# Counter-GAP: Counterfactual Bias Evaluation through Gendered Ambiguous Pronouns

**Zhongbin Xie[1], Vid Kocijan[2], Thomas Lukasiewicz[3,1], Oana-Maria Camburu[4]**
[1] University of Oxford  [2] Kumo.ai  [3] TU Wien
[4] University College London
`zhongbin.xie@cs.ox.ac.uk, vid@kumo.ai,`
`thomas.lukasiewicz@tuwien.ac.at, o.camburu@ucl.ac.uk`

## Abstract

Bias-measuring datasets play a critical role in detecting biased behavior of language models and in evaluating progress of bias mitigation methods. In this work, we focus on evaluating gender bias through coreference resolution, where previous datasets are either hand-crafted or fail to reliably measure an explicitly defined bias. To overcome these shortcomings, we propose a novel method to collect diverse, natural, and minimally distant text pairs via counterfactual generation, and construct Counter-GAP, an annotated dataset consisting of 4008 instances grouped into 1002 quadruples. We further identify a bias cancellation problem in previous group-level metrics on Counter-GAP, and propose to use the difference between inconsistency across genders and within genders to measure bias at a quadruple level. Our results show that four pre-trained language models are significantly more inconsistent across different gender groups than within each group, and that a name-based counterfactual data augmentation method is more effective to mitigate such bias than an anonymization-based method.

## 1 Introduction

It is a common practice to train state-of-the-art natural language processing (NLP) models by unsupervised pre-training and supervised fine-tuning (e.g., Devlin et al., 2019; Joshi et al., 2020), both of which rely heavily on large corpora of real-world text. However, these corpora often reflect societal stereotypes and may lead to models exhibiting biased behaviors (Bender et al., 2021). Hence, much research effort has been put to reveal and mitigate unintended biases (Meade et al., 2022).

While early work focuses on detecting and mitigating gender bias in the space of word embeddings (e.g., Bolukbasi et al., 2016), recent approaches turn to design bias-measuring datasets on specific NLP tasks (Nangia et al., 2020; Nadeem et al., 2021; Barikeri et al., 2021). In this work, we

focus on gender bias in coreference resolution and adopt a kind of representational harm (Blodgett et al., 2020) to define gender fairness: a gender-neutral model should rely on the semantic information, rather than on the gender information contained in the texts, to make predictions. Otherwise, a model should be considered gender-biased. In line with this definition, WinoBias (Zhao et al., 2018) and WinoGender (Rudinger et al., 2018) leverage pairs of minimally distant sentences, i.e., two sentences that contain the same semantic information but different gender information, to measure models' performance difference in resolving pronouns of different genders under the same context. This minimally distant setting enables us to isolate the influence of gender information on model predictions.

A limitation of WinoBias and WinoGender is that they are made up of *hand-crafted sentences*, which prevents us from measuring gender bias in the more diverse real-world scenarios. An alternative to overcome this shortcoming is the GAP dataset (Webster et al., 2018), which exploits linguistic patterns to automatically extract instances from a *real-world corpus*. However, since GAP's masculine and feminine instances cannot be grouped into minimally distant pairs, we are not sure whether a difference in model performance is due to different gender information or to different semantic information. For example, compared to masculine instances, GAP's feminine instances have more candidate entities serving as distractors, and longer distance between the correct name and the pronoun (Kocijan et al., 2021). So, it is *not equally hard to resolve* the masculine and feminine instances in GAP. Hence, the performance difference between masculine and feminine instances on GAP is not a reliable measure of gender bias according to the above definition of gender fairness.[1]

---

[1]In GAP (Webster et al., 2018), the authors did not explicitly describe the fairness definiton that they adopt.

Given these observations, we propose a novel method to construct coreference-resolution-based bias-measuring datasets consisting of minimally distant text pairs that originate from real-world corpora. Specifically, we leverage the method from GAP (Webster et al., 2018) to extract original instances containing gendered ambiguous pronouns, and generate minimally distant instances by asking the counterfactual question "How would the prediction change if we swapped the roles of masculine and feminine people in this context?" (Garg et al., 2019). The resulting instances are grouped into quadruples, each of which consists of an original, a gender-controlled, and two gender-swapped instances. An example is shown in Table 1.

Furthermore, we find that bias in different directions may be canceled out if we aggregate the results by performance difference across groups of instances, and we call this problem *bias cancellation*. To alleviate it, we propose a new metric, inconsistency across genders, to measure bias at the quadruple level. We also leverage the gender-controlled instances to disentangle inconsistency *within* genders from inconsistency *across* genders, so that we can eliminate the impact of name perturbations.

Our contributions are as follows: (i) We propose a novel method to construct coreference resolution datasets consisting of diverse, natural, and minimally distant instances to reliably detect gender bias. (ii) We apply our method to online books and collect Counter-GAP, an annotated dataset with 4008 instances grouped into 1002 quadruples. (iii) We propose a new metric, the difference ($\Delta I$) between inconsistency across genders and within genders, to alleviate the bias cancellation problem of previous metrics. (iv) We use Counter-GAP to empirically evaluate four pre-trained language models and two debiasing methods based on Counterfactual Data Augmentation (CDA, Zhao et al., 2018; Webster et al., 2020). Our results show that $\Delta I$ can detect significant gender bias hidden by group-level performance difference, and that name-based CDA is more effective than vanilla CDA in mitigating such bias.[2]

## 2 Dataset Construction

The Counter-GAP dataset is derived from 1575 fictional books in Project Gutenberg[3] and BookCorpus (Zhu et al., 2015). It is constructed through a generic multi-stage process, as described below. Here, we follow the GAP dataset (Webster et al., 2018) and focus only on the English language, as well as adopting a notion of binary gender.

### 2.1 Original Instance Extraction

First, we detect all the occurrences of personal names and pronouns in a book with a dependency parser and a named entity recognizer (NER).[4] For each occurrence of a gendered non-reflexive pronoun (*he*, *him*, *his*, *she*, *her*, *hers*), we extract a surrounding context that consists of a maximum of five sentences and contains exactly two masculine and two feminine personal names. Personal names are identified by NER tags, and the gender specification of a name is determined by statistics from a gender-guesser.[5] Genders for titled names (e.g., *Mr. Smith*) are assumed from the traditional gender associations of those titles.

