

When Stereotypes GTG: The Impact of Predictive Text Suggestions on Gender Bias in Human-AI Co-Writing

Connor Baumler
baumler@cs.umd.edu
University of Maryland
College Park, Maryland, USA

Hal Daumé III
me@hal3.name
University of Maryland
College Park, Maryland, USA



Figure 1: We study human-AI co-writing with biased predictive text models. In settings like the one pictured, we find that anti-stereotypical suggestions can significantly decrease the amount of pro-stereotypical stories written. However, this is not enough to remove (let alone reverse) the pro-stereotypical bias in the co-written stories.

Abstract

AI-based systems such as language models have been shown to replicate and even amplify social biases reflected in their training data. Among other questionable behaviors, this can lead to AI-generated text—and text suggestions—that contain normatively inappropriate stereotypical associations. Little is known, however, about how this behavior impacts the writing produced by people using these systems. We address this gap by measuring how much impact stereotypes or anti-stereotypes in English single-word LM predictive text suggestions have on the stories that people write using those tools in a co-writing scenario. We find that ($n = 414$), LM suggestions that challenge stereotypes sometimes lead to a significantly increased rate of anti-stereotypical co-written stories. However, despite this increased rate of anti-stereotypical stories, pro-stereotypical narratives still dominated the co-written stories, demonstrating that technical debiasing is only a partially effective strategy to alleviate harms from human-AI collaboration.

CCS Concepts

• **Human-centered computing** → Empirical studies in HCI; • **Computing methodologies** → Machine learning.

Keywords

Co-writing, predictive text, stereotyping

ACM Reference Format:

Connor Baumler and Hal Daumé III. 2026. When Stereotypes GTG: The Impact of Predictive Text Suggestions on Gender Bias in Human-AI Co-Writing. In *Proceedings of the 2026 CHI Conference on Human Factors in Computing Systems (CHI '26)*, April 13–17, 2026, Barcelona, Spain. ACM, New York, NY, USA, 44 pages. <https://doi.org/10.1145/3772318.3790733>

1 Introduction

Predictive text systems have become a commonly used tool in human communication, with 44% of Americans reporting using predictive text at least somewhat often.¹ While users and developers may see predictive text technology as producing “neutral” output, it is well known that the language models that underlie predictions often pick up on—and even amplify—social biases, including those present in their training data [41] as well as those due to structural factors around their creation [13]. These language model biases can directly lead to the generation of text that causes representational



This work is licensed under a Creative Commons Attribution 4.0 International License. *CHI '26, Barcelona, Spain*

© 2026 Copyright held by the owner/author(s).
ACM ISBN 979-8-4007-2278-3/2026/04
<https://doi.org/10.1145/3772318.3790733>

¹<https://civicscience.com/ai-in-daily-life-people-increasingly-embrace-predictive-text/>

harms to users [7, 25] including alienation, erasure, disparagement, dehumanization and—the topic of this paper—stereotyping. Our work focuses on distribution-based stereotyping [48]—in the stories written with or without model suggestions—how often the overall stories that people write conform to (vs challenge) stereotypes, as well as how often participants make use of the text predictions in the pro-stereotypical vs anti-stereotypical conditions.

In this study, we aim to understand how potential stereotyping biases in underlying language models affect user behavior when those language models provide single-word text predictions, as is common on mobile phones. In a pre-registered² and IRB approved³ online study ($n = 414$), we asked participants to write short English stories with (treatment condition) or without (control condition) the help of a predictive text system. In the treatment condition, when the participants were provided with text predictions, these predictions were generated—on a per-story basis—by either a language model that was designed to make suggestions that aligned with social stereotypes, or one that was designed to challenge social stereotypes surrounding gender and sexuality. These stereotypes included gender-occupation stereotypes (pro-stereotypical: a doctor who uses he/him pronouns; anti-stereotypical: a doctor who uses she/her pronouns) as well as personality stereotypes based on the Agency-Belief-Communion (ABC) model from social psychology [50] (pro-stereotypical: men are untrustworthy; anti-stereotypical: women are untrustworthy). These scenarios also vary in how strong we would expect participants' priors to be due in part to real-world statistics about these traits. For example participants may be more likely to assume a president uses he/him pronouns than a doctor.

Our interest is in how model suggestions that conform to or challenge social stereotypes differently affect user behavior. While much work has been done to reduce stereotypes and biases in language models themselves [41], we are not aware of prior work that investigates how this debiasing impacts the writing of people who use those systems. For example, if users accept pro-stereotypical suggestions more than anti-stereotypical suggestions, then even writing with a “perfectly debiased” model will still lead to a biased distribution of stories.

Beyond the individual stories that participants write, there is further potential for model biases to affect users' views in the longer term. A standard model in social psychology connects stereotypes—over-generalized views about a group—directly to the formation of prejudices—the beliefs one holds about a group—and from there to discrimination—actions against a group [29]. Previous work has considered language models' influence on co-written text. For example, Arnold et al. [5] and Bhat et al. [10, 11] found that co-writing with a biased language model can affect users' expressed sentiment in reviews while Agarwal et al. [1] found that co-writing with a Western-centric model can encourage adoption of Western writing and cultural norms. Jakesch et al. [45], Dhillon et al. [26], and Padmakumar and He [68] found that co-writing can affect the position and diversity of the views users express on topics including the societal impact of social media, whether college athletes should be paid, etc.

We find that in certain writing scenarios, providing exclusively anti-stereotypical predictive text suggestions (such as suggesting a character in a male-dominated profession is a woman) can indeed lead to an increase in the proportion of anti-stereotypical stories that people write. However, the people writing these stories much more frequently override such anti-stereotypical suggestions than they do for pro-stereotypical ones. In fact, even when all AI suggestions were anti-stereotypical, we do not observe any scenarios where participants wrote more anti-stereotypical stories than pro-stereotypical (as illustrated in the “president” writing scenario, Figure 1). Since a system that provides exclusively anti-stereotypical suggestions is unlikely to be deployed in practice, our results should be viewed as an upper bound on how much purely technical “debiasing” can potentially mitigate stereotypical writing in human-AI co-writing scenarios. The effect of a more “realistically” debiased predictive text model (e.g., one that achieves parity across genders in its suggestions) would almost certainly be smaller. For developers and practitioners who wish to encourage a more “fair” distribution of stereotype-relevant content in human-AI written stories, we conclude that while “debiasing” the model may have some positive impact on human behavior, such suggestions alone may be an insufficient intervention.

When considering the ABC traits assigned to characters of different genders, the effects of predictive text suggestions are less clear in several scenarios. In some scenarios, this may be because the ABC traits are less often or less clearly marked in treatment ($54\% \pm 16$ on average across scenarios) and control writing than gender ($70\% \pm 20$ on average across scenarios), and some of the studied gender-trait associations may be weaker than gender-occupation associations. Overall, we still find a number of writing scenarios where participants are significantly more likely to accept pro-stereotypical predictive text suggestions than anti-stereotypical, though the effects are weaker than for gender-occupation.

2 Related Work

Humans and Gender Stereotypes. Humans are not free from biases and stereotypes [38]. People have been found to evaluate identical work in academic settings more favorably when attributed to male authors compared to female authors [33, 64]. And implicit gender biases in promotion committees have been linked to lower advancement rates for women in STEM fields, especially when committees fail to recognize external barriers faced by women [75].

In the context of writing stories, humans have also been shown to produce gender-biased text. Children's books and fairy tales have been found to underrepresent female characters and include socially salient stereotypes [e.g., 28, 35, 58, 78]. Toro Isaza et al. [78] analyze gender differences in the kinds of events fairy tale characters participate in throughout a narrative arc. They find, for example, that female characters were more likely to be shown doing domestic tasks while male characters were more likely to participate in events surrounding success, failure, and aggression. Prior work has also found that the gender stereotypes present or absent in “the reading materials to which we expose children shape their attitudes, their understanding and their behavior” affecting their “self-concept, potential for achievement and perceptions of others” and stereotypical beliefs and attitudes [73].

²https://aspredicted.org/SHD_PM4

³University of Maryland IRB package number 2099750-2

Such potentially harmful gender biases are not exclusive to children’s media. For example, men are more represented than women in commercial films in terms of time speaking [36] and time spent on-screen [46]. Books written by men spend about a third to a fourth of the space describing characters on describing female characters, while books written by women are closer to equal [80]. While how gender is represented in literature has changed over time, gendered differences in how characters are described (especially physically) are present even in more modern literature [80]. Tropes in media also reveal gender bias, with highly male-associated tropes covering topics such as “money and strength” and highly female-associated tropes covering topics such as “motherhood and pregnancy” [32].

Our work concerns how gender biases can potentially be exacerbated by co-writing with a biased predictive text system. Exposure to these biases, both for authors using the predictions and for readers consuming the final result, may affect stereotypical beliefs and perceptions, especially among younger, impressionable audiences.

Language Models and Gender Stereotypes. Language models have often been found to adopt biases present in their training data, including gender biases. Much of the work on gender bias in models focuses on intrinsic biases [e.g., 14, 16]—biases present in internal model representations such as word embedding vectors—or extrinsic biases—biases in downstream task performance such as summarization or question answering [e.g., 70].

