

# Harms of Gender Exclusivity and Challenges in Non-Binary Representation in Language Technologies

**Sunipa Dev**  
she/her  
UCLA

**Masoud Monajatipoor\***  
he/him  
UCLA

**Anaelia Ovalle\***  
they/he/she  
UCLA

**Arjun Subramonian\***  
they/them  
UCLA, Queer in AI

**Jeff M Phillips**  
he/him  
University of Utah

**Kai-Wei Chang**  
he/him  
UCLA

## Abstract

*Content Warning: This paper contains examples of stereotypes and associations, misgendering, erasure, and other harms that could be offensive and triggering to trans and non-binary individuals.*

Gender is widely discussed in the context of language tasks and when examining the stereotypes propagated by language models. However, current discussions primarily treat gender as binary, which can perpetuate harms such as the cyclical erasure of non-binary gender identities. These harms are driven by model and dataset biases, which are consequences of the non-recognition and lack of understanding of non-binary genders in society. In this paper, we explain the complexity of gender and language around it, and survey non-binary persons to understand harms associated with the treatment of gender as binary in English language technologies. We also detail how current language representations (e.g., GloVe, BERT) capture and perpetuate these harms and related challenges that need to be acknowledged and addressed for representations to equitably encode gender information.

## 1 Introduction

As language models are more prolifically used in language processing applications, ensuring a higher degree of fairness in associations made by their learned representations and intervening in any biased decisions they make has become increasingly important. Recent work analyzes, quantifies, and mitigates language model biases such as gender, race or religion-related stereotypes in static word embeddings (GloVe (Pennington et al., 2014)) and contextual (e.g., BERT (Devlin et al., 2019)) representations (Bolukbasi et al., 2016; DeArteaga et al., 2019; Ravfogel et al., 2020; Dev et al., 2020b).

\* Equal contribution  
{sunipa,anaelia,arjunsub,kwchang}@cs.ucla.edu,  
monajati@g.ucla.edu, jeffp@cs.utah.edu

A bulk of social bias studies on language models have focused on binary gender and the stereotypes associated with masculine and feminine attributes (Bolukbasi et al., 2016; Webster et al., 2018; Dev et al., 2020b). Additionally, models often rely on gendered information for decision making, such as in named entity recognition, coreference resolution, and machine translation (Mehrabi et al., 2020; Zhao et al., 2018; Stanovsky et al., 2019), but the purview of gender in these tasks and associated measures of performance focus on binary gender. While discussing binary gender bias and improving model performance are important, it is important to reshape our understanding of gender in language technologies in a more accurate, inclusive, non-binary manner.

Current language models can perpetrate harms such as the cyclical erasure of non-binary gender identities (Uppunda et al., 2021; Sap, 2021; Lakoff; Fiske, 1993; Fast et al., 2016; Behm-Morawitz and Mastro, 2008). These harms are driven by model and dataset biases due to tainted examples, limited features, and sample size disparities (Wang et al., 2019; Barocas et al., 2019; Tan and Celis, 2019), which are consequences of the non-recognition and a lack of understanding of non-binary genders in society (MAP, 2016; Rajunov and Duane, 2019).

Some recent works attempt to mitigate these harms by building task-specific datasets that are not restricted to binary gender and building metrics that on extension, could potentially measure biases against all genders (Cao and Daumé III, 2020; Rudinger et al., 2018). While such works that intentionally inject real-world or artificially-created data of non-binary people into binary-gendered datasets are well-intentioned, they could benefit from a broader perspective of harms as perceived by non-binary persons to avoid mischaracterizing non-binary genders as a single gender (Sun et al., 2021) or perpetuating biases through non-intersectional training examples, i.e. examples that

do not capture the interconnected nature of social identities (Crenshaw, 1989).

In this paper, we conduct a detailed investigation into the representational and allocational harms (Barocas et al., 2017; Blodgett et al., 2020) related to the treatment of gender with binary predilections in English language technologies. We do so by explaining the complexity of gender and language around it, and surveying non-binary persons with some familiarity with AI on potential harms in common NLP tasks. While the challenges associated with limited or tainted data are loosely hypothesized, they are not well understood.

We study the extent of these data challenges and detail how they manifest in the resultant language representations and downstream tasks. We examine both static embeddings (GloVe) and contextual representations (BERT) with respect to the quality of representations (Section 4.2) of non-binary-associated words and pronouns. We highlight how the disparity in representations cyclically propagates the biases of underrepresentation and misrepresentation and can lead to the active misgendering and erasure of non-binary persons in language technologies.

## 2 Gender, Language, and Bias

We first discuss the complex concepts of gender and bias and their expression in English language.

### 2.1 Gender

In this paper, *gender* refers to *gender identity*, as opposed to *gender expression* or *sex*. *Gender identity* concerns how individuals experience their own gender. In contrast, *gender expression* concerns how one expresses themselves, through their “hair length, clothing, mannerisms, makeup” and *sex* relates to one’s “genitals, reproductive organs, chromosomes, hormones, and secondary sex characteristics” (Rajunov and Duane, 2019). Gender identity, gender expression, and sex do not always “align” in accordance with Western cisnormativity (Rajunov and Duane, 2019). However, people are conditioned to erroneously believe otherwise, which leads to “societal expectations and stereotypes around gender roles” and the compulsive (*mis*)gendering of others (Cao and Daumé III, 2020; Serano, 2007).

**Gender in Western Society** In Western society, *discourse* around one’s gender identity can, but does not always, comprise two intersecting aspects:

(i) a description of how it is similar to or different from the binary genders, i.e. male and female. For instance, genderfluid persons do not identify with a single gender, and agender individuals do not subscribe to gender at all (Rajunov and Duane, 2019). It is important to note that gender may fluctuate over an individual’s lifetime, and it is extremely problematic to assume a biologically essentialist view of it (Weber, 2019), and (ii) whether it is the same as or differs from the individual’s gender assigned at birth, i.e. cisgender or transgender, respectively. Many individuals who are not cis, including non-binary people, identify as trans.

Non-binary genders encompass *all* the genders that do not conform to the Western gender binary (Rajunov and Duane, 2019). There are many non-Western non-cis identities, like the Jogappas of Karnataka, Muxes of Oaxaca, and Mahuwahines of Hawai’i (Desai, 2018; Mirandé, 2016; Clarke, 2019). However, non-Western non-cis identities cannot be accurately described by the Western-centric, English-based gender framework afore established (Mirandé, 2016; Thorne et al., 2019). Hence, as this paper focuses on the English language, its treatment of non-binary genders does *not* adequately include non-Western non-cis identities.

**Pronouns and Gendered Names** In societies where language has referential gender, i.e., when an entity is referred to, and “their gender (or sex) is realized linguistically” (Cao and Daumé III, 2020), it is difficult to escape gendering others. In English, pronouns are gendered; hence, pronouns can be central to English speakers’ gender identity. However, pronouns cannot be bijectively mapped to gender. For example, not all non-binary persons use *they/them/theirs* pronouns, nor do all persons who use *they/them/theirs* pronouns identify as non-binary (Clarke, 2019). Furthermore, the use of binary pronouns, *he* and *she*, is not exclusive to cis individuals; trans and non-binary individuals also use them. English pronouns are always evolving (McCulloch and Gawne, 2016). Singular *they* has become widely adopted by trans and non-binary persons (McCulloch and Gawne, 2016; Feraday, 2016; Clarke, 2019). Neopronouns like *xe/xem/xyr* and *ze/hir/hirs* are also in use by non-cis individuals (Feraday, 2016).

Not everyone who speaks English chooses to use pronouns, and some individuals use multiple sets of pronouns (e.g. *she/her/hers* and *they/them/theirs*) (Feraday, 2016). Many non-

binary people use different pronouns depending on the space in which they are, especially if they are not publicly out; for example, a non-binary person may accept *she/her* pronouns at work but use *they/them* pronouns outside of work. Additionally, non-binary people can find multiple sets of pronouns affirming; for instance, non-binary men may use a combination of *they/them/theirs* and *he/him/his*. Furthermore, genderfluid individuals can use different sets of pronouns based on their “genderfeels” at a certain time (Gautam, 2021). This may also lead individuals to be open to being referenced by “all pronouns” or “any pronouns.” Ultimately, individuals use the pronouns that allow them to feel gender euphoria in a given space, at a given time (Gautam, 2021).

In languages without referential gender or where pronouns are seldom used (e.g. Estonian), pronouns can be less central to one’s gender identity (Crouch, 2018).

Another form of referential gender is gendered names, which are assumed for binary gender, even in language technologies, which itself can be inaccurate and problematic. Additionally, trans and non-binary persons may choose a new name that matches their gender identity to replace their *dead-name*, i.e. name assigned at birth (Rose, 2020). Many Western non-binary chosen names are creative and diverse, overlapping with common nouns or nature words, having uncommon orthographic forms, and/or consisting of a single letter (Rose, 2020).

**Lexical Gender** Lexical gender in English language is gender (or sex) conveyed in a non-referential manner (Cao and Daumé III, 2020). Examples include “mother” and “Mr.” Non-binary persons have adopted honorifics like “Mx.” to eliminate gendering (Clarke, 2019), and often use gender-neutral terms like “partner” to refer to their significant other. However, their adoption into written text and narratives is recent and sparse.

**Implications in Language Technologies** Given the complex and evolving nature of gender and the language around it, for language technologies to truly equitably encode gender, they would need to capture the full diversity and flexibility therein.

## 2.2 Biases

There has been an increase in awareness of the social biases that language models carry. In this paper, we use the term *bias* to refer to a skewed and

undesirable association in language representations which has the potential to cause representational or allocational harms (Barocas et al., 2017). There have been multiple attempts to understand social biases in language processing (Sheng et al., 2021; Caliskan et al., 2017), quantify them (Rudinger et al., 2018; Webster et al., 2018; De-Arteaga et al., 2019), and mitigate them (Zhao et al., 2019; Ravfogel et al., 2020; Sun et al., 2019). A primary focus has been on gender bias, but the narrative has been dominated by biases associated with binary gender, primarily related to occupations and adjectives. However, the biases faced by non-binary persons can be distinct from this. Non-binary genders are severely underrepresented in textual data, which causes language models to learn meaningless, unstable representations for non-binary-associated pronouns and terms. Furthermore, there are derogatory adjectives associated with non-binary-related terms (as seen in Appendix B.1). Thus, analyzing and quantifying biases associated with non-binary genders cannot be treated merely as a corollary of those associated with binary gender.

## 3 Harms

Utilizing and perpetuating the binary construction of gender in English in language technologies can have adverse impacts. We focus on specific tasks within language processing and associated applications in human-centered domains where harms can be perpetrated, motivated by their frequent mention in a survey we conduct (Section 3.1). The primary harms we discuss are misgendering and erasure.

**Misgendering:** Misgendering is the act of accidentally or intentionally addressing someone (oneself or others) using a gendered term that does not match their gender identity. Misgendering persons and the associated harms have been studied in contexts of computer vision (Keyes, 2018) and human-computer interaction (Keyes et al., 2021), which highlight its adverse impact on the mental health of non-binary individuals. Language applications and their creators can also perpetrate misgendering. For instance, language applications that operationally ask non-binary users to choose between *male* and *female* as input force non-binary users to misgender themselves (Keyes, 2018; Spiel et al., 2019). Furthermore, language models which do not explicitly collect gender information are capable of both accidental and intentional misgendering. Specifically, language models accidentally misgen-

der non-binary persons when there is insufficient information to disambiguate the gender of an individual, and so they default to binary pronouns and binary-gendered terms, potentially based on stereotypes. However, as shown in Section 4.2, language models can also misgender non-binary individuals even when their pronouns are provided.

