# MISGENDERMENDER:
# A Community-Informed Approach to Interventions for Misgendering

**Tamanna Hossain**
University of California, Irvine
tthossai@uci.edu

**Sunipa Dev**
Google Research
sunipadev@google.com

**Sameer Singh**
University of California, Irvine
sameer@uci.edu

## Abstract

*Content Warning:* This paper contains examples of misgendering and erasure that could be offensive and potentially triggering.

Misgendering, the act of incorrectly addressing someone's gender, inflicts serious harm and is pervasive in everyday technologies, yet there is a notable lack of research to combat it. We are the first to address this lack of research into interventions for misgendering by conducting a survey of gender-diverse individuals in the US to understand perspectives about automated interventions for text-based misgendering. Based on survey insights on the prevalence of misgendering, desired solutions, and associated concerns, we introduce a misgendering interventions task and evaluation dataset, MISGENDERMENDER. We define the task with two sub-tasks: (i) detecting misgendering, followed by (ii) correcting misgendering where misgendering is present in domains where editing is appropriate. MISGENDERMENDER comprises 3790 instances of social media content and LLM-generations about non-cisgender public figures, annotated for the presence of misgendering, with additional annotations for correcting misgendering in LLM-generated text. Using this dataset, we set initial benchmarks by evaluating existing NLP systems and highlighting challenges for future models to address. We release the full dataset, code, and demo at https://tamannahossainkay.github.io/misgendermender/.

## 1 Introduction

Misgendering is the act of referring to someone using a word, e.g. a pronoun or title, that does not correctly reflect the gender with which they identify (Dictionary, 2023). While there is growing awareness about the adverse impacts of misgendering on peoples' lives (Dev et al., 2021), there is insufficient scholarship or resources that iden-

| Linguistic Gender Profile | |
|---|---|
| **Name:** | Elliot Page |
| **Gender identity:** | Trans man, Non-binary |
| **Pronouns:** | he/him/his/his/himself, they/them/their/theirs/themselves |
| **Gendered Terms:** | masculine, neutral |
| **Deadname:** | Ellen Grace Philpotts-Page |

**Annotated Content**

Detect-Only
**X Post:** John Wayne was a man and Elliot Page is a woman...
**Detect Label:** Misgendering

**X Post:** ..."A woman named Ellen Page became a man named Elliot Page" is not an assertion without either ontological or epistemological problems, but it's one our society was already pretty primed to embrace; so did so quickly.
**Detect Label:** No Misgendering

Detect+Correct
**LLM-generation:** Ellen Grace credits her mother with her success, and she is eternally grateful for her love and support.
**Detect Label:** Misgendering
**Corrected:** ~~Ellen~~ →Elliot credits ~~her~~ →his mother with ~~her~~ →his success, and ~~she~~ →he is eternally grateful for her love and support.

Figure 1: MISGENDERMENDER examples consisting of a gender linguistic profile and corresponding annotated content for detecting and correcting misgendering.

tify and attempt to mitigate misgendering in these various daily use platforms and technologies.

Efforts to measure and mitigate gender bias in natural language processing primarily focus on cisgender and binary gender categories (Guo et al., 2022; Choubey et al., 2021). Few efforts to address non-traditional gender categories have evaluated LLMs' abilities to use non-binary pronouns (Hossain et al., 2023), coreference resolution using neo-pronouns (Cao and Daumé III, 2020), and representational biases in word embeddings (Dev et al., 2021). Furthermore, even though misgendering is both a factual inaccuracy and a toxic act of identity erasure, research on factuality and toxicity has largely ignored it (Gao and Emami, 2023; Lees et al., 2022).

Our contribution is two-fold: (i) we conduct a community survey to understand opinions about automated interventions for text-based misgendering, and (ii) based on the survey, we define a task and evaluation dataset for addressing misgendering in text-based content. Our survey of gender-diverse[1] individuals revealed a prevalent issue of misgendering, especially on social media, but also in other areas like AI-generated content, news articles, and academic journals (§ 2). While there was a general preference for automatic detection of misgendering across domains, opinions diverged on measures such as correcting or hiding misgendered content (§ 2.1). Participants were more receptive to the idea of auto-correction in AI-generated content than social media, citing concerns over limiting freedom of speech and creating a false sense of allyship. Importantly, there were significant apprehensions regarding the implementation of any automated systems to address misgendering, encompassing issues like the fundamental infeasibility of these systems, privacy, the risk of profiling or targeting based on gender linguistic preferences databases, and doubts about the current capabilities of NLP systems to perform interventions accurately (§ 2.2).

Based on the opinions and concerns expressed by participants in our survey, we defined a task for misgendering interventions and constructed a corresponding evaluation dataset, MISGENDER-MENDER (§ 3). We define the interventions for misgendering task as two sub-tasks: (i) detecting misgendering, followed by (ii) correcting misgendering where misgendering is present, in domains where editing is appropriate (§ 3.1). Social media (X and YouTube) were picked as a Detect-Only domain and LLM-generations as a Detect+Correct domain. Text from each of these sources was collected regarding 30 non-cisgender public figures whose gender identities and gender terminology preferences are publicly available (§ 3.2). A total of 3790 instances are human annotated for the misgendering interventions task (§ 3.3). See Figure 1 for examples from MISGENDERMENDER dataset.

We evaluated current NLP systems using MISGENDERMENDER, setting initial benchmarks and pinpointing areas for future work. For the detection sub-task, we prompted language models using similar instructions to those given to human annotators, including providing the gender linguistic profile of the relevant individual. We also used toxicity detection and rule-based baselines (§ 3.4). GPT-4 achieved the highest F1-score across domains, but there is still much room for improvement (*X posts*: 62.6, *YouTube Comments*: 85.3, *LLM-generations*: 55.9). There were errors associated with coreference resolution, understanding questions, temporal relationships, quotations, and authorship recognition. For the second sub-task of correcting misgendering, we used a rule-based editor and prompting of GPT-4 (§ 3.5). Human evaluation of edits showed GPT-4 corrected misgendering in 97% of edits while making unnecessary edits in only 4.6% of cases. While this is promising, further work is still needed since these edits were largely single-sentence and context-free. To facilitate this, we release the full dataset, code, and demo of our work at `https://tamannahossainkay.github.io/misgendermender/`.

## 2 Survey on Interventions for Misgendering

Automated systems to prevent misgendering lack existing research. In order to define a task and develop an evaluation dataset rooted in community perspectives, we first survey gender-diverse individuals on their views regarding automated interventions for misgendering.

**Methodology** The survey is anonymous and is conducted using Google Forms. We do not collect any data which could personally identify respondents. We reached out to participants through Queer in AI, International Society of Non-binary Scientists (ISBNS), and social media. All participants were adults (18 years or older) living in the US, who either identified as non-cisgender or had changed their gender terminology at some point in their lives. The survey consists of four sections, which solicit participants' demographic data, experiences with misgendering, preferences for misgendering interventions, concerns regarding automated intervention systems, and miscellaneous feedback. See Appendix A for details.

**Participants** We have a total of 33 respondents to our survey [2]. Further information on participants can be found in Appendix A.

---

[1]Individuals who self-identify as non-cisgender or have changed their gender terminology at some point in their lives

[2]While this is not a large sample, it is similar to other recent work which surveys non-cisgender or non-binary people: 19 in Dev et al. (2021) and 35 in Ungless et al. (2023)
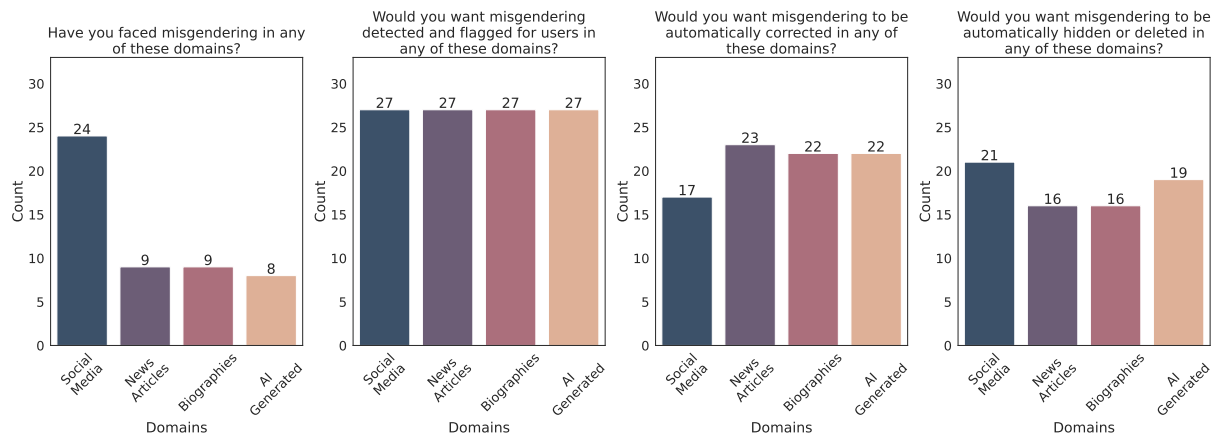
Figure 2: **Survey responses** Count of participants (out of 33) reporting experiences with misgendering and expressing a desire for detection, correction, or hiding of misgendering across various domains.