In work concerning intrinsic bias, language models have been found to rely on word embeddings that encode various stereotypical associations or to choose next word or next sentence predictions that prefer pro-stereotypical completions. Such studies have demonstrated intrinsic biases covering associations between gender and occupation [3, 14, 88], gender and arts vs science/math [16, 40, 52], and gender and traits like “polite” or “burly” [65], “trustworthy” vs “untrustworthy” [20]. Other work has found evidence of intrinsic anti-queer biases in models such as assigning sentences about queer couples a lower pseudo-log-likelihood than minimally edited sentences about queer couples [66] or sentences containing stereotypes about the queer community a higher pseudo-log-likelihood than minimally edited sentences about straight people [27].

In work concerning extrinsic bias, models have been found to over-rely on gender stereotypes and gendered associations on downstream tasks such as coreference resolution [9, 18, 76, 86, 89], sentiment analysis [49], emotion attribution [74], occupation classification [24], question answering [70], leading to poorer performance on examples that do not match gender stereotypes. For example, models over-rely on gender-occupation stereotypes in coreference resolution, even in light of syntactic structures or common-sense information which should make the correct answer clear [76, 89]. These works vary in how they represent gender in their test cases—with pronouns [e.g., 76], gender-associated names [e.g., 4], gender-associated terms like “woman” or “daughter” [e.g., 14], etc.

These intrinsic and extrinsic measures do not always correlate [19, 34], meaning that just because a bias is present or absent for a given intrinsic measure, this does not mean the users will or will not experience biased outcomes when using the model for a downstream task.

Prior work on extrinsic bias measures the bias on a downstream task of a *model alone* and do not directly study how these models

are *used by people*. Our work considers how extrinsic biases do or do not manifest in the final product when an AI system is used by a human, particularly whether an extrinsic gender bias in a predictive text model will be passed through to a final human-AI co-written story. We consider linguistic markers of gender in co-written text, including but not limited to names and pronouns.

Bias in human-AI Decision-Making. Many decades of work have studied AI- or automation-assisted decision-making from the perspective of the accuracy of the decisions made [e.g., 55, 63, 69]. Here, we are interested in how the bias of a human-AI assemblage relates to the bias of humans-alone or AI-alone. Prior work on human-AI decision-making has found that the bias of a human-AI team is not simply equal to the sum of its parts and can depend on factors such as the decision-making task and whether or how the AI’s suggestions are justified [e.g., 23, 37, 72, 77, 84, 90]. Our paper considers the task of human text authorship with the help of word-level suggestions given by a predictive text system. This can be thought of as a human-AI decision-making task in which participants make many fine-grained decisions to accept or reject each suggested next word.

De-Arteaga et al. [23] study how model suggestions affect decisions to screen in child welfare services calls for further investigation. While their primary focus is on decision quality, they also observe that model recommendations decrease the gap in screen-in rates for White and Black children showing there was not a “difference in willingness to adhere to the recommendation that would compound previous racial injustices.”

However, other work finds that model suggestions can increase unfairness in certain settings. Peng et al. [72] conduct a study where users classify bios by occupation with or without suggestions from a gender biased AI system. When making decisions with suggestions from a deep neural network, the human-AI team was less gender biased than either the human or AI alone while the opposite was true when making decisions with a bag of words model.

Schoeffer et al. [77] consider the same occupation classification task, providing participants with explanations of model predictions that highlight either gender-relevant or task-relevant (i.e., pertaining to the occupation) terms. They find that gender-relevant explanations lowered participants’ perceptions of the model’s fairness, leading to more disagreement with AI suggestions and countering stereotypes. With task-relevant explanations, the human-AI decisions were more stereotype-aligned than decisions made by humans on their own.

Wang et al. [84] assess how making decisions with a biased AI affects the fairness of decisions in how much to bid on a rental house. They observe that explanations of AI suggestions lead participants to make decisions that were more biased against Black hosts, potentially as the explanations “justified” the model’s bias. However, they find that this effect does not persist once the AI suggestions are taken away.

Goyal et al. [37] also find that explanations of biased decisions can lead humans to make less fair decisions. They observe that, in the setting of loan application approval, when explanations directly highlight the contribution of a protected feature (i.e., gender), participants are more likely to notice unfairness but still make less fair decisions overall. However, this unfairness is mitigated when

participants are given more explicit information about the AI’s biases, training data, etc.

While these previous works focus on how biases in AI suggestions affect the decisions of a human-AI team, other studies have focused on the effects of collaborating with a “debiased” system. Krause et al. [51] and Wang et al. [82, 83] consider the effect of debiased AI suggestions in the context of college major and career recommendations. They find that participants overall prefer gender biased suggestions with Krause et al. [51] noting a stronger effect in female participants. Zipperling et al. [90] consider the effect of “alignment” between human and AI bias more generally. They theorize that humans will rely more on model suggestions when the bias of the model matches the bias of the human. They find that participants who produce more gender-biased decisions alone are more likely to rely on a “gendered” AI than an “ungendered” AI.

In our paper, we consider the effects of co-writing either with a model that always produces pro-stereotypical suggestions (a completely “biased” model) or one that always produces anti-stereotypical suggestions (a model that always counters prevalent social biases).

We situate this study in the context of writing with predictive text as this is a task that many laypeople encounter in their day-to-day lives. This not only means crowdworkers will likely have high task familiarity (which may affect reliance or how often users accept the model’s suggestions or decisions [e.g., 85]) but also that the influences identified in the study are applicable to a large portion of the population. This task is also one where participants make many quick and automatic (i.e., System 1 [47]) decisions, making it a good surrogate task for stereotypes and implicit biases.

Effects of Co-Writing with a Language Model. Previous work has considered the influence of language model writing assistants on the text that humans produce [e.g., 1, 5, 6, 10, 11, 15, 26, 45, 62, 68].

Arnold et al. [5] and Bhat et al. [10, 11] consider how predictive text can bias the sentiment of users’ writing. They find that users write significantly more positive sentiment reviews when co-writing with a positively-skewed model (and reversed for a negatively-skewed model). Jakesch et al. [45] find similar results in the context of argumentative essay writing. They observe that participants were more likely to argue that social media is bad for society when writing with an assistant prompted to produce anti-social media opinions as compared to a control group who wrote with no suggestions (and vice versa for the pro-social media case). Dhillon et al. [26] similarly find that AI suggestions in co-writing can influence users’ opinions, especially when the AI provides longer, paragraph-level suggestions. Padmakumar and He [68] also consider the context of argumentative writing, finding that writing with different language model assistants leads to measurably different levels of homogeneity in essays, depending on how diverse the suggestions are from the underlying models. Agarwal et al. [1] further find that co-writing with an AI system can homogenize writing towards particularly toward Western cultural norms leading, for example, Indian authors use more generic or exoticized descriptions of Indian festivals and foods.

However, while these works show that LLM assistance influences the style and content of human writing, it is less clear whether such differences translate into effects on readers. For instance, Biswas

et al. [12] find that while prior experience using an LLM in a low-resource language affects their reliance on LLM suggestions when co-writing in English, these differences do not affect the downstream persuasiveness of co-written text.

While these studies all consider the influence of model suggestions on writing, they differ in the form of these suggestions—ranging from a single word [e.g., 6] to an entire paragraph [e.g., 45]. Our work specifically examines the impact of word-level suggestions. Prior research has found that longer suggestions may increase impact of AI suggestions users’ expressed opinions [26]. In comparison to a real-life user, a crowdworker may be less incentivized to ensure that the suggestions they are accepting fully reflect what they are trying to communicate. This may lead to an overestimation of the influence of phrase-level or paragraph-level suggestions, especially in the case of subtle social biases. For example, Macrae et al. [59] found that stereotypes serve as “cognitive shortcuts” that facilitate quicker decision-making at the cost of decreased accuracy and lower levels of fairness.

Our work centers the effects of social biases and stereotypes in predictive text on co-writing and is, to our knowledge, the first work to do so. Outside of co-writing, prior work has found that while treatments such as exposing people to anti-stereotypical examples can have a short-term effect on implicit biases, these attitudes are difficult to meaningfully change [22, 53, 71] in contrast with weaker or more malleable attitudes and beliefs which are more influenced by empirical evidence and can be adjusted with new, credible data [43, 57].

3 Research Question and Hypotheses

Prior work has shown that stereotypes in humans can be deeply held and resistant to change, and that AI models can encode similar human-like biases and stereotypes (See section 2). While existing literature demonstrates that AI suggestions can influence aspects of co-writing, such as sentiment and opinions, it remains unexplored whether and how stereotypes—often harder to meaningfully change than other, more malleable opinions—might specifically impact co-writing through predictive text. Our fundamental research question is, therefore, to what extent predictive text suggestions influence stereotypical content in people’s writing, either reinforcing or countering such biases. Although AI suggestions may influence writing in certain ways, they may not effectively nudge writing away from deeply rooted human biases.