**Erasure:** In one sense, erasure is the accidental or intentional invalidation or obscuring of non-binary gender identities. For example, the language technology Genderify, which purportedly “identif[ied] someone’s [binary] gender based on their name, email address or username” erased non-binary people by reductively distributing individuals into binary “gender bins” by their name, based on the assumption that they were cisgender (Lauer, 2020; Spiel et al., 2019; Serano, 2007). Another sense of erasure is in how stereotypes about non-binary communities are portrayed and propagated (see Appendix Table 12). Since non-binary individuals are often “denied access to media and economic and political power,” individuals in power can paint negative narratives of non-binary persons or erase the diversity in gender communities (Serano, 2007; Rajunov and Duane, 2019).

Language applications are capable of automating erasure, in a cyclical fashion (Hashimoto et al., 2018; Sap, 2021). We posit the cycle of non-binary erasure in text, in which: (i) language applications, trained on large, binary-gendered corpora, reflect the misgendering and erasure of non-binary communities in real life (Lakoff; Fiske, 1993) (ii) this reflection is viewed as a “source of truth and scientific knowledge” (Keyes et al., 2021) (iii) consequently, authors buy into these harmful ideas and other language models encode them, leading them to stereotypically portray non-binary characters in their works or not include them at all, and (Fast et al., 2016) (iv) this further amplifies non-binary erasure, and the cycle continues.

### 3.1 Survey on Harms

To understand harms associated with skewed treatment of gender in English NLP tasks and applications, the perspective of those facing the harms is essential. We conduct a survey for the same.

**Survey Respondents** We focused this survey on non-binary persons who have familiarity with AI. We acknowledge that this indeed is a limitation, as it narrows our focus to non-binary persons of

specific socioeconomic status, ethnicity, and English fluency. However, we field this survey as the first in a series to gain foray into harms experienced by non-binary individuals who build AI and know its effects. Furthermore, it allows us to gather what tasks could potentially cause harm without asking leading questions with explicit examples of tasks that exhibit stereotypes or skews against non-binary genders. We distributed the survey through channels like social media and mailing lists at universities and organizations. We had 19 individuals respond to our survey. While existing research has surveyed non-binary individuals on the harms of gendered web forms (Scheuerman et al., 2021), there is no precedent for our survey on language technology harms, so our primary intent with this sample of respondents was to assess the efficacy of our survey design.

**Survey Structure** The survey was anonymous, with no financial compensation, and questions were kept *optional*. Further ethical considerations are presented in Section 6. In the following subsections, we briefly summarize our survey design and survey responses. We provide the full survey, our rationale for each question, and qualitative analysis of all responses received in Appendix A.

#### 3.1.1 Demographic information

We asked survey respondents for demographic information to better understand the intersections of their identities. Demographic information included gender identity, ethnicity, AI experience, etc. 84.2% of respondents use pronouns *they/them*, 26.3% use *she/her*, 15.8% use *he/him*, and 5.3% use *xe/xem*. 31.6% use multiple sets of pronouns. Additionally, an overwhelming majority (all but two) of our respondents identified as white and/or Caucasian. No respondents were Black, Indigenous, and/or Latinx, and two respondents were people of color. Furthermore, 52.6% of respondents are originally from the US, 63.2% current live in the US, and the majority of others are originally from or currently live in Canada and countries in Western Europe. This limits the conclusions we can reach from this sample’s responses. All respondents were familiar with AI, through their occupation, coursework, books, and social media (more details in Appendix A.1).

#### 3.1.2 Harms in Language Tasks

This segment first defined representational and allocational harms (Barocas et al., 2017) and intro-

duced three common NLP tasks (Named Entity Recognition (NER), Coreference Resolution, and Machine Translation) using publicly-available AllenNLP demos (Gardner et al., 2018), which survey respondents engaged with to experiment with potential harms. The demos were accompanied by *non-leading* questions about representational and allocational harms, if any, that non-binary communities could face as a result of these tasks. The questions were intentionally phrased to ask about the harms that could occur rather than imply likely harms (see Appendix A). We summarize the responses to these questions in Table 1, where we see that misgendering of persons is a common concern across all three tasks. We found that, for all tasks, above 84% of respondents could see/think of undesirable outcomes for non-binary genders. Furthermore, the severity of harms, as perceived by subjects of the survey, is the highest in machine translation, which is also a task more commonly used by the population at large. We provide descriptions of the tasks and in-depth analyses of all the responses in Appendix A.2.

### 3.1.3 Broader Concerns with Language Technologies

This segment was purposely kept less specific to understand the harms in different domains (healthcare, social media, etc.) as perceived by different non-binary individuals. We first list some domains to which language models can be applied along with summarized explanations of respondents regarding undesirable outcomes (see Appendix A.3 for in-depth analyses).

- *Social Media*: LGBTQ+ social media content is automatically flagged at higher rates. Ironically, language models can fail to identify hateful language targeted at non-binary people. Further, if social media sites attempt to infer gender from name or other characteristics, this can lead to incorrect pronouns for non-binary individuals. Additionally, “language models applied in a way that links entities across contexts are likely to out and/or deadname people, which could potentially harm trans and non-binary people”. Moreover, social media identity verification could incorrectly interpret non-binary identities as fake or non-human.
- *Healthcare*: Respondents said that “healthcare requires engaging with gender history as well as identity”, which current language models are not capable of doing. Additionally, language models could “deny insurance claims, e.g. based on a ‘mis-

match’ between diagnosis and gender/pronouns”.

- *Education*: Language models in automated educational/grading tools could “automatically mark things wrong/‘ungrammatical’ for use of non-standard language, singular *they*, neopronouns, and other new un- or creatively gendered words”.

Additionally, respondents discussed some language applications that could exacerbate misgendering, non-binary erasure, transphobia, and the denial of cisgender privilege. Some examples were how automated summarization could fail to recognize non-binary individuals as people, language generation cannot generate text with non-binary people or language, speech-to-text services cannot handle neopronouns, machine translation cannot adapt to rapidly-evolving non-binary language, and automated gender recognition systems only work for cis people (Appendix A.3).

The barriers (Barocas et al., 2019) to better including non-binary persons in language models, as explained in the responses, are as follows (definitions and in-depth analyses in Appendix A.3).

- *Tainted Examples*: Since the majority of training data are scraped from sources like the Internet, which represent “hegemonic viewpoints”, they contain few mentions of non-binary people; further, the text is often negative, and positive gender non-conforming content is not often published.
- *Limited Features*: Data annotators may not recognize or pay attention to non-binary identities and may lack situational context.
- *Sample Size Disparities*: Non-binary data may be “discarded as ‘outliers’” and “not sampled in training data”, non-binary identities may not be possible labels, developer/research teams tend to “want to simplify variables and systems” and may not consider non-binary persons prevalent enough to change their systems for.

## 3.2 Limitations and Future Directions

We found that our survey, without any leading questions, was effective at getting respondents to recount language technology harms they had experienced on account of their gender, and brainstorm harms that could affect non-binary communities. However, our survey reaches specific demographics of ethnicity, educational background, etc. The responses equip us to better reach out to diverse groups of persons, including those without familiarity with AI and/or not fluent in English. Some respondents also indicated that language models could be used violently or to enable existing dis-

|                                       | Named Entity Recognition (NER)  | Coreference Resolution  | Machine Translation   |
|---------------------------------------|---|---|---|
| <b>Example representational harms</b> | <ul style="list-style-type: none"> <li>systematically mistags neopronouns and singular <i>they</i> as non-person entities</li> <li>unable to tag non-binary chosen names as <i>Person</i>, e.g. the name “A Boyd” is not recognized as referring to a <i>Person</i></li> <li>tags non-binary persons as <i>Person – male</i> or <i>Person – female</i></li> </ul>   | <ul style="list-style-type: none"> <li>may incorrectly links <i>s/he</i> pronouns with non-binary persons who do not use binary pronouns</li> <li>does not recognize neopronouns</li> <li>cannot link singular <i>they</i> with individual persons, e.g. In “Alice Smith plays for the soccer team. They scored the most goals of any player last season.”, <i>they</i> is linked with <i>team</i> instead of with <i>Alice</i></li> </ul>  | <ul style="list-style-type: none"> <li>translates from a language where pronouns are unmarked for gender and picks a gender grounded in stereotypes associated with the rest of the sentence, e.g. translates “(3SG) is a nurse” (in some language) to “She is a nurse” in English</li> <li>translates accepted non-binary terms in one language to offensive terms in another language, e.g. <i>kathoey</i>, which is an accepted way to refer to trans persons in Thailand, translates to <i>ladyboy</i> in English, which is derogatory</li> </ul> |
| <b>Example allocational harms</b>     | <ul style="list-style-type: none"> <li>NER-based resume scanning systems throw out resumes from non-binary persons for not having a recognizable name</li> <li>non-binary persons are unable to access medical and government services if NER is used as a gatekeeping mechanism on websites</li> <li>non-binary people with diverse and creative names are erased if NER is employed to build a database of famous people</li> </ul> | <ul style="list-style-type: none"> <li>a coref-based ranking system undercounts a non-binary person’s citations (including pronouns) in a body of text if the person uses <i>xe/xem</i> pronouns</li> <li>a coref-based automated lease signing system populates referents with <i>s/he</i> pronouns for an individual who uses <i>they/them</i> pronouns, forcing self-misgendering</li> <li>a coref-based law corpora miner undercounts instances of discrimination against non-binary persons, which delays more stringent anti-discrimination policies</li> </ul> | <ul style="list-style-type: none"> <li>machine-translated medical and legal documents applies incorrectly-gendered terms, leading to incorrect care and invalidation, e.g. a non-binary AFAB person is not asked about their pregnancy status when being prescribed new medication if a translation system applies masculine terms to them</li> <li>machine-translated evidence causes non-binary persons to be denied a visa or incorrectly convicted of a crime</li> </ul>  |

Table 1: Summary of survey responses regarding harms in NLP tasks.

crimutory policies, which should be explored in future related work. Ultimately, we hope our survey design serves as a model for researching the harms technologies pose to marginalized communities.

## 4 Data and Technical Challenges

As a consequence of historical discrimination and erasure in society, narratives of non-binary persons are either largely missing from recorded text or have negative connotations. Language technologies also reflect and exacerbate these biases and harms, as discussed in Section 3.1, due to tainted examples, limited features, and sample size disparities. These challenges are not well understood. We discuss the different fundamental problems that need to be acknowledged and addressed to strategize and mitigate the cyclical erasure and misgendering of persons as a first step towards building language models that are more inclusive.

### 4.1 Dataset Skews

The large text dumps often used to build language representations have severe skews with respect to gender and gender-related concepts. Just observing pronoun usage, English Wikipedia text (March 2021 dump), which comprises 4.5 billion tokens, has over 15 million mentions of the word *he*, 4.8 million of *she*, 4.9 million of *they*, 4.5 thousand

of *xe*, 7.4 thousand of *ze*, and 2.9 thousand of *ey*. Furthermore, the usages of non-binary pronouns were mostly not meaningful with respect to gender (Appendix B). *Xe*, as we found by annotation and its representation, is primarily used as the organization *Xe* rather than the pronoun *xe*. *Ze* was primarily used as the Polish word *that*, as indicated by its proximity to mostly Polish words like *nie*, i.e. *no*, in the GloVe representations of the words, and was also used for characterizing syllables. Additionally, even though the word *they* occurs comparably in number to the word *she*, a large fraction of the occurrences of *they* is as the plural pronoun, rather than the singular, non-binary pronoun *they*. Some corpora do exist such as the Non-Binary Wiki which contain instances of meaningfully used non-binary pronouns. However, with manual evaluation, we see that they have two drawbacks: (i) the narratives are mostly short biographies and lack the diversity of sentence structures as seen in the rest of Wikipedia, and (ii) they have the propensity to be dominated by Western cultures, resulting in further sparsification of diverse narratives of non-binary persons.