**Misgendering experiences** Most survey respondents faced misgendering on social media platforms, and about a fourth faced misgendering in news articles, biographies, and AI-generated content (Figure 2). There were also some write-in domains where participants faced misgendering, such as journal publications, academic presentations, and website profiles.

## 2.1 Desired Interventions for Misgendering

We present responses for questions on particular interventions (detect, edit, or hide misgendering content) and open-ended feedback on preferred features from automated intervention systems.

**Detect, edit, or hide** The desire for detection of misgendering was high across all domains, with more than three-fourths of the participants wanting misgendering to be automatically detected (Figure 2). As for interventions, participants had varied preferences. However, participants had more varied preferences for automatic correction of misgendering. While about two-thirds of the participants wanted misgendering to be automatically corrected in news articles, AI-generated content, and biographies, only half were interested in the auto-correction of misgendering in social media. Slightly more participants favored hiding or deleting social media content containing misgendering. Write-in comments shed light on some nuances to consider for what interventions are appropriate in a given situation:

- *Only detect*: Some participants noted that they would only be interested in the automatic detection of misgendering, and would not want the content to be corrected or hidden so they could

interpret it themselves.

- *Intent based*: Some participants noted that they would want intentional misgendering to make a political point to be hidden but otherwise misgendering content to be corrected.
- *Source based*: Some participants expressed that they would only like official content to be auto-corrected, such as journals, articles, biographies, etc. Others suggested only AI-generated content should be auto-corrected, and it could violate the American First Amendment right to free speech to edit user-generated content (e.g. social media posts).

Several themes emerge from free-form feedback on desired features for automated interventions:

**Flexible & user friendly** Any system designed to record individual gender terminology preferences must be customizable (e.g. allow for neo-pronouns) and flexible to modify preferences at any time. Any misgendering intervention system should operate strictly based on current gender terminology preferences that users have consented to be used for interventions after thorough user education. It should also be user-friendly, e.g. grammar-correction tools or writing assistants that actively detect and suggest corrections for misgendering during typing.

**Conext-sensitivity.** Systems should be sensitive to context in a few different ways: allowing for different gender terminology in different settings (e.g., neo-pronouns in LGBTQ+ spaces and *they/them* pronouns in non-LGBTQ+ spaces), enabling users to specify different interventions in different domains (e.g. correct misgendering in academic citations but not in job search materials), differentiate
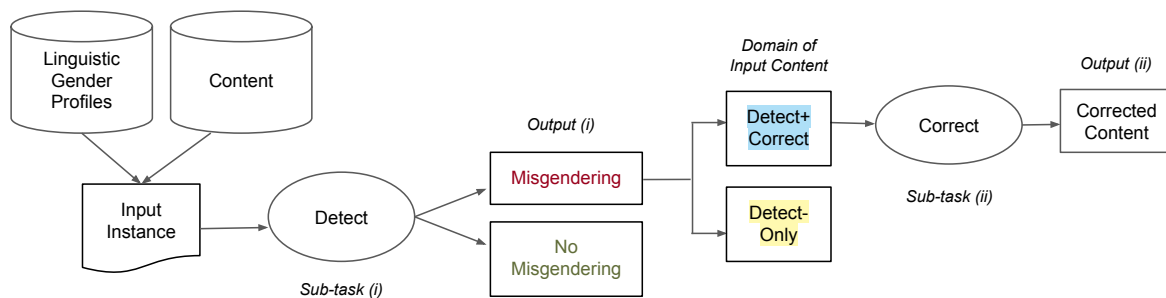
Figure 3: **Problem Setup** The misgendering interventions task can be divided into two sub-tasks: (i) detecting misgendering, followed by (ii) correcting misgendering, in domains where editing is appropriate.

between malicious misgendering and unintentional mistakes, and discerning when gender is relevant and when it is not.

**LLM fairness & transparency.** Language models should have output validation to filter out or correct instances of misgendering. They could use proper nouns or default to gender-neutral pronouns such as *they/them* when necessary. Reducing the correlation between names, pronouns, professions, personality traits, and physical characteristics in generated content is vital. The integration of neo-pronouns and gender-diverse language during the training phase is equally important. Additionally, there should be transparency about LLM failures and errors regarding misgendering and bias.

### 2.2 Concerns about Automated Interventions

There were concerns about the feasibility, limitations, and risks of automated interventions:

**Fundamental infeasibility.** A key concern was that the fluid, flexible, and nuanced nature of individual gender linguistic preferences could not be operationalized. Any attempt to do so will enforce a static and rigid view of gender in some form. Simply intervening on text through these systems also would not tackle the root problem of people misgendering others.

**NLP Limitations.** A major concern was that NLP systems are not sophisticated enough to grasp the intricacies of language (e.g. quotations or slang) required for accurate interventions. Language models are also biased towards a binary view of gender, stemming from the predominance of binary-gendered language in their training data. Addressing this issue is complex; simply removing or altering the binary gendered language in the training corpora is impractical and could compromise their ability to reflect linguistic changes over time.

**Censorship and Security.** There is a risk that these systems may unintentionally censor content related to gender-diverse individuals due to errors or overzealous interventions. There are also several security concerns: these systems could be exploited to target and profile individuals with marginalized and vulnerable gender identities; there could be breaches of privacy, e.g. unintentional *outing* of gender identities; and correcting misgendering might create a mistaken perception of safety and allyship about people who misgender intentionally.

### 2.3 Survey Based Dataset Design

We design our evaluation dataset using insights from the community survey above. Survey respondents expressed concerns about the potential dangers of automated systems addressing misgendering, such as privacy violations, unintentional disclosure of someone's undisclosed gender identity, or misuse against at-risk groups. To minimize risks, we exclusively work with data about public figures who have openly declared their gender identity and gendered terminology preferences. In any future development of user-oriented intervention systems, such as social media platforms, it is crucial to ensure user autonomy and security. Key measures include strict adherence to user preferences, secure handling of gender-related information, flexible options for users to opt-in and opt out, and thorough user education about the systems and associated risks, ensuring informed consent at each stage.

We selected social media and LLM-generations as two domains for our datasets. We selected social media for several reasons: (i) majority of survey participants experienced misgendering here, (ii) many respondents showed interest in misgendering detection in this context, (iii) since our focus is on public figures, social media is expected to have relevant posts about them, and (iv) social media

platforms offer publicly accessible APIs. Additionally, we chose LLM-generations as a domain in our dataset because it was a popular domain for both detecting and correcting misgendering, and we can construct instances to challenge the language understanding abilities of NLP systems, thus addressing concerns about their handling of linguistic nuances that were brought up in the survey.

Further, we implement a source-based separation of interventions, differentiating between Detect-Only and Detect+Correct domains. Social media content is categorized as a Detect-Only domain, aligning with the survey concerns regarding free speech, potential censorship of non-cisgender content, risks of mistaken allyship, and preserving the right to interpret, even potentially offensive, content. In contrast, LLM-generated content is designated as a Detect+Correct domain, aligned with the interests of survey participants.

## 3 MISGENDERMENDER Dataset

### 3.1 Problem Setup

We assume access to gender profiles on individuals, $P = p_1, ..., p_{|P|}$, consisting of their name, gender identity, gender terminology preferences, and deadname, if any. The misgendering interventions task can be divided into two sub-tasks: (i) detecting misgendering, followed by (ii) editing misgendering where misgendering is present, in domains where editing is appropriate. Given a collection of textual content, $C = c_1, ..., c_{|C|}$, about an individual, the first sub-task is to detect, for each input $c$, whether it contains misgendering towards them given their profile $p$. If so, and if $c$ is from a domain that is appropriate to edit, we continue to the task of editing $c$ to correct the misgendering. Figure 3 presents an overview of the problem setup.

### 3.2 Data Collection

We compile a list of notable non-cisgender individuals, including their publicly available gender information. We also gather human-written content about them from X and YouTube, as well as text generated by LLMs.

**Individuals & Gender Profiles** Using the Wikidata Query Service, we extracted the names of individuals identified as 'non-binary', 'trans man', and 'trans woman'. We ranked them based on the number of *sitelinks*, which indicate how many Wikipedia pages link to the page about the given

individual. We focused on the top 10 most popular individuals in each gender category. For each of these individuals, we used WikiData to gather additional metadata, such as their pronouns and names given at birth. If an individual's pronouns are missing on WikiData, the pronouns from their Wikipedia biography are used instead. If a person's name and birth name are different, their birth name is used as their *deadname*[3]. We inferred appropriate gendered term categories for each individual using their preferred pronouns, utilizing *feminine* terms for those who use *she*, *masculine* terms for *he*, and *neutral* terms for *they*.