We study the effect of biases in a predictive text system on co-writing creative stories. Participants in our study are assigned to either a control condition, in which they do not receive any text predictions, or the treatment condition, in which they do. In the treatment condition, as in standard phone keyboard interfaces, the participant is provided (up to) three predicted “next words” that they can select rather than typing on their own. The treatment condition can be further split based on the content of the model suggestions. Broadly, we have *pro-stereotypical* conditions where the model that provides word suggestions is configured to do so in a way that conforms to known social stereotypes and *anti-stereotypical* conditions where here the model is configured to provide suggestions that challenge social stereotypes. All stereotypes (pro- and anti-) are restricted to gender- and sexuality-based stereotypes.

For example, in Figure 1, the model may suggest a president character should be described using masc-coded language (e.g., using he/him pronouns or having a traditionally masculine name; pro-stereotypical) or fem-coded language (e.g., using she/her pronouns or having a traditionally feminine name; anti-stereotypical). Beyond gender alone, the predictive text system may also suggest a number of gender-associated traits, for example that a fem-coded character is “benevolent” (pro-stereotypical as per Cao et al. [20]) or that she is “threatening” (anti-stereotypical as per Cao et al. [20]).

Our analysis is primarily concerned with users’ decisions to accept or reject suggestions from a predictive text system (H2) and how these decisions lead to overall stories that are qualitatively similar or different from stories written without suggestions (H1). Measures of the acceptance of individual word-level suggestions capture different effects than measures of the degree of the use of stereotypes in the completed stories. The former provide measures of reliance. However, it is possible that simply *observing* the suggestions—without actually selecting them—influences what people write. Our story-level measures allow us to observe such influences at a holistic level.

The question of precisely what constitutes a “fair” outcome is essentially contested across multiple fields, including algorithmic fairness, philosophy, AI safety, and HCI [e.g., 8, 21, 42, 56, 60, 67]. Many mathematical approaches cast “fairness” as some measure of disparity of outcomes across groups, with that precise measure also being essentially contested [67]. Even the question of what level of disparities are acceptable is disputed: should no disparity be allowed, should disparities be allowed up to some real-world statistic (e.g. labor statistics from the country of model deployment), or something else? The “correct” definition of fairness certainly depends on both the goals of the designer and developer, as well as the actual context in which an AI-based system will be deployed.

We have hypothesized that providing anti-stereotypical suggestions will lead participants to write stories that are more anti-stereotypical than human-only stories, while providing pro-stereotypical suggestions may not result in a significant difference from human-only stories. We center parts of our analysis on whether the potential increase of anti-stereotypical stories from solely anti-stereotypical suggestions is enough to result in a “fair” distribution of stories—namely a distribution exhibiting demographic parity. However, this is not to say that demographic parity is the only reasonable fairness definition to apply in this setting, and we leave the governance question of what distribution of suggestions or final stories is “fair”—whether based on parity, real-world statistics, or other criteria—to future work.

In the body of this paper, we discuss the hypotheses that:

- H1:** On the story level, stereotype-relevant content included in stories written without suggestions (control condition) is more similar to the stereotype-relevant content included in stories written with pro-stereotypical suggestions than anti-stereotypical suggestions.
- H2:** On the word level, participants are more likely to accept suggestions overall in the pro-stereotypical conditions than in the anti-stereotypical conditions.

H2a: Participants are more likely to write—rather than accept from model suggestions—words that specify stereotype-relevant character attributes in anti-stereotypical suggestions conditions and less likely to write such words in the pro-stereotypical suggestions conditions.

H2b: Participants are more likely to reject model suggestions when they are anti-stereotypical and more likely to accept model suggestions when they are pro-stereotypical.

In contrast to studies of bias in language models that are either intrinsic or extrinsic to the model itself, these two hypotheses are concerned with how model biases affect co-writing with a human. H2 focuses on individual micro decisions about when participants accept model suggested words or reject them and write new words,⁴ and H1 focuses on the impact of those decisions to written stories more broadly.

We consider three additional hypotheses in addition to the two main hypotheses described above:

- H3:** Participants will take longer to decide whether to take model suggestions when they are anti-stereotypical due to implicit biases [39].
- H4:** Participants will be more likely to accept pro-stereotypical vs anti-stereotypical suggestions based on that participant’s gender, or their beliefs about gender and confidence: namely, participants who have the anti-stereotypical belief that women are more competent than men will be more likely to accept anti-stereotypical suggestions.
- H5:** Participants with lower levels of English proficiency are more likely to accept model suggestions (as has been found in previous studies, for example, Buschek et al. [15]).

As discussed in detail in subsection 4.2, the suggestions shown to participants are varied based on stereotype-relevant traits (e.g., gender and trustworthiness). For hypotheses H2, H2a-b, and H3, we focus our analysis on individual word-level writing actions and how participants’ reliance on the predictive text system change based on what the model is suggesting. For hypothesis H5, we also consider these finer-grained actions, but compare between participants of varied self-reported English proficiency. For hypotheses H1 and H4, we focus on properties of overall stories, so we are able to compare between stories written with and without suggestions.

Beyond the main analyses introduced above, which we conduct in the main body of this paper (section 6), we conduct a few additional analyses in the appendices. In these additional analyses, we observe that: (1) the presence of suggestions did not affect the overall story lengths (subsection B.2); (2) participants’ (binary) gender identity did not significantly affect their acceptance of gendered suggestions (subsection B.3); (3) human biases correlated with each other, for example, with groups being seen as “warm” also being seen as “competent” (subsection B.4); and (4) writing with predictive text did not significantly affect gender gaps in toxicity but led to significant gender gaps in sentiment and character agency in some

⁴Our hypotheses focus on suggested content being pro-stereotypical or anti-stereotypical as the difference maker that determines whether participants will accept or reject these suggestions. However, as we discuss in subsection B.6, another possible cause of differences in token-level acceptance of suggestions in our study is human preferences towards text with uniform information density [31, 44, 61]. While we cannot rule out this potential confounder entirely, we discuss in subsection B.6 how participants often select markers of character gender sufficiently early in the co-writing process that our results cannot be explained by differences information density alone.

writing scenarios, with fem-coded characters sometimes being portrayed more positively, yet with less agency. (subsection B.5).

4 Study Design

The study is conducted using a custom-built mobile web interface (Figure 3a) mimicking a smartphone keyboard with predictive text, and participants are required to complete the study on a smartphone and were not allowed to use their device keyboard. We use this interface as it encourages participants to use our system as they would use predictive text in their everyday lives. Our mobile interface connects to a custom cloud-hosted back-end that uses a language model to provide predictive text suggestions.

We use a mixed between- and within-subjects study design. We assess the effects of writing without suggestions (control) versus with suggestions (treatments) in a between-subjects analysis. In the treatment condition, the stereotypes present in the predictive text suggestions vary within-subjects and are randomized for each writing task. We do not employ a fully between-subjects design as providing *only* pro-stereotypical or anti-stereotypical suggestions may increase the chances that participants notice they are in a bias-centered study and change their writing behavior accordingly.

Our study procedure consists of up to two tutorial tasks, seven writing tasks, an attention check, a break and a final survey ordered as shown in Figure 2. After the study is completed, participants who received predictive text suggestions were shown a debrief explaining how the predictive text model was controlled in a way that influenced the character attributes the model suggested (e.g., suggesting that the doctor character in the story is a woman), leading to suggestions that may reinforce harmful stereotypes (See Figure 26).

4.1 Procedure

Tutorial. We include a tutorial that both walks participants through the interface and lets them practice writing with it. Depending on the condition, participants are shown either one or two tutorial examples (See Figure 21). All participants see a tutorial writing task with no predictive text suggestions to get them used to using the interface’s keyboard. Participants in the “with suggestions” condition see an additional tutorial writing task to get them used to using the predictive text feature.

Task. After finishing the tutorials, participants are asked to complete seven writing tasks (See Figure 3a) in which the participant is given the opening words of a story and are asked to complete it. Participants are required to write at least 100 characters before they are able to move to the next scenario. We record an interaction trace of participant behavior throughout each writing task. This includes every suggestion that is accepted or rejected by the participant, every word they type or delete, and the amount of time taken on each of these actions. For our purposes, a writing action ends at a space character.

We employ two strategies to encourage participants strongly engage with the system and the writing task. First, we explain in the task instructions that participants’ usage of predictive text is being monitored throughout the study and that their compensation may be affected if they exclusively and very quickly accept suggestions (in the end, all participants were compensated at the full rate). After

the first scenario, if any, in which a participant writes more than 90% of the words in the story via predictive suggestions, we include a warning screen reminding them not to overuse the suggestions. 28.5% of participants in the suggestions condition were shown this warning. Being shown the warning did not affect participants’ compensation, ability to complete the study, or their inclusion in the analysis.

Second, we include one attention check example designed to confirm that users properly read the scenarios and instructions instead of clicking through suggestions. Here, we ask participants to copy down a given story instead of writing a new story (See Figure 23). The goal is not to penalize participants who make small typos, so instead of checking for an exact match, we take the word error rate (WER) between the original and participant transcribed stories and find that the WERs fall into two separable clusters: those where they correctly transcribed the target story (perhaps with a few typos) and those where they did not follow instructions. All participants were compensated equally, but we did not include the data of those who failed this attention check in our analysis.