Neopronouns and gendered pronouns not “he” or “she”  
[https://nonbinary.wiki/wiki/Main\\_Page](https://nonbinary.wiki/wiki/Main_Page)

| Pronoun | Top 5 Neighbors                           |
|---------|---|
| He      | 'his', 'man', 'himself', 'went', 'him'    |
| She     | 'her', 'woman', 'herself', 'hers', 'life' |
| They    | 'their', 'them', 'but', 'while', 'being'  |
| Xe      | 'xa', 'gtx', 'xf', 'tl', 'py'             |
| Ze      | 'ya', 'gan', 'zo', 'l'ovic', 'kan'        |

Table 2: Nearest neighbor words in GloVe for binary and non-binary pronouns.

## 4.2 Text Representation Skews

Text representations have been known to learn and exacerbate skewed associations and social biases from underlying data (Zhao et al., 2017; Bender et al., 2021; Dev, 2020), thus propagating representational harm. We examine representational skews with respect to pronouns and non-binary-associated words that are extremely sparsely present in text.

**Representational erasure in GloVe.** Table 2 shows the nearest neighbors of different pronouns in their GloVe representations trained on English Wikipedia data. The singular pronouns *he* and *she* have semantically meaningful neighbors as do their possessive forms (Appendix B.1). The same is not true for non-binary neopronouns *xe* and *ze* which are closest to acronyms and Polish words, respectively. These reflect the disparities in occurrences we see in Section 4.1 and show a lack of meaningful encodings of non-binary-associated words.

**Biased associations in GloVe.** Gender bias literature primarily focuses on stereotypically gendered occupations (Bolukbasi et al., 2016; De-Arteaga et al., 2019), with some exploration of associations of binary gender and adjectives (Dev and Phillips, 2019; Caliskan et al., 2017). While these associations are problematic, there are additional, significantly different biases against non-binary genders, namely misrepresentation and under-representation. Furthermore, non-binary genders suffer from a sentiment (positive versus negative) bias. Gender-occupation associations are not a dominant stereotype observed across all genders (Table 13), where non-binary words like *transman* and *nonbinary* are not dominantly associated with either stereotypically male or female occupations. In fact, most occupations exhibit no strong correlation with words and pronouns associated with non-binary genders (see Appendix B.1).

To investigate sentiment associations with binary versus non-binary associated words, we use the WEAT test (Caliskan et al., 2017) with respect to pleasant and unpleasant attributes (listed in Appendix B.2). Since neopronouns are not

| Word       | Doctor | Engineer | Nurse  | Stylist |
|------------|--------|----------|--------|---------|
| man        | 0.809  | 0.551    | 0.616  | 0.382   |
| woman      | 0.791  | 0.409    | 0.746  | 0.455   |
| transman   | -0.062 | -0.152   | -0.095 | 0.018   |
| transwoman | -0.088 | -0.271   | 0.050  | 0.062   |
| nonbinary  | 0.037  | -0.243   | 0.129  | 0.015   |

Table 3: Cosine similarity: gendered words vs common occupations.

well-embedded, we compare disparate sentiment associations between binary versus non-binary pronouns, gendered words and proxies (e.g., *male*, *female* versus *transman*, *genderqueer*, etc.). The WEAT score is 0.916, which is non-zero, i.e. ideal, significantly large (detailed analysis in Appendix B.2), and indicates disparate sentiment associations between the two groups. For *man* and *woman*, the top nearest neighbors include *good*, *great* and *good*, *loving*, respectively. However, for *transman* and *transwoman*, top words include *dishonest*, *careless* and *unkind*, *arrogant*. This further substantiates the presence of biased negative associations, as seen in the WEAT test. Furthermore, the nearest neighbors of words associated with non-binary genders are derogatory (see Appendix Table 12). In particular, *agender* and *genderfluid* have the neighbor *negrito*, meaning “little Black”, while *genderfluid* has *Fasiq*, which is an Arabic word used for someone of corrupt moral character.

**Representational erasure in BERT.** Pronouns like *he* or *she* are part of the word-piece embedding vocabulary that composes the input layer in BERT. However, similar length neo-pronouns *xe* or *ze* are deemed as out of vocabulary by BERT, indicating infrequent occurrences of each word and a relatively poor embedding.

BERT’s contextual representations should ideally be able to discern between singular mentions of *they* (denoted  $they(s)$ ) and plural mentions of *they* (denoted  $they(p)$ ), and to some extent it indeed is able to do so, but not with high accuracy. For this, we train BERT as a classifier to disambiguate between singular and plural pronouns. Given a sentence containing a masked pronoun along with two preceding sentences, it predicts whether the pronoun is singular or plural. We build two separate classifiers  $C_1$  and  $C_2$ . Both are first trained on a dataset containing sentences with *i* or *we* (singular versus plural; details on this experi-

Code and supporting datasets can be found at <https://github.com/uclanlp/harms-challenges-non-binary-representation-NLP>

ment in Appendix B.3). Next,  $C_1$  is trained on classifying  $they(s)$  vs  $they(p)$  while  $C_2$  is trained on classifying  $he$  vs  $they(p)$ . This requires balanced, labeled datasets for both classifiers. The text spans for  $they(p)$  are chosen randomly from Wikipedia containing pairs of sentences such that the word  $they$  appears in the second sentence (with no other pronoun present) and the previous sentence has a mention of two or more persons (determined by NER). This ensures that the word  $they$  in this case was used in a plural sense. Since Wikipedia does not have a large number of sentences using  $they(s)$ , for such samples, we randomly sample them from Non-Binary Wiki (Section 4.1). The sentences are manually annotated for further confirmation of correct usage of each pronoun. We follow the procedure of data collection for  $they(s)$  to create datasets for sentences using the pronoun  $he$  from Wikipedia. Therefore, while  $C_1$  sees a dataset containing samples with  $they(s)$  or  $they(p)$ ,  $C_2$  sees samples with  $he$  or  $they(p)$ . In each dataset, however, we replace the pronouns with the  $[MASK]$  token. We test  $C_1$  and  $C_2$  on their ability to correctly classify a new dataset for  $they(p)$  (collected the same way as above). If  $C_1$  and  $C_2$  learn the difference between the singular and plural representations, each should be able to classify all sentences as plural with net accuracy 1. While the accuracy of  $C_2$  is 83.3%,  $C_1$ 's accuracy is only 67.7%. This indicates  $they(s)$  is not as distinguishable from  $they(p)$  as a binary-gendered pronoun (further experiments are in Appendix B.3).

**Biased representations with BERT.** To understand biased associations in BERT, we must look at representations of words with context. For demonstrating skewed associations with occupations (as shown for GloVe), we adopt the sentence template “[pronoun] is/are a [target]”. We iterate over a commonly-used list of popular occupations (Dev et al., 2020a), broken down into stereotypically female and male (Bolukbasi et al., 2016). We get the average probability for predicting each gendered pronoun (Table 4)  $P([pronoun] | [target] = occupation)$  over each group of occupations. The results in Table 4 demonstrate that the occupation biases in language models with respect to binary genders is not meaningfully applicable for all genders.

**BERT and Misgendering.** Misgendering is a harm experienced by non-binary persons, as empha-

| Pronouns | Occupations Categories |            |           |
|----------|------------------------|------------|-----------|
|          | Male                   | Female     | All       |
| he       | 0.5781                 | 0.1788     | 0.5475    |
| she      | 0.1563                 | 0.4167     | 0.2131    |
| they     | 0.1267                 | 0.1058     | 0.1086    |
| xe       | 2.1335e-05             | 1.9086e-05 | 1.6142e-5 |
| ze       | 7.4232e-06             | 6.0601e-06 | 5.6769e-6 |

Table 4: Pronoun associations with (i) stereotypically male, (ii) stereotypically female, and (iii) extensive list of 180 popular occupations. Values are aggregated probabilities (higher value implies more associated; see main text for more details).

sized in the survey (see Section 3.1). Further, misgendering in language technologies can reinforce erasure and the diminishing of narratives of non-binary persons. We propose an evaluation framework here that demonstrates how BERT propagates this harm. We set up sentence templates as such:

[Alex] [went to] the [hospital] for [PP] [appointment]. [MASK] was [feeling sick].

Every word within [] is varied. The words in bold are varied to get a standard set of templates (Appendix B.3). These include the verb, the subject, object and purpose. We iterate over 919 names available from SSN data which were unisex or least statistically associated with either males or females (Flowers, 2015). We choose this list to minimize binary gender correlations with names in our test. Next, we vary the underlined words in pairs. The first of each pair is a possessive pronoun (PP) which we provide explicitly (thus indicating correct future pronoun usage) and use BERT to predict the masked pronoun in the second sentence in each template. The ability to do so for the following five pairs is compared: (i) *his, he* (ii) *her, she* (iii) *their, they* (iv) *xir, xe* and (v) *zir, ze* in Table 5, where *Accuracy* is the fraction of times the correct pronoun was predicted with highest probability and the score *Probability* is the average probability associated with the correct predictions. The scores are high for predicting *he* and *she*, but drop for *they*. For *xe* and *ze* the amount by which the accuracy drops is even larger, but we can attribute this to the fact that these neopronouns are considered out of vocabulary by BERT. This demonstrates how models like BERT can explicitly misgender non-binary persons even when context is provided for correct pronoun usage.

## 5 Discussion and Conclusion

This work documents and demonstrates specific challenges towards making current language modeling techniques inclusive of all genders and re-



| Pronouns pairs | Accuracy | Probability |
|----------------|----------|-------------|
| his-he         | 0.861    | 0.670       |
| her-she        | 0.785    | 0.600       |
| their-they     | 0.521    | 0.391       |
| xir-xe         | 0.0      | 1.137e-05   |
| zir-ze         | 0.0      | 1.900e-04   |

Table 5: BERT performance for gendered pronoun predictions. Accuracy is the fraction of times the correct pronoun was predicted and probability is the aggregated probability associated with correct prediction.

ducing the extent of discrimination, misgendering, and cyclical erasure it can perpetrate. In particular, our survey identifies numerous representational and allocational harms, as voiced from individuals affected, and we demonstrate and quantify several cases where the roots of these concerns arise in popular language models. Some efforts in the NLP community have worked towards countering problems in task-specific data sets with skewed gender tags, due to underrepresentation of non-binary genders. Notably, [Cao and Daumé III \(2020\)](#) and [Cao and Daumé III \(2021\)](#) introduce a gender-inclusive dataset GICoref for coreference resolution and [Sun et al. \(2021\)](#) propose rewriting text containing gendered pronouns with *they* as the substituted pronoun to obtain more gender-neutral text. The challenges still remain that (i) not all neopronouns will have sufficient data in real-world text, and (ii) considering non-binary genders as a monolithic third category (i.e. male, female, and gender-neutral) is counter-productive and perceived as harmful (Section 3.1). While these efforts are a start in moving away from binary gender, it is questionable if gender should be defined as discrete quantities in language modeling, when in reality, it is of a fluid nature. Furthermore, models currently do not account for the mutability of gender and the language around it, and even if they did, they would likely assume there exist well-defined points at which individuals and words transition, which too is detrimental (as documented in our survey, see Section 3.1). Representing gender is as complex as the concept of gender itself. Bucketing gender in immutable, discrete units and trying to represent each, would inevitably result in marginalization of sections of the population to varied extents. As our survey catalogs how pronounced the harms of being consistently misgendered and diminished are, we encourage future work to carefully examine how (*and if*) to define and model gender in language representations and tasks.