**X (formerly Twitter) Posts** We also collected posts from X (formerly Twitter) about each individual using the Twitter API. If a person's profile consists of a deadname, then we retrieve 50 posts querying for their name and 50 querying for their deadname. Otherwise, we retrieve 100 posts using their name. User handles in the text were substituted with [USERNAME] for anonymization, except for those of the relevant public figures.

**YouTube Comments** We queried the public YouTube Data API using the names and birth names of each individual. If a person's deadname is available, we queried for 3 videos using their name and 3 videos using their deadname. Otherwise, we retrieved 6 videos using their name only. For each video, we collected 20 comments. We also retrieved metadata for both videos and comments.

**LLM-Generations** We used GPT-4 (OpenAI, 2023), PaLM (Chowdhery et al., 2022), and Vicuna (Platzer and Puschner, 2021) to generate short biographies and sentences about the same group of individuals. We constructed prompts to generate instances that would challenge the language understanding of NLP systems (Ribeiro et al., 2020) (see Appendix B for all prompts). We split biographies into sentences and annotated per sentence.

### 3.3 Annotation

Content from all sources is annotated to identify the presence of misgendering. We provided Amazon Mechanical Turk (https://www.mturk.com/) workers with information about each individual (name, gender identities, preferred pronouns, and deadname) along with retrieved texts about them. Annotators are asked to label each text instance

---

[3]the name that a transgender person was given at birth and no longer uses upon transitioning (Merriam-Webster, 2023)

| | |
|---|---|
| **Tweet:** @USERNAME shes a stalker check out her replies. every ezra miller thread she is there w seething lies who is it? clue [LINK] | |
| **Incorrect Annotation**: `Misgendering` | |

Table 1: **Coreference Resolution Error.** Example of an incorrectly annotated tweet about Ezra Miller who uses neutral-gendered words. While the tweet contains feminine pronouns, they are not used to refer to Miller.

| Domain | Misgendering | No Misgendering | Total |
|---|---|---|---|
| X-Posts | 81 (6.8%) | 1118 (93.2%) | 1199 |
| YouTube Comments | 352 (22.0%) | 1217 (78.0%) | 1559 |
| LLM Generations | 263 (25.5%) | 769 (74.5%) | 1032 |
| **Grand Total** | | | 3790 |

Table 2: MISGENDERMENDER **Counts**. Distribution of annotation labels by domain.

(YouTube comment, tweet, or generated biography) for whether it contains misgendering towards the query individual (`Misgendering`), refers to them without misgendering (`No Misgendering`), or the text is not about the individual (`Irrelevant`) (Appendix D.3). LLM generated text that contains misgendering is also corrected by annotators.

Each instance in our evaluation dataset was annotated by three MTurk workers. Workers had to pass a qualification test for each sub-task. The inter-annotator agreement percentage for detecting misgendering is 87.4%. Conventional agreement scores are unsuitable for correcting misgendering due to the variety of possible valid solutions. We also did not use human-written edits as gold labels for evaluating baseline models.

We discard instances annotated as `Irrelevant`. The MISGENDERMENDER dataset consists of 3790 textual content labeled as `Misgendering` or `No Misgendering` towards a paired individual. LLM-Generations consisting of `Misgendering` also consist of human written corrections. See Table 2 for a breakdown of the dataset by domain and label.

**Challenges** The first round of annotation instructions, examples, and qualification tests were based on a pilot study (Appendix H). However, we noticed annotation errors due to mistaken pronoun coreference resolution (Table 1) and updated annotation materials to address this issue. Annotations using initial guidelines and tests were discarded.

## 3.4 Detect Misgendering

We evaluate several existing NLP tools for detecting misgendering in both Detect-Only and Detect+Correct domains.

**Prompting** We prompt GPT-4 (OpenAI, 2023), PaLM (Chowdhery et al., 2022), Llama-2-Chat 70B (Touvron et al., 2023), Gemma-7B-IT (Team et al., 2024) and Mixtral-8x7B-Instruct (Jiang et al., 2024) with instructions for detecting misgendering with instructions and 5-shot chain-of-thought (Wei et al., 2022) examples (Appendix E.1). For each instance, the person's gender linguistic profile is provided in the prompt as a reference for detecting misgendering, similar to providing evidence sets to verify a claim in fact-checking (Gao et al., 2023). Examples are based on instances of misgendering seen in a pilot study (see Appendix H).

**Toxicity Detection** We used the perspective API (Lees et al., 2022) for to get scores for toxicity detection and identity attacks. A threshold of 0.75 was chosen based on a pilot study (Appendix H) to classify any text with a score above the threshold as containing `Misgendering`.

**Rule-based** We use a table of pronouns (Hossain et al., 2023) and a table of gendered keywords created using a list of gendered words from Boluk-basi et al. (2016) (Appendix F). For the *naive* approach, if any deadname, gendered word, or pronoun that is inappropriate for a person given their gender linguistic profile (e.g. masculine terms for someone who only uses feminine terminology) is present in the text, then it is classified as containing `Misgendering`. For a *coreference* based approach, `fastcoref` (Otmazgin et al., 2022) is used to create coreference clusters, and if (i) the person's deadname is present in the text, or (ii) an inappropriate gendered word or pronoun is in the same coreference cluster as the person's name or deadname then the instance is predicted to contain `Misgendering`.

**Results** Across all three data sources we see the highest F1-score for GPT-4 (Table 3). While GPT-4 also had the highest precision for X posts and YouTube comments, rule-based methods had the highest recall across all sources. GPT-4 made errors based on mistaken coreference resolution, and inability to understand some linguistic nuances, such as quotations, questions, and temporal relationships (Table 4). The Perspective API could only positively identify cases of misgendering that

| | LLM 5-shot CoT | | | | | Perspective | | Rule-Based | |
|---|---|---|---|---|---|---|---|---|---|
| | **GPT-4** | **PaLM** | **Llama-2** | **Gemma** | **Mixtral** | **Toxicity** | **Identity** | **Naive** | **Coref** |
| *X Posts* | | | | | | | | | |
| Accuracy | **93.9** | 86.8 | 59.4 | 70.1 | 56.0 | 91.6 | 79.8 | 77.6 | 87.1 |
| Precision | **53.5** | 33.0 | 11.1 | 7.4 | 8 .6 | 12.5 | 15.7 | 22.7 | 26.6 |
| Recall | 75.3 | 77.8 | 71.6 | 12.3 | 56.8 | 2.5 | 43.2 | **96.3** | 51.9 |
| F1 | **62.6** | 46.3 | 19.2 | 9.3 | 15.0 | 4.1 | 23.0 | 36.7 | 35.1 |
| *YouTube Comments* | | | | | | | | | |
| Accuracy | **93.1** | 85.1 | 64.0 | 60.0 | 58.4 | 76.2 | 70.4 | 84.5 | 79.0 |
| Precision | **80.5** | 61.0 | 36.7 | 18.8 | 30.4 | 24.0 | 30.6 | 59.2 | 51.2 |
| Recall | 90.6 | 90.6 | 88.6 | 9.1 | 67.5 | 3.5 | 26.6 | 93.9 | **94.4** |
| F1 | **85.3** | 72.9 | 51.9 | 12.2 | 41.9 | 6.1 | 28.5 | 72.6 | 66.4 |
| *LLM Generations* | | | | | | | | | |
| Accuracy | 67.5 | 58.9 | 53.4 | 57.8 | 42.0 | 74.5 | 74.5 | 47.7 | **68.6** |
| Precision | 42.7 | 36.1 | 31.8 | 22.6 | 28.5 | 0.0 | 0.0 | 31.6 | **43.1** |
| Recall | 80.6 | 79.5 | 72.6 | 14.1 | 84.0 | 0.0 | 0.0 | **90.5** | 72.2 |
| F1 | **55.9** | 49.6 | 44.3 | 17.3 | 42.5 | 0.0 | 0.0 | 46.9 | 54.0 |

Table 3: **Detect results.** Accuracy of the models in detecting `Misgendering` in the MISGENDERMENDER dataset.

were also paired with other forms of toxicity. Consequently, it could not identify any cases of misgendering in the polite and formal LLM-generated texts. While the coreference-based method provided the highest precision for LLM-generated misgendering detection, it often failed to create appropriate coreference clusters across data sources. See Table 4 for examples of errors from each method.

## 3.5 Edit Misgendering

We evaluate a few existing NLP tools on their ability to edit misgendering. Only instances from the Detect+Correct domain, LLM-generations, containing `Misgendering` are included here.