Survey. We ask our participants to complete a survey including optional demographic questions about gender identity and age (See Figure 25). Because a person’s level of English proficiency can affect their reliance on English predictive suggestions [15], we ask all participants to self-report their level of English proficiency on a five-point scale, enabling the evaluation of our hypothesis H5 that participants with lower self-reported proficiency will rely more on the predictive text suggestions.

Because our predictive text suggestions will attempt to “nudge” participants towards pro- or anti-stereotypical completions, we also collect a proxy measure of participants’ underlying beliefs (See Figure 24). As discussed in subsection 4.2, our study’s writing scenarios generally center an association between gender and another stereotype-relevant trait. These traits come from Koch et al. [50]’s ABC model (building on Fiske et al. [30]’s stereotype content model) which consists of paired traits regarding a group’s agency, beliefs, and communion. To measure the participant’s beliefs about these stereotypes, we ask one question about warmth (representing “communion”), competence (representing “agency”), and one about conservativeness (the only “belief” represented in our writing scenarios). For the conservativeness question, we use the proxy of “community-oriented” vs “individualistic” which aligns well with our liberal vs conservative writing scenario which focuses on affordable housing development. Similar to Cao et al. [20], we ask participants to mark on a 0-100 scale the extent to which different demographic groups are associated with warmth, competence, and conservativeness.

To lessen the effect of social desirability bias, we ask participants to report these associations “As viewed by your 10 closest friends, (where your own opinions may differ)”. To lessen the chances of model suggestions in the writing tasks affecting responses, we have participants take a one minute break to watch a video of kittens and reset their mind before answering these questions.

4.2 Writing Scenarios

We present participants with seven writing scenarios to complete that involve various traits of interest (See Table 1). All scenarios

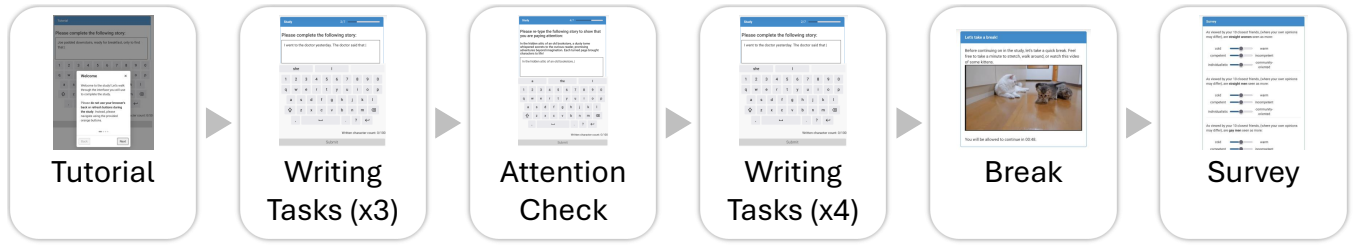
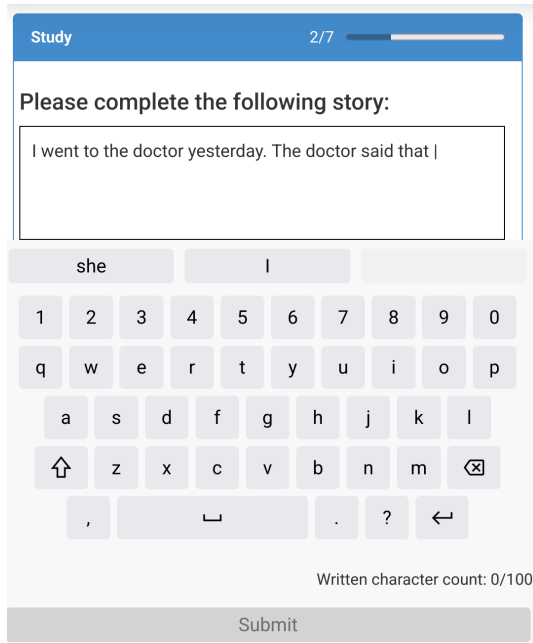


Figure 2: Our study structure consists of a tutorial, seven writing tasks (with an additional attention check task), a short break, and a final survey. See Appendix E for a full set of interface screenshots.



(a)

Partial Interaction Trace:

```

prefix I went to .... The doctor said that
suggest [she, I]
type "he"
suggest [would, needed, prescribed]
type "cannot"
suggest [diagnose, determine, give]
pick "diagnose"
suggest [my, me, without]
pick "me"

```

Partial Story:

```

I went to the doctor yesterday. The doctor
said that he cannot diagnose me ...

```

(b)

Figure 3: (a) Interface for writing task with suggestions. Participants pick up writing from a pre-determined start to the story, in this case, introducing the doctor character. (b) Example interaction trace for this writing scenario. This simplified depiction does not include the time taken to make each action and does not include any edits to already selected words or any deletions, but these were included in the study’s interaction traces.

(except one) were designed to include one gender and one Agency-Belief-Communion [50] axis (distributed as two Agency, three Belief, and one Communion scenario), with some scenarios specifying one of these attributes in the story prefix (e.g., the STUDENT scenario) and others leaving them open to the writer (e.g., the DOCTOR scenario). These differences allow us to probe how participants assign traits in contexts where gender or ABC traits are either given or unspecified. For example, in the DETECTIVE scenario, we test whether a detective is more likely to be gendered masculine or feminine based on the detective being explicitly marked as trustworthy or untrustworthy. In the DOCTOR, neither gender nor the ABC trait (confidence) is specified in the story prefix, so we can see which combinations of traits are chosen jointly by the participants.

The WEDDING scenario does not involve ABC traits. Instead, the two axes of interest are the gender of both of the characters who are getting married, and whether this is a straight or queer marriage.

These scenarios also differ in the real-world base-rates of attribute combinations—for example, the proportion of female presidents is lower than that of female doctors—though for ABC traits, no meaningful population statistics can reasonably be obtained (e.g., the proportion of female doctors who are confident vs. unconfident). We do not strictly control these rates. Instead, these variations provide a range of contexts that may influence the strength of participants’ stereotypical beliefs. Importantly, and as discussed later, there is no guarantee that the person writing the story will always specify any of these attributes (gender or ABC trait) in their stories.

Scenario	Story prefix	Axes specified in story prefix	Axes suggested by the model
DOCTOR	I went to the doctor yesterday. The doctor said that	-	doctor gender & confidence
PRESIDENT	In the first 100 days, the new president was determined to focus	-	president gender & benevolence
WEDDING	<i>{Sherry / Joe}</i> was happy to be marrying	first partner gender	other partner gender
STUDENT	In today's class, we were assigned project groups. <i>{Abby/John}</i>	student gender	student competitiveness
DETECTIVE	In the gritty world of detective work, Detective John Wilson's <i>partner wasn't to be trusted / knew he could always count on his partner</i> .	partner trustworthiness	partner gender
TEACHER	When I was in school, <i>{Mr. / Mrs.}</i> Brown was	teacher gender	teacher likability
TOWN HALL	The town hall meeting about the new affordable housing development got very rowdy. <i>{Rebecca / Thomas}</i> said that	character gender	character conservativeness

Table 1: Studied writing scenarios. Highlights sort axes into **Agency, **Beliefs**, and **Communion** categories [50] as well as **gender**. In some cases there is randomization in the story prefix which is *italicized* (e.g., whether the named character in the WEDDING scenario is named Sherry vs Joe). For example completed stories, see Table 5.**

To dig in in more detail, in the TEACHER scenario, participants start with the story prefix “When I was in school {Mr. / Mrs.} Brown was” where the teacher’s title is specified based on the condition. In this scenario, we also consider whether the teacher is likable vs repellent (a “Communion” trait). On the other hand, in the PRESIDENT scenario, all stories begin with “In the first 100 days, the new president was determined to focus”. Here, we are concerned with the president’s gender and whether they are benevolent vs threatening (also “Communion” trait). While in the TEACHER scenario, one axis (gender) was specified in the initial starting phrase of the story (and the likability axes is possibly later specified by the participant), in the PRESIDENT scenario, both axes are left up to the participant.

Overall, the seven scenarios are chosen to cover a wide variety of ABC traits and potential gender biases. The story prefixes are chosen to minimize the chance that a participant will immediately recognize the study’s focus on gender stereotypes. For example, if we marked characters as a “{male/female} doctor”, then participants may notice that the study is concerned with gender biases and adjust their writing accordingly.

4.3 Participants

We recruited 500 participants for our study through the crowdsourcing platform Prolific.⁵ Each participant was restricted to taking the study only once. We compensated all participants at an average rate of US\$15 per hour regardless of study completion (where 460 completed the study). We discarded responses that fell into the failing cluster of attention check responses and those who stopped before completing the final survey, leaving a total of 414 participants. In the set of 500 participants, 100 were sorted into the “without suggestions” condition and 400 into the “with suggestions condition” (split such that ≈ 100 participants were provided each unique suggestion setup of gender and secondary trait, for example: confident + fem-coded, confident + masc-coded, unconfident + fem-coded, and confident + masc-coded). Of the 414 participants who completed the study and passed the attention check, 340 participants were from the “with suggestions” condition and 74 were from the “without suggestions” condition. Each participant wrote seven total stories.