This work sets the interdisciplinary stage for rethinking and addressing challenges with inclusively

modeling gender in language representations and tasks. Any viable solution cannot simply be a quick fix or patch, but must rely on a bottom-up approach involving affected persons system-wide, such as in annotation and human-in-the-loop mechanisms. Simultaneously, research into monitoring language technologies over time to detect harms against non-binary individuals is critical. It is further paramount to transparently communicate the performance of language technologies for non-binary persons and possible harms. In the case of harm, non-binary individuals must be able to obtain valid recourse to receive a more favorable outcome, as well as have the opportunity to provide feedback on the model’s output and have a human intervene.

**Acknowledgements** We would like to thank Emily Denton and Vasundhara Gautam for their extensive feedback on our survey and our anonymous reviewers for detailed feedback on the paper. Furthermore, we would like to thank and deeply appreciate all the survey respondents for the time and effort they invested into drawing from their lived experiences to provide innumerable informative insights, without which this project would not have been possible. We would also like to thank Vasundhara Gautam, Gauri Gupta, Krithika Ramesh, Sonia Katyal, and Michael Schmitz for their feedback on drafts of this paper.

This work was supported by NSF grant #2030859 to the Computing Research Association for the CIFellows Project. Additionally, we also thank the support from Alfred P. Sloan foundation, NSF IIS-1927554, NSF CCF-1350888, CNS-1514520, CNS-1564287, IIS-1816149, and Visa Research.

## 6 Broader Impact and Ethics

Our survey was reviewed by an Institutional Review Board (IRB), and the IRB granted the survey Exempt status, requiring only signed informed consent for the entire study, as the survey posed minimal risk, was conducted online, and did not involve treating human subjects. In the survey, we ask non-leading questions to record perceptions of participants and not preemptively impose possibilities of harm. Questions were also optional to enable participants to control the amount of emotional labor they incurred while taking the survey. We were cautious to protect identities of persons who took the survey; we analyzed aggregated data and any quotes of text analyzed or mentioned cannot be

traced back to an individual.

Inclusivity and fairness are important in NLP and its wide ranging applications. Gender, when modeled in these applications has to reflect fairly the concepts of gender identity and expression. Failure to do so leads to severe harms, especially for persons not subscribing to binary gender. In this work, we attempt to broaden the awareness of gender disparities and motivate future work to discuss and further address the harms propagated by language technologies. We emphasize the importance of centering the lived experiences of marginalized communities therein.

## References

- Solon Barocas, Kate Crawford, Aaron Shapiro, and Hanna Wallach. 2017. [The problem with bias: Allocative versus representational harms in machine learning](#). In *SIGCIS Conference*.
- Solon Barocas, Moritz Hardt, and Arvind Narayanan. 2019. *Fairness and Machine Learning*. fairmlbook.org. <http://www.fairmlbook.org>.
- Elizabeth Behm-Morawitz and Dana Mastro. 2008. Mean girls? the influence of gender portrayals in teen movies on emerging adults' gender-based attitudes and beliefs. *Journalism and Mass Communication Quarterly*, 85:131 – 146.
- Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big. *Proceedings of FAccT*.
- 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 *Advances in Neural Information Processing Systems*, volume 29.
- Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. 2017. [Semantics derived automatically from language corpora contain human-like biases](#). *Science*, 356(6334):183–186.
- Yang Trista Cao and Hal Daumé III. 2020. [Toward gender-inclusive coreference resolution](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4568–4595, Online. Association for Computational Linguistics.
- Yang Trista Cao and Hal Daumé III. 2021. [An Analysis of Gender and Bias Throughout the Machine Learning Lifecycle\\*](#). *Computational Linguistics*, pages 1–46.
- Jessica Clarke. 2019. [They, them, and theirs](#). *132 Harvard Law Review*, page 894.
- Kimberle Crenshaw. 1989. [Demarginalizing the intersection of race and sex: A black feminist critique of antidiscrimination doctrine, feminist theory and antiracist policies](#). *University of Chicago Legal Forum*, 1989(1):139–167.
- Erin Crouch. 2018. [Being non-binary in a language without gendered pronouns – estonian](#). *Deep Baltic*.
- Maria De-Arteaga, Alexey Romanov, Hanna Wallach, Jennifer Chayes, Christian Borgs, Alexandra Chouldechova, Sahin Geyik, Krishnaram Kenthapadi, and Adam Tauman Kalai. 2019. Bias in bios: A case study of semantic representation bias in a high-stakes setting. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 120–128.
- Rishikesh Bahadur Desai. 2018. [Karnataka’s jogappas can now live a gender-fluid life](#). *The Hindu*.
- Sunipa Dev. 2020. [The geometry of distributed representations for better alignment, attenuated bias, and improved interpretability](#).
- Sunipa Dev, Tao Li, Jeff M. Phillips, and Vivek Srikumar. 2020a. [On measuring and mitigating biased inferences of word embeddings](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):7659–7666.
- Sunipa Dev, Tao Li, Jeff M Phillips, and Vivek Srikumar. 2020b. [OSCaR: Orthogonal subspace correction and rectification of biases in word embeddings](#). *arXiv*.
- Sunipa Dev and Jeff M. Phillips. 2019. [Attenuating bias in word vectors](#). In *The 22nd International Conference on Artificial Intelligence and Statistics, AISTATS 2019, 16-18 April 2019, Naha, Okinawa, Japan*, volume 89 of *Proceedings of Machine Learning Research*, pages 879–887. PMLR.
- 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.
- Ethan Fast, Tina Vachovsky, and Michael Bernstein. 2016. Shirtless and dangerous: Quantifying linguistic signals of gender bias in an online fiction writing community. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 10.

- Christine Feraday. 2016. [For lack of a better word: Neo-identities in non-cisgender, non-straight communities on tumblr](#). *Ryerson University*.
- Susan T Fiske. 1993. Controlling other people: The impact of power on stereotyping. *American psychologist*, 48(6):621.
- Andrew Flowers. 2015. [The most common unisex names in america: Is yours one of them?](#) *FiveThirtyEight*.
- Matt Gardner, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson F. Liu, Matthew Peters, Michael Schmitz, and Luke S. Zettlemoyer. 2018. AllenNLP: A deep semantic natural language processing platform. In *ACL workshop for NLP Open Source Software*.
- Vasundhara Gautam. 2021. [Guest lecture in pronouns: Vasundhara](#). In Kirby Conrod, editor, *Pronoun Studies*.
- Tatsunori Hashimoto, Megha Srivastava, Hongseok Namkoong, and Percy Liang. 2018. [Fairness without demographics in repeated loss minimization](#). In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1929–1938. PMLR.
- Os Keyes. 2018. The misgendering machines: Trans/hci implications of automatic gender recognition. *Proc. ACM Hum.-Comput. Interact.*, 2(CSCW).
- Os Keyes, Zoë Hitzig, and Mwenza Blell. 2021. [Truth from the machine: artificial intelligence and the materialization of identity](#). *Interdisciplinary Science Reviews*, 46(1-2):158–175.
- Robin Lakoff. Language and woman’s place. *Language in society*, 2(1).
- Dave Lauer. 2020. [You cannot have ai ethics without ethics](#). In *AI and Ethics*.
- MAP. 2016. [Unjust: How the broken criminal justice system fails transgender people](#). *Movement Advancement Project and Center for American Progress*.
- Gretchen McCulloch and Lauren Gawne. 2016. [Episode 2: Pronouns. little words, big jobs](#). *Linguicism*.
- Ninareh Mehrabi, Thamme Gowda, Fred Morstatter, Nanyun Peng, and Aram Galstyan. 2020. [Man is to person as woman is to location: Measuring gender bias in named entity recognition](#). In *Proceedings of the 31st ACM Conference on Hypertext and Social Media*, HT ’20, page 231–232, New York, NY, USA. Association for Computing Machinery.
- Alfredo Mirandé. 2016. [Hombres mujeres: An indigenous third gender](#). *Men and Masculinities*, 19(4):384–409.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. [GloVe: Global vectors for word representation](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.
- Micah Rajunov and Scott Duane. 2019. *Nonbinary: Memoirs of Gender and Identity*. Columbia University Press.
- Shauli Ravfogel, Yanai Elazar, Hila Gonen, Michael Twiton, and Yoav Goldberg. 2020. [Null it out: Guarding protected attributes by iterative nullspace projection](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7237–7256, Online. Association for Computational Linguistics.
- The Blunt Rose. 2020. [Nonbinary name list](#).
- 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.
- Maarten Sap. 2021. Positive ai with social commonsense models. *Allen Institute for Artificial Intelligence*.
- Morgan Klaus Scheuerman, Aaron Jiang, Katta Spiel, and Jed R. Brubaker. 2021. [Revisiting gendered web forms: An evaluation of gender inputs with \(non-\)binary people](#). In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, CHI ’21, New York, NY, USA. Association for Computing Machinery.
- Julia Serano. 2007. *Whipping Girl: A Transsexual Woman on Sexism and the Scapegoating of Femininity*. Seal Press.
- Emily Sheng, Kai-Wei Chang, Prem Natarajan, and Nanyun Peng. 2021. Societal biases in language generation: Progress and challenges. 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)*.
- Robyn Speer, Joshua Chin, Andrew Lin, Sara Jewett, and Lance Nathan. 2018. [Luminosinsight/wordfreq: v2.2](#).
- Katta Spiel, Os Keyes, and Pinar Barlas. 2019. [Patching gender: Non-binary utopias in hci](#). In *Association for Computing Machinery, CHI EA ’19*, page 1–11, New York, NY, USA.
- 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*.
- Tony Sun, Kellie Webster, Apu Shah, William Yang Wang, and Melvin Johnson. 2021. [They, them, theirs: Rewriting with gender-neutral english](#).
- Yi Chern Tan and L. Elisa Celis. 2019. [Assessing social and intersectional biases in contextualized word representations](#). In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 13209–13220.
- Nat Thorne, Andrew Kam-Tuck Yip, Walter Pierre Bouman, Ellen Marshall, and Jon Arcelus. 2019. [The terminology of identities between, outside and beyond the gender binary - a systematic review](#). *International Journal of Transgenderism*.
- Ankith Uppunda, Susan D. Cochran, Jacob G. Foster, Alina Arseniev-Koehler, Vickie M. Mays, and Kai-Wei Chang. 2021. [Adapting coreference resolution for processing violent death narratives](#).
- Tianlu Wang, Jieyu Zhao, Mark Yatskar, Kai-Wei Chang, and Vicente Ordonez. 2019. [Balanced datasets are not enough: Estimating and mitigating gender bias in deep image representations](#). In *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*, pages 5309–5318. IEEE.
- Shannon Weber. 2019. *Queer Media Images: LGBT Perspectives (Born This Way: Biology and Sexuality in Lady Gaga's Pro-LGBT Media)*. Lexington Books.
- 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.
- Jieyu Zhao, Tianlu Wang, Mark Yatskar, Ryan Cotterell, Vicente Ordonez, and Kai-Wei Chang. 2019. [Gender bias in contextualized word embeddings](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 629–634, Minneapolis, Minnesota. Association for Computational Linguistics.
- Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2017. [Men also like shopping: Reducing gender bias amplification using corpus-level constraints](#). *Proceedings of the EMNLP 2017*.
- 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.

# Appendix: Harms of Gender Exclusivity and Challenges in Non-Binary Representation in Language Technologies

## A Survey

Below, we provide the full Survey on Harms, our rationale for each question, and qualitative analysis of all responses received.

### A.1 Demographic information

**Q:** What pronouns do you use? (checkboxes)

**Options:** they/them, she/her, he/him, xe/xem, e/em, ze/hir, I don't use pronouns, I am questioning about my pronouns, I don't want to answer this question, Other (option to specify in text field)

| Pronouns  | Percentage of Total Respondents |
|-----------|---------------------------------|
| they/them | 84.2%                           |
| she/her   | 26.3%                           |
| he/him    | 15.8%                           |
| xe/xem    | 5.3%                            |

Table 1: Survey Pronouns Distribution

For this question, we ensured to allow respondents to check multiple options, as many non-binary persons use more than one set of pronouns. Table 1 summarizes the distribution of pronouns. 31.6% of respondents used more than one set of pronouns (e.g. *she/her* and *xe/xem*).