Second, we select the subset of contexts that contain gendered ambiguous pronouns as defined by the following three patterns from GAP (Webster et al., 2018) (henceforth, the gendered ambiguous pronoun is called *target pronoun*, and the two names that are gender-consistent with the target pronoun are called *candidate names*):

- **FINALPRO.** Both candidate names must be in the same sentence, and the target pronoun may appear in the same or directly following sentence.

- **MEDIALPRO.** The first candidate name must be in the sentence directly preceding the target pronoun and the second candidate name, both of which must be in the same sentence. The target pronoun must be in an initial subordinate clause or be a possessive in an initial prepositional phrase.

- **INITIALPRO.** Both the candidate names and the target pronoun must be in the same sentence, and the target pronoun must be in an

---

[2]The dataset and code are available at `https://github.com/x-zb/Counter-GAP`.

[3]`https://www.gutenberg.org/`

[4]We use Spacy (`https://spacy.io/`).

[5]`https://pypi.org/project/gender-guesser/`. As we focus on English books, we use the default setting where the gender of a name is first considered according to its use in English-speaking countries.

| original | Tom did not appear to hear this, but tried to keep up the conversation with Julia, desiring to have it appear that they were intimate friends; but the young lady gave brief replies, and finally, turning away, devoted herself once more to Herbert, much to Tom's disgust. In fact, what he saw made Tom pass a very unpleasant evening, and when, on their return home, Maria suggested that Julia had taken a fancy to Herbert, **he** told her to mind her own business. |
|---|---|
| gender-controlled | Herbert did not appear to hear this, but tried to keep up the conversation with Maria, desiring to have it appear that they were intimate friends; but the young lady gave brief replies, and finally, turning away, devoted herself once more to Tom, much to Herbert's disgust. In fact, what he saw made Herbert pass a very unpleasant evening, and when, on their return home, Julia suggested that Maria had taken a fancy to Tom, **he** told her to mind her own business. |
| gender-swapped-1 | Maria did not appear to hear this, but tried to keep up the conversation with Herbert, desiring to have it appear that they were intimate friends; but the young *gentleman* gave brief replies, and finally, turning away, devoted *himself* once more to Julia, much to Maria's disgust. In fact, what *she* saw made Maria pass a very unpleasant evening, and when, on their return home, Tom suggested that Herbert had taken a fancy to Julia, **she** told *him* to mind *his* own business. |
| gender-swapped-2 | Julia did not appear to hear this, but tried to keep up the conversation with Tom, desiring to have it appear that they were intimate friends; but the young *gentleman* gave brief replies, and finally, turning away, devoted *himself* once more to Maria, much to Julia's disgust. In fact, what *she* saw made Julia pass a very unpleasant evening, and when, on their return home, Herbert suggested that Tom had taken a fancy to Maria, **she** told *him* to mind *his* own business. |

Table 1: Counterfactual generation of a quadruple in Counter-GAP. Personal names and their genders are depicted in colors: masculine names are in blue and cyan; feminine names are in violet and orange. The target pronoun is in **bold** and underlined; also underlined is the true coreferent name. Other words constitute the context, and words in *italic* are gendered words swapped according to the gendered words list.

initial subordinate clause or a possessive in an initial prepositional phrase.

After filtering, we get 2585 contexts and adopt them as original instances.

## 2.2 Counterfactual Generation

We generate minimally distant instances in Counter-GAP through two counterfactual generation functions. An example is illustrated in Table 1. Formally, we denote an original instance as $x_o = s(P, C_1, C_2, O_1, O_2)$, where $P$ is the target pronoun to be resolved, $C_1$ and $C_2$ are the two candidate names that are gender-consistent with $P$, $O_1$ and $O_2$ are two personal names of the opposite gender, and $s(\cdot)$ denotes the context around these mentions.

**Gender-controlled generation.** We swap all the occurrences of $C_1$ and $C_2$, and of $O_1$ and $O_2$, to generate a gender-controlled instance $x_c = s(P, C_2, C_1, O_2, O_1)$. We choose to swap names within an instance instead of introducing new names, so that the candidate names naturally occur in the same real-world context.

**Gender-swapped generation.** We first substitute all gendered words with their opposite gendered words (e.g., man→woman, he→she),[6] and swap all the occurrences of $C_1$ and $O_1$ (or $O_2$), $C_2$ and

$O_2$ (or $O_1$). As a result, we obtain two gender-swapped instances $\widetilde{x_o} = \tilde{s}(\tilde{P}, O_1, O_2, C_1, C_2)$ and $\widetilde{x_c} = \tilde{s}(\tilde{P}, O_2, O_1, C_2, C_1)$, where $\tilde{s}(\cdot)$ is the context with all the gendered words substituted in $s(\cdot)$, and $\tilde{P}$ is the opposite-gendered pronoun for $P$. We call $x_o, x_c, \widetilde{x_o}, \widetilde{x_c}$ minimally distant instances, in that the words at the same position in the context ($s(\cdot)$ or $\tilde{s}(\cdot)$) are either the same (for gender-neutral words) or have the same role but opposite gender (for gendered words).

We consider a generated counterfactual instance to be invalid if (i) it contradicts commonsense knowledge, e.g., historical people being of the opposite gender; or (ii) the meaning of the counterfactual is different from the original, resulting in the gold coreference labels changing or becoming undetermined. To tackle these, we take three measures. First, we extract original instances mainly from fictional books, whose content is less likely to involve real-world people. Second, during human annotation (Section 2.3), we explicitly ask annotators to validate whether an instance contradicts commonsense knowledge, and discard such instances. Third, we discard the whole quadruple $(x_o, x_c, \widetilde{x_o}, \widetilde{x_c})$ if not all of its four instances get the same majority labels from annotators.[7] Here, *same*

---

[6] We adopt an augmented list of gendered words from (Zhao et al., 2018).

[7] Note that discarding inconsistent quadruples can also cover some error cases caused by the gender-guesser's incorrect prediction. For example, if an incorrect gender prediction occurs for a TRUE coreferent name, the target pronoun and the TRUE coreferent name will not be gender-consistent, and

*label* means the coreferent names' positions are the same in the context (e.g., in Table 1, the position of "Tom" in the original instance and that of "Maria" in the gender-swapped-1 instance).

## 2.3 Human Annotation

We use Amazon Mechanical Turk to collect coreference labels for all the 2585 original instances and their counterfactual counterparts (hence, $2585 \times 4$ instances in total). Each instance was assigned to three annotators. Annotation instructions and a sample task interface are presented in Appendix B. Specifically, we ask annotators to highlight token spans that are coreferent with the target pronoun. We adopt majority vote to aggregate the collected annotations, and generate a `TRUE/FALSE` label for each of the two candidate names indicating if it is coreferent with the target pronoun.

