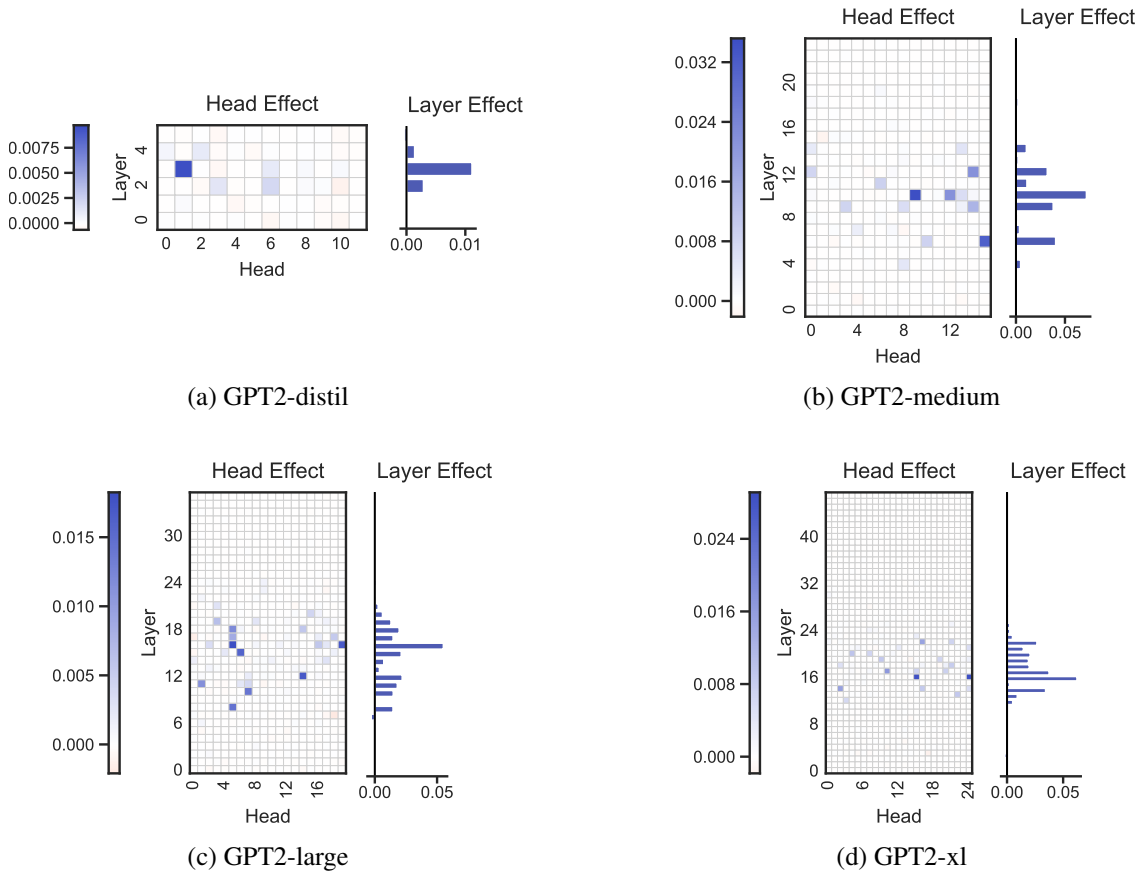


Figure 29: Mean indirect effect on Winobias for heads (the heatmap) and layers (the bar chart) over additional GPT2 variants as measured using the alternate metric, normalized difference.