**Q:** What are your pronouns in other languages? (text field)

We collected this information because in languages without referential gender or where pronouns are seldom used, pronouns can be less central to one's gender identity; thus, we wanted to discover pronouns non-binary persons use in languages other than English. These data could be useful to future research on the harms of non-English language technologies to non-binary communities.

Many respondents listed their pronouns in other languages. Unmodified responses included:

hen (swedish), iel/any pronoun (french)

| Sexual Orientation                   | Percentage of Total Respondents |
|--------------------------------------|---------------------------------|
| Queer                                | 57.9%                           |
| Bisexual                             | 26.3%                           |
| Asexual                              | 26.3%                           |
| Pansexual                            | 26.3%                           |
| Straight                             | 10.3%                           |
| Gay                                  | 5.3%                            |
| Questioning                          | 5.3%                            |
| I don't want to answer this question | 5.3%                            |

Table 2: Survey Sexual Orientation Distribution

Ta (Mandarin), sie (German)

nin/nim (German)

“Hen” in Swedish, “hän” in Finnish, none in Japanese (pronouns are seldom used at all). Unfortunately, many languages still do not have more commonly accepted gender-neutral pronouns, or I'd use them.

**Q:** What is your sexual orientation? (checkboxes)

**Options:** Lesbian, Gay, Bisexual, Asexual, Pansexual, Queer, Straight, Questioning, I don't want to answer this question, Other (option to specify in text field)

For this question, we ensured to allow respondents to check multiple options. Results may be found in Table 2.

**Q:** What is your gender? (checkboxes)

**Options:** Man, Woman, Non-binary, Genderqueer, Third-gender, Genderfluid, Gender non-conforming, Pangender, Two-Spirit, Agender, Questioning, I don't want to answer this question, Other (option to specify in text field)

For this question, we ensured to allow respondents to check multiple options. Results may be found in Table 3. Table 3 demonstrates that many individuals identify with multiple genders and that we achieved a Western gender-diverse sample.

| Gender                | Percentage of Total Respondents |
|-----------------------|---------------------------------|
| Non-binary            | 73.7%                           |
| Genderqueer           | 31.6%                           |
| Agender               | 31.6 %                          |
| Gender non-conforming | 21.1%                           |
| Questioning           | 21.1%                           |
| Woman                 | 15.8%                           |
| Man                   | 10.5%                           |
| Genderfluid           | 5.3%                            |
| demi-boy              | 5.3%                            |

Table 3: Survey Gender Distribution

**Q:** In a few words, how would you describe your gender and sexual orientation? If you feel that the above questions are not able to capture your gender and sexual orientation, feel free to use this question to enter your response. (text field)

We achieved a Western gender and sexuality-diverse sample. Unmodified responses included:

agender gray-asexual, falling under the broader categories of nonbinary, trans, asexual, and queer

I am a bisexual demi-boy existing somewhere between male and a third gender space

gender non-conforming ace; sex-averse; at my happiest with intimate- or queer-platonic relationships which are generally not labeled “romantic”

I identify as pansexual and queer. I’m questioning my gender but I’m likely somewhere between nonbinary woman and agender.

I present as somewhat femme to “tomboy” and generally live life as a woman, but internally am very agender. I wouldn’t care about my pronouns, but I feel obligated to state them to support various groups of people and generally support the visibility of women. And I’m a pan-romantic sex-disinterested grey-asexual.

I’m panromantic asexual and I’m coming to the conclusion that I’m also agender.

| Trans Identification                 | Percentage of Total Respondents |
|--------------------------------------|---------------------------------|
| Yes                                  | 63.2%                           |
| I am questioning about my gender     | 21.1%                           |
| No                                   | 10.5%                           |
| I don’t want to answer this question | 5.3%                            |

Table 4: Survey Trans Identification Distribution

**Q:** Are you trans? (radio buttons)

**Options:** Yes, No, I am questioning about my gender, I don’t want to answer this question

Results may be found in Table 4.

**Q:** In a few words, how would you describe your ethnicity? (text field)

**Q:** Are you Black, Latinx, and/or Indigenous? (radio buttons)

**Options:** Yes, No, I don’t want to answer this question

**Q:** Are you a person of color? (radio buttons)

**Options:** Yes, No, I don’t want to answer this question

**Q:** Which country are you from originally? (text field)

**Q:** Which country do you live in? (text field)

We intentionally made many of the above questions free-response to allow respondents to explain their ethnicity and nationality.

An overwhelming majority (all but two) of our respondents identified as white, white British, Western European, and/or Caucasian. No respondents were Black, Indigenous, and/or Latinx, and two respondents were people of color.

Furthermore, 52.6% of respondents are originally from the US, and 63.2% current live in the US. Two respondents grew up in India and Singapore, and Taiwan. The remaining respondents are originally from or currently live in Canada and countries in Western Europe, like France, the UK, Germany, and Sweden.

This severely limits the conclusions we can reach from this sample’s responses. In the future, we will improve survey outreach to diversify our sample.

**Q:** How would you describe your occupation?

| Occupation        | Percentage of Total Respondents |
|-------------------|---------------------------------|
| Researcher        | 42.1%                           |
| Software Engineer | 31.6%                           |
| Student           | 26.3%                           |
| Unemployed        | 5.3%                            |

Table 5: Survey Occupation Distribution

(text field)

**Q:** How would you describe your familiarity with AI? (text field)

Occupation results may be found in Table 5. All respondents were familiar with AI, through their occupation, coursework, books they had read, and/or social media. We recognize that these occupations and level of familiarity of AI are correlated with privilege and socioeconomic status; in the future, we will expand our sample beyond those who work in/are familiar with AI.

## A.2 Harms in Language Tasks

This segment first defined representational and allocational harms, and then introduced three tasks (Named Entity Recognition, Coreference Resolution, and Machine Translation) using publicly-available demos (Gardner et al., 2018) which survey subjects engaged with. The demos were accompanied with non-leading questions about potential representational and allocational harms that non-binary communities could face as a result of these tasks.

**Named Entity Recognition (NER)** involves taking an unannotated block of text, such as this one: “Microsoft acquired GitHub in 2018”, and producing an annotated block of text that highlights the names of entities: “[Microsoft]<sub>Organization</sub> acquired [GitHub]<sub>Organization</sub> in [2018]<sub>Time</sub>”. We provided survey participants with the AllenNLP Named Entity Recognition demo (Gardner et al., 2018).

**Q:** Can you see/think of cases where Named-Entity Recognition with current language models could have an undesirable outcome for non-binary genders? (radio buttons)

<https://demo.allennlp.org/named-entity-recognition/fine-grained-ner>

**Options:** Yes, No

84.2% of respondents indicated Yes, while 15.8% indicated No.

**Q:** Does it cause representational harm? Please provide an example(s). (text field)

**Q:** What’s the severity of the representational harm? (1-5 scale)

Respondents argued that NER “systematically mistags neopronouns, which reinforces the stereotype that neopronouns/alternate pronouns are ‘hard’ or ‘complicated’ and is thus directly harmful to non-binary people”. Additionally, NER can assume singular “they” refers to multiple people, and it may label those who use “it/its” pronouns as objects, which is “dehumanizing and reinforces a negative stereotype of non-binary persons”.

Another concern respondents raised was NER’s inability to recognize the names of non-binary persons and correctly tag the Person entity, since many Western non-binary chosen names are creative and diverse, “overlapping with common nouns” (especially nature-related nouns), having “uncommon orthographic forms”, and/or consisting of a single letter. For example, the AllenNLP NER demo cannot correctly tag the full name of a person in the case of a single-letter first name. NER performing badly on these names would “reinforce that non-binary names are ‘weird’”.

Finally, respondents mentioned that NER systems that classify human entities as ‘Person-male’ or ‘Person-female’ and reinforce the gender binary can be psychologically harmful.

Overall, on a scale of 1-5, where 1 indicates “No impact on lives”, 3 indicates “Noticeably negatively affects lives”, and 5 indicates “Significantly hinders lives”, 47.1% of respondents said that the severity of NER’s representational harm is a 3, 23.5% said 2, 17.6% said 4, 5.9% said 1, and 5.9% said 5.

**Q:** Can it cause allocational harm? Please provide an example(s) of use cases and the resultant allocational harm. (text field)

**Q:** What’s the severity of the allocational harm? (1-5 scale)

Respondents said that NER systems can be devastating when they are unable to recognize non-binary chosen names. For example, if organizations scan resumes using NER systems, job and fellowship applications from non-binary persons may be thrown out for “not having a name”. Additionally, if NER systems are used for identity verification, non-binary persons could be “systemically incorrectly labeled by these systems, which could come into play when a system that wants to verify the identity of an account concludes the account does not belong to a human”. Similarly, non-binary people may be unable to access medical and government-administered services if NER is used as a gatekeeping mechanism on healthcare and government websites. NER systems may also be used to automatically build a database of famous people from text data, and if non-binary names are less likely to be correctly recognized, they will be excluded from the database, which could exacerbate erasure.

Overall, on a scale of 1-5, where 1 indicates “No impact on lives”, 3 indicates “Noticeably negatively affects lives”, and 5 indicates “Significantly hinders lives”, 25% of respondents said that the severity of NER’s allocational harm is a 3, 25% said 2, 18.8% said 4, 18.8% said 5, and 12.5% said 1.

**Coreference Resolution** is the task of finding all expressions that refer to the same entity in a block of text. For example, a coreference resolution system would determine that in “[UCLA] is a public university. [It] offers courses in Computer Science.”, “UCLA” corefers with “it”. It is an important step for a lot of higher level NLP tasks that involve natural language understanding such as document summarization, question answering, and information extraction. We provided survey participants with the AllenNLP Coreference Resolution demo (Gardner et al., 2018).

**Q:** Can you see/think of cases where Coreference Resolution with current language models could have an undesirable outcome for non-binary genders? (radio buttons)

**Options:** Yes, No

84.2% of respondents indicated Yes, while 15.8%

<https://demo.allennlp.org/coreference-resolution>

indicated No.

**Q:** Does it cause representational harm? Please provide an example(s). (text field)

**Q:** What’s the severity of the representational harm? (1-5 scale)

Respondents argued that “the potential of accidental misgendering is high”. For instance, coreference resolution systems could “apply ‘s/he’ to individuals who might not identify that way”. Furthermore, coreference resolution systems might “incorrectly apply non-binary pronouns to people who do not use them, like applying ‘it’ to a trans woman who uses ‘she’ pronouns”; this would “echo the societal harm in which people with nonstandard gender presentations are treated as less than human”.

Additionally, respondents mentioned that “neopronoun users as a group are diminished when software does not work on language referencing them”, especially since neopronouns are often underrepresented or even non-existent in textual data. Erasing and neglecting neopronouns contribute to queer erasure, and “when we build coreference systems that cannot handle neopronouns, we reinforce the stereotype that neopronouns/alternate pronouns are ‘hard’ or ‘complicated’, which is directly harmful to non-binary people”.

Similarly, a non-binary person referred to by name and then subsequently by “they/them” or “it/its” pronouns “might fail to be identified as referring to the same person”, because coreference resolution systems could erroneously assume the person is multiple people or an object. For example, the AllenNLP coreference resolution demo cannot correctly handle singular “they” pronouns. One respondent found that in the example, “Alice Smith plays for the soccer team. They scored the most goals of any player last season,” the model connects “they” with “team”; however, English speakers would be able to disambiguate and understand that “they” actually refers to “Alice”.

Furthermore, respondents emphasized that coreference resolution systems can reinforce the idea that names/occupations/roles are gendered and that there are only two genders, e.g. ‘doctor’ is much more likely to link to ‘he’ than ‘they’, ‘she’, ‘xe’, etc.