After discarding quadruples containing invalid counterfactuals as discussed in Section 2.2, we further filter out quadruples containing real-world people to avoid grounding. Next, we randomly downsample the remaining quadruples to balance the number of original masculine and feminine instances. The final Counter-GAP dataset consists of 1002 quadruples with an inter-annotator agreement[8] of 86.5%.

## 3 Evaluation Metrics on Counter-GAP

We use $X = (x_o, x_c, \widetilde{x_o}, \widetilde{x_c}) \in \mathcal{X}$ to denote a quadruple, and lowercased $x$ to denote an arbitrary instance, which could be each of $x_o, x_c, \widetilde{x_o}, \widetilde{x_c}$ from a quadruple. Given a model $f(\cdot)$, assume that $f(x) \in \{0, 1\}$ indicates whether $f$'s prediction on instance $x$ is correct (1) or not (0).

### 3.1 Bias Cancellation in Accuracy Difference

A so far commonly used metric is to directly compare the model's performance difference (or ratio) between different gender groups (Webster et al., 2018; Sun et al., 2019; Blodgett et al., 2020). For example, if we divide a test set $\mathcal{X}$ into a group of masculine instances $\mathcal{D}^{(m)}$ and a group of feminine instances $\mathcal{D}^{(f)}$ according to the gender information contained in the instances (e.g., the gender of the target pronoun), gender bias can be measured by model $f$'s accuracy difference ($Acc_{Diff}$) on $\mathcal{D}^{(m)}$

and $\mathcal{D}^{(f)}$:

$$Acc_{Diff} = \frac{\sum_{x \in \mathcal{D}^{(m)}} f(x)}{|\mathcal{D}^{(m)}|} - \frac{\sum_{x \in \mathcal{D}^{(f)}} f(x)}{|\mathcal{D}^{(f)}|}. \quad (1)$$

However, the above metric may suffer from *bias cancellation* on Counter-GAP. Consider two quadruples from Counter-GAP. In the first, the model makes correct predictions on the two masculine instances and incorrect predictions on the two feminine ones. The model should be deemed gender-biased (towards masculine), since it makes different predictions for instances containing the same semantic information. If, in the second quadruple, the model makes reversed predictions, i.e., correct on the two feminine instances and incorrect on the two masculine ones, it should also be deemed gender-biased, yet in the opposite direction towards feminine. However, the model's accuracies on the masculine and feminine groups are both $2/4 = 50\%$, making Eq. (1) equal to zero. In short, biases in opposite directions may be canceled out in some cases if we use Eq. (1) to aggregate them.

## 3.2 Measuring Bias via Inconsistencies

Given the bias cancellation problem of accuracy difference, we propose to measure gender bias through inconsistencies, i.e., whether a model's prediction is consistent on a pair of minimally distant instances. Specifically, we adopt two metrics, *inconsistency across genders* ($I_{across}$):

$$\frac{1}{4|\mathcal{X}|} \sum_{X \in \mathcal{X}} \Big( |f(x_o) - f(\widetilde{x_o})| + |f(x_c) - f(\widetilde{x_c})| \\ + |f(x_o) - f(\widetilde{x_c})| + |f(x_c) - f(\widetilde{x_o})| \Big), \quad (2)$$

and *inconsistency within genders* ($I_{within}$):

$$\frac{1}{2|\mathcal{X}|} \sum_{X \in \mathcal{X}} \Big( |f(x_o) - f(x_c)| + |f(\widetilde{x_o}) - f(\widetilde{x_c})| \Big). \quad (3)$$

Inconsistency across genders ($I_{across}$) measures inconsistency in instance pairs containing two instances of different genders, while inconsistency within genders ($I_{within}$) measures inconsistency in instance pairs containing two instances of the same gender. The two instances in a pair should be minimally distant (i.e., from the same quadruple) to guarantee that they contain the same semantic information. Since Counter-GAP adopts personal names as proxies for person entities, we need to

---

this may confuse the annotators, leading to inconsistent labels.

[8]Average percentage of agreed annotations on each instance.

disentangle the part of inconsistency caused by different names ($I_{within}$) from that caused by different genders ($I_{across}$). Therefore, our final metric to measure gender bias is

$$\Delta I = I_{across} - I_{within}. \qquad (4)$$

In practice, a positive $\Delta I$ indicates biased behaviors of the model, while a zero or negative $\Delta I$ means that the measured inconsistency across genders are mostly noises from name perturbations, thus no bias can be detected.

## 4 Bias Evaluation on Counter-GAP

For evaluation, we adopt the coreference resolution system `c2f-coref`[9] (Lee et al., 2018) based on four pre-trained language models: BERT-base/large and SpanBERT-base/large (Joshi et al., 2020). All four models are fine-tuned on OntoNotes (Pradhan et al., 2012),[10] and training details are shown in Appendix A. In our evaluation, no candidate names are provided as input to the models, and models are responsible to detect candidate names in the text by themselves. A model's prediction on an instance is considered correct if the candidate name with gold label TRUE and none of those with gold label FALSE are in the target pronoun's coreferent cluster.

### 4.1 Results

Results for gender bias measured on Counter-GAP by accuracy difference (Eq. (1)) and $\Delta I$ (Eq. (4)) are shown in Tables 2 and 3, respectively. In Table 3, we also report the inconsistency metrics on each gender group (M, F) and each swapping direction (M2F, F2M), together with their differences.

Results in Table 3 show that for all four models, not only $I_{across}$ is larger than $I_{within}$ ($\Delta I$ being positive), but also the difference is statistically significant, which indicates biased behaviors in these models. Note that the absolute values of accuracy difference ($Acc_{Diff}$) in Table 2 are in general smaller than the corresponding values of $\Delta I$ in Table 3, and $Acc_{Diff}$ for BERT-large even becomes statistically insignificant, which is contrary to the well-known conclusion that BERT encodes social bias (Nadeem et al., 2021). This brings evidence

---

[9]We use the implementations from `https://github.com/mandarjoshi90/coref`.

[10]Since the annotation conventions of OntoNotes are a little different from those of Counter-GAP, we omitted the abbreviation period "." in titles like "Mr.", "Mrs.", and "Dr." in Counter-GAP during evaluation.

towards the bias cancellation problem, i.e., bias measured by accuracy difference (Eq. (1)) may be canceled out compared to that measured by inconsistency difference ($\Delta I$).