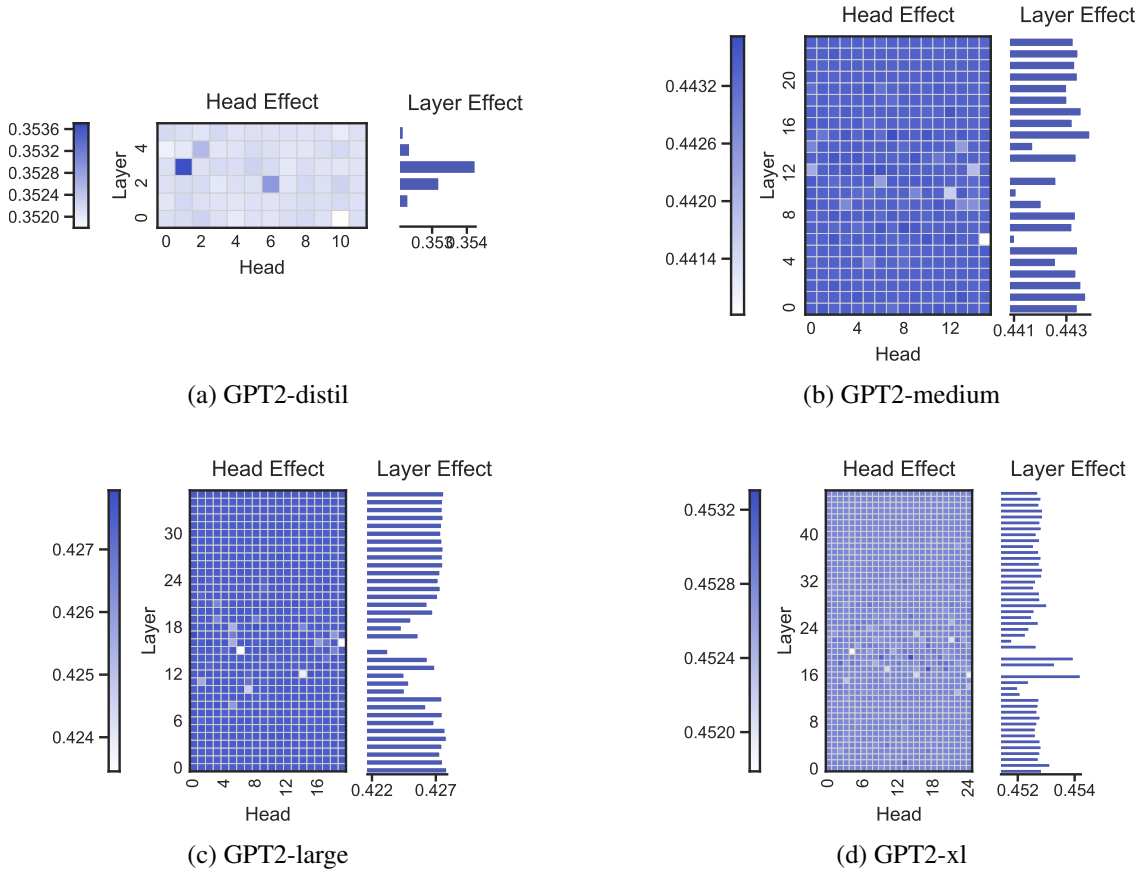


Figure 30: Mean indirect effect on Winobias for heads (the heatmap) and layers (the bar chart) over additional GPT2 variants as measured using the alternate metric, TV-distance.

Model	Winobias Dev Filtered			
	Original	Relative ℓ_∞	Norm. Diff.	TV-distance
GPT2-distil	0.118	0.105	0.045	0.358
GPT2-small	0.249	0.183	0.079	0.405
GPT2-medium	0.774	0.409	0.164	0.438
GPT2-large	0.751	0.383	0.166	0.425
GPT2-xl	1.049	0.446	0.205	0.456

Table 5: Total effects on Winobias Dev filtered for all metrics using different GPT2 variants.

Model	Winogender Bergsma Filtered			
	Original	Relative ℓ_∞	Norm. Diff.	TV-distance
GPT2-distil	0.075	0.068	0.029	0.341
GPT2-small	0.135	0.105	0.050	0.318
GPT2-medium	0.384	0.226	0.126	0.360
GPT2-large	0.350	0.222	0.120	0.322
GPT2-xl	0.362	0.224	0.110	0.364

Table 6: Total effects on Winogender Bergsma filtered for all metrics using different GPT2 variants.

Appendix G. Attention intervention experiments with masked language models

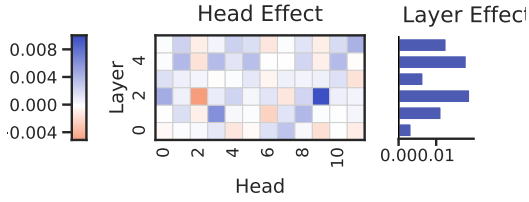
Scoring a multi-token continuation, say $[x_3, x_4]$, given a prefix, say $[x_1, x_2]$, is straightforward with autoregressive models: we simply need to combine the individual token-level probabilities $p_\theta(x_3 \mid x_1, x_2)$ and $p_\theta(x_4 \mid x_1, x_2, x_3)$ (which we do by taking a geometric mean). The problem becomes less trivial when we are faced with masked LMs, as there is not a single obvious way to compute token-level probabilities anymore. In this work, we try out three different ways of doing so, which in our running example would correspond to defining token-level probabilities as:

1. $p_\theta(x_3 \mid x_1, x_2, \text{mask})$ and $p_\theta(x_4 \mid x_1, x_2, x_3, \text{mask})$;
2. $p_\theta(x_3 \mid x_1, x_2, \text{mask}, x_4)$ and $p_\theta(x_4 \mid x_1, x_2, x_3, \text{mask})$;
3. $p_\theta(x_3 \mid x_1, x_2, \text{mask}, \text{mask})$ and $p_\theta(x_4 \mid x_1, x_2, x_3, \text{mask})$.

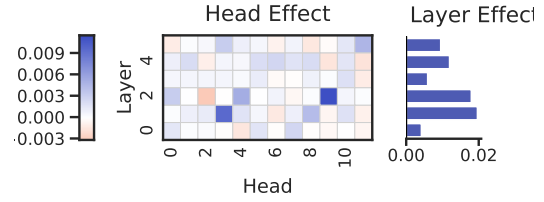
We propose scheme 1 because it seems to give the closest formulation to the autoregressive setting, while schemes 2 and 3 are derived from existing literature (Shin, Lee, & Jung, 2019; Wang & Cho, 2019; Salazar, Liang, Nguyen, & Kirchhoff, 2020). We also have a choice of whether or not to include the special `cls` and `sep` tokens (used during the pre-training of the models; e.g., `[CLS]`, and `[SEP]` in the case of BERT) when feeding examples to the masked LMs. Consequently, we define three additional scoring schemes which are exact copies of the original ones but with the special tokens included:

- 1'. $p_\theta(x_3 \mid \text{cls}, x_1, x_2, \text{mask}, \text{sep})$ and $p_\theta(x_4 \mid \text{cls}, x_1, x_2, x_3, \text{mask}, \text{sep})$;
- 2'. $p_\theta(x_3 \mid \text{cls}, x_1, x_2, \text{mask}, x_4, \text{sep})$ and $p_\theta(x_4 \mid \text{cls}, x_1, x_2, x_3, \text{mask}, \text{sep})$;
- 3'. $p_\theta(x_3 \mid \text{cls}, x_1, x_2, \text{mask}, \text{mask}, \text{sep})$ and $p_\theta(x_4 \mid \text{cls}, x_1, x_2, x_3, \text{mask}, \text{sep})$.

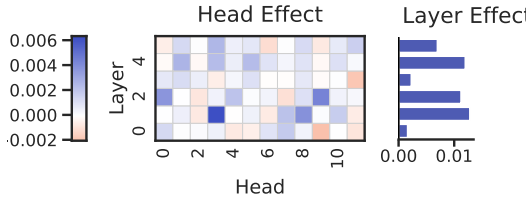
As mentioned in Section 5.5, we find considerable variation in the results with different masked LMs and scoring schemes. Figures 31, 32, 33, 34, and 35 show the indirect effects for each head and layer in DistilBERT, BERT-base-uncased, BERT-large-uncased, RoBERTa-base, and RoBERTa-large, respectively, using all six of the scoring schemes, on the filtered Winobias Dev dataset. We do observe the general trend of total effects being larger for larger variants of the same model, however. This is illustrated in tables 7 and 8, which list the total effects for different models and scoring schemes on the filtered Winobias Dev, Winogender Bergsma, and Professions datasets.