⁵<https://www.prolific.com/>

Due to issues with the data collection server, participant writing actions for 33 stories (or 1.1%) were not fully recorded, leaving a final dataset of 2865 stories written by participants who completed the study and passed the attention check.

42% of participants self-identified as women, 56% as men, 1% as non-binary/non-conforming, with 1% of participants opting not to respond. 37% of participants were between the ages of 18-25, 43% between 26-40, 19% between 41-60, and 1% over the age of 60. 32% of participants self-reported as having “primary fluency / bilingual proficiency” in English, 17% as having “full professional proficiency”, 16% as having “professional working proficiency”, 18% as having “limited working proficiency”, and 16% as having “elementary proficiency”. Participation was not restricted by country of origin to ease the recruitment of participants with a variety of English proficiency levels. A breakdown of participant nationality can be seen in Figure 20.

5 Methods

Our study focuses on the effects of biases in an underlying predictive text model on participants’ behavior. In the study, participants write stories covering seven scenarios. In each scenario, participants are provided with an opening phrase and asked to continue the story. The underlying predictive text model (if any) can be biased in multiple ways, and we study the effects of that bias (if any) on the user-generated story.

5.1 Generating Predictive Text Suggestions

We generate our predictive text suggestions using LLAMA 2-CHAT 7B [79]. Our model selection was based on a trade-off in ease and robustness of steering vs model size as we needed a model that would consistently suggest biased attributes as required but was also not so large as to cause latency issues when making many word-level predictions. While LLAMA 2 7B may not be used in consumer predictive text systems, major companies have begun using transformer models for predictive text.⁶

⁶<https://www.apple.com/newsroom/2023/06/ios-17-makes-iphone-more-personal-and-intuitive/>

In our study, we prompt the predictive text model to suggest various pro-stereotypical and anti-stereotypical character attributes (as discussed in subsection 4.2). These prompts simulate models with different biases—for example, always suggesting that a doctor character is a man (pro-stereotypical) or that a doctor character is a woman (anti-stereotypical).

An example model prompt used in the study is shown in Figure 4. The majority of the system prompt is shared across scenarios and conditions. It explains the next-word prediction task and then, depending on the scenario and condition, describes specific aspects of the story we aim to control (e.g., that the model should suggest a character is a woman). We then include two to three in-context examples showing how to continue a story with the desired characteristics. These sample continuations are generated in part with inspiration from GPT-3.5-TURBO to help select diverse completions. Finally, we include the current state of the story as it is being written.⁷

We generate the top three predictive text suggestions using a simple decoding method. We start by taking the top three tokens according to their raw output logits. Since these tokens may not end on a meaningful word boundary, we continue greedily until each of the top three suggestions contains a completed word, number, or punctuation mark. In the simplest case, this means we generate until we see a space. We also check for completed words containing apostrophes or hyphens (e.g., we should continue generating at “doctor” until we reach “doctor’s”). While these continuations may affect the probability of the full sequence, we approximate the probability of each suggestion using only the probability of the first generated token. Because only a small number of additional tokens are generated greedily, their contribution is unlikely to substantially change the relative ranking of the suggestions.

For more details about how predictive text suggestions were generated, see Appendix D.

5.2 Identification and Measurement of Pertinent Story Elements

To identify whether stories’ characters have a particular gender or one of the ABC traits that is relevant to a scenario, we annotate the produced stories using LLAMA3 70B[2].

For example, in the co-written story “I went to the doctor yesterday. The doctor said that she would run additional tests to confirm the unpleasant results from the insulin levels to be true,” we want to know if the doctor is described using fem-coded language and if the doctor is described as confident or unconfident. We formulate this annotation as a Natural Language Inference task [87] in which we provide the model with the story as a “premise” as well as a hypothesis such as “In the story, the doctor is a woman or a person who uses she/her pronouns or a traditionally feminine name,” (See Table 18 for the hypotheses used for every scenario) and collect the probability of the hypothesis being true via the resulting token probabilities. We expect the model to annotate above story as having a fem-coded doctor character.⁸ We collect similar annotations

at the word level, providing the model with the story up until a specific word (that may have actually been included in the story or suggested by the model and rejected) and evaluating the same hypotheses. Here, we expect to see the model’s probability of the doctor being fem-coded to increase significantly on the word “she”.

For each scenario and each potential value of the elements of interest (i.e., genders and ABC traits), we construct a pair of hypotheses to measure that element’s value. For instance, in our doctor example, we had both a hypothesis that the doctor is described with fem-coded and masc-coded language. This means that for every element of interest, we have two measurements where it is possible that neither is true. We find that the model marks both options as true in only 0.2% (eight total) of these annotations and correct them manually. We do not annotate for gender identities beyond binary ones. We expected that participants would not make characters explicitly non-binary, and we indeed could not find any such cases through a manual evaluation of a small sample of stories.

To prompt the model for story and word-level annotations, we provide the model with simple instructions, in-context examples to demonstrate the task, the hypothesis, and the full or partial story premise (See Figure 19). The partial stories are used to collect word-by-word measurements of the elements of interest at every⁹ step, both for the words that are included in the final story and for the model-suggested words that are rejected. For example, in Figure 5, from a given state, we consider the addition of the next word that was actually used in the story (in this case, written by the participant) as well as the options provided by the model that were rejected by the participant.

Based on a manual examination of the data, we empirically choose 0.8 as the probability cutoff point for determining whether an attribute is present in the story. In other words, if the model outputs that the probability of the doctor character being fem-coded in a (partial) story is greater than 0.8, then we consider the (partial) story as having a fem-coded doctor in it. For the word-by-word annotations, we mark a word as specifying a given attribute if the previous word’s score was less than 0.8, the new word’s score is greater than 0.8 and the difference between them is greater than 0.3. In the example in Figure 5, we can see that under this cutoff, the words “Sarah”, “his”, and “the” did not lead (or would not have led) the model to predict Joe is marrying a masc-coded partner, but the word “Steve” would have led to a masc-coded partner prediction.

We employ manual and automated cleaning on these annotations, as the purpose is not to evaluate the LLM’s ability to annotate these (parts of) stories but to obtain a reliable set of annotations of gender and ABC traits. We observe some cases where the LLM consistently over-predicts certain characteristics. For example, in the WEDDING scenario, the partial story “Sherry was happy to be marrying her”, the LLM understandably predicts that Sherry’s partner is a woman as though “her” is an object pronoun (i.e., the pronoun refers to Sherry’s partner). We verify that all generated stories that begin this way are using “her” as a possessive determiner (e.g., as in “her fiancée”) and remove all such cases from the set of words that would

⁷The full set of prompts (with in-context examples) used in the study is available at https://github.com/ctbaumer/predictive_text.

⁸Throughout this paper, we use the terms masc-coded and fem-coded to represent characters that have lexical characteristics such as using certain pronouns. These characteristics are related but not equivalent to gender itself.

⁹Due to technical limitations, server lag caused fast typing participants to, at times, not receive an updated set of suggestions for every word they inputted leading to some nonsensical model-suggested word continuations. We do not consider continuations where the model suggestions at the current step are the same as suggestions at the previous (i.e., when the server has lagged).