Regarding the effect of model size on gender bias, results from both metrics show that larger models seem to be less biased than smaller models. Note that both our large and base models are (pre-)trained on the same datasets, but in general larger language models are pre-trained on larger amount of data, so they are still at a higher risk of exhibiting biased behaviors (Bender et al., 2021).

Regarding the detected bias direction, different metrics provide information from different perspectives. We can learn from the sign of $Acc_{Diff}$ in Table 2 that the overall bias directions of these models are all towards masculine. In Table 3, all of the Diff. for $I_{within}$ are negative, indicating a larger inconsistency within the feminine group. All of the Diff. for $I_{across}$ being negative indicates that inconsistency will increase when we change genders in an originally feminine context.

| Models | $Acc_M$ | $Acc_F$ | $Acc_{Diff}$ |
|---|---|---|---|
| BERT-base | 63.12% | 59.53% | +3.59%* |
| BERT-large | 72.60% | 72.11% | +0.50% |
| SpanBERT-base | 71.36% | 69.06% | +2.30%* |
| SpanBERT-large | 77.25% | 75.40% | +1.85%* |

Table 2: Gender bias measured by Eq. (1) on Counter-GAP. We report accuracy on masculine instances ($Acc_M$), feminine instances ($Acc_F$), and their difference ($Acc_{Diff} = Acc_M - Acc_F$). A "∗" means that the difference is statistically significant ($p < 0.01$) under one-sided bootstrap resampling (Graham et al., 2014).

### 4.2 No Bias Between the Original and Counterfactual Instances

Since the gender-swapped instances in Counter-GAP are generated automatically, although they have been validated by annotators, we still check whether there is a systematic bias towards the original or counterfactual instances. To investigate this, we measure two statistics. First, we measure a model's accuracy on instances with the original gender ($x_o, x_c$) and the counterfactual gender ($\widetilde{x_o}, \widetilde{x_c}$), and report their difference. From Table 4, we see that the differences are very small and not statistically significant. Second, we measure the correlation between the inconsistency across genders score ($|f(x_o) - f(\widetilde{x_o})| + |f(x_c) - f(\widetilde{x_c})| + |f(x_o) - f(\widetilde{x_c})| + |f(x_c) - f(\widetilde{x_o})|$) and the original gender of a quadruple $X$. From Table 4, we

| Models | inconsistency within genders | | | | inconsistency across genders | | | | $\Delta I = I_{across}$ $-I_{within}$ |
|---|---|---|---|---|---|---|---|---|---|
| | M | F | Diff. | $I_{within}$ | M2F | F2M | Diff. | $I_{across}$ | |
| BERT-base | 15.47% | 16.47% | -1.00% | 15.97% | 18.26% | 23.25% | -4.99% | 20.76% | +4.79%* |
| BERT-large | 10.28% | 10.28% | 0.00% | 10.28% | 10.88% | 14.27% | -3.39% | 12.57% | +2.30%* |
| SpanBERT-base | 9.98% | 12.18% | -2.20% | 11.08% | 12.18% | 15.07% | -2.89% | 13.62% | +2.54%* |
| SpanBERT-large | 5.79% | 6.29% | -0.50% | 6.04% | 6.89% | 8.18% | -1.30% | 7.53% | +1.50%* |

Table 3: Gender bias measured by the inconsistency metrics on Counter-GAP. We report inconsistency within genders on masculine instances (M), feminine instances (F), and their difference (Diff. = M - F), as well as inconsistency within genders on all the instances ($I_{within}$). We also report inconsistency across genders on quadruples generated by transforming masculine instances to feminine instances (M2F), transforming feminine instances to masculine instances (F2M), and their difference (Diff. = M2F - F2M), as well as inconsistency across genders on all the quadruples ($I_{across}$). $\Delta I = I_{across} - I_{within}$ measures gender bias, where a "*" means that the difference is statistically significant ($p < 0.01$) under one-sided bootstrap resampling (Graham et al., 2014).

| Models | Accuracy | | | Spear-man's $\rho$ |
|---|---|---|---|---|
| | Orig. | Counter. | Diff. | |
| BERT-base | 61.58% | 61.08% | +0.50% | -0.083 |
| BERT-large | 72.06% | 72.65% | -0.60% | -0.065 |
| SpanBERT-base | 70.21% | 70.21% | 0.00% | -0.060 |
| SpanBERT-large | 76.55% | 76.10% | +0.45% | -0.030 |

Table 4: Results on systematic bias evaluation. We report accuracy on instances with the original gender (Orig.), with the counterfactual gender (Counter.), and their difference (Diff.= Orig. - Counter.). All the differences are not statistically significant ($p > 0.01$) under one-sided bootstrap resampling (Graham et al., 2014). We also report Spearman's $\rho$ between inconsistency across genders and the original gender of a quadruple.

see that the values of Spearman's $\rho$ are all close to zero, indicating no significant correlations. Hence, we conclude that the counterfactual instances in Counter-GAP do not introduce systematic bias.

### 4.3 Comparison with GAP

We further compare Counter-GAP with two GAP-like datasets: the original GAP test set (Webster et al., 2018) and a subset of Counter-GAP where only the original instances $x_o$ are kept (we call this dataset our-GAP). The results of SpanBERT-large on the above datasets are shown in Table 5. We see that the accuracy differences between masculine and feminine instances are much smaller on Counter-GAP than on the original GAP and our-GAP. This empirically verifies that datasets without minimally distant instances cannot reliably measure bias (they amplify bias in this case) due to the different semantic information contained in its masculine and feminine instances. Moreover, the overall direction of detected gender bias (the sign of $Acc_{Diff}$) is different on the original GAP and our-GAP, which shows that different source corpora (Wikipedia for GAP vs. fictions for our-GAP)

may detect different bias in the model. This highlights the importance of domain diversity when using data-centric methods for bias detection.

## 5 Bias Mitigation

We evaluate two debiasing methods based on counterfactual data augmentation (CDA) (Zhao et al., 2018; Webster et al., 2020): (i) anonymization-based CDA (a-CDA), where the training set (OntoNotes) is augmented by substituting all gendered words with their opposite gendered words, while the gold coreference labels are kept unchanged. Personal names in the training set are anonymized using place holders such as "E1, E2, . . ."; (ii) name-based CDA (n-CDA), where, in addition to the substitution between gendered words, masculine and feminine names also substitute each other according to their frequencies (Hall Maudslay et al., 2019). See Appendix A for more details. Performance on bias mitigation is measured by $Acc_{Diff}$ and $\Delta I$, while performance on coreference resolution is measured by overall accuracy on Counter-GAP and $F_1$ score on OntoNotes' dev set.