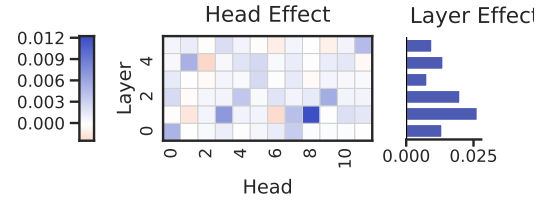
(a) Scoring scheme 1



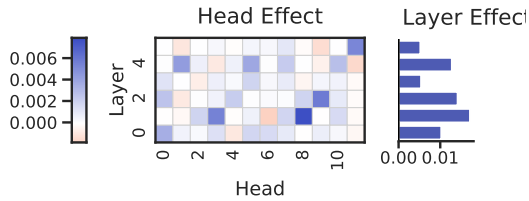
(b) Scoring scheme 2



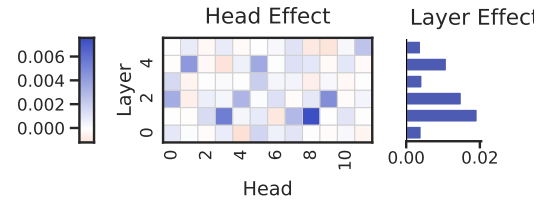
(c) Scoring scheme 3



(d) Scoring scheme 1'

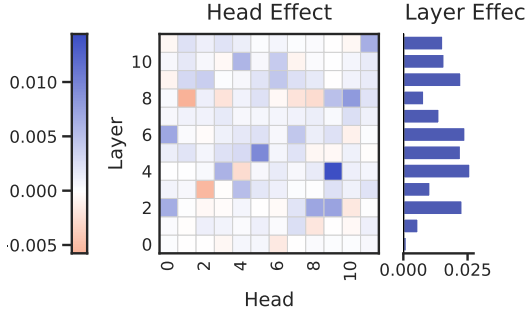


(e) Scoring scheme 2'

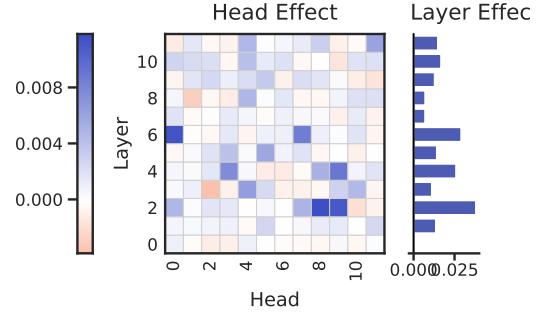


(f) Scoring scheme 3'

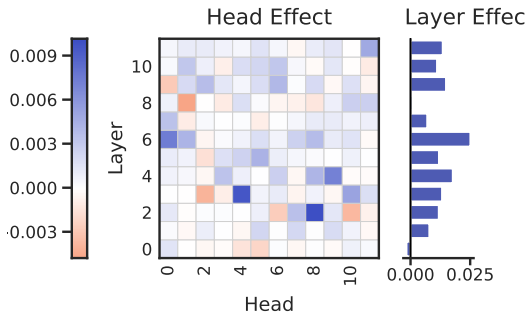
Figure 31: Mean indirect effect on Winobias for heads (the heatmap) and layers (the bar chart) in DistilBERT over different scoring schemes.



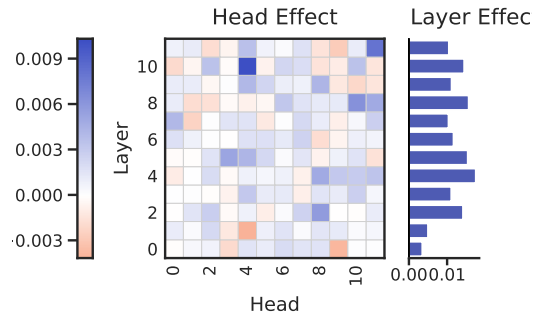
(a) Scoring scheme 1



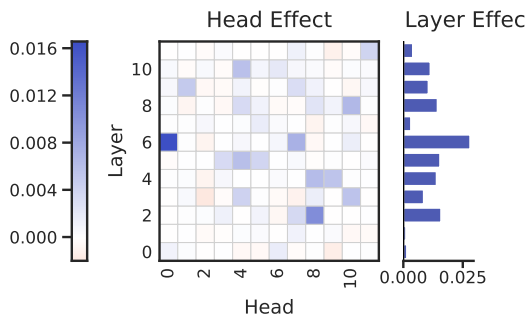
(b) Scoring scheme 2



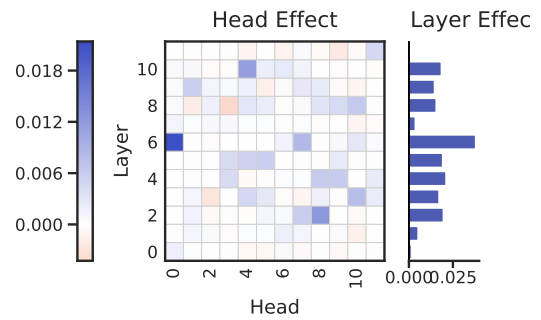
(c) Scoring scheme 3



(d) Scoring scheme 1'



(e) Scoring scheme 2'



(f) Scoring scheme 3'

Figure 32: Mean indirect effect on Winobias for heads (the heatmap) and layers (the bar chart) in BERT-base-uncased over different scoring schemes.