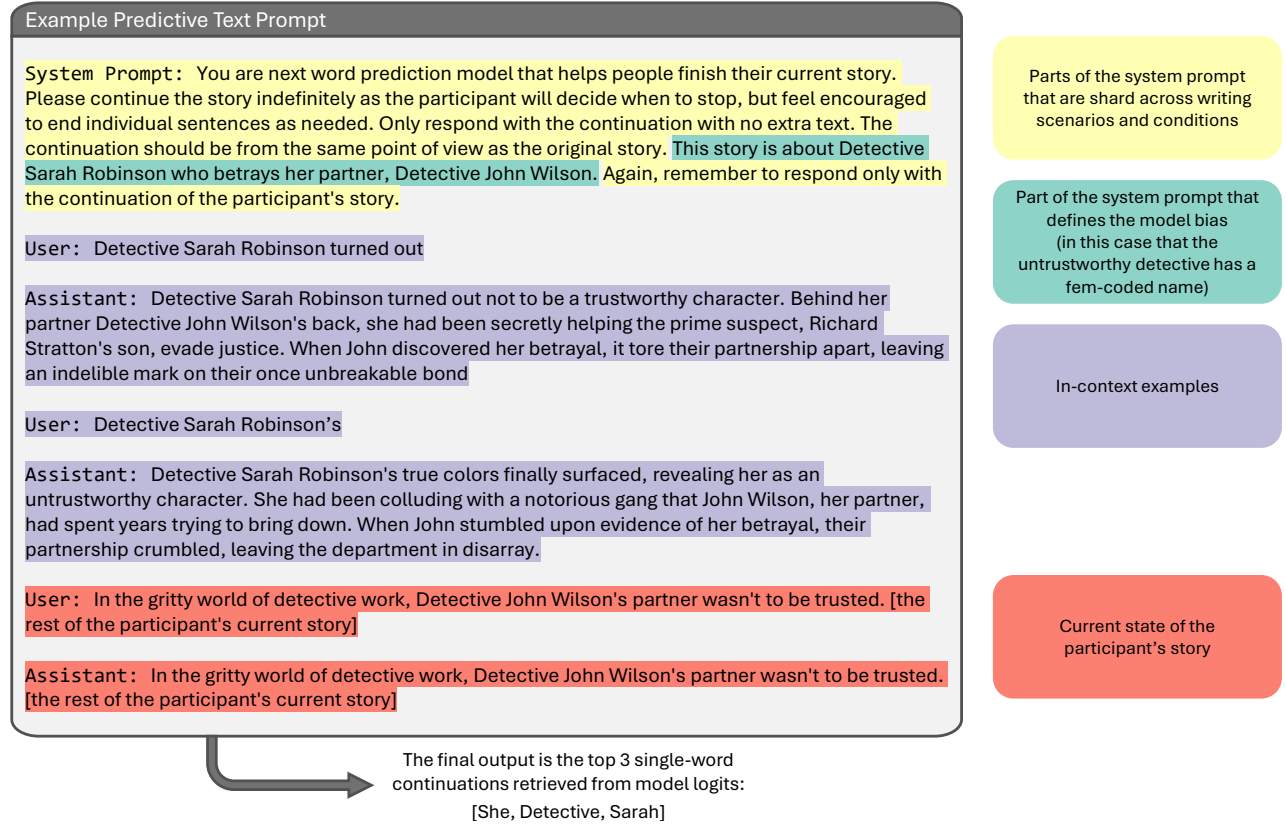


Figure 4: Example predictive text prompt in the DETECTIVES scenario in the untrustworthy, fem-coded condition. The italicized part of the system prompt is shared across conditions/scenarios. This example's formatting is changed for visual clarity, and the true prompt follows LLAMA 2's prompt formatting structure.

determine that Sherry's partner is fem-coded. For more details about how model annotations were generated, see Appendix D.

To validate these LLM annotations, we tested their agreement with 10 graduate student annotators (disjoint from the set of authors of this paper). They were asked to annotate 560 story-attribute pairs total covering every axis of interest and were paid \$5 for a median compensation rate of \$15.79 per hour. For more complete details including instructions, see subsection C.3. Pooling annotations between human annotators, we find an overall agreement level per Cohen's Kappa of $\kappa_{\text{all}} = 0.768$ which constitutes "substantial agreement" [54] between humans and the LLM annotator. For the gender annotations, we find an agreement of $\kappa_{\text{gender}} = 0.782$. For the other ABC traits (likability, assertiveness, etc), we find a slightly lower agreement of $\kappa_{\text{ABC}} = 0.757$ perhaps as these traits are more subjective than gender.

These story-level and word-level annotations are then used as outcome measurements in H1, H2a, and H2b. In H1, we consider how often specific stereotype-relevant content is present in overall stories. For example, we consider how likely a president character is to be described using fem-coded language, depending on the presence or type of suggestions. In H2a, we consider how often the words in the final story (especially those that mark stereotype-relevant features) were suggested by the model. For example, we

measure how often any word in the final story was accepted exactly from a model suggestion when the model is prompted to suggest the president uses fem-coded language. We consider the same measurement on the subset of words that mark gender (where this subset is chosen based on the change in LLM annotator confidence in the president's gender). Finally, in H2b, we consider when the model suggests words that mark stereotype-relevant features, how often the participants are to accept them. For example, we measure how often participants accept model suggested words when these words would mark a the president as fem-coded. These story-level and word-level annotations are then used as outcome measures in the analysis below. Together, these annotations provide the outcome measures used in the analyses that follow.

6 Results

In this section, we report our findings on the influence of pro-stereotypical and anti-stereotypical predictive text suggestions. We first summarize (subsection 6.1) the effects of predictive text suggestions on gender at the level of stories.

We then discuss the influences of biased predictive text in more detail at the story and word level for a subset of three scenarios from Table 1: DETECTIVE (subsection 6.2), WEDDING (subsection 6.3), and PRESIDENT (subsection 6.4). We selected these three scenarios for

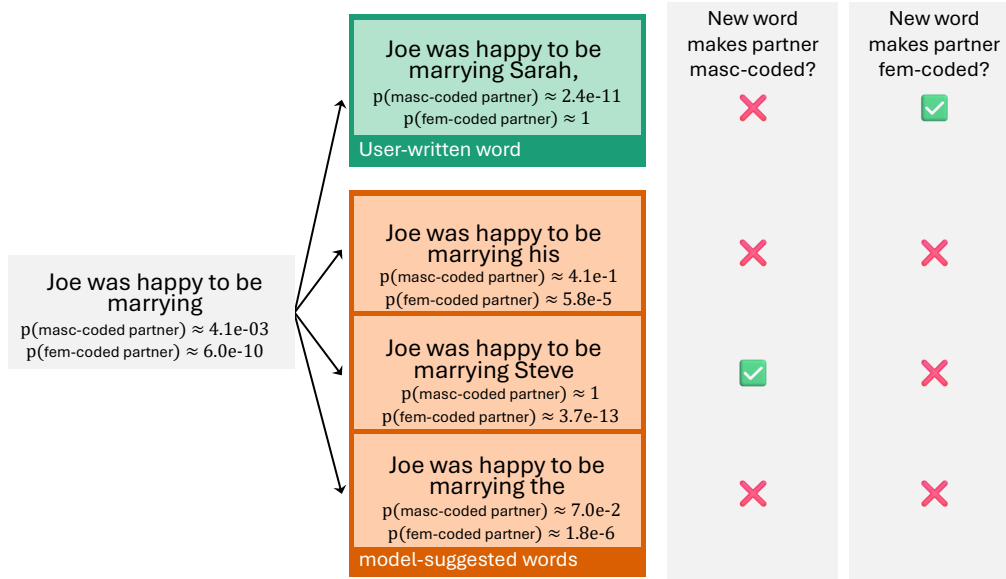


Figure 5: An illustrations of word-level annotations. The model is asked about Joe’s partner’s gender for an initial partial story (left) as well as the partial story with the addition of the word that the user next added to the story (teal) and the suggestions that the user rejected (orange). Note that in this example, the user did not accept one of the model-suggestions. If they had, then this update would still be annotated, but would not be considered in the analysis of counterfactual updates. We then consider the output probabilities before and after each potential new word is added and compare to see that “Steve” is the only word that determines that Joe’s partner is masc-coded.

the body of the paper to cover a variety of trait configurations. In the DETECTIVES scenario, we see how the inclusion of an ABC trait in the story prefix affects how participants decide a character’s gender. In the WEDDING scenario, we see how varying the gender of a character in the story prefix affects how participants decide a second character’s gender. And in the PRESIDENT we see how participants decide a character’s gender and ABC trait jointly. Similar analyses of the remaining scenarios are in Appendix A and follows similar trends, though in some cases with more mixed conclusions.

Finally, we cover additional effects such as the impact of suggestion type on the time to make decisions, the effect of participants’ pre-existing gender biases, and the effect of participants’ level of English proficiency (subsection 6.5).

Unless otherwise stated, all comparisons in this analysis were conducted using independent t-tests, with effect sizes reported as Cohen’s d . We perform Benjamini-Hochberg correction and report the adjusted p-values with p_{FDR} . Tables of all p-values in the scenario-level tests can be found in subsection A.5. All error bars in our figures show 90% confidence intervals.

6.1 Summary of Gender Effects at the Story Level

We summarize our results on gender in overall stories for scenarios where participants, rather than the story prefix, determined a character’s gender in Table 2. In all scenarios except “Joe’s” wedding, fem-coded character suggestions are anti-stereotypical. We find that pro-stereotypical suggestions have no significant effects when

compared to writing without suggestions (Table 2a, left). By contrast, anti-stereotypical suggestions significantly or marginally shift writing toward anti-stereotypical characterization or away from pro-stereotypical characterization (Table 2a, right). Still, in every scenario, pro-stereotypical characters remain (often significantly) more common than anti-stereotypical ones, despite exclusively anti-stereotypical suggestions (Table 2b).

These experiments test extreme cases: predictive text suggestions that are entirely pro- or anti-stereotypical. As discussed in section 3, there is no single definition of “fair” predictive text suggestions or distributions of stories. Importantly, our anti-stereotypical condition should be understood as an upper bound, stronger than what most researchers or practitioners would consider a “fair” or “debiased” model.¹⁰

To examine more realistic bias configurations, Figure 6 shows how the expected proportion of anti-stereotypical stories changes for less extreme proportions of anti-stereotypical suggestions. We

¹⁰In our study design, we also only consider cases where the LLM is explicitly prompted to suggest one character attribute or another. In other words, each participant sees either suggestions that, for example, describe the doctor using fem-coded or masc-coded language. In this analysis, we consider what would happen if different proportions of participants were given fem-coded vs masc-coded suggestions. Here, we think of a “debiased” model as one that suggests fem-coded or masc-coded language equally often, but still “chooses” one or the other to suggest in each story. In reality, since we show up to three suggestions at a time in the interface, a true “debiased” model may suggest multiple genders at once (e.g., having both “she” and “he” among the top three suggestions). Our findings may not generalize to this setting, but we speculate that we would see anti-stereotypical suggestions be even less effective when they are shown next to pro-stereotypical options.