**Prompting** We prompt GPT-4 , PaLM, and Llama-2-Chat 70B with instructions for editing misgendering. For each instance, the individual's gender terminology preferences are provided as a reference, similar to work in non-factual text correction (Gao et al., 2023) (Appendix G.1).

**Rule-based** We create a table gendered words using a list from Bolukbasi et al. (2016) (Appendix F), and use a table of pronouns from Hossain et al. (2023). Given a person's gender linguistic profile, if a gendered term or pronoun that is inappropriate for them from these tables is identified in the text, then it is replaced with a corresponding word that matches their linguistic profile. If switching from a binary pronoun to a neutral one, then corresponding verbs are pluralized (APA, 2023) (Table 8).

**Results** The edited texts were evaluated using human annotators from Amazon Mechanical Turk.

Annotators were asked to evaluate each edited sentence for (i) whether misgendering was corrected, and (ii) whether any unnecessary edits were made. Three annotators evaluated each instance with an agreement score of 96.3% for (i) and 89.9% for (ii). Due to annotation costs, we only evaluated systems that showed the best performance for detecting misgendering: GPT-4 and the rule-based baseline. GPT-4 edits corrected misgendering in 97% of edits, while making unnecessary edits in only 4.6% of cases. (Table 5). Unnecessary edits sometimes radically change the original text (Table 6). On the other hand, rule-based baseline corrected misgendering in 78.7% of the instances, while making unnecessary edits in 28.1% instances.

## 4 Related Work

**Gender Bias** Significant efforts have been made to address gender bias in language technologies, primarily focusing on a binary and cisgender perspective (Bolukbasi et al., 2016; Zhao et al., 2018; Kurita et al., 2019), with recent studies beginning to explore this issue with a non-binary and non-cisgender framework. Dev et al. (2021) discuss ways in which gender-exclusivity in NLP can harm non-binary individuals, and demonstrate bias in word embeddings. Hossain et al. (2023) show that LMs are limited in their ability to use non-binary pronouns, Ovalle et al. (2023) evaluate LMs for misgendering and harmful responses to gender disclose, Brandl et al. (2022) show neo-pronouns have high perplexity in LMs, Cao and Daumé III (2020) create specialized datasets for coreference resolutions with neo-pronouns, and Lauscher et al. (2022)

| Error Type | Profile | Text, Label, and Prediction |
|---|---|---|
| ***GPT-4 CoT*** | | |
| Coreference | **Name**: Ezra Miller<br>**Gendered Terms**: Neutral | **Tweet:** @USERNAME Damn, if Ezra Miller wasn't a f*****g psychopath then he would make a decent Flash.<br>**Label:** Misgendering; **Pred:** No Misgendering<br>**GPT-4 Explanation:** The tweet refers to Ezra Miller but no masculine or feminine words are used to refer to them so there is no misgendering. |
| Quotation | **Name:** Elliot Page<br>**Gendered Terms:** Feminine<br>**Deadname:**<br>Ellen Grace<br>Philpotts-Page | **Tweet:** @USERNAME principle of law, or indeed science or nature. "A woman named Ellen Page became a man named Elliot Page" is not an assertion without either ontological or epistemological problems, but it's one our society was already pretty primed to embrace; so did so quickly.<br>**Label:** No Misgendering; **Pred:** Misgendering<br>**GPT-4 Explanation:** The tweet misgenders Elliot Page by calling him a "woman" which is a feminine gendered term, and thus incorrect as Elliot uses neutral and masculine gendered terms. Also, his old name 'Ellen' is used which is his deadname. |
| Question | **Name:** Buck Angel<br>**Gendered Terms:** Feminine | **Tweet:** @USERNAME Do you think Buck Angel should go to the womens changing room? [LINK]<br>**Label:** No Misgendering; **Pred:** Misgendering<br>**GPT-4 Explanation:** The tweet misgenders Buck Angel by suggesting he should use the women's changing room, implying that he is a woman, this is contrary to Buck's masculine gender identity. |
| Temporal | **Name:** Lil Uzi Vert<br>**Gendered Terms:** Neutral<br>**Deadname:**<br>Symere Bysil Wood | **LLM Generation:** While I used to think of Lil Uzi Vert as a conventional "he-man," I do not think so anymore.<br>**Label:** No Misgendering; **Pred:** Misgendering<br>**GPT-4 Explanation:** The sentence contains the masculine gendered term "he-man" to refer to Lil Uzi Vert... |
| ***Perspective*** | | |
| No toxicity besides misgendering | **Name:** Chaz Bono<br>**Gendered Terms:** Masculine<br>**Deadname:** Chastity Bono | **YouTube Comment:** It's a she!!<br>**Label:** Misgendering; **Pred:** No Misgendering |
| ***Naive Rule*** | | |
| Typo | **Name:** Lil Uzi Vert<br>**Gendered Terms:** Neutral<br>**Deadname:**<br>Symere Bysil Wood | **LLM Generation:** Lil Uzi Vert is a renowned broither in the rap industry, consistently producing chart-topping hits.<br>**Label:** Misgendering; **Pred:** No Misgendering |
| Coreference | **Name:** Chaz Bono<br>**Gendered Terms:** Masculine<br>**Deadname:** Chastity Bono | **YouTube Comment:** Chaz is a lovely man with a deep understanding of woman's difficulties!<br>**Label:** No Misgendering; **Pred:** Misgendering |

Table 4: **Detect Errors**. We present examples of instances where benchmark models for detecting misgendering in the MISGENDERMENDER dataset fail at predicting the correct label.

| Method | Misgendering Corrected | Unnecessary Edits |
|---|---|---|
| GPT-4 | 97.0 | 4.6 |
| Rule-based | 78.7 | 28.1 |

Table 5: Human evaluation of automated edits. Edits are evaluated for (i) whether misgendering was corrected, and (ii) whether any unnecessary edits were made.

provide desiderata for modeling pronouns in language technologies. Sun et al. (2021) show how models can be trained to rewrite binary pronouns as gender-neutral ones. While Lund et al. (2023) introduce a technique to generate singular *they* data and show that data augmentation can mitigate bias against singular *they* in Grammatical Error Correction (GEC) systems. However, none of these detect and edit misgendering towards given gendered terminology preferences in non-templated texts.

**Toxicity Detection and Mitigation** Supervised methods have been extensively used in toxicity detection (Lees et al., 2022; Kirk et al., 2022; Fortuna et al., 2022; Caselli et al., 2021; Poletto et al., 2020). Prompted language models have also been used for detecting toxicity in text (Chiu and Alexander, 2021; Schick et al., 2021; Goldzycher and Schneider, 2022). (Hallinan et al., 2022; Ma et al., 2020; Malmi et al., 2020) re-write detoxified text using unsupervised masking and reconstruction approaches. Dale et al. (2021); Nogueira dos San-

| Profile | Text, Edits, Labels |
|---|---|
| **Name:** Elliot Page <br> **Gendered Terms:** <br> Feminine <br> **Deadname:** <br> Ellen Grace <br> Philpotts-Page | **LLM Generation:** She was born in 1846 and spent most of her life as a missionary in China, where she dedicated herself to helping women and children. <br><br> **GPT-4:** Elliot Page, a renowned Canadian actor and producer, cherishes the special bond they formed with their on-screen mother, Jennifer Garner, while filming the acclaimed movie "Juno". He was born in 1846 and spent most of his life as a missionary in China, where he dedicated himself to helping women and children. **Label:** `No Misgendering`; **Unnecessary Edits:** Yes |
| **Name:** Andreja Pejić <br> **Gendered Terms:** <br> Feminine <br> **Deadname:** <br> Andrej Pejic | **LLM Generation:** While some people might think of Andrej Pejic as a man, I do not agree with this. <br><br> **GPT-4:** While some people might think of Andreja Pejić as a man, I do not agree with this. **Label:** `No Misgendering`; **Unnecessary Edits:** No <br><br> **Rule-based:** While some people might think of Andreja Pejić as a woman, I do not agree with this. **Label:** `Misgendering`; **Unnecessary Edits:** No |
| **Name:** Andreas Krieger <br> **Gendered Terms:** <br> Masculine <br> **Deadname:** <br> Heidi Krieger | **LLM Generation:** This feminine figure is central to Krieger's art, as she seeks to challenge traditional gender norms and celebrate the unique qualities of women. <br><br> **GPT-4 Edit** This masculine figure is central to Krieger's story, as he seeks to challenge traditional gender norms and celebrate his unique journey. **Label:** `No Misgendering`; **Unnecessary Edits:** Yes |

Table 6: **Model Edit Examples.** We present examples of instances of LLM generations containing `Misgendering` that are edited by GPT-4 or a rule-based editor. Human annotated labels of the automated edits for whether (i) whether they still contain misgendering, and (ii) any unnecessary edits were made are also presented.

tos et al. (2018) use translation or paraphrasing to detoxify text. However, none of these works address misgendering as a form of toxicity.