Overall, on a scale of 1-5, where 1 indicates



“No impact on lives”, 3 indicates “Noticeably negatively affects lives”, and 5 indicates “Significantly hinders lives”, 47.1% of respondents said that the severity of coreference resolution’s representational harm is a 4, 23.5% said 2, 23.5% said 3, 5.9% said 5, and 0% said 1.

**Q:** Can it cause allocational harm? Please provide an example(s) of use cases and the resultant allocational harm. (text field)

**Q:** What’s the severity of the allocational harm? (1-5 scale)

Respondents provided numerous realistic harmful use cases of coreference resolution. For instance, a “ranking system where you count citations of a person from a body of text (including references to their pronouns which you would resolve through coreference resolution) could miss a lot of instances of people being cited with ‘xe/xem’ pronouns, which would give them a lower ranking”. Another respondent conceived an example in which “a person using singular ‘they’ pronouns who was required to sign a lease that populated referents with ‘s/he’ pronouns instead is forced to sign an incorrect acknowledgement or not obtain housing”. Furthermore, coreference resolution systems might cause applications for financial aid targeted solely at individuals from non-binary persons who use “they/them” pronouns to be automatically flagged as ineligible. Finally, a respondent described a situation in which “if coreference resolution is used to sort through large law corpora to find instances of non-binary people being discriminated against to see if more stringent policy should be put in place to stop discrimination, it may erroneously find that there are not many cases of this since ‘they’ is not often linked to a specific person”.

Overall, on a scale of 1-5, where 1 indicates “No impact on lives”, 3 indicates “Noticeably negatively affects lives”, and 5 indicates “Significantly hinders lives”, 35.7% of respondents said that the severity of coreference resolution’s allocational harm is a 4, 21.4% said 2, 21.4% said 3, 14.3% said 5, and 7.1% said 1.

**Machine Translation** systems translate text from one language to another. When translating from languages with pronouns that do not carry gender information (e.g. Tagalog) to those that have

gendered pronouns (e.g. English), translation systems may impose incorrect binary pronouns on individuals. This can be problematic in several ways such as reinforcing gender stereotypes, and misgendering and excluding non-binary persons. We provided survey participants with Google Translate .

**Q:** Can you see/think of cases where Machine Translation with current language models could have an undesirable outcome for non-binary genders? (radio buttons)

**Options:** Yes, No

89.5% of respondents indicated Yes, while 10.4% indicated No.

**Q:** Does it cause representational harm? Please provide an example(s). (text field)

**Q:** What’s the severity of the representational harm? (1-5 scale)

Respondents overwhelmingly discussed the harm of machine translation systems “translating from a language where pronouns are unmarked for gender, and picking a gender grounded in stereotypes associated with the rest of the sentence” in the translation in the target language. Many respondents raised the example of translating “‘3SG is a nurse’ (in some language) to ‘She is a nurse’ in English and ‘3SG is a doctor’ (in some language) to ‘He is a doctor’ in English”. Another example, based in heteronormativity, is Google Translate French-to-English “translates ‘sa femme’ (his/her/their wife) as ‘his wife’ and ‘son mari’ (his/her/their husband) as ‘her husband’ even in sentences with context, e.g. ‘Elle et sa femme se sont mariées hier’ (‘she and her wife got married yesterday’) is translated as ‘she and his wife got married yesterday’”.

Furthermore, the long-established gender-neutral pronouns “‘hen’ and ‘hän’ from Swedish and Finnish” and “strategies to mix gendered inflections” all often “automatically translate to ‘her’ or ‘him’” in English. In addition, machine translation systems can “misinterpret non-binary names and pronouns as referring to objects, thereby dehumanizing non-binary people”. This can lead to nonbinary people being misgendered if their pronouns do not align with the ones that the machine

<https://translate.google.com/>

translation system imposed upon them. Moreover, “if neopronouns are not even represented, then this also contributes to erasure of queer identity”; it is likely that neopronouns “are represented as unknown tokens,” which can be problematic.

Additionally, many grammatically-gendered languages lack non-binary gender options, so a person may have their gender incorrectly “binary-ified” in the target language, which constitutes misgendering and “is hurtful”.

Finally, differences in how languages talk about non-binary people are “extremely nuanced”, which can lead to “extremely disrespectful” translations. One respondent explained that, while “a common and accepted way to refer to trans people in Thailand is the word ‘kathoe’, which translates to ‘ladyboy’”, if someone called this respondent a “lady-boy” in English, the respondent would be extremely offended.

Overall, on a scale of 1-5, where 1 indicates “No impact on lives”, 3 indicates “Noticeably negatively affects lives”, and 5 indicates “Significantly hinders lives”, 47.1% of respondents said that the severity of machine translation’s representational harm is a 4, 29.4% said 3, 17.6% said 5, 5.9% said 1, and 0% said 2.

**Q:** Can it cause allocational harm? Please provide an example(s) of use cases and the resultant allocational harm. (text field)

**Q:** What’s the severity of the allocational harm? (1-5 scale)

Respondents argued that if machine translation is used in medical or legal contexts, a translation that automatically applies incorrectly gendered terms can result in incorrect care or invalidation. An example provided was “a nonbinary AFAB person might not be asked about their pregnancy status when being prescribed a new medication if a cross-lingual messaging system assigned ‘male’ terms to them”. Furthermore, non-binary persons might be “denied a visa or convicted of a crime due to mistranslation of evidence”. Another very real consequence of machine translation systems misgendering that respondents brought up is that they can deny non-binary persons “gender euphoria” (i.e. the joy of having one’s gender affirmed) and cause psychological harm.

Overall, on a scale of 1-5, where 1 indicates

“No impact on lives”, 3 indicates “Noticeably negatively affects lives”, and 5 indicates “Significantly hinders lives”, 33.3% of respondents said that the severity of machine translation’s allocational harm is a 5, 20% said 2, 20% said 3, 20% said 4, and 6.7% said 1.

**Q:** Rank the representational harms caused by the aforementioned tasks by severity for the worst realistic use case.

Results may be found in Table 6.

**Q:** Rank the allocational harms caused by the aforementioned tasks by severity for the worst realistic use case.

Results may be found in Table 7.

### A.3 Broader Concerns with Language Models

This segment was purposely kept less specific to understand the harms in different domains (healthcare, social media, etc.) and their origins, as perceived by different non-binary individuals.

**Q:** Can you see/think of domains (e.g. healthcare, social media, public administration, high-tech devices, etc.) to which language models can/could be \*applied\* in a way that produces undesirable outcomes for non-binary individuals? If so, please list such domains below. (text field)

**Q:** For each domain you listed above, please provide an example(s) of harmful applications and use cases and evaluate the severity of the resultant harms.

### Social Media

LGBTQ+ social media content is automatically flagged at higher rates. Ironically, language models can fail to identify hateful language targeted at nonbinary people. Furthermore, if social media sites attempt to infer gender from name or other characteristics, this can lead to incorrect pronouns for non-binary individuals. Additionally, “language models applied in a way that links entities across contexts are likely to out and/or deadname people, which could harm trans and nonbinary people”. Moreover, social media

| Severity Ranking \ Task | NER   | Coreference Resolution | Machine Translation |
|-------------------------|-------|------------------------|---------------------|
| Lowest                  | 52.6% | 21.1%                  | 21.1%               |
| In-Between              | 21.1% | 42.1%                  | 31.6%               |
| Highest                 | 21.1% | 31.6%                  | 42.1%               |

Table 6: Language Task Rankings by Severity of Representational Harm

| Severity Ranking \ Task | NER   | Coreference Resolution | Machine Translation |
|-------------------------|-------|------------------------|---------------------|
| Lowest                  | 36.8% | 26.3%                  | 26.3%               |
| In-Between              | 21.1% | 42.1%                  | 26.3%               |
| Highest                 | 31.6% | 21.1%                  | 36.8%               |

Table 7: Language Task Rankings by Severity of Allocational Harm

identity verification could incorrectly interpret non-binary identities as fake or non-human.

### Productivity Technologies

Autocomplete could suggest “only binary pronouns, or make predictions that align with gender stereotypes”.

### Healthcare

Respondents said that “healthcare requires engaging with gender history as well as identity”, which language models are not capable of doing, and “even humans intending to do well and using the best terms they know often struggle with the limitations of our language for nonbinary people and their bodies”. Language models could further “misgender patients”. Additionally, language models could “deny insurance claims, e.g. based on a ‘mismatch’ between diagnosis and gender/pronouns”.

### Policing

Respondents said that “any system which incorrectly handles singular ‘they’ might result in communications being flagged as false, self-contradictory, or incomplete”.

### Marketing and Customer Service

Language models could enable “predatory or adversarial advertising” for non-binary persons.

### Hiring

Respondents explained that “a system which incorrectly handles singular ‘they’ might result in non-binary people’s achievements being misattributed to group work or to organizations they worked for”.

### Finance

Finance-related identity verification could incorrectly interpret non-binary identities as fake or non-human.

### Government-Administrated Services

Government services could misgender non-binary persons or reject their applications based on language analysis.

### Education

Language models employed in automated educational/grading tools could “automatically mark things wrong/‘ungrammatical’ for use of non-standard language, singular ‘they’, neopronouns, and other ‘new’ un- or creatively gendered words”.

**Q:** Can you see/think of applications of language models that can/could exacerbate non-binary erasure? If so, please list such applications below. (text field)

**Q:** For each application you listed above, please provide an example(s) of harmful use cases and evaluate the severity of the resultant harms using the 1-5 scale below. (text field)

Automated summarization (e.g. used in Google's information boxes) could erase non-binary persons. For example, non-binary people are "more likely to be tagged as non-human and thus less likely to have their achievements accurately summarized, which makes them make invisible".

Moreover, current language models cannot generate text with nonbinary people or language (e.g. "it never generates sentences with 'they/them' pronouns or 'ze/hir' pronouns or a sentence like 'She is nonbinary', but it regularly generates examples with 'he/him' and 'she/her' pronouns and sentences like 'He is a man'); this decidedly contributes to nonbinary erasure.

Screen readers and speech-to-text services that cannot handle neopronouns may also erase non-binary individuals. Similarly, "neopronouns are almost always listed as 'wrong'" by spelling and grammar checkers.

Machine translation is particularly prone to erasing non-binary gender because "nonbinary people often create new ways of using language out of necessity, and these usages are rare/new enough to not be reflected in machine translation".

One respondent said that "any model that attempts to classify gender" contributes to nonbinary erasure because "'nonbinary' is not a single entity or identity type"; further, "treating 'nonbinary' as a distinct third gender or as some 'ambiguous' category" is also erasing.

**Cisgender Privilege** is the unearned benefits you receive when your gender identity matches your sex assigned at birth.

**Q:** Can you see/think of applications of language models that can/could exacerbate transphobia or denial of cisgender privilege? If so, please list such applications below. (text field)

**Q:** For each application you listed above, please provide an example(s) of harmful use cases and evaluate the severity of the resultant harms. (text field)

Respondents said that language models can exacerbate transphobia "by incorrectly flagging non-toxic content from trans persons as toxic at higher rates, or by not recognizing transphobic comments as toxic".

Furthermore, "any system that attempts to ascertain gender or pronouns from a name or other attributes" can enable cisgender privilege. A respondent explained that "if a model misgenders someone because it accounts for a history of them having another name, or does not allow for flexibility in gender signifiers to change over time, it reinforces the idea that gender is or should be immutable, that individuals have a 'true' gender and then 'change it,' that gender can only 'change' once if it happens at all, and that there is some clear point of demarcation where you 'transition' and only ever in one direction (binary transition)"; furthermore, "any corrections to those assumptions in the model would necessarily be post-hoc, marking oneself as 'other' for not fitting into the binary construction" of gender.