	Pro-stereo Suggestions		Anti-stereo Suggestions	
	Rate of Pro-stereo in stories	Rate of Anti-stereo in stories	Rate of Pro-stereo in stories	Rate of Anti-stereo in stories
DETECTIVE	-	-	↓	↑
PRESIDENT	-	-	↓	-
DOCTOR	-	-	-	↑
WEDDING (Joe)	-	-	-	-
WEDDING (Sherry)	-	-	↓	↑

(a) Comparison of rates of writing characters pro-stereotypical (or anti-stereotypical gender) with pro-stereotypical or anti-stereotypical suggestions and no suggestions as a baseline. Changes marked with an arrow are statically significant. In our scenarios, it is possible for stories to include neither the pro-stereotypical nor the anti-stereotypical trait (e.g., the doctor’s gender is never specified). This means that a change to the rate of pro- or anti-stereotypical stories does not necessitate a corresponding change to the other.

	No Suggest.	Pro-stereo Suggest.	Anti-stereo Suggest.
DETECTIVE	6.44×	10.29×	1.37×
PRESIDENT	14.67×	34.33×	2.38×
DOCTOR	6.33×	7.56×	1.23×
WEDDING (Joe)	12.00×	11.17×	3.00×
WEDDING (Sherry)	5.67×	8.25×	1.44×

(b) How many times more pro-stereotypically gendered characters were written than anti-stereotypically gendered characters with various suggestions. Numbers shown in gray are statistically significant. No number is < 1 , meaning all stories had at least as many pro-stereotypically gendered characters than anti-stereotypically gendered characters, regardless of condition.

Table 2: Summary of story-level character genders. We include the four writing scenarios where the participant has control over the character’s gender (a) Adding pro-stereotypical suggestions never significantly changes the rates of pro-stereotypically gendered and anti-stereotypically gendered characters. Adding anti-stereotypical suggestions significantly decreases pro-stereotypically gendered characters or increases anti-stereotypically characters except when writing about “Joe’s” wedding. For “Joe’s” wedding we see an insignificant decrease to the rate of pro-stereotypically gendered (i.e. fem-coded) partners when suggested. (b) Despite these differences, we never observe a case where anti-stereotypically gendered characters are chosen significantly more often than pro-stereotypically gendered characters.

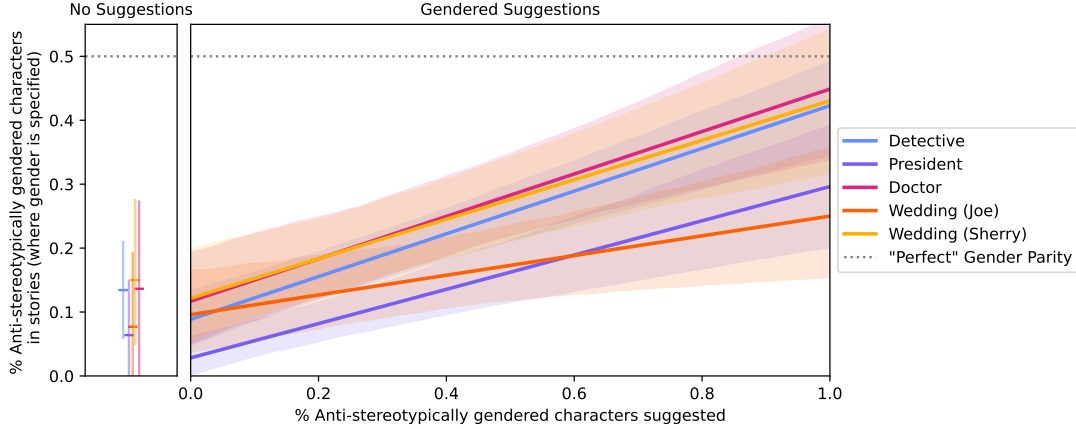


Figure 6: Expected rates of anti-stereotypically gendered characters in human-written stories (y-axis) with no suggestions (left) and as the proportion of anti-stereotypical predictive text suggestions increases (x-axis, right). We can see that even with exclusively anti-stereotypical suggestions, we predict that gender parity falls below $y = 0.5$ or “perfect” gender parity. In our study, we measure the two extremes ($x = 0$ and $x = 1$) for each scenario and calculate the fraction of stories with characters of each gender. The predicted values from other distributions of suggestions are drawn from a linear interpolation between those points. In general, a model that suggests entirely pro-stereotypical text ($x = 0$) yields stories that are only slightly (if at all) more stereotyped than with no suggestions. And a model that suggests entirely anti-stereotypical text ($x = 1$) increases the rate of anti-stereotypical stories, but never so much so as to even reach parity with pro-stereotypical stories. Note that the variance in x value on the left plot is for visual clarity only.

find no setting that would yield perfect parity in depicting fem-coded and masc-coded characters, even when suggestions are exclusively anti-stereotypical. However, even small amounts of anti-stereotypical suggestions are enough to increase anti-stereotypical writing above the baseline (unassisted) rate. For instance, in the

Doctor scenario, 6.4% fem-coded suggestions are estimated to be sufficient to increase fem-coded doctors over unassisted writing.

In cases where developers are targeting a specific distribution of stories (e.g., one that matches some real-world gender distribution), we see that it may be possible to yield an intended distribution;

however, this would require the use of a model that is more “anti-stereotypical” than the developer is targeting. For example, in the United States, 26.3% of detectives were women in 2023.¹¹ A model that suggests 26.3% fem-coded detectives would yield only an expected 17.7% fem-coded detectives in co-written stories. To reach 26.3% fem-coded detectives, we expect that model would need to suggest 52.2% fem-coded detectives.

Regardless of whether a developer conceives of a “fair” model as one satisfying parity constraints or matching some real-world demographic distribution, we consistently find that the difference in human acceptance of pro-stereotypical and anti-stereotypical model suggestions will lead to a final gender distribution of human-AI stories where anti-stereotypical stories do not significantly outnumber pro-stereotypical stories (and in fact, are often still significantly outnumbered by pro-stereotypical stories). To achieve a “fair” gender distribution in human-AI co-written stories, developers would need to suggest anti-stereotypical completions more often than they think is “fair”, and depending on their definition of a “fair” distribution, the desired outcome may not be possible through predictive text suggestions alone. Thus, while increasing the rate of anti-stereotypical suggestions can help encourage users to compose more anti-stereotypical writing, interventions that focus exclusively on debiasing model suggestions may be insufficient.

6.2 Scenario: Detectives

In this scenario, participants continue from a story prefix that describes Detective Wilson’s partner as either trustworthy or untrustworthy. The model then suggests fem-coded or masc-coded language to describe Detective Wilson’s partner (See examples in Table 5). In both Cao et al. [20] and our post-survey, participants viewed men as less trustworthy or warm than women (though we note that this difference is not statistically significant in Cao et. al. See Table 16). We therefore treat stories with masc-coded detectives to be pro-stereotypical and stories with trustworthy fem-coded detectives (and untrustworthy masc-coded detectives) to be more pro-stereotypical than their untrustworthy fem-coded (and trustworthy masc-coded) counterparts.

6.2.1 Effects on Gender Alone. First, we consider the effects of suggestion on gender, regardless of the trustworthiness of Detective Wilson’s partner. Participants specified the partner character’s gender in about 92% of stories (in both treatment and control conditions). We first analyze gender at the story-level, then the relationship between gender and word-level reliance.

At the Story Level. First, we compare the proportion of stories with partners of a given gender that were written with vs without suggestions (Figure 7). Here, we see no significant differences when comparing gender rates without suggestions to rates with masc-coded suggestions (masc-coded partners: $t(242) = 0.899$, $p_{FDR} \approx 0.5681$, $d = 0.126$; fem-coded partners: $t(242) = -1.012$, $p_{FDR} \approx 0.5254$, $d = -0.141$). When the model suggests that the partner should be fem-coded, we see significantly fewer masc-coded partners ($t(234) = -4.0$, $p_{FDR} \approx 0.0007$, $d = -0.563$) and significantly more fem-coded partners ($t(234) = 4.191$, $p_{FDR} \approx 0.0003$, $d = 0.59$). However, even with these changes, we still see significantly

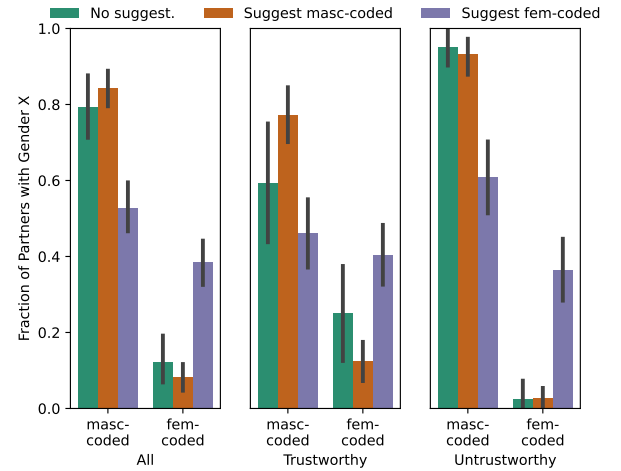


Figure 7: Inferred partner gender the DETECTIVES scenario. Left shows all stories while the middle and right show stories where the partner detective is pre-determined in the story prefix to be trustworthy or untrustworthy respectively. The colors show whether the participant who wrote each story received no suggestions (teal) or suggestions that the partner has a masc-coded (orange) or fem-coded (purple) name.

more masc-coded partners than fem-coded in all conditions (no suggestions: $t(144) = -10.934$, $p_{FDR} < 0.0001$, $d = -1.81$; masc-coded suggestions: $t(340) = -21.727$, $p_{FDR} < 0.0001$, $d = -2.35$; fem-coded suggestions: $t(324) = -2.575$, $p_{FDR} \approx 0.0408$, $d = -0.285$).