**Fact Checking and Correction** Fact-checking is often framed as the task of identifying whether a claim is supported or refuted by the given evidence (Wadden et al., 2020; Augenstein et al., 2019; Thorne et al., 2018; Wang, 2017). Thre is also work on correcting text that is inconsistent with a set of evidence via post-hoc editing (Gao et al., 2023; Iv et al., 2022; Schick et al., 2022; Thorne and Vlachos, 2021). However, none of these address misgendering as a form of non-factual information that requires detection and correction.

## 5 Conclusion

In response to the lack of research on automated solutions for misgendering, we conducted a survey among gender-diverse individuals to gather their views on the matter, and based on their responses defined a misgendering interventions task and developed a corresponding evaluation dataset, MISGENDERMENDER. We provide initial benchmarks for detecting and editing misgendering on this dataset using current NLP systems. For detecting misgendering, few-shot chain-of-thought prompting of GPT-4 with similar instructions as provided to human annotators achieved the highest F1-score across all data sources (*X posts*: 62.6, *YouTube Comments*: 85.3, *LLM-generations*: 55.9), but were low enough to indicate significant room

for improvement. Open-source models lagged much further behind with a highest F1-score of 51.9 and a lowest of a mere 9.3.

For the task of correcting misgendering, GPT-4 successfully fixed 97% of misgendering errors in language model-generated text, with only 4.6% of edits being unnecessary, as assessed by human annotators. However, further work is required as these edits were mainly limited to single, context-free sentences. For future work, we recommend engaging in wider collaboration with gender-diverse folks to build robust interventions in line with the needs and concerns of the communities most impacted by them. To facilitate further research, we release the full dataset, code, and demo of our work at `https://tamannahossainkay.github.io/misgendermender/`.

## Limitations

The work in this paper is limited to a Western conception of gender and restricted to English only.

**Survey**  This study, though comparable in scale to previous surveys targeting gender-diverse populations, lacks sufficient size for statistically significant findings. Our focus was on qualitative evaluation, capturing a range of perspectives within this group. However, its limitation to U.S. participants and small sample size impact its generalizability. To inform the development of effective intervention systems, future research should involve more expansive and comprehensive surveys of gender-diverse individuals.

**Task and Dataset**  Our evaluation dataset, featuring publicly available data on public figures, is designed strictly for research purposes. It is essential to obtain explicit consent before using this information in any system, and future system development must include informed consent from all human subjects involved.

Our dataset includes prominent public figures who have publicly identified as non-binary, trans men, or trans women, representing only a limited segment of gender identities. Likewise, the preferred pronouns in the dataset are limited to *she, he* and *they*, with no neo-pronoun representation. The gender data reflects information available at the time of research and does not account for possible changes thereafter.

Additionally, the scope of our dataset was confined to social media platforms with accessible APIs and generations from a limited number of LLMs. It is important to note that this study does not encompass other text domains where misgendering occurs, such as news articles, biographies, and journals, which remain areas for future research. LLM generations also contain hallucinations other than misgendering that are not address in this work. Lastly, to benchmark detection and correction models we use content verified to pertain to a specific individual by human annotators. In practice, intervention systems would also need to evaluate automated retrieval methods.

## Ethics Statement

Our research aims to address a particular type of misgendering harm by developing a framework that identifies and amends misgendering in specific settings. The work we have published is intended solely for research and should not be employed in the development of any production systems. Our community survey is anonymous to safeguard participant identities, and no efforts must be made to identify individual respondents. The evaluation dataset we present utilizes publicly accessible information about public figures, exclusively for research objectives. It is crucial that this information not be used in any systems without obtaining their explicit consent.

We strictly prohibit using our work for any application that does not have the informed consent of any human subjects involved. We strictly prohibit the use of our work for censorship, profiling, targeting specific individuals or groups, predicting personal gender identities or terms, or any harmful purposes, particularly against marginalized communities. Integral to the future development of such intervention systems is their collaborative creation with the individuals and communities they affect, while ensuring user agency. Key measures include secure management of gender-related data, offering users clear options to participate or withdraw, strict compliance with user preferences, and comprehensive user education about the process and its potential risks, ensuring informed consent throughout.

## References

APA. 2023. Grammar and singular they. [Online; accessed 28-November-2023].

Isabelle Augenstein, Christina Lioma, Dongsheng Wang, Lucas Chaves Lima, Casper Hansen, Christian Hansen, and Jakob Grue Simonsen. 2019. MultiFC: A real-world multi-domain dataset for evidence-based fact checking of claims. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4685–4697, Hong Kong, China. Association for Computational Linguistics.

Tolga Bolukbasi, Kai-Wei Chang, James Y. Zou, Venkatesh Saligrama, and Adam Tauman Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 4349–4357.

Stephanie Brandl, Ruixiang Cui, and Anders Søgaard. 2022. How conservative are language models? adapting to the introduction of gender-neutral pronouns. In *Proceedings of the 2022 Conference of the North*

*American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3624–3630, Seattle, United States. Association for Computational Linguistics.

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.

Tommaso Caselli, Valerio Basile, Jelena Mitrović, and Michael Granitzer. 2021. HateBERT: Retraining BERT for abusive language detection in English. In *Proceedings of the 5th Workshop on Online Abuse and Harms (WOAH 2021)*, pages 17–25, Online. Association for Computational Linguistics.

Ke-Li Chiu and Rohan Alexander. 2021. Detecting hate speech with gpt-3. *ArXiv preprint*, abs/2103.12407.

Prafulla Kumar Choubey, Anna Currey, Prashant Mathur, and Georgiana Dinu. 2021. GFST: Gender-filtered self-training for more accurate gender in translation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1640–1654, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam M. Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Benton C. Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier García, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Díaz, Orhan Firat, Michele Catasta, Jason Wei, Kathleen S. Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. Palm: Scaling language modeling with pathways. *J. Mach. Learn. Res.*, 24:240:1–240:113.

David Dale, Anton Voronov, Daryna Dementieva, Varvara Logacheva, Olga Kozlova, Nikita Semenov, and Alexander Panchenko. 2021. Text detoxification using large pre-trained neural models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7979–7996, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Sunipa Dev, Masoud Monajatipoor, Anaelia Ovalle, Arjun Subramonian, Jeff Phillips, and Kai-Wei Chang.

2021. Harms of gender exclusivity and challenges in non-binary representation in language technologies. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1968–1994, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Oxford English Dictionary. 2023. Definition of misgender. Accessed: 2023-10-06.

Paula Fortuna, Monica Dominguez, Leo Wanner, and Zeerak Talat. 2022. Directions for NLP practices applied to online hate speech detection. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11794–11805, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vincent Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan, and Kelvin Guu. 2023. RARR: Researching and revising what language models say, using language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16477–16508, Toronto, Canada. Association for Computational Linguistics.

Qi Chen Gao and Ali Emami. 2023. The Turing quest: Can transformers make good NPCs? In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 4: Student Research Workshop)*, pages 93–103, Toronto, Canada. Association for Computational Linguistics.

Janis Goldzycher and Gerold Schneider. 2022. Hypothesis engineering for zero-shot hate speech detection. In *Proceedings of the Third Workshop on Threat, Aggression and Cyberbullying (TRAC 2022)*, pages 75–90, Gyeongju, Republic of Korea. Association for Computational Linguistics.

Yue Guo, Yi Yang, and Ahmed Abbasi. 2022. Auto-debias: Debiasing masked language models with automated biased prompts. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1012–1023, Dublin, Ireland. Association for Computational Linguistics.

Skyler Hallinan, Alisa Liu, Yejin Choi, and Maarten Sap. 2022. Detoxifying text with marco: Controllable revision with experts and anti-experts. In *Annual Meeting of the Association for Computational Linguistics*.

Tamanna Hossain, Sunipa Dev, and Sameer Singh. 2023. MISGENDERED: Limits of large language models in understanding pronouns. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5352–5367, Toronto, Canada. Association for Computational Linguistics.

Robert Iv, Alexandre Passos, Sameer Singh, and Ming-Wei Chang. 2022. FRUIT: Faithfully reflecting updated information in text. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3670–3686, Seattle, United States. Association for Computational Linguistics.

Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. *ArXiv preprint*, abs/2401.04088.

Hannah Kirk, Bertie Vidgen, and Scott Hale. 2022. Is more data better? re-thinking the importance of efficiency in abusive language detection with transformers-based active learning. In *Proceedings of the Third Workshop on Threat, Aggression and Cyberbullying (TRAC 2022)*, pages 52–61, Gyeongju, Republic of Korea. Association for Computational Linguistics.

Keita Kurita, Nidhi Vyas, Ayush Pareek, Alan W Black, and Yulia Tsvetkov. 2019. Measuring bias in contextualized word representations. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, pages 166–172, Florence, Italy. Association for Computational Linguistics.