Results are shown in Table 6. In terms of $\Delta I$, both a-CDA and n-CDA can effectively reduce gender bias, while n-CDA is more effective than a-CDA in that its $\Delta I$ values are smaller and less significant. Comparing the results of $Acc_{Diff}$ and $\Delta I$, we discover that bias measured by $Acc_{Diff}$ tends to be more easily mitigated by the debiased methods. For example, a-CDA fails to reduce $\Delta I$ to an insignificant level for all the four models, but it succeeds to do so for $Acc_{Diff}$ on BERT-base and SpanBERT-large; n-CDA can reduce BERT-base's $Acc_{Diff}$ to an insignificant level, but fails to do so under the measurement of $\Delta I$.

Regarding the trade-off between bias mitigation and overall performance, both a-CDA and n-CDA

| Datasets | Accuracy | | | | $\Delta I$ |
|---|---|---|---|---|---|
| | $Acc_M$ | $Acc_F$ | $Acc_{Diff}$ | Overall | |
| original GAP | 85.10% | 79.80% | +5.30% | 82.45% | – – |
| our-GAP | 75.25% | 78.44% | -3.19% | 76.85% | – – |
| Counter-GAP | 77.25% | 75.40% | +1.85% | 76.32% | +1.50% |

Table 5: Results from SpanBERT-large on three datasets.

| Models | Debiasing Method | $Acc_{Diff} =$ $Acc_M - Acc_F$ | $\Delta I =$ $I_{across} - I_{within}$ | Overall Accuracy | $F_1$ on OntoNotes |
|---|---|---|---|---|---|
| BERT-base | none | +3.59%* | +4.79%* | 61.33% | 74.39% |
| | a-CDA | +1.90% | +2.30%* | 65.17% | 73.91% |
| | n-CDA | **+0.20%** | +1.85%* | 66.82% | 73.60% |
| BERT-large | none | +0.50% | +2.30%* | 72.36% | 77.35% |
| | a-CDA | +2.99%* | +1.75%* | 72.95% | 76.96% |
| | n-CDA | +0.95% | +1.30% | 73.53% | 77.13% |
| SpanBERT -base | none | +2.30%* | +2.54%* | 70.21% | 77.71% |
| | a-CDA | +5.14%* | +2.54%* | 69.49% | 78.04% |
| | n-CDA | +0.95% | +1.25% | 71.03% | 77.70% |
| SpanBERT -large | none | +1.85%* | +1.50%* | 76.32% | 80.06% |
| | a-CDA | +0.35% | +1.45%* | 76.72% | **80.07%** |
| | n-CDA | +0.65% | **+0.15%** | **77.92%** | 79.93% |

Table 6: Bias mitigation results. For $Acc_{Diff}$ and $\Delta I$, lower is better; for overall accuracy and $F_1$ on OntoNotes, higher is better. Best results are in bold. A "*" on $\Delta I$ indicates that the difference is statistically significant ($p < 0.01$) under one-sided bootstrap resampling (Graham et al., 2014).

can maintain or even increase the overall accuracy on Counter-GAP. This indicates that they do not sacrifice model performance for fairness, which is a favorable characteristic of debiasing methods. However, n-CDA achieves decreased $F_1$ scores on OntoNotes for all the four models, indicating that it is more suitable for tasks involving mostly personal names.

## 6 Qualitative Analysis

In Table 7, we show some Counter-GAP examples with predictions from SpanBERT-large. In Example 1, SpanBERT-large makes correct decisions for both the original and gender-controlled instance. But for the two gender-swapped instances, it either refers the target pronoun to the incorrect feminine person (Denise), or includes both feminine names (Roxanne and Denise) in the coreferent cluster, which indicates a worse performance on resolving feminine pronouns under the same context. This kind of gender bias can only be detected when we counterfactually augment the original instance. In Example 2, SpanBERT-large correctly finds "Alice" in the gender-swapped-2 instance, but confuses "Alice" and "Dora" in the gender-swapped-1 instance. This illustrates the inconsistency brought by name perturbations within the same gender, and we take this into account by subtracting inconsis-

tency within genders from inconsistency across genders ($\Delta I = I_{across} - I_{within}$).

## 7 Related Work

**Measuring Bias in NLP models.** Human-like biases are first detected and measured in word embeddings (Bolukbasi et al., 2016; Caliskan et al., 2017; Garg et al., 2018; Gonen and Goldberg, 2019; Manzini et al., 2019). For pre-trained language models, May et al. (2019) adopt bleached sentence templates to contextualize target words, while most recent works leverage crowd-sourced benchmark datasets on NLP tasks such as language modeling (Nangia et al., 2020; Nadeem et al., 2021), sentiment analysis (Kiritchenko and Mohammad, 2018), dialog generation (Barikeri et al., 2021), natural language inference (Dev et al., 2020), and machine translation (Stanovsky et al., 2019). Our work follows GAP (Webster et al., 2018), WinoBias (Zhao et al., 2018), and WinoGender (Rudinger et al., 2018) to measure gender bias in coreference resolution, with a specific focus on collecting diverse, natural, and minimally distant instances.

**Counterfactual Bias Evaluation.** Kusner et al. (2017) propose the notion of *counterfactual fairness*, which requires similar predictions before and after counterfactual interventions in casual graphs. Garg et al. (2019) apply this notion to text classifi-