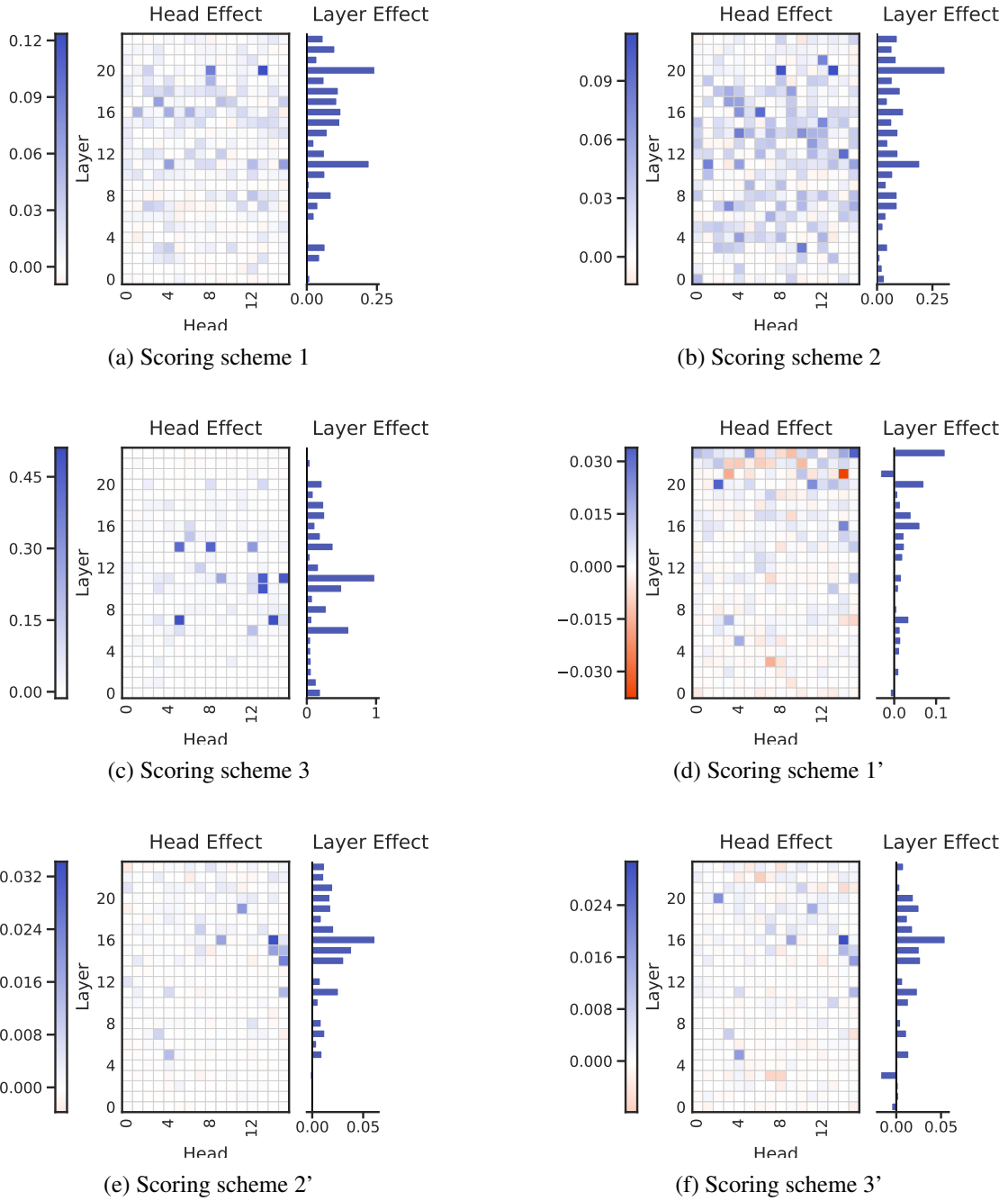


Figure 33: Mean indirect effect on Winobias for heads (the heatmap) and layers (the bar chart) in BERT-large-uncased over different scoring schemes.