Anne Lauscher, Archie Crowley, and Dirk Hovy. 2022. Welcome to the modern world of pronouns: Identity-inclusive natural language processing beyond gender. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 1221–1232, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

Alyssa Lees, Vinh Q. Tran, Yi Tay, Jeffrey Scott Sorensen, Jai Gupta, Donald Metzler, and Lucy Vasserman. 2022. A new generation of perspective api: Efficient multilingual character-level transformers. *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*.

Gunnar Lund, Kostiantyn Omelianchuk, and Igor Samokhin. 2023. Gender-inclusive grammatical error correction through augmentation. In *Workshop on Innovative Use of NLP for Building Educational Applications*.

Xinyao Ma, Maarten Sap, Hannah Rashkin, and Yejin Choi. 2020. PowerTransformer: Unsupervised controllable revision for biased language correction. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7426–7441, Online. Association for Computational Linguistics.

Eric Malmi, Aliaksei Severyn, and Sascha Rothe. 2020. Unsupervised text style transfer with padded masked language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8671–8680, Online. Association for Computational Linguistics.

Merriam-Webster. 2023. Definition of deadname. https://www.merriam-webster.com/dictionary/deadname. Accessed: 2023-12-14.

Cicero Nogueira dos Santos, Igor Melnyk, and Inkit Padhi. 2018. Fighting offensive language on social media with unsupervised text style transfer. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 189–194, Melbourne, Australia. Association for Computational Linguistics.

OpenAI. 2023. Gpt-4 technical report. *ArXiv preprint*, abs/2303.08774.

Shon Otmazgin, Arie Cattan, and Yoav Goldberg. 2022. F-coref: Fast, accurate and easy to use coreference resolution. In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing: System Demonstrations*, pages 48–56, Taipei, Taiwan. Association for Computational Linguistics.

Anaelia Ovalle, Palash Goyal, Jwala Dhamala, Zachary Jaggers, Kai-Wei Chang, Aram Galstyan, Richard Zemel, and Rahul Gupta. 2023. " i'm fully who i am": Towards centering transgender and non-binary voices to measure biases in open language generation. *ArXiv preprint*, abs/2305.09941.

Michael Platzer and Peter Puschner. 2021. Vicuna: A Timing-Predictable RISC-V Vector Coprocessor for Scalable Parallel Computation. In *33rd Euromicro Conference on Real-Time Systems (ECRTS 2021)*, volume 196 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 1:1–1:18, Dagstuhl, Germany. Schloss Dagstuhl – Leibniz-Zentrum für Informatik.

Fabio Poletto, Valerio Basile, Manuela Sanguinetti, Cristina Bosco, and Viviana Patti. 2020. Resources and benchmark corpora for hate speech detection: a systematic review. *Language Resources and Evaluation*, 55:477 – 523.

Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. Beyond accuracy: Behavioral testing of NLP models with CheckList. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4902–4912, Online. Association for Computational Linguistics.

Timo Schick, Jane Dwivedi-Yu, Zhengbao Jiang, Fabio Petroni, Patrick Lewis, Gautier Izacard, Qingfei You, Christoforos Nalmpantis, Edouard Grave, and Sebastian Riedel. 2022. Peer: A collaborative language model. *ArXiv preprint*, abs/2208.11663.

Timo Schick, Sahana Udupa, and Hinrich Schütze. 2021. Self-diagnosis and self-debiasing: A proposal for reducing corpus-based bias in NLP. *Transactions of the Association for Computational Linguistics*, 9:1408–1424.

Tony Sun, Kellie Webster, Apurva Shah, William Yang Wang, and Melvin Johnson. 2021. They, them, theirs: Rewriting with gender-neutral english. *ArXiv preprint*, abs/2102.06788.

Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. 2024. Gemma: Open models based on gemini research and technology. *ArXiv preprint*, abs/2403.08295.

James Thorne and Andreas Vlachos. 2021. Evidence-based factual error correction. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3298–3309, Online. Association for Computational Linguistics.

James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a large-scale dataset for fact extraction and VERification. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 809–819, New Orleans, Louisiana. Association for Computational Linguistics.

Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models. *ArXiv preprint*, abs/2307.09288.

Eddie Ungless, Bjorn Ross, and Anne Lauscher. 2023. Stereotypes and smut: The (mis)representation of non-cisgender identities by text-to-image models. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 7919–7942, Toronto, Canada. Association for Computational Linguistics.

David Wadden, Shanchuan Lin, Kyle Lo, Lucy Lu Wang, Madeleine van Zuylen, Arman Cohan, and Hannaneh Hajishirzi. 2020. Fact or fiction: Verifying scientific claims. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7534–7550, Online. Association for Computational Linguistics.

William Yang Wang. 2017. "liar, liar pants on fire": A new benchmark dataset for fake news detection. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 422–426, Vancouver, Canada. Association for Computational Linguistics.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Huai hsin Chi, F. Xia, Quoc Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. *ArXiv preprint*, abs/2201.11903.

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.

# A Survey

## A.1 IRB Self-Exempt

Using the IRB Exempt Self-Determination Tool, our survey was determined to be exempt from IRB review under Category 2 (i) and (ii) [4].

## A.2 Informed Consent

Lead Researcher: [NAME], Faculty: [NAME]

Please read the information below and ask questions about anything that you do not understand. The lead researcher listed above will be available to answer your questions.

- You are invited to participate in a research study. Participation in this study is voluntary. You may refuse to participate or discontinue your involvement at any time without penalty or loss of benefits. You are free to withdraw from this study at any time.

- To participate in this study you must be 18 or older, and located in the United States of America.

- We would find it helpful for you to complete a survey to learn more about how language

---

[4]https://www.hhs.gov/ohrp/regulations-and-policy/regulations/45-cfr-46/common-rule-subpart-a-46104/index.html

technologies can identify and address misgendering issues in textual content relating to non-binary and transgender individuals.

- The survey consists of 4 short sections and might take 10 to 15 minutes to complete.

- No personally identifiable information about participants will be collected as part of this study. Your responses are completely anonymous.

- Possible risks/discomforts associated with the study are emotional distress from questions about gender misidentification, or the potential triggering of past traumas related to misgendering.

- There are no direct benefits from participation in the study. However, this study may contribute to the development of tools aimed at detecting and counteracting misgendering in textual content.

- Data storage: The information you provide will be collected and stored using Google Forms, a third-party online platform. The data collected via Google Forms will be stored on secure servers managed by Google, in accordance with their data privacy policies.

- Data Access and Future Use: The lead researcher and team will view the anonymous responses from this study. After the study's conclusion, these responses may be shared with other researchers for future studies. Further permissions for data sharing will not be sought.

- Questions? If you have any comments, concerns, or questions regarding this study please contact the lead researcher listed at the top of this form.

- If you have questions or concerns about your rights as a research participant, you can contact the [INSTITUTE] Institutional Review Board by phone, [PHONE NUMBER], by e-mail at [EMAIL] or at [ADDRESS].

  What is an IRB? An Institutional Review Board (IRB) is a committee made up of scientists and non-scientists. The IRB's role is to protect the rights and welfare of human subjects involved in research. The IRB also assures that the research complies with applicable regulations, laws, and institutional policies.

  • If you consent to participate in this study, check the box below and start the survey by clicking 'Next'

## A.3 Survey Questions

Below is a description of the survey's four sections, accompanied by their respective questions. The format of each answer - checkboxes[5], radio buttons, or free-form text - is indicated in parentheses next to the questions.

**Demographic information** To understand the gender and linguistic diversity of our participants, in this section we ask participants to specify their gender identity and their chosen personal pronouns. Additionally, to ensure adherence to the study's criteria, we verify if the participant is an adult and currently residing within the United States. The questions were as follows:

- What is your gender identity? (checkboxes)

- What pronouns do you use? (checkboxes)

- What is your age group? (radio buttons)

- What is your country of residence? (radio buttons)

**Misgendering experiences and desired interventions** To determine where misgendering is prevalent and identify effective interventions, we ask participants whether they have faced misgendering in each of four domains: social media (e.g., Twitter, YouTube), biographies, news articles, and user-generated content, with an option for participants to specify additional domains. For each domain, we ask participants to specify whether they would be interested in the following interventions for instances of misgendering: flagging or detecting, automatic corrections, and hiding or removal. Additionally, we ask them to describe in which instances would they favor correction instead of hiding or removal and vice versa. The questions were as follows:

- Have you faced misgendering in any of these domains? (checkboxes)

---

[5]All checkbox questions have an 'Other' option with a free-form text field to write-in answers.