Additionally, there are dangerous applications, like bots on social media, that systematically harass non-binary people.

Language models can also be used to empower the enforcement mechanisms of transphobic policies. This could occur in "visual/textual systems for things like passport verification or other legal processing like getting driver's licenses, applying for loans, attempting to obtain citizenship".

While developing and testing language model-based systems, developers may find that the language nonbinary persons have created for themselves is not compatible with their systems. Hence, developers may "blame nonbinary people" for the difficulty associated with including them and "decide that the systems will only serve binary-aligned people". However, "this both increases cis/binary privilege (by making those systems inaccessible to nonbinary people) and increases transphobia (by creating or strengthening feelings of resentment towards people who do not fit conveniently into the gender binary)".

**Q:** What are the top reasons that there exist limited data concerning non-binary individuals in Natural Language Processing (NLP) in your opinion?

Most respondents cited a lack of diversity in developer/research teams. They said there exists "limited trans/non-binary representation so knowledge gaps exist", and many developers and researchers have a "lack of knowledge about nonbinary identities, transness, queerness, etc." Further, devel-

oper/research teams tend to “want to simplify variables and systems after and in spite of learning about the complexity” of gender identity, and may not consider non-binary persons “important enough to change their systems for”.

Respondents also discussed the sources of training data. They explained that “most training data pulls from large scale internet sources and ends up representing hegemonic viewpoints”. Additionally, “lots of our models are built on Wikipedia which has few non-binary people, Reddit which is hardly a queer utopia, and books written when gender non-conforming content did not get published much”. Moreover, non-binary data may be “discarded as ‘outliers’” and “not sampled in training data”. Data annotators may also “not recognize non-binary identities”, “non-binary identities may not be possible labels”, and/or annotators “may not be paying attention to non-binary identities due to minimal wages, lack of situational context, and lack of epistemic uncertainty in the models”.

Other major reasons included “historic erasure, active discrimination, invisibility of some non-binary identities, small non-binary population”.

From Barocas et al. (2019):

**Skewed sampling:** model-influenced, positive feedback loops in data collection

**Tainted examples:** historically-biased training data

**Limited features:** missing or erroneous features for training examples

**Sample size disparities:** insufficient training examples for minority classes

**Proxies:** correlated features leak undesirable information

**Q:** Score the following barriers to better including non-binary individuals in language models using this 1-5 scale: 1 (Easy solutions that could be deployed immediately), 3 (Could eventually achieve solutions), 5 (Impossible to remedy)

Results may be found in Table 8.

**Q:** Can you see/think of cases where harms (representational or allocational) are compounded for non-binary individuals with particular intersecting identities? (radio buttons)

**Options:** Yes, No

89.5% of respondents indicated Yes, while 10.4% indicated No.

**Q:** If Yes, could you give examples of such intersecting identities? (text field)

**Q:** For each intersecting identity, please provide an example(s) of harmful use of language models and evaluate the severity of the resultant harm. (text field)

Issues with coreference resolution could be compounded for non-binary persons with non-Western names. In addition, machine translation can fail for someone when translating “from a culture where their nonbinary identity does not fit neatly into the nonbinary boxes we’ve devised in English”. Non-binary racial minorities in Western societies will also be misrepresented and underrepresented in data samples. And, non-white non-binary individuals are more susceptible to misgendering and related harms in policing that employs language models.

Furthermore, medical harms can be worsened for non-binary persons who already have limited access to healthcare due to other aspects of their identity, like race, immigration status, fluency in English, etc. Moreover, non-binary persons with limited fluency in a language who have more interactions with machine translation systems are more likely to regularly incur the aforementioned representational and allocational harms posed by the systems.

Additionally, “some neurodivergent people refer to themselves with traditionally dehumanizing language, which could compound the issue of models not recognizing their identities as real and human” if they’re also non-binary. Further, non-binary persons with certain disabilities who rely on language model-based speech-to-text services may not have their pronouns recognized.

Other examples of intersecting identities included: class, body size, religious affiliation, nationality, sexual or romantic orientation, age, and education level.

## B Dataset Skews

Usage of non-binary pronouns in text is not always meaningful with respect to gender, as seen in Table 9.

| Source of Bias          | Feasibility Ranking |       |       |       |       |
|-------------------------|---------------------|-------|-------|-------|-------|
|                         | 1                   | 2     | 3     | 4     | 5     |
| Skewed Sampling         | 0%                  | 36.8% | 47.4% | 5.3%  | 0%    |
| Tainted Examples        | 5.3%                | 26.3% | 47.4% | 5.3%  | 5.3%  |
| Limited Features        | 10.5%               | 21.1% | 21.1% | 47.4% | 10.5% |
| Sample Size Disparities | 5.3%                | 36.8% | 26.3% | 26.3% | 0%    |
| Proxies                 | 5.3%                | 5.3%  | 31.6% | 36.8% | 10.5% |

Table 8: Feasibility of Mitigation Rankings for Sources of Bias

Table 9: Example sentences containing nonbinary pronouns

| Pronoun | Sentence  |
|---------|---|
| Ey      | “The difference in the alphabets comes only in the Faroese diphthongs (ei being 26, ey 356, oy 24...)”  |
| Em      | Approximating the em dash with two or three hyphens.  |
| Xem     | “‘Em ði xem hoi trang ram’”, establishing her icon for Vietnamese women as well as earning the title of the “‘Queen of Folk’”   |
| Ze      | “He taught himself to write with his left hand and described his experiences before, during, and after the accident in a deeply moving journal, later published under the title ‘Pogodzic sie ze swiatem’ (‘To Come to Terms with the World’).” |
| Zir     | “The largest operation in the Struma Valley was the capture by 28th Division of Karajakoi Bala, Karajakoi Zir and Yenikoi in October 1916.”   |

Further, the distribution of different pronouns is also not equal across genders. Overall, using the Python library *wordfreq* <https://pypi.org/project/wordfreq/> which samples over diverse data to give an approximate usage of different words in all of the text curated from the web, we observe how vastly different the frequencies of different gendered words are per billion words in English (Speer et al., 2018). While ‘he’ and ‘she’ occur 0.49% and 0.316% per billion words respectively, the percent for ‘xe’ and ‘ze’ is only 0.0005% and 0.0011% respectively. We list these percentages for a larger set of gendered words in Appendix 10 to highlight this disparity.

### B.1 Representation Skews

Glove was trained on English Wikipedia articles with a window size of 15, a dimension of 50 for each word, and a minimum word frequency of 5. Skews as seen in GloVe representations are seen here with respect to nearest neighbors in Table 11 and often even with derogatory associations reflecting social biases (Table 12).

**Biases With Respect to Occupations** Binary gender and their stereotypically associated occupations is a bias widely discussed. We see in Table 13 that is not very relevant for non-binary-gendered persons and the biases faced by them.

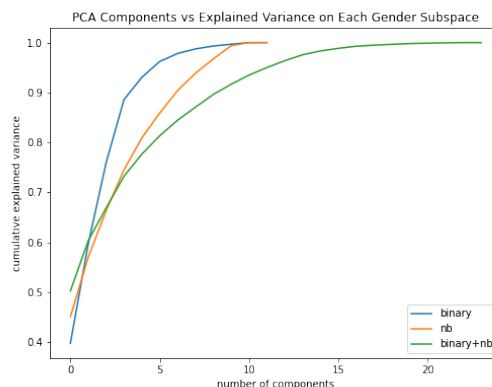


Figure 1: PCA Components for each Gender Subspace

**Subspace analyses** Capturing a gender subspace has been useful in techniques of bias analysis and techniques in subsequent debiasing in binary gender (Dev and Phillips, 2019), especially in context-free or static representations like GloVe or word2vec. These methods postulate expanding this to nonbinary gender by determining a general subspace for gender which captures both binary and non-binary genders. We test if we can approach capturing the all-gender subspace by extending one such general subspace capturing method (Bolukbasi et al., 2016) - principal component analysis (PCA) - on 3 groups of words:

1. **Binary set:** ‘he’, ‘she’, ‘man’, ‘woman’, ‘hers’, ‘his’, ‘herself’, ‘himself’, ‘girl’, ‘boy’,

<https://dumps.wikimedia.org/>

Table 10: Per Billion Word Frequency in the English Language

| Word        | Frequency (%) |
|-------------|---------------|
| he          | 0.49          |
| his         | 0.324         |
| they        | 0.316         |
| she         | 0.182         |
| them        | 0.155         |
| man         | 0.0661        |
| girl        | 0.024         |
| woman       | 0.0224        |
| himself     | 0.0178        |
| boy         | 0.0148        |
| female      | 0.01          |
| male        | 0.00776       |
| herself     | 0.00603       |
| two-spirit  | 0.00588       |
| em          | 0.00372       |
| hers        | 0.00093       |
| transgender | 0.00081       |
| queer       | 0.00057       |
| ey          | 0.00019       |
| ze          | 0.00012       |
| xe          | 5e-05         |
| nonbinary   | 2e-05         |
| cisgender   | 2e-05         |
| genderqueer | 1e-05         |
| zem         | 1e-05         |
| genderfluid | 1e-05         |
| xem         | 0             |
| zey         | 0             |
| zir         | 0             |
| bigender    | 0             |
| xir         | 0             |
| cisman      | 0             |
| ciswoman    | 0             |
| xey         | 0             |

Table 11: Five Nearest neighbors for binary and non-binary possessive pronouns

| Pronoun | Top 5 Neighbors   |
|---------|---|
| His     | 'he', 'him', 'who', 'after', 'himself'                          |
| Hers    | 'somehow', 'herself', 'thinks', 'someone', 'feels'              |
| Theirs  | 'weren', 'tempted', 'couldn', 'gotten', 'willingly'             |
| Xers    | "yogad", 'doswelliids', 'hlx', 'cannibalize', 'probactrosaurus' |
| Zers    | 'ditti', 'bocook', 'kurikkal', 'felimy', 'hifter'               |
| Eirs    | 'cheor', 'yha', 'mnetha', 'scalier', 'paynet'                   |

Table 12: Ten Nearest neighbors of non-binary terms highlighting derogatory Terms

| Term        | 10 Nearest Neighbors  |
|-------------|---|
| agender     | bigender, genderfluid, genderqueer, tosin, cisgender, nonbinary, laia, muhafazat, <b>negrito</b> , farmgirl                 |
| bigender    | pangender, agender, genderfluid, overcontact, pnong, genderqueer, nonbinary, eczemas, gega                                  |
| queer       | lesbian, lgbtq, feminism, lgbt, lesbians, feminist, racism, sexuality, stereotypes, gay                                     |
| nonbinary   | genderqueer, <b>transsexual</b> , cisgender, transsexuals, bigender, genderfluid, chorti, referents, pansexual, hitchhikers |
| transgender | lesbian, lgbt, lgbtq, bisexual, intersex, gender, <b>transsexual</b> , lesbians, heterosexual, discrimination               |
| genderfluid | agender, bigender, genderqueer, transwoman, nonbinary, pansexual, montserratian, <b>negrito</b> , supercouple, <b>fasiq</b> |
| genderqueer | pansexual, nonbinary, lgbtqia, <b>transsexual</b> , genderfluid, agender, bisexuality, bigender, diasporic, multiracial     |

| occupation | man   | woman | transman | transwoman | cisgender | transgender | nonbinary | genderqueer | genderfluid | bigender |
|------------|-------|-------|----------|------------|-----------|-------------|-----------|-------------|-------------|----------|
| doctor     | 0.809 | 0.791 | -0.062   | -0.088     | 0.094     | 0.388       | 0.037     | 0.022       | 0.069       | -0.107   |
| engineer   | 0.551 | 0.409 | -0.152   | -0.271     | -0.227    | 0.043       | -0.243    | -0.176      | -0.084      | -0.298   |
| nurse      | 0.616 | 0.746 | -0.095   | 0.050      | 0.206     | 0.527       | 0.129     | 0.083       | 0.182       | 0.022    |
| stylist    | 0.382 | 0.455 | 0.018    | 0.062      | 0.117     | 0.318       | 0.015     | 0.126       | 0.207       | -0.017   |