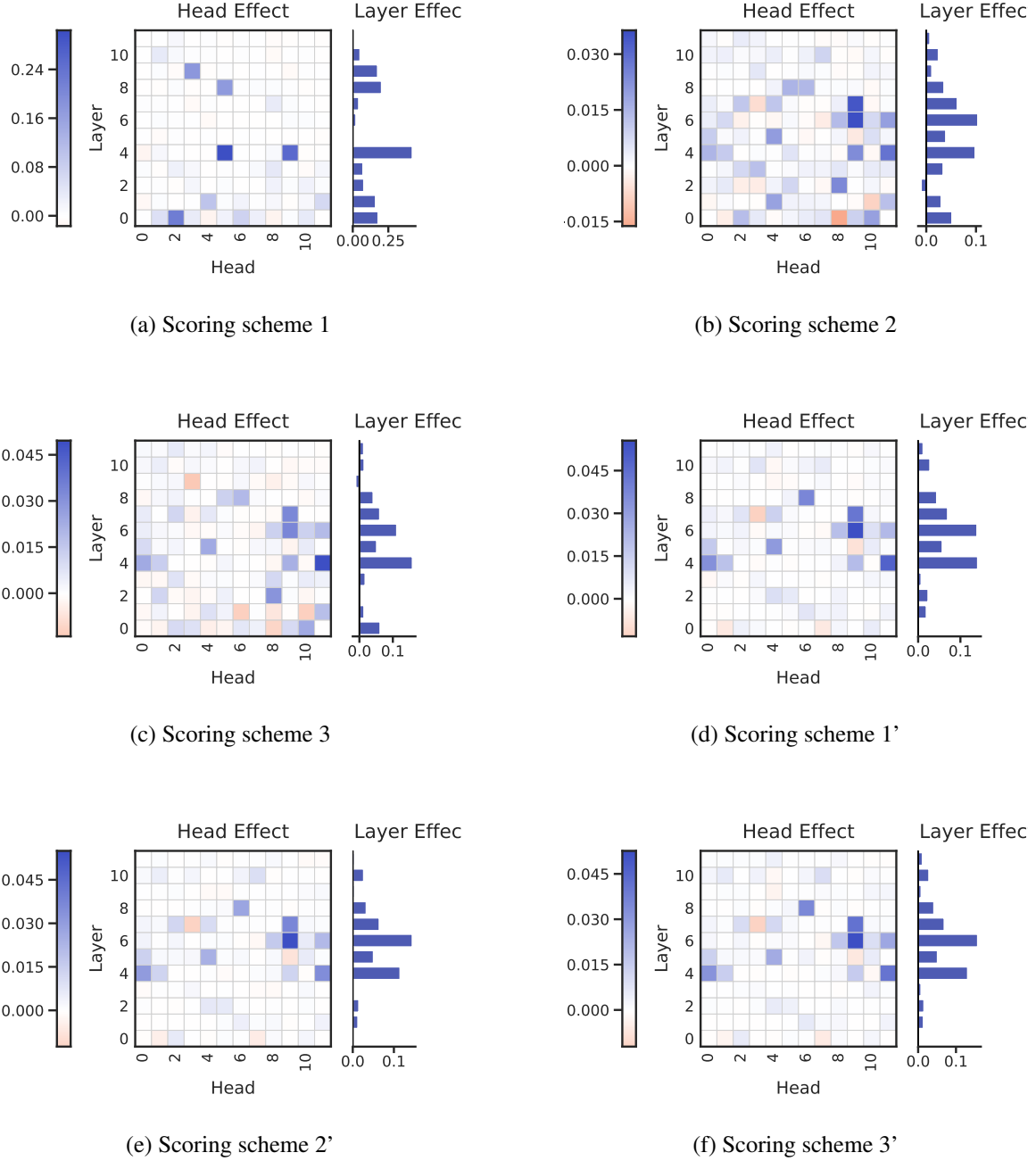


Figure 34: Mean indirect effect on Winobias for heads (the heatmap) and layers (the bar chart) in RoBERTa-base over different scoring schemes.

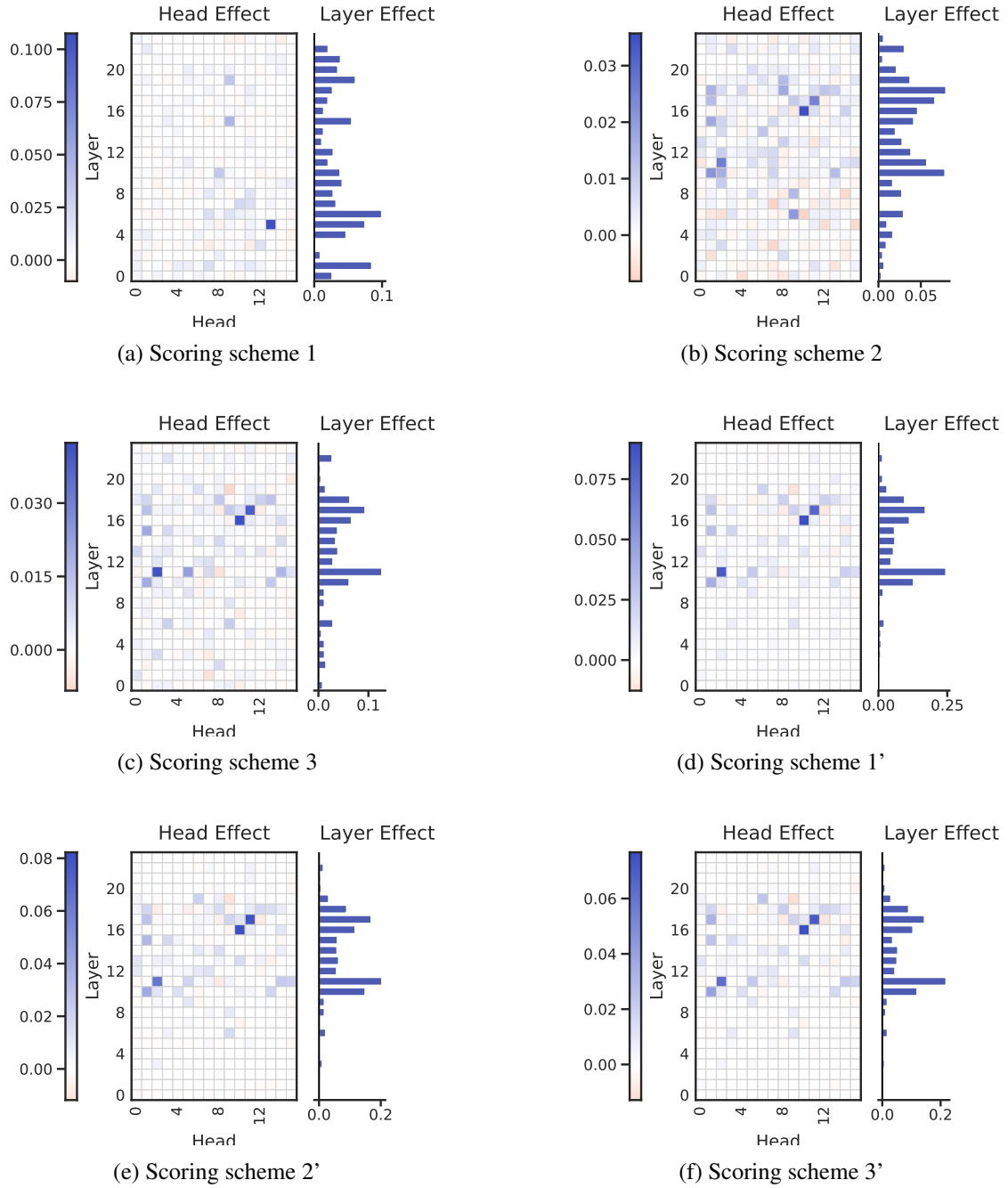


Figure 35: Mean indirect effect on Winobias for heads (the heatmap) and layers (the bar chart) in RoBERTa-large over different scoring schemes.