| Example | 1 |
|---|---|
| original (correct) | "In fact, Roxanne told me that she had scheduled an interview with a source tonight." Denise sipped at her lemonade through her straw until she found the bottom of her glass at last. Scotty told Chris that Denise didn't give **him** any particulars about why she needed to hire a private detective when she sought **his** advice. |
| gender-controlled (correct) | "In fact, Denise told me that she had scheduled an interview with a source tonight." Roxanne sipped at her lemonade through her straw until she found the bottom of her glass at last. Chris told Scotty that Roxanne didn't give **him** any particulars about why she needed to hire a private detective when she sought **his** advice. |
| gender-swapped-1 (incorrect) | "In fact, Chris told me that he had scheduled an interview with a source tonight." Scotty sipped at his lemonade through his straw until he found the bottom of his glass at last. Denise told Roxanne that Scotty didn't give **her** any particulars about why he needed to hire a private detective when he sought **her** advice. |
| gender-swapped-2 (incorrect) | "In fact, Scotty told me that he had scheduled an interview with a source tonight." Chris sipped at his lemonade through his straw until he found the bottom of his glass at last. Roxanne told Denise that Chris didn't give **her** any particulars about why he needed to hire a private detective when he sought **her** advice. |
| Example | 2 |
| original (correct) | Dora said, "You ought not to bet, especially on Sunday," and Alice altered it to "You may be sure." "Well, but what then?" Oswald asked Denny. "Out with it," for **he** saw that **his** youthful friend had got an idea and couldn't get it out. |
| gender-controlled (correct) | Alice said, "You ought not to bet, especially on Sunday," and Dora altered it to "You may be sure." "Well, but what then?" Denny asked Oswald. "Out with it," for **he** saw that **his** youthful friend had got an idea and couldn't get it out. |
| gender-swapped-1 (incorrect) | Oswald said, "You ought not to bet, especially on Sunday," and Denny altered it to "You may be sure." "Well, but what then?" Dora asked Alice. "Out with it," for **she** saw that her youthful friend had got an idea and couldn't get it out. |
| gender-swapped-2 (correct) | Denny said, ""You ought not to bet, especially on Sunday,"" and Oswald altered it to ""You may be sure."" ""Well, but what then?"" Alice asked Dora. ""Out with it,"" for **she** saw that **her** youthful friend had got an idea and couldn't get it out. |

Table 7: Examples from Counter-GAP. In each instance, the predicted coreference cluster from the model is highlighted in yellow, and the correctness of the prediction is annotated in the first column. The target pronoun is in **bold** and underlined; also underlined is the true coreferent name. Other notations follow those in Table 1.

cation and propose the metric of Counterfactual Token Fairness, which is similar to our inconsistency metrics, but we further distinguish inconsistency within genders from inconsistency across genders in our quadruple setting. Counterfactual Data Augmentation (CDA) (Webster et al., 2020; Zmigrod et al., 2019) is a widely adopted method for bias evaluation (Cao et al., 2020; Zhang et al., 2021), and we additionally focus on personal names during counterfactual generation of coreference resolution instances.

**Name Artifacts in NLP Models.** Since neural language models do not treat personal names as interchangeable, there are various biases in the learned representations of personal names (Shwartz et al., 2020; Prabhakaran et al., 2019; Wolfe and Caliskan, 2021; Wang et al., 2022). Counter-GAP considers name biases as the source of gender bias, and exhibits these name biases through the task of coreference resolution.

# 8 Summary and Outlook

In this work, we proposed a method to construct minimally distant bias-measuring datasets for coreference resolution, and exemplified it in the collection of Counter-GAP. We proposed the inconsistency metric $\Delta I$ to overcome the bias cancellation problem and noise from name perturbations. We showed that four pre-trained language models exhibit significant gender bias, and name-based CDA is most effective in mitigating the detected bias.

Limitations of Counter-GAP include that around half of the instances are from historical fictions in Project Gutenburg, making the dataset less representative of contemporary bias; the rules in our method for constructing Counter-GAP are specific for English, and might not be easily adapted to languages with more complex morphology; while we recognize that gender is non-binary, we adopt the simplifying setup of binary gender construct, which prevents us from detecting gender bias against minority groups with non-binary genders.

In future work, we will apply our method to different domains and more contemporary corpora such as news articles. Leveraging the data augmentation method for languages with grammatical genders (Zmigrod et al., 2019), as well as linguistic resources for non-binary genders (Cao et al., 2020) is also an important future direction to construct more gender- and language-inclusive datasets.

## Limitations

Counter-GAP adopts the setup of a binary gender construct, which restricts it from detecting bias against non-binary gender groups. Future work may extend Counter-GAP using non-binary gendered word lists, and correspondingly extend our metric (inconsistency across and within binary gender groups) for multiple gender groups.

Our method relies on specific characteristics of the English language. Directly applying it to other languages may be non-trivial. For example, languages like French or Italian adopt grammatical genders that need extra rules in our counterfactual generation method, while Chinese names are, in principle, gender neutral, which makes it impossible to identify genders from personal names. Therefore, adaptation efforts are required for researchers working on multilingual problems.

Like many other bias-measuring datasets, Counter-GAP only serves as a diagnostic dataset. This means that, if our dataset and metric detect significant bias, we could deem a model biased; but if little or no bias is detected, we cannot guarantee that the model is unbiased. Practitioners may adopt diverse bias benchmarks before reaching a conclusion.

## Acknowledgements

## References

Soumya Barikeri, Anne Lauscher, Ivan Vulić, and Goran Glavaš. 2021. RedditBias: A real-world resource for bias evaluation and debiasing of conversational language models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1941–1955, Online. Association for Computational Linguistics.

Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, page 610–623, New York, NY, USA. Association for Computing Machinery.

Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. Language (technology) is power: A critical survey of "bias" in NLP. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5454–5476, Online. Association for Computational Linguistics.

Tolga Bolukbasi, Kai-Wei Chang, James Y. Zou, Venkatesh Saligrama, and Adam T. Kalai. 2016. Man is to computer programmer as woman is to homemaker? Debiasing word embeddings. In *Proceedings of Neural Information Processing Systems (NeurIPS)*.

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

Qingqing Cao, Aruna Balasubramanian, and Niranjan Balasubramanian. 2020. Towards accurate and reliable energy measurement of NLP models. In *Proceedings of SustaiNLP: Workshop on Simple and Efficient Natural Language Processing*, pages 141–148, Online. Association for Computational Linguistics.

Sunipa Dev, Tao Li, Jeff Phillips, and Vivek Srikumar. 2020. On measuring and mitigating biased inferences of word embeddings. In *Proceedings of AAAI*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Nikhil Garg, Londa Schiebinger, Dan Jurafsky, and James Zou. 2018. Word embeddings quantify 100 years of gender and ethnic stereotypes. *Proceedings of the National Academy of Sciences*, 115(16):E3635–E3644.

Sahaj Garg, Vincent Perot, Nicole Limtiaco, Ankur Taly, Ed H. Chi, and Alex Beutel. 2019. Counterfactual fairness in text classification through robustness. In

*Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society, AIES '19.*

Hila Gonen and Yoav Goldberg. 2019. Lipstick on a pig: Debiasing methods cover up systematic gender biases in word embeddings but do not remove them. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 609–614, Minneapolis, Minnesota. Association for Computational Linguistics.