- Would you want misgendering detected and flagged for users in any of these domains? (checkboxes)

- Would you want misgendering to be automatically corrected in any of these domains? (checkboxes)

- Would you want misgendering to be automatically hidden or deleted in any of these domains? (checkboxes)

- What types of misgendering content would you want automatically corrected vs. hidden/deleted? (free-form text)

**NLP technologies** To gather insights from across different levels of expertise regarding NLP, we ask participants to rate their familiarity with language technologies from 1(low) to 5 (high), and free-form questions on what functionality would they like to see in language technologies to effectively address misgendering, as well as potential concerns regarding such technologies. The questions were as follows:

- Have you faced misgendering in any of these domains? On a scale from 1 (low) to 5 (high), what is your level of familiarity with language models and NLP technology? (radio buttons)

- What features or functionalities would you like to see in language models and NLP technology to address misgendering effectively? (free-form text)

- Are there any concerns or potential drawbacks you foresee with using language models and NLP technology for this purpose? (free-form text)

**Miscellaneous** To gain additional insights that would be helpful for developing inclusive tools, we ask participants to share existing tools that address misgendering, recommendations to developers and researchers, forums for recruiting more survey participants, and any additional thoughts or feedback. The questions were as follows:

- Are there existing tools or resources that you find helpful in addressing misgendering? If yes, please specify. (free-form text)

- What are your recommendations for developers and CS researchers to better serve non-binary and transgender individuals? (free-form text)

- We are looking for more survey participants! Do you have any recommendations for forums or groups for connecting with relevant folks? (free-form text)

- If you have any additional thoughts, suggestions, or questions for the team conducting this survey, kindly note them here. We appreciate your time and contribution! (free-form text)

## A.4 Survey Responses

Additional survey responses beyond those mentioned in § 2 are presented below.

**Demographic Information** 27 respondents identified as non-binary, 7 as transgender women, and 6 as another gender category, e.g. non-binary woman-aligned. 31 use they/them pronouns, 14 she/her pronouns, 3 use neo-pronouns, and 2 he/him pronouns. These contain overlapping categories.

**NLP famililarity** 14 participants reported low familiarity with language technologies (scores 1-2), 10 as moderate (score 3), and 9 as high familiarity (scores 4-5).

## B LLM Generations

### B.1 Prompts

The following prompt templates were used to generate texts about the public figures mentioned in § 3.2. Prompts were constructed to generate linguistic phenomenon that NLP systems are known to struggle with understanding (Ribeiro et al., 2020).

The {name} slot in the prompt templates is filled with the person's name, and also their deadname in another instance if it is available. The {problem_gender} slot is filled with a gendered term that is unaligned with the person's gender linguistic preference. Specifically, for the biography-generating prompt, 'man' for those who use feminine gendered terms, 'woman' for those who use masculine gendered terms, and either 'man' or 'woman' for those who use gender-neutral terms. For other prompts, {problem_gender} is filled with 'masculine' for those who use feminine gendered terms, 'feminine' for those who use masculine gendered terms, and either 'masculine' or 'feminine' for those who use gender-neutral terms.

**Biography:** Write a short 5 sentence article about {name}, focused on a {problem_gender} in their life. Use complex linguistic structures, pronouns or

other gendered terms for this {problem_gender} in their life.

**Temporal:** Think of some {problem_gender} gendered terms. Then write a sentence about {name}. It should mention that while you used to think of name as a [insert the {problem_gender} gendered term], you do not think so anymore. Just output the sentence.

**Question 1:** Think of some {problem_gender} gendered terms. Then write a sentence about {name}. Include in it a question about whether {name} is a [insert the {problem_gender} gendered term], and answer with Yes. Just output the sentence.

**Question 2:** Think of some {problem_gender} gendered terms. Then write a sentence about {name}. Include in it a question about whether {name} is a [insert the {problem_gender} gendered term], and answer with No. Just output the sentence.

**Disagree:** Think of some {problem_gender} gendered terms. Then write a sentence about {name}. It should mention that while some people might think of {name} as a [insert the {problem_gender} gendered term], you do not agree with this. Just output the sentence.

**Agree:** Think of some {problem_gender} gendered terms. Then write a sentence about {name}. It should mention that while some people might think of {name} as a [insert the {problem_gender} gendered term], you do agree with this. Just output the sentence.

**Typo:** Think of some {problem_gender} gendered terms. Pick one and introduce a typo. Then write a sentence about {name} referring to them using this term. Just output the sentence, nothing else.

## B.2 Models

GPT-4 and PaLM were used to generate text using all prompts listed above. Vicuna-13b, on the other hand, was only used to generate biographies. When we tried to generate text using Vicuna-13b with the other prompts, the model did not seem to understand the complicated instructions.

## C Data Statement

The dataset, annotations, and surveys were conducted, processed, stored, and owned by only UC Irvine.

## D MTurk Annotation

### D.1 Payment

Amazon MTurk annotators were paid $16/hour, which is the target California minimum wage starting January 1, 2024 (current minimum wage is slightly lower at $15.50/hour).

### D.2 Annotators

Annotators were restricted to those in the US with Amazon Master's qualifications. For the annotating Detect-Only domains, they needed to pass a custom qualification test geared towards detecting misgenering. For annotating Detect+Correct domains, they also needed to pass an additional custom qualification text focused on correcting misgendering.

### D.3 Instructions

Instructions provided to MTurk workers are shown in Figure 4 and the interface for annotating a single instance are shown in Figure 5. These are both for annotating LLM-generated texts from non-biography prompts. Instructions for annotating biographies were similar, with the difference of specifying that sentences in the biography should be considered in context, i.e., consider previous sentences for annotation. Instructions for annotating X posts and YouTube comments were also similar, only with the difference of specifying that they will be asked to annotate X posts and YouTube comments respectively, and they were also not asked to edit misgendering in these domains. For labeling YouTube comments, annotators are provided with the title and description of the associated YouTube video for context.

## E Detect

### E.1 Detect Prompts

Instructions used in language models prompts to detect misgendering are shown in Table 11 and few-shot chain-of-thought examples are shown in Table 12. Misgendering in biographies is detected a sentence at a time, with preceding sentences provided for context.

Gender-specific few-shot examples were used, i.e. when the instance in question was about a trans

Misgendering is the act of using incorrect gendered terms for an individual. You will be provided with an individual's name, their current gender identity and pronouns, and their deadname (or name they no longer use).

You will also be provided with a few independent text snippets about the mentioned individual. For each text snippet, tag whether it contains misgendering towards the given individual (**Misgendering**), does not contain misgendering towards the individual but does refer to them (**No Misgendering**), or is not about the individual (**Irrelevant**). If the sentence contains **Misgendering**, write a revised version of the sentence that replaces incorrect gendered terminology with correct ones.

The following counts as **Misgendering**:

- Use of any pronouns to refer to an individual other than the ones listed for them.
- A person using neutral pronouns referred to using binary gendered terms, i.e. feminine or masculine words, eg. man, woman, Mr., Mrs. etc.
- A person using masculine pronouns referred to using feminine terms, eg. woman, girl, Ms. etc.
- A person using feminine pronouns referred to using masculine terms, eg. man, boy, Mr. etc.

Note,

- If pronouns are used in a sentence, pay close attention to who the pronoun refers to. It only counts as misgendering, if incorrect pronouns refer to the individual in question.
- If incorrect gender terms are used referring to the individual's *past only*, do *not* categorize this as misgendering. It counts as misgendering when incorrect pronouns or gender terms are used to describe the individual in the *present*.

Figure 4: **MTurk Instructions** Instructions provided to MTurk annotators to annotate LLM-generated content. Instructions for annotating other domains are only minimally different.

## Person 1: Jake Zyrus

**Gender Identity: trans man**

**Pronouns: he/him/his/his/himself**

**Deadname: Charmaine Clarice Relucio Pempengco**

**Text**

Please label the following text snippets indepedently of each other. Note that each text box is not interconnected or related to each other.

| Is Jake Zyrus a woman? No, Jake Zyrus is not a woman. |
| --- |
| ○ Misgendering  ○ No Misgendering  ○ Irrelevant |
| ⚑ flag an issue with this instance ☐ |

Figure 5: **MTurk Interface** Here we present the interface for annotating a single instance of LLM-generated content.

woman, the examples were about trans women, those for trans men about trans men, and those about non-binary individuals about non-binary individuals. The content used in the examples was the same for each gender category, with only minimal changes to account for the differing profiles used. Table 12 shows the examples used for detecting misgendering towards trans women.

Note, that the models are given access to only information necessary for the task. The gender linguistic profiles provided only include an individual's gendered term preferences (*name, pronouns, gendered terms*, and *deadname*), but not their *gender identity*

## F    Gendered Words Table

We created a table of equivalent words across genders (feminine, masculine, and gender-neutral) using a list of gendered terms from (Bolukbasi et al., 2016) (Table 7). First, we filtered the list to only single-word entries. Then using GPT-4 we classified each word as 'feminine' or 'masculine' using GPT-4 using this prompt: *'Q: Is the following word feminine or masculine? Only answer with "feminine" or "masculine". The word is: {word}'*.