Model/Scoring scheme	Winobias Dev Filtered					
	1	2	3	1'	2'	3'
DistilBERT	0.224	0.191	0.158	0.215	0.189	0.161
BERT-base-uncased	0.213	0.179	0.178	0.262	0.240	0.245
BERT-large-uncased	0.935	0.744	0.950	0.430	0.303	0.304
RoBERTa-base	0.812	0.478	0.469	0.399	0.457	0.436
RoBERTa-large	0.402	0.610	0.614	0.994	1.033	0.859

Table 7: Total effects on Winobias Dev filtered for all masked LMs using different scoring schemes.

Model	WB	WG	Prof.
Transformer-XL	0.356	0.353	36.168
XLNet-base-cased	0.327	0.201	0.729
XLNet-large-cased	1.140	0.405	
DistilBERT	0.189	0.082	16.278
BERT-base-uncased	0.240	0.076	3.675
BERT-large-uncased	0.303	0.165	1.775
RoBERTa-base	0.457	0.191	29.625
RoBERTa-large	1.033	0.230	1.789

Table 8: Total effects (TE) of gender bias in models other than GPT2 evaluated on Winobias (WB), Winogender (WG), and the professions dataset (Prof.). Results for masked LMs on WB and WG are with scoring scheme 2'.

References

- Adi, Y., Kermany, E., Belinkov, Y., Lavi, O., & Goldberg, Y. (2017). Fine-grained analysis of sentence embeddings using auxiliary prediction tasks. In *Proceedings of the International Conference for Learning Representations (ICLR)*.
- Arras, L., Horn, F., Montavon, G., Müller, K.-R., & Samek, W. (2017). What is relevant in a text document?: An interpretable machine learning approach. *PLOS ONE*, 12(8), 1–23.
- Avin, C., Shpitser, I., & Pearl, J. (2005). Identifiability of path-specific effects. In *Proceedings of the 19th International Joint Conference on Artificial Intelligence*, pp. 357–363. Morgan Kaufmann Publishers Inc.
- Balkanski, E., Breuer, A., & Singer, Y. (2018). Non-monotone submodular maximization in exponentially fewer iterations. In *Advances in Neural Information Processing Systems*, pp. 2359–2370.
- Balkanski, E., & Singer, Y. (2018a). The adaptive complexity of maximizing a submodular function. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2018*, pp. 1138–1151, New York, NY, USA. ACM.
- Balkanski, E., & Singer, Y. (2018b). Approximation guarantees for adaptive sampling. In Dy, J., & Krause, A. (Eds.), *Proceedings of the 35th International Conference on Machine Learning*, Vol. 80 of *Proceedings of Machine Learning Research*, pp. 384–393, Stockholmsmassan, Stockholm Sweden.
- Barocas, S., Hardt, M., & Narayanan, A. (2019). *Fairness and Machine Learning*. fairmlbook.org. <http://www.fairmlbook.org>.
- Barrett, M., Kementchedjieva, Y., Elazar, Y., Elliott, D., & Søgaard, A. (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)*, pp. 6331–6336.
- Basta, C., Costa-jussà, M. R., & Casas, N. (2019). Evaluating the underlying gender bias in contextualized word embeddings. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, pp. 33–39, Florence, Italy. Association for Computational Linguistics.
- Belinkov, Y., & Glass, J. (2019). Analysis methods in neural language processing: A survey. *Transactions of the Association for Computational Linguistics*, 7, 49–72.
- Bergsma, S., & Lin, D. (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*, pp. 33–40, Sydney, Australia. Association for Computational Linguistics.
- Bolukbasi, T., Chang, K.-W., Zou, J. Y., Saligrama, V., & Kalai, A. T. (2016). Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In Lee, D. D., Sugiyama, M., Luxburg, U. V., Guyon, I., & Garnett, R. (Eds.), *Advances in Neural Information Processing Systems 29*, pp. 4349–4357. Curran Associates, Inc.

- Buchbinder, N., Feldman, M., Naor, J. S., & Schwartz, R. (2014). Submodular maximization with cardinality constraints. In *Proceedings of the Twenty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 1433–1452. Society for Industrial and Applied Mathematics.
- Caliskan, A., Bryson, J. J., & Narayanan, A. (2017). Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334), 183–186.
- Cao, Y. T., & Daumé III, H. (2020). Toward gender-inclusive coreference resolution. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 4568–4595, Online. Association for Computational Linguistics.
- Clark, K., Khandelwal, U., Levy, O., & Manning, C. D. (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*, pp. 276–286, Florence, Italy. Association for Computational Linguistics.
- Conneau, A., Kruszewski, G., Lample, G., Barrault, L., & Baroni, M. (2018). What you can cram into a single $\$&!#*$ vector: Probing sentence embeddings for linguistic properties. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2126–2136, Melbourne, Australia. Association for Computational Linguistics.
- Dai, Z., Yang, Z., Yang, Y., Carbonell, J., Le, Q., & Salakhutdinov, R. (2019). Transformer-XL: Attentive language models beyond a fixed-length context. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 2978–2988, Florence, Italy. Association for Computational Linguistics.
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (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.
- Dwork, C., Hardt, M., Pitassi, T., Reingold, O., & Zemel, R. (2011). Fairness through awareness..
- Elazar, Y., & Goldberg, Y. (2018). Adversarial removal of demographic attributes from text data. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 11–21.
- Elazar, Y., Ravfogel, S., Jacovi, A., & Goldberg, Y. (2020). When bert forgets how to POS: Amnesic probing of linguistic properties and MLM predictions. ArXiv e-prints.
- Ene, A., & Nguyen, H. L. (2019). *Submodular Maximization with Nearly-optimal Approximation and Adaptivity in Nearly-linear Time*, pp. 274–282.
- Fahrbach, M., Mirrokni, V., & Zadimoghaddam, M. (2019a). Non-monotone submodular maximization with nearly optimal adaptivity and query complexity. In Chaudhuri, K., & Salakhutdinov, R. (Eds.), *Proceedings of the 36th International Conference on Machine Learning*, Vol. 97 of *Proceedings of Machine Learning Research*, pp. 1833–1842, Long Beach, California, USA. PMLR.
- Fahrbach, M., Mirrokni, V., & Zadimoghaddam, M. (2019b). *Submodular Maximization with Nearly Optimal Approximation, Adaptivity and Query Complexity*, pp. 255–273.