Table 13: Cosine similarity between occupations and words

| Set      | Adjectives   |
|----------|--|
| positive | <i>smart, wise, able, bright, capable, ambitious, calm, attractive, great, good, caring, loving, adventurous</i> |
| negative | <i>dumb, arrogant, careless, cruel, coward, boring, lame, incapable, rude, selfish, dishonest, lazy, unkind</i>  |

Table 14: Positive and Negative Adjective Sets

'female', 'male'

2. **Nonbinary set:** 'they', 'them', 'xe', 'ze', 'xir', 'zir', 'xey', 'zey', 'xem', 'zem', 'ey', 'em'

3. **Binary + Non-Binary set**

If we truly captured the gender subspace, we could safely assume that the difference between the binary subspace and the all-gender subspace, along with the non-binary subspace and the all-gender subspace, is somewhat negligible. We make the following observations leveraging the cosine distance, defined as  $1 - c$ , where  $c$  is the cosine similarity between two vectors. We observe, opposite to what we expected, that the distance was quite different in these respective pairs. Between the binary and all-gender subspace was a cosine distance of 1.48, while the distance between the non-binary and all-gender subspace was larger, at 1.93. This tells us that the binary subspace is much less dissimilar than the nonbinary subspace with respect to the all-gender subspace, i.e., extending the approach of subspace capture to all genders would result in a subspace more dominantly aligned with binary gender than non-binary gender. Further, due to the poor representation of non-binary pronouns, the subspace is likely representing the difference in frequency of terms rather than the concept of gender as a whole. Due to weaker alignment with the non-binary gender, any tasks performed using this new 'gender' subspace would not be very effective or applicable to non-binary genders, thus indicating towards further skews and harm.

## B.2 Words for WEAT and Similarity Tests

WEAT tests require groups of words. We list herein Tables 16, 17, 20 the pronouns and words used as well as the pleasant and unpleasant words we compared them against.

## B.3 BERT experiments

We create a balanced, labeled dataset containing sentences using either *they(s)* or *they(p)*. The text spans for *they(p)* are chosen randomly from Wikipedia containing pairs of sentences such that the word *they* appears in the second sentence (with no other pronoun present) and the previous sentence has a mention of two or more persons (determined by NER). This ensures that the word *they* in this case was used in a plural sense. For the samples with singular usage of *they*, since Wikipedia does not have a large number of such sentences in general, we randomly sample them from the Non-Binary Wiki . The sentences are manually annotated for further confirmation of correct usage of each pronoun. We follow the procedure of data collection for *they(s)* to create datasets for sentences with words *he* and *she* from Wikipedia. In each dataset, we replace the pronouns with the [MASK] token and use BERT to predict the token's representation. Figure 5 demonstrates that the representations of *he*, *she* and *they(s)* acquired in this manner are similarly separable from the representations of *they(p)*.

**BERT templates for Misgendering** We used templates of sentences to evaluate if misgendering is propagated by BERT. The sentence templates used are:

[Name] went to the hospital for [pronoun] appointment. [MASK] was feeling sick.

[Name] went to the salon for [pronoun] appointment. [MASK] got a haircut.

[Name] went to the university for [pronoun] classes. [MASK] had a morning class.

[Name] went to the restaurant near [pronoun] house. [MASK] bought lunch.

[Name] drove to the hospital for [pronoun]

[https://en.wikipedia.org/wiki/Category:Non-binary\\_gender](https://en.wikipedia.org/wiki/Category:Non-binary_gender)



| Proxy       | Average |
|-------------|---------|
| man         | 0.462   |
| woman       | 0.494   |
| transman    | -0.043  |
| transwoman  | 0.088   |
| cisgender   | 0.101   |
| transgender | 0.228   |
| nonbinary   | 0.025   |
| genderqueer | 0.054   |
| genderfluid | 0.126   |
| bigender    | -0.052  |

Table 15: Average similarity between occupations and words

| Set                | Words  |
|--------------------|--|
| binary pronouns    | <i>he, him, his, she, her, hers</i>  |
| binary words       | <i>man, woman, herself, himself, girl, boy, female, male, cisman*, ciswoman*</i>         |
| binary all         | <i>binary pronouns + binary words</i>  |
| nonbinary pronouns | <i>zey, ey, em, them, xir, they, zem, ze, their, zir, zers, eirs, xey, xers, xe, xem</i> |
| nonbinary words    | <i>transgender, queer, nonbinary, genderqueer, genderfluid, bigender, two-spirit</i>     |
| nonbinary all      | <i>nonbinary pronouns + nonbinary words</i>  |

Table 16: Word set definitions for binary and non-binary concepts

| Set        | Words   |
|------------|---|
| pleasant   | <i>joy, love, peace, wonderful, pleasure, friend, laughter, happy</i> |
| unpleasant | <i>agony, terrible, horrible, nasty, evil, war, awful, failure</i>    |

Table 17: Set of unpleasant and pleasant words

| Words | Average | Absolute Average |
|-------|---------|------------------|
| he    | 0.509   | 0.509            |
| him   | 0.465   | 0.465            |
| his   | 0.498   | 0.498            |
| she   | 0.495   | 0.495            |
| her   | 0.473   | 0.473            |
| xir   | -0.197  | 0.207            |
| they  | 0.395   | 0.395            |
| xey   | -0.007  | 0.094            |
| them  | 0.389   | 0.390            |
| ey    | 0.086   | 0.111            |
| zey   | -0.056  | 0.108            |
| xe    | -0.054  | 0.111            |
| their | 0.378   | 0.379            |
| xers  | -0.088  | 0.105            |
| em    | 0.185   | 0.209            |
| zir   | -0.035  | 0.092            |
| zem   | -0.068  | 0.091            |
| eirs  | -0.158  | 0.185            |
| zers  | -0.104  | 0.116            |
| ze    | 0.123   | 0.143            |
| xem   | -0.169  | 0.180            |

Table 18: Average cosine similarity between occupations and pronouns

| Words       | Average | Absolute Average |
|-------------|---------|------------------|
| man         | 0.469   | 0.469            |
| woman       | 0.473   | 0.473            |
| herself     | 0.400   | 0.400            |
| himself     | 0.483   | 0.483            |
| girl        | 0.421   | 0.421            |
| boy         | 0.457   | 0.457            |
| female      | 0.393   | 0.394            |
| male        | 0.350   | 0.353            |
| transman    | -0.052  | 0.098            |
| transwoman  | 0.023   | 0.125            |
| transgender | 0.262   | 0.262            |
| queer       | 0.184   | 0.192            |
| nonbinary   | 0.010   | 0.088            |
| genderqueer | 0.048   | 0.099            |
| genderfluid | 0.079   | 0.113            |
| bigender    | -0.085  | 0.118            |

Table 19: Average cosine similarity between occupations and words

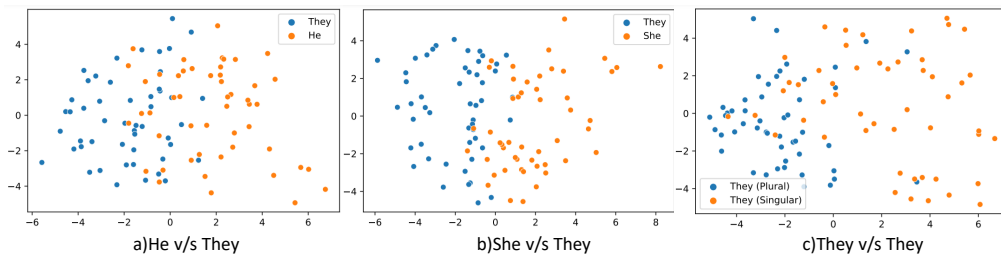


Figure 2: BERT representation analysis: a) he v/s they(p) b) she v/s they(p) c) they(s) v/s they(p)

| Sets  | Weat Score |
|---|------------|
| Random Vectors  | -0.02      |
| Binary Pronouns vs. Non-Binary Pronouns                   | 0.2        |
| Binary Words vs. Non-Binary Proxies                       | 0.718      |
| Binary Pronouns + Words vs. Non-Binary Pronouns + Proxies | 0.916      |

Table 20: WEAT Scores (vs. pleasant and unpleasant attributes)

| Classifier | Accuracy |
|------------|----------|
| $C_1$      | 67.7%    |
| $C_2$      | 83.3%    |
| $C_3$      | 83.1%    |

Table 21: The performance of BERT classifier

appointment. [MASK] was feeling sick.

[Name] drove to the salon for [pronoun] appointment. [MASK] got a haircut.

[Name] drove to the university for [pronoun] classes. [MASK] had a morning class.

[Name] drove to the restaurant near [pronoun] house. [MASK] bought lunch.

[Name] walked to the hospital for [pronoun] appointment. [MASK] was feeling sick.

[Name] walked to the salon for [pronoun] appointment. [MASK] got a haircut.

[Name] drove to the university for [pronoun] classes. [MASK] had a morning class.

[Name] fed [pronoun] dog. [MASK] had to leave for work.

[Name] met [pronoun] friend at the cafe. [MASK] ordered a coffee.

[Name] attached a file to [pronoun] email. [MASK] sent the email.

[Name] realized [pronoun] left [pronoun] keys at home. [MASK] ran back to get the keys.

[Name] found [pronoun] drivers license on the pavement. [MASK] picked it up.

[Name] checks [pronoun] phone constantly. [MASK] is expecting an important email.

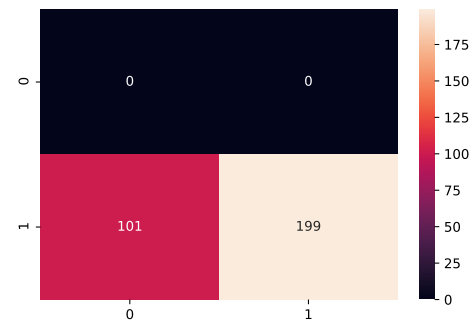


Figure 3: Confusion matrix of Classifier C1

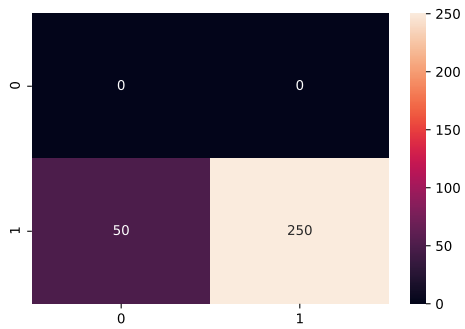


Figure 4: Confusion matrix of Classifier C2

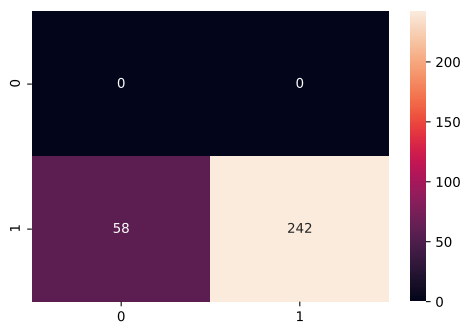


Figure 5: Confusion matrix of Classifier C3

[Name] said that [pronoun] child was just born.  
 [MASK] is excited for the future.  
 [Name] is in a rush to attend [pronoun] lecture.  
 [MASK] eats lunch quickly.  
 [Name] enjoys riding [pronoun] bike. [MASK] is  
 able to get anywhere.

We vary the name (over 900 names) and pro-  
 nouns as described in the paper in Section 4.2.