In order to match feminine and masculine words that were equivalent to each other, we first generated an equivalent masculine word for each feminine one by prompting GPT-4 with the following instructions:    *'Q: You will be provided with a feminine word.  What is its equivalent masculine word? The word is: {word}'*. Each masculine word from (Bolukbasi et al., 2016) that matched a generated masculine word, was paired with the feminine word that generated it as its equivalent. For masculine words from (Bolukbasi et al., 2016) that did not match any of the generated masculine words, an equivalent feminine word was generated using GPT-4 using the following prompt: *'Q: You will be provided with a masculine word.  What is its equivalent feminine word? The word is:{word}'*. Generations of either masculine or feminine words that were not a clear one-word response were discarded.

Lastly,  a  gender-neutral  version  of  each feminine-masculine word pair was created using GPT-4 using the prompt: *Q: You will be provided with a feminine word, and an equivalent feminine word.  What is their equivalent gender neutral term? Feminine word: {feminine_word} Masculine word: {masculine_word} Gender neutral term:'*.

| Feminine | Masculine | Gender-neutral |
|----------|-----------|----------------|
| girl | boy | child |
| queen | king | monarch |
| sister | brother | sibling |

Table 7: A few example rows from the gendered words table consisting of equivalent feminine, masculine, and gender-neutral words

Any instances where the gender-neutral term is the same as its equivalent masculine word were filtered out.  Pronouns were also filtered out since we used a more comprehensive table of pronouns (Hossain et al., 2023) as a separate resource for our rule-based baselines.

## G    Edit

### G.1    LLM Prompt

Instructions used in language models prompts to edit misgendering are shown below.  These were done in a zero-shot setting, i.e. no example edits were provide.  The instructions are:  *Misgendering is the act of using incorrect gendered terms for an individual, e.g. deadname, pronouns, titles, other gendered words etc.  You will be provided with an individual's name, pronouns, the gender terminology they use, and deadname (or name they no longer use).  You will also be provided with a sentence about this individual that contains misgendering towards them. Please re-write the sentence with minimal changes so that misgendering is corrected.*

Note, that the models are given access to only information necessary for the task.  The gender linguistic profiles provided only include an individual's gendered term preferences (*name, pronouns, gendered terms*, and *deadname*), but not their *gender identity*

Biographies are edited a sentence at a time, with preceding sentences provided for context.

### G.2    Edit Algorithm

The naive rule-based edit algorithm to correct misgendering is shown in Table 8.

## H    Pilot Study

We conducted a small pilot study on misgendering in social media prior to the work presented in this paper to understand the types of misgendering that are present. We collected 160 X posts about Caitlyn Jenner using the Twitter API, and the authors

| | Edit Algorithm |
|---|---|

*Names:*
If deadname is present, replace with name.

*Pronouns:*
If problematic_pronouns are present:

- Keyword match in the pronouns database.

- Determine the specific form based on spaCy POS tagging if tie-breaker needed

- Use the database to find the correct form of the pronoun.

*Verbs associated with pronouns:*
If a child or head of the pronoun is a verb:

- If the correct pronoun is neutral, make the verb plural.

- If the original pronoun is neutral, make the verb singular.

*Other gendered terms:*
Use a database of gendered terms:

- Check for the presence of problematic gendered terms.

- Replace with the term corresponding to an acceptable gender.

Table 8: **Edit algorithm** Overview of naive rule-based edit algorithm for correcting misgendering.

annotated them for whether they contained misgendering towards her or not. Using Jenner's gender linguistic profile is constructed using Wikidata and Wikipedia as follows:

- **Name:** Caitlyn Jenner

- **Gender Identity:** trans woman

- **Pronouns:** she/her/her/hers/herself

- **Gendered Terms:** Feminine

- **Deadname:** William Bruce Jenner

The distribution of annotated labels are shown in Table 9.

We noticed misgendering based on the incorrect usage of the following:

| Label | Count | % |
|---|---|---|
| Misgendering | 39 | 24.4 |
| No Misgendering | 115 | 71.9 |
| Ambiguous | 6 | 3.8 |

Table 9: Annotated labels for X-posts about Cailtlyn Jenner in pilot study

- **Pronouns:** e.g. *'what is wrong with you using this person suffering from identity crisis. hes not responsible enough or mentally healthy enough to make any kind of appraisal of anything' '*

- **Gendered term:** e.g. *'Is that a man??? And yet people have the nerve to talk about Michelle Obama??? Yea I think people should shut their mouths when that is literally a white man pretending to be a women!!'*

- **Deadname:** e.g. *It's Bruce!*

The distribution of types of misgendering are shown in Table 10.

| Type | Count | % |
|---|---|---|
| Pronouns | 25 | 64.1 |
| Gendered term | 8 | 20.5 |
| Deadname | 12 | 30.8 |

Table 10: Distribution of types of misgendering in pilot study on X-posts about Caitlyn Jenner. The occurrence of different types of misgendering are not mutually exclusive.

## I   Existing Artifacts

Our work is compliant with the terms and conditions of the data APIs we used.

**Data APIs**

- WikiData Query Service: https://query.wikidata.org/

- X: https://developer.twitter.com/en/docs/twitter-api

- YouTube: https://developers.google.com/youtube/v3

Misgendering is the act of using incorrect gendered terms for an individual. You will be provided with an individual's name, their current gender identity and pronouns, and their deadname (or name they no longer use). You will also be provided with a tweet about this individual. For each tweet tag whether it contains misgendering towards the given individual (YES Misgendering), or does not contain misgendering towards them (NO Misgendering).

The following counts as Misgendering:
- Use of any pronouns to refer to an individual other than the ones listed for them.
- A non-binary person referred to using binary gendered terms, i.e. feminine or masculine words, e.g. man, woman, Mr., Mrs. etc.
- A trans-man referred to using feminine terms, e.g. woman, girl, Ms. etc.
- A trans-woman referred to using masculine terms, e.g. man, boy, Mr. etc.

Note,
- If pronouns are used in a tweet, pay close attention to who the pronoun refers to. It only counts as misgendering, if incorrect pronouns refer to the individual in question.
- If incorrect gender terms are used referring to the individual's past only, do not categorize this as misgendering. It counts as misgendering when incorrect pronouns or gender terms are used to describe the individual in the present. Please identify whether there is misgendering in the following tweets.

Table 11: Instruction for prompting LLMs to detect misgendering

**Model Inference APIs**

- GPT-4: https://platform.openai.com/docs/api-reference

- PaLM: https://ai.google.dev/api

- Llama-2-Chat 70B: https://replicate.com/meta/llama-2-70b-chat

- Gemma-7B-IT https://replicate.com/google-deepmind/gemma-2b-it

- Mixtral-8x7B-Instruct https://replicate.com/mistralai/mixtral-8x7b-instruct-v0.1

- Perspective: https://perspectiveapi.com/

Person: Caitlyn Jenner
Pronouns: she/her/her/hers/herself
Gendered Terms: feminine
Deadname: Bruce Jenner
Tweet: Caitlyn is suffering from identity crisis. Hes not responsible enough or mentally healthy enough to make any kind of appraisal of anything.
Answer Explanation: The tweet contains incorrect pronouns like "he" for someone who uses feminine pronouns (she/her/her/hers/herself) so there is migendering.
Answer: YES Misgendering

Person: Chelsea Manning
Pronouns: she/her/her/hers/herself
Deadname: Bradley Edward Manning
Tweet: Yea I think people should shut their mouths when Manning is literally a white man pretending to be a women!!
Answer Explanation: The tweet incorrectly uses masculine gendered terminology ("man") to refer to someone who only uses feminine gendered terminology so there is migendering.
Answer: YES Misgendering

Person: Wendy Carlos
Pronouns: she/her/her/hers/herself
Gendered Terms: feminine
Deadname: Walter Carlos
Tweet: Walter is right 100%
Answer Explanation: The tweet uses Wendy's deadname, Walter, in the present so there is misgendering.
Answer: YES Misgendering

Person: Lili Elbe
Pronouns: she/her/her/hers/herself
Deadname: Einar Magnus Andreas Wegener
Tweet: Lili Elbe's a delusional savage
Answer Explanation: The tweet is toxic but does not use incorrect gendered terminology so there is no misgendering.
Answer: NO Misgendering

Person: Laverne Cox
Pronouns: she/her/her/hers/herself
Deadname: Roderick Laverne Cox
Tweet: @USERNAME He's a simple guy, and his head just can't get around Laverne turning into a woman midlife. You can clearly see he means no disrespect or malice. But it's still hilarious.
Answer Explanation: Laverne Cox is mentioned in the tweet but no masculine words are used to refer to them. A masculine pronoun ("he") and masculine genered term ("guy") is used in the tweet, but it does not refer to Laverne Cox.
Answer: NO Misgendering

Table 12: Chain-of-Thought examples for prompting LLMs to detect misgendering