- Feder, A., Oved, N., Shalit, U., & Reichart, R. (2020). Causalm: Causal model explanation through counterfactual language models. ArXiv e-prints.
- Gehrmann, S., Dernoncourt, F., Li, Y., Carlson, E. T., Wu, J. T., Welt, J., Foote Jr, J., Moseley, E. T., Grant, D. W., Tyler, P. D., et al. (2018). Comparing deep learning and concept extraction based methods for patient phenotyping from clinical narratives. *PloS one*, 13(2), e0192360.
- Gehrmann, S., Strobel, H., Krüger, R., Pfister, H., & Rush, A. M. (2019). Visual interaction with deep learning models through collaborative semantic inference. *IEEE Transactions on Visualization and Computer Graphics*, 26.
- Giulianelli, M., Harding, J., Mohnert, F., Hupkes, D., & Zuidema, W. (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*, pp. 240–248, Brussels, Belgium. Association for Computational Linguistics.
- Gonen, H., & Goldberg, Y. (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)*, pp. 609–614, Minneapolis, Minnesota. Association for Computational Linguistics.
- Hall Maudslay, R., Gonen, H., Cotterell, R., & Teufel, S. (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)*, pp. 5267–5275, Hong Kong, China. Association for Computational Linguistics.
- Hewitt, J., & Liang, P. (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)*, pp. 2733–2743, Hong Kong, China. Association for Computational Linguistics.
- Holland, P. W. (1986). Statistics and causal inference. *Journal of the American statistical Association*, 81(396), 945–960.
- Hoover, B., Strobel, H., & Gehrmann, S. (2019). exbert: A visual analysis tool to explore learned representations in transformers models. ArXiv e-prints.
- Huang, P.-S., Zhang, H., Jiang, R., Stanforth, R., Welbl, J., Rae, J., Maini, V., Yogatama, D., & Kohli, P. (2020). Reducing sentiment bias in language models via counterfactual evaluation. ArXiv e-prints.
- Hupkes, D., Veldhoen, S., & Zuidema, W. (2018). Visualisation and ‘diagnostic classifiers’ reveal how recurrent and recursive neural networks process hierarchical structure. *Journal of Artificial Intelligence Research*, 61, 907–926.
- Imai, K., Keele, L., & Tingley, D. (2010a). A general approach to causal mediation analysis. *Psychological methods*, 15(4), 309.
- Imai, K., Keele, L., & Yamamoto, T. (2010b). Identification, inference and sensitivity analysis for causal mediation effects. *Statistical Science*, 25(1), 51–71.

- Isabelle, P., Cherry, C., & Foster, G. (2017). A challenge set approach to evaluating machine translation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pp. 2486–2496, Copenhagen, Denmark. Association for Computational Linguistics.
- Kaushik, D., Hovy, E., & Lipton, Z. C. (2019). Learning the difference that makes a difference with counterfactually-augmented data. ArXiv e-prints.
- Kiritchenko, S., & Mohammad, S. (2018). Examining gender and race bias in two hundred sentiment analysis systems. In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, pp. 43–53.
- Kovaleva, O., Romanov, A., Rogers, A., & Rumshisky, A. (2019). Revealing the dark secrets of bert. ArXiv e-prints.
- Kurita, K., Vyas, N., Pareek, A., Black, A. W., & Tsvetkov, Y. (2019). Quantifying social biases in contextual word representations. In *1st ACL Workshop on Gender Bias for Natural Language Processing*.
- Kusner, M. J., Loftus, J., Russell, C., & Silva, R. (2017). Counterfactual fairness. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., & Garnett, R. (Eds.), *Advances in Neural Information Processing Systems 30*, pp. 4066–4076. Curran Associates, Inc.
- Li, J., Chen, X., Hovy, E., & Jurafsky, D. (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*, pp. 681–691, San Diego, California. Association for Computational Linguistics.
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., & Stoyanov, V. (2019). Roberta: A robustly optimized bert pretraining approach. ArXiv e-prints.
- Lu, K., Mardziel, P., Wu, F., Amancharla, P., & Datta, A. (2020). *Gender Bias in Neural Natural Language Processing*, pp. 189–202. Springer International Publishing, Cham.
- Madras, D., Creager, E., Pitassi, T., & Zemel, R. (2019). Fairness through causal awareness: Learning causal latent-variable models for biased data. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pp. 349–358. ACM.
- Montavon, G., Binder, A., Lapuschkin, S., Samek, W., & Müller, K.-R. (2019). Layer-wise relevance propagation: an overview. In *Explainable AI: interpreting, explaining and visualizing deep learning*, pp. 193–209. Springer.
- Murdoch, W. J., Liu, P. J., & Yu, B. (2018). Beyond word importance: Contextual decomposition to extract interactions from LSTMs. In *International Conference on Learning Representations*.
- Naik, A., Ravichander, A., Sadeh, N., Rose, C., & Neubig, G. (2018). Stress test evaluation for natural language inference. In *Proceedings of the 27th International Conference on Computational Linguistics*, pp. 2340–2353, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Nemhauser, G. L., & Wolsey, L. A. (1978). Best algorithms for approximating the maximum of a submodular set function. *Mathematics of operations research*, 3(3), 177–188.

