

1 We thank the reviewers for their thoughtful feedback and helpful suggestions. We address specific points below.

2 **R1: missing assumptions.** The standard assumptions made in the causality literature are required when we observe only
3 one outcome per unit and cannot observe the counterfactual outcomes. These assumptions are needed for identification
4 of causal effects from observed outcomes. However, in our case, we perform two interventions on every unit (and on the
5 mediators) and observe all counterfactual outcomes. Therefore, these assumptions are not needed for our calculations
6 of mediation effects, and no statistical bias is expected from the analysis.

7 **R1: broader applicability.** While we focused in this work on binary interventions and outcomes, the existing literature
8 on causal mediation analysis enables the study of more general scenarios, including a different combination of variable
9 types (binary, categorical or continuous) for interventions, mediators and the outcomes.

10 **R1: defining bias.** Dwork (2012) defines an algorithm to be fair if it gives similar predictions to similar individuals.
11 The formalization of this definition was extended into Counterfactual Fairness (Kusner, 2017). We will explicitly define
12 bias as the extent to which an algorithm is not counterfactually fair and draw the connection to our outcome variable y .

13 **R1: dataset differences.** We believe the difference in NIE between the professions dataset (NIE concentrated in initial
14 layers) and the Winograd-style datasets (NIE concentrated in middle layers) reflects the fact that the former relates to
15 bias in word embeddings (lexical semantics), while the latter relates to bias in coreference, a higher-level phenomenon.

16 **R1: harms of gender binarization** We acknowledge that our current discussion of the unintended harms of treating
17 grammatical gender as binary variable is insufficient. Experimental results on he/they show very similar total effects to
18 he/she ($\pm 15\%$), although with a lower variance. We will add these results and a discussion of the measuring difficulties
19 of this effect under the singular/plural “they” confounder, as well as suggestions for mitigation, to the main body, and
20 will extend the impact statement.

21 **R1/R3: other model variants.** We now have additional results for Transformer-XL, BERT, DistilBERT, RoBERTa, and
22 XLNet, which are consistent with the results from GPT-2.

23 **R2: insights and takeaways.** Debiasing is an important research direction. Although we feel it is beyond the scope of
24 this paper, we believe our insights point to promising applications in evaluating and developing debiasing techniques.
25 One could envision manipulating mediators found through our method to reduce gender bias, e.g., setting them to a
26 null/neutral value. Further study is needed to evaluate how this approach impacts model bias and general performance.

27 **R2: heads targeting anti-stereotypical candidates.** Attention may capture negative as well as positive relationships,
28 depending on the head-specific value vectors to which the attention weights are applied. We hypothesize that attention
29 towards an anti-stereotypical candidate may decrease the probability of it being treated as the antecedent.

30 **R2: concentration of attention in specific heads.** Past work has shown that attention in middle layers correlates with
31 coreference (as you allude to), which is tightly related to our analysis of gender bias. Specialization of attention heads
32 has also been observed more generally, e.g., for various types of dependency relations (Clark et al., 2019). We will
33 expand the discussion of this point in the camera-ready version.

34 **R2: other types of biases.** For this novel adaptation of mediation analysis, we perform an extensive analysis of a
35 specific case study rather than a broader study of multiple phenomena, which would be a great area for future work.

36 **R2: correctness.** Would R2 point out the methodological problems that warrant a “no” answer to the correctness
37 question, such that we may address them?

38 **R2: missing related work.** Thank you for pointing this out. We will add this to the related work.

39 **R2/R3: limited scale and setup.** We are constrained by available resources. However, we find consistent results across
40 multiple models/datasets. In addition, the Winograd-style datasets are fairly nuanced in the linguistic phenomena.

41 **R3: only the reporting clause is reported.** In the professions dataset, we deliberately use verbs that are as neutral as
42 possible to focus on the profession word and the bias it leads to. In the Winograd-style datasets, the examples are much
43 more nuanced, in fact containing similar examples to the second one suggested by the reviewer. In this case, the bias
44 depends on the entire context in the prompt. We agree that our method doesn’t take into account potential bias in the
45 continuation itself.

46 **R3: only examines nouns.** This is true for the professions dataset, though in the Winograd-style datasets the verbs
47 play a role as well. We feel that a focused analysis of bias in verbs, while valuable, would warrant a separate study.

48 **R3: inter-sentential context.** Indeed, we have only looked at intra-sentential context. We note that some contexts are
49 rather nuanced in these Winograd-style datasets. Our methodology may be applied to inter-sentential contexts as well.

50 **R3: direct effects.** Figure 5 is representative of the direct effects we observed in all models: direct effects approximate
51 the difference between the total effects and the indirect effects. We will include additional results on direct effects.

52 **R3: difference from other studies on gender bias.** The main difference is that our research questions and methodology
53 focus specifically on mediators in bias by performing interventions. We will clarify this in the related work.

54 **R4: computational cost.** We recognize that that computational cost is non-trivial. We discuss computational complexity
55 in Appendix D with respect to the subset selection algorithm, but we will also discuss more generally in the main body.