Yvette Graham, Nitika Mathur, and Timothy Baldwin. 2014. Randomized significance tests in machine translation. In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 266–274, Baltimore, Maryland, USA. Association for Computational Linguistics.

Rowan Hall Maudslay, Hila Gonen, Ryan Cotterell, and Simone Teufel. 2019. It's all in the name: Mitigating gender bias with name-based counterfactual data substitution. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5267–5275, Hong Kong, China. Association for Computational Linguistics.

Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, and Omer Levy. 2020. Span-BERT: Improving pre-training by representing and predicting spans. *Transactions of the Association for Computational Linguistics*, 8:64–77.

Mandar Joshi, Omer Levy, Luke Zettlemoyer, and Daniel Weld. 2019. BERT for coreference resolution: Baselines and analysis. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5803–5808, Hong Kong, China. Association for Computational Linguistics.

Svetlana Kiritchenko and Saif Mohammad. 2018. Examining gender and race bias in two hundred sentiment analysis systems. In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, pages 43–53, New Orleans, Louisiana. Association for Computational Linguistics.

Vid Kocijan, Oana-Maria Camburu, and Thomas Lukasiewicz. 2021. The gap on GAP: Tackling the problem of differing data distributions in bias-measuring datasets. In *Proceedings of AAAI*.

Harold T. Kuhn. 1955. The Hungarian method for the assignment problem. *Naval Research Logistics Quarterly*, 2(1-2):83–97.

Matt Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. 2017. Counterfactual fairness. In *Proceedings of Neural Information Processing Systems (NeurIPS)*.

Kenton Lee, Luheng He, and Luke Zettlemoyer. 2018. Higher-order coreference resolution with coarse-to-fine inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 687–692, New Orleans, Louisiana. Association for Computational Linguistics.

Thomas Manzini, Lim Yao Chong, Alan W Black, and Yulia Tsvetkov. 2019. Black is to criminal as caucasian is to police: Detecting and removing multiclass bias in word embeddings. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 615–621, Minneapolis, Minnesota. Association for Computational Linguistics.

Chandler May, Alex Wang, Shikha Bordia, Samuel R. Bowman, and Rachel Rudinger. 2019. On measuring social biases in sentence encoders. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 622–628, Minneapolis, Minnesota. Association for Computational Linguistics.

Nicholas Meade, Elinor Poole-Dayan, and Siva Reddy. 2022. An empirical survey of the effectiveness of debiasing techniques for pre-trained language models. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1878–1898, Dublin, Ireland. Association for Computational Linguistics.

Moin Nadeem, Anna Bethke, and Siva Reddy. 2021. StereoSet: Measuring stereotypical bias in pretrained language models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5356–5371, Online. Association for Computational Linguistics.

Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R. Bowman. 2020. CrowS-pairs: A challenge dataset for measuring social biases in masked language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1953–1967, Online. Association for Computational Linguistics.

Vinodkumar Prabhakaran, Ben Hutchinson, and Margaret Mitchell. 2019. Perturbation sensitivity analysis to detect unintended model biases. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5740–5745, Hong Kong, China. Association for Computational Linguistics.

Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Olga Uryupina, and Yuchen Zhang. 2012. CoNLL-2012 shared task: Modeling multilingual unrestricted

coreference in OntoNotes. In *Joint Conference on EMNLP and CoNLL - Shared Task*, pages 1–40, Jeju Island, Korea. Association for Computational Linguistics.

Rachel Rudinger, Jason Naradowsky, Brian Leonard, and Benjamin Van Durme. 2018. Gender bias in coreference resolution. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 8–14, New Orleans, Louisiana. Association for Computational Linguistics.

Vered Shwartz, Rachel Rudinger, and Oyvind Tafjord. 2020. "you are grounded!": Latent name artifacts in pre-trained language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6850–6861, Online. Association for Computational Linguistics.

Gabriel Stanovsky, Noah A. Smith, and Luke Zettlemoyer. 2019. Evaluating gender bias in machine translation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1679–1684, Florence, Italy. Association for Computational Linguistics.

Tony Sun, Andrew Gaut, Shirlyn Tang, Yuxin Huang, Mai ElSherief, Jieyu Zhao, Diba Mirza, Elizabeth Belding, Kai-Wei Chang, and William Yang Wang. 2019. Mitigating gender bias in natural language processing: Literature review. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1630–1640, Florence, Italy. Association for Computational Linguistics.

Jun Wang, Benjamin Rubinstein, and Trevor Cohn. 2022. Measuring and mitigating name biases in neural machine translation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2576–2590, Dublin, Ireland. Association for Computational Linguistics.

Kellie Webster, Marta Recasens, Vera Axelrod, and Jason Baldridge. 2018. Mind the GAP: A balanced corpus of gendered ambiguous pronouns. *Transactions of the Association for Computational Linguistics*, 6:605–617.

Kellie Webster, Xuezhi Wang, Ian Tenney, Alex Beutel, Emily Pitler, Ellie Pavlick, Jilin Chen, Ed Chi, and Slav Petrov. 2020. Measuring and reducing gendered correlations in pre-trained models. *CoRR*, cs.CL/2010.06032v2.

Robert Wolfe and Aylin Caliskan. 2021. Low frequency names exhibit bias and overfitting in contextualizing language models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 518–532, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Chong Zhang, Jieyu Zhao, Huan Zhang, Kai-Wei Chang, and Cho-Jui Hsieh. 2021. Double perturbation: On the robustness of robustness and counterfactual bias evaluation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3899–3916, Online. Association for Computational Linguistics.

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. Gender bias in coreference resolution: Evaluation and debiasing methods. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 15–20, New Orleans, Louisiana. Association for Computational Linguistics.

Yukun Zhu, Ryan Kiros, Richard S. Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *International Conference on Computer Vision (ICCV)*.

Ran Zmigrod, Sabrina J. Mielke, Hanna Wallach, and Ryan Cotterell. 2019. Counterfactual data augmentation for mitigating gender stereotypes in languages with rich morphology. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1651–1661, Florence, Italy. Association for Computational Linguistics.