- Pearl, J. (1995). Causal diagrams for empirical research. *Biometrika*, 82(4), 669–688.
- Pearl, J. (2001). Direct and indirect effects. In *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*, UAI’01, pp. 411–420, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Pearl, J. (2009). *Causality*. Cambridge university press.
- Qian, S., & Singer, Y. (2019). Fast parallel algorithms for feature selection..
- Qian, Y., Muaz, U., Zhang, B., & Hyun, J. W. (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*, pp. 223–228, Florence, Italy. Association for Computational Linguistics.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8).
- Robins, J. M. (2003). *Semantics of causal DAG models and the identification of direct and indirect effects*. OXFORD UNIV PRESS.
- Robins, J. M., & Greenland, S. (1992). Identifiability and exchangeability for direct and indirect effects. *Epidemiology*, 3, 143–155.
- Romanov, A., De-Arteaga, M., Wallach, H., Chayes, J., Borgs, C., Chouldechova, A., Geyik, S., Kenthapadi, K., Rumshisky, A., & Kalai, A. (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)*, pp. 4187–4195, Minneapolis, Minnesota. Association for Computational Linguistics.
- Rudinger, R., May, C., & Van Durme, B. (2017). Social bias in elicited natural language inferences. In *Proceedings of the First ACL Workshop on Ethics in Natural Language Processing*, pp. 74–79.
- Rudinger, R., Naradowsky, J., Leonard, B., & Van Durme, B. (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)*, pp. 8–14, New Orleans, Louisiana. Association for Computational Linguistics.
- Salazar, J., Liang, D., Nguyen, T. Q., & Kirchhoff, K. (2020). Masked language model scoring. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 2699–2712, Online. Association for Computational Linguistics.
- Sanh, V., Debut, L., Chaumond, J., & Wolf, T. (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)*.
- Sennrich, R. (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*, pp. 376–382, Valencia, Spain. Association for Computational Linguistics.
- Shin, J., Lee, Y., & Jung, K. (2019). Effective sentence scoring method using bert for speech recognition. In Lee, W. S., & Suzuki, T. (Eds.), *Proceedings of The Eleventh Asian Conference*

- on *Machine Learning*, Vol. 101 of *Proceedings of Machine Learning Research*, pp. 1081–1093, Nagoya, Japan. PMLR.
- Spirtes, P., Glymour, C. N., Scheines, R., & Heckerman, D. (2000). *Causation, prediction, and search*. MIT press.
- Strobelt, H., Gehrmann, S., Pfister, H., & Rush, A. M. (2017). Lstmvis: A tool for visual analysis of hidden state dynamics in recurrent neural networks. *IEEE transactions on visualization and computer graphics*, 24(1), 667–676.
- Tan, Y. C., & Celis, L. E. (2019). Assessing social and intersectional biases in contextualized word representations. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., & Garnett, R. (Eds.), *Advances in Neural Information Processing Systems 32*, pp. 13230–13241. Curran Associates, Inc.
- Tenney, I., Das, D., & Pavlick, E. (2019). BERT rediscovers the classical NLP pipeline. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4593–4601, Florence, Italy. Association for Computational Linguistics.
- VanderWeele, T. J., & Vansteelandt, S. (2009). Conceptual issues concerning mediation, interventions and composition. *Statistics and its Interface*, 2(4), 457–468.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, I. (2017). Attention is all you need. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., & Garnett, R. (Eds.), *Advances in Neural Information Processing Systems 30*, pp. 5998–6008. Curran Associates, Inc.
- Vig, J. (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*, pp. 37–42, Florence, Italy. Association for Computational Linguistics.
- Vig, J., & Belinkov, Y. (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*, pp. 63–76, Florence, Italy. Association for Computational Linguistics.
- Vig, J., Gehrmann, S., Belinkov, Y., Qian, S., Nevo, D., Singer, Y., & Shieber, S. (2020). Investigating gender bias in language models using causal mediation analysis. In *Advances in Neural Information Processing Systems*.
- Wang, A., & Cho, K. (2019). BERT has a mouth, and it must speak: BERT as a Markov random field language model. In *Proceedings of the Workshop on Methods for Optimizing and Evaluating Neural Language Generation*, pp. 30–36, Minneapolis, Minnesota. Association for Computational Linguistics.
- Webster, K., Recasens, M., Axelrod, V., & Baldrige, J. (2018). Mind the gap: A balanced corpus of gendered ambiguous pronouns. *Transactions of the Association for Computational Linguistics*, 6, 605–617.
- Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., & Brew, J. (2019). HuggingFace’s transformers: State-of-the-art natural language processing. *ArXiv, abs/1910.03771*.
- Yang, Z., & Feng, J. (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)*.

- Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R. R., & Le, Q. V. (2019). Xlnet: Generalized autoregressive pretraining for language understanding. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., & Garnett, R. (Eds.), *Advances in Neural Information Processing Systems* 32, pp. 5753–5763. Curran Associates, Inc.
- Zhao, J., Wang, T., Yatskar, M., Cotterell, R., Ordonez, V., & Chang, K.-W. (2019a). 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, J., Wang, T., Yatskar, M., Cotterell, R., Ordonez, V., & Chang, K.-W. (2019b). 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.
- Zhao, J., Wang, T., Yatskar, M., Ordonez, V., & Chang, K.-W. (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)*, pp. 15–20, New Orleans, Louisiana. Association for Computational Linguistics.
- Zhao, J., Zhou, Y., Li, Z., Wang, W., & Chang, K.-W. (2018b). Learning gender-neutral word embeddings. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 4847–4853.
- Zhao, Q., & Hastie, T. (2019). Causal interpretations of black-box models. *Journal of Business & Economic Statistics*, 0(0), 1–10.