## A Training Details

Training details follow those by Joshi et al. (2019, 2020), and the hyperparameters adopted for each model during fine-tuning on OntoNotes are shown in Table 8. Specifically, each document in OntoNotes is divided into non-overlapping segments of length `max_segment_len`. The segments are then encoded independently by BERT/SpanBERT to contextualized word embeddings and fed to the `c2f-coref` model (Lee et al., 2018). The models are fine-tuned for 20 epochs with a dropout rate of 0.3, `bert_learning_rate` on parameters in BERT/SpanBERT, and `task_learning_rate` on parameters in `c2f-coref`. The learning rates are linearly decayed. A batch size of one document is used, where each document is randomly truncated to contain `max_training_sentences` segments. All the experiments are conducted on one Tesla-V100 GPU with 32 GB memory.

For the debiasing method n-CDA, we adjust the bipartite graph matching method from Hall Maudslay et al. (2019) to fit the name list in our gender-guesser. Specifically, in our name list, each first name is assigned a label in {"male", "mostly male", "female", "mostly female"} indicating its gender specification, as well as a 55-dimensional frequency vector. The value in each dimension is an integer in $[0, 13]$ that indicates the name's relative frequency in one of the 55 countries. Below, we only describe how we match "male" with "female" names; the method to match "mostly male" with "mostly female" names is the same. We build a bipartite graph where "male" and "female" names are nodes in distinct parts, and define the weight of an edge between a "male" and a "female" name as $w_{i,j} = \|v_i - v_j\|_2 \cdot (\alpha - cos\langle v_i, v_j \rangle)$, where $v_i$ and $v_j$ are the frequency vectors of the two names, and $\alpha > 1$ is a hyperparameter balancing the $\ell_2$ and the cosine distance. Our motivation is to encourage a rare name to be matched with even a popular name in the same country other than another rare name in a different country, so we choose $\alpha = 12/11$. Finally, we leverage a minimum weight full matching algorithm (Kuhn, 1955) to compute the matches between the names in the two parts.

## B Amazon Mechanical Turk Annotation Details

Our annotation instructions and a sample Human Intelligence Task (HIT) interface are shown in Figure 2. To ensure the annotation quality, we implement a series of on-submission checks including checks on whether the selected span is a personal name, whether multiple entities are selected, whether the "no names are coreferent" box is misused, and so on. We require annotators to have at least a 95% approval rate with more than 50 approved HITs. The average time an annotator spent on one HIT is around 30 minutes.

## C Illustration of the Inconsistency Metrics

A conceptual illustration of the proposed inconsistency metrics is shown in Figure 1: inconsistency across genders ($I_{across}$) measures inconsistency in instance pairs containing two instances of different genders, while inconsistency within genders ($I_{within}$) measures inconsistency in instance pairs containing two instances of the same gender.
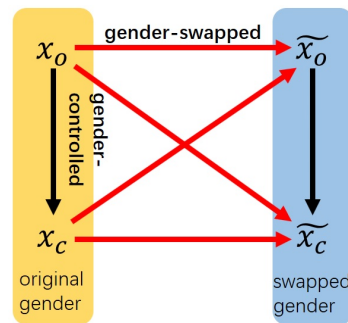


Figure 1: Illustration of the two inconsistency metrics. When computing inconsistency within genders, we use instance pairs linked by the two black arrows; when computing inconsistency across genders, we use instance pairs linked by the four red arrows.

| | BERT-base | BERT-large | SpanBERT-base | SpanBERT-large |
|---|---|---|---|---|
| $bert\_learning\_rate$ | 1e-5 | 1e-5 | 2e-5 | 1e-5 |
| $task\_learning\_rate$ | 2e-4 | 2e-4 | 1e-4 | 3e-4 |
| $max\_segment\_len$ | 128 | 384 | 384 | 512 |
| $max\_training\_sentences$ | 11 | 3 | 3 | 3 |

Table 8: Hyperparameters for fine-tuning.

**Instructions for the Pronoun Resolution task**

Please answer 10 **independent** Questions.

For each Question:

1. You are given a piece of text and a pronoun. The given pronoun is already highlighted in **yellow**.
2. You will have to **find whom the given pronoun refers to** by highlighting the names (including titles) of the correct person.
3. To **highlight** a name, select it in the text and press the green "Highlight" button.
4. If you need to **clear** your highlighted names, press the "Clear" button.
5. If the correct person's name is not mentioned in the given text, or the correct person doesn't have a name, **check** the box before "No names are coreferent with the given pronoun".
6. If the given text contains anything that **contradicts commonsense knowledge** (e.g. a man gives birth), please specify in the following textbox.

Example-1

> Then he went away, regained his automobile and drove straight to the Alexandria Hotel. Mr. Cumberford had insisted on the Kanes taking rooms at the hotel during the meet, and all three were now established there, Mrs. Kane having decided to go each day to Dominguez, where Stephen and Sybil could tell **her** of the events as they occurred.

Example-2

> "Oh," said Julia, arching her eyebrows, "I thought you were both in papa's counting-room." "We shall know each other better by and by," said Herbert, smiling. Tom did not appear to hear this, but tried to keep up the conversation with Julia, desiring to have it appear that they were intimate friends; but the young lady gave brief replies, and finally, turning away, devoted herself once more to Herbert, much to Tom's disgust. In fact, what **he** saw made Tom pass a very unpleasant evening, and when, on their return home, Maria suggested that Julia had taken a fancy to Herbert, he told her to mind her own business, which Maria justly considered a piece of rudeness wholly uncalled for.

Example-3

> "You might have been quite sure that my father's house would have been open to Polly," said Jack quite warmly, "or Mr Wilkins's, for the matter of that." "I know it lad, I know it" returned the captain, slapping **his** friend on the shoulder, "but after all, this Aunt Maria--this lady-like individual--is the most natural protector."
> **You should check the box before "No names are coreferent with the given pronoun"** (because the given pronoun "his" refers to "the captain", whose name doesn't appear in the text)

Quality checks and known answers are placed throughout the questionnaire. You will not be able to submit the HIT if those checks are not passed!

**Question-1**

> Luis's rapid breathing began to slow in tempo, his rage slowly becoming calmer. He knew that Luna was right; he would go after her and try to stop her from doing anything that would risk her life. Since day one, Luna had been his best friend; someone that he could always find comfort in. Chelsea and Phil followed them down into the ditch, checking to see if Luis had gone through any more damage. Miraculously, **he** didn't seem to have damaged his wound any further than it already was.

[Highlight] [Clear] ☐ **No names are coreferent with the given pronoun**

If the text contains anything that **contradicts commonsense knowledge**, please specify below:

[                                                                                    ]

**Question-2**

> "Only the judge knows, but it's probable unless a suitable home is found. The county is usually given

Figure 2: A sample HIT interface and annotation instructions.