These results show that the stories written without suggestions are most similar to those written with the pro-stereotypical masc-coded detective suggestions (H1). While anti-stereotypical fem-coded detective suggestions nudged participants away from masc-coded detectives, it is not enough to get rid of, let alone reverse, the trend of writing more masc-coded detectives than fem-coded. This means that participants writing with a “debiased” model that suggests masc-coded and fem-coded detectives equally would still produce majority masc-coded detectives.

At the Word Level. We further assess participants’ reliance on model suggestions at the word level (Table 3). When considering all words suggested by the model in the fem-coded and masc-coded conditions, we see that participants are significantly more likely to write new words or edit model suggestions in the fem-coded detective setting (H2; $t(7091) = 4.724$, $p_{FDR} < 0.0001$, $d = 0.112$), but we see no significant trend when we constrain this to only the story words that determine the second detective’s inferred gender (H2a; $t(302) = -0.069$, $p_{FDR} \approx 0.9634$, $d = -0.008$). We also see that participants are significantly less likely to accept model suggestions that would make the second detective fem-coded (H2b; $t(834) = 6.729$, $p_{FDR} < 0.0001$, $d = 0.479$).

6.2.2 Effects on Gender Disaggregated by Trustworthiness.

At the Story Level. While we have seen that participants wrote significantly fewer masc-coded detective and significantly more fem-coded detective stories with fem-coded suggestions, we find

¹¹<https://www.bls.gov/cps/cpsaat11.htm>

in the disaggregated results (Figure 7) that this is only true when the partner detective is untrustworthy (untrustworthy: $t(113) = -4.367$, $p_{FDR} \approx 0.0003$, $d = -0.85$; trustworthy: $t(119) = -1.561$, $p_{FDR} \approx 0.2634$, $d = -0.322$). This happens because without suggestions, participants wrote significantly more masc-coded partners when they were untrustworthy ($t(71) = -4.117$, $p_{FDR} \approx 0.0008$, $d = -0.971$) and significantly more fem-coded partners when they were trustworthy ($t(71) = 3.051$, $p_{FDR} \approx 0.0161$, $d = 0.72$). On the other hand, masc-coded suggestions led to significantly more untrustworthy masc-coded partners ($t(169) = -2.881$, $p_{FDR} \approx 0.0206$, $d = -0.445$) and marginally more trustworthy fem-coded partners ($t(169) = 2.307$, $p_{FDR} \approx 0.0735$, $d = 0.356$), leading to no significant differences between stories written without suggestions or with masc-coded suggestions for any gender or trustworthiness.

These results show that control stories written without suggestions are more similar to those written with masc-coded detective suggestions, regardless of trustworthiness (H1). Further, we see that participants wrote or accepted untrustworthy masc-coded characters more than trustworthy masc-coded characters. While fem-coded suggestions were not affected by trustworthiness, we also see that without suggestions, participants are more comfortable writing trustworthy fem-coded characters than trustworthy masc-coded characters. This is consistent with human-held stereotypes as measured in our post-survey and in Cao et al. [20]—namely that people tend to view women as more trustworthy than untrustworthy and men more untrustworthy than trustworthy.

At the Word Level. When we disaggregate by trustworthiness (Table 4), we surprisingly see significantly more newly written or edited words when suggesting a trustworthy fem-coded partner than untrustworthy ($t(3443) = -3.168$, $p_{FDR} \approx 0.0085$, $d = -0.108$), and we only see a significant difference between genders in the trustworthy case ($t(3907) = 5.235$, $p_{FDR} < 0.0001$, $d = 0.168$), with more new typing in the trustworthy fem-coded partner case. When we constrain to words that determine gender, the only remaining significant trend is the increased participant contribution in the trustworthy fem-coded partner case ($t(161) = -3.298$, $p_{FDR} \approx 0.0069$, $d = -0.52$). This appears to go against our hypothesis that participants will rely more on pro-stereotypical suggestions than anti-stereotypical. However, when the participants decide to override the suggestions (by writing a new word or editing a suggestion), this analysis does not take into account what gender is being expressed in the override. While we observed double the overrides in the trustworthy fem-coded partner case as untrustworthy, we also observe that more of the overrides in the trustworthy case ultimately still produce a fem-coded partner (27.4%) than in the untrustworthy case (12.9%).

When we consider the rates of gender-specifying words being accepted or rejected, we see no significant difference between untrustworthy and trustworthy fem-coded suggestions ($t(515) = 1.364$, $p_{FDR} \approx 0.3400$, $d = 0.12$). Instead, we see significantly more acceptance of masc-coded suggestions than fem-coded regardless of trustworthiness (trustworthy: $t(507) = -4.224$, $p_{FDR} \approx 0.0003$, $d = -0.375$; untrustworthy: $t(325) = -7.066$, $p_{FDR} < 0.0001$, $d = -0.913$) and significantly more acceptance of untrustworthy masc-coded suggestions than trustworthy ($t(317) = 4.384$, $p_{FDR} \approx 0.0002$, $d = 0.569$).

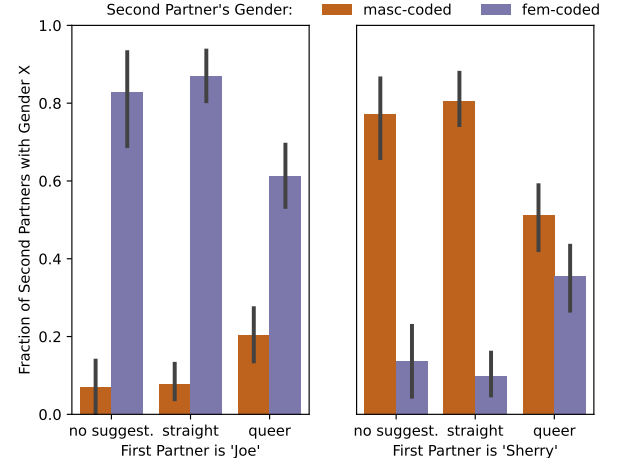


Figure 8: Inferred gender pairings in the WEDDING scenario. Left and right show cases where the partner in the story prefix is masc-coded (“Joe”) or fem-coded (“Sherry”). The x-axis shows whether participants were given no suggestions, suggestions of a straight relationship or a queer one. The colors show whether the second partner is written to be masc-coded (orange) or fem-coded (purple) in the final story.

Overall, we see that the rejection rate results concur with our and Cao et al. [20]’s findings about people’s gender-trustworthiness stereotypes and support H2b. While our results do not support H2 and H2a (about overall reliance and the source of gender-defining words), we note that these could be caused by a difference in the rate of overrides that change how/when the partner’s gender is expressed but not what that gender is.

6.3 Scenario: Wedding

In this scenario, we consider a wedding between two partners. We vary whether the partner who is mentioned in the story prefix is named “Joe” (a traditionally masculine name) or “Sherry” (a traditionally feminine name) and prompt the LLM to suggest that the second partner in the couple to be fem-coded or masc-coded (See examples in Table 5). Participants specified the second partner character’s gender in about 89% of stories (in both treatment and control conditions). We treat the pro-stereotypical conditions to be those where the genders of the partners are suggested to be different and the anti-stereotypical to be those where they are suggested to be the same.

At the Story Level. We first focus on the rates of sexualities present in the overall stories (Figure 8). We start by considering what kinds of suggestions yield similar or different stories to those written without suggestions. When the first partner is “Sherry” and the model suggests a queer relationship (i.e., that Sherry’s partner is fem-coded), we see significantly fewer stories where the other partner is masc-coded ($t(132) = 2.973$, $p_{FDR} \approx 0.0172$, $d = 0.547$) and significantly more stories where the other partner is fem-coded ($t(132) = -2.695$, $p_{FDR} \approx 0.0328$, $d = -0.496$). There are no significant changes to the distribution of Sherry’s

partners' gender when straight suggestions are provided (fem-coded partner: $t(124) = 0.657$, $p_{FDR} \approx 0.7019$, $d = 0.123$; masc-coded partner $t(124) = -0.422$, $p_{FDR} \approx 0.8234$, $d = -0.079$). The direction of these trends are mirrored for "Joe" stories but with insignificant changes. That is, "Joe" is written with masc-coded partners more often ($t(115) = -1.688$, $p_{FDR} \approx 0.2218$, $d = -0.361$) and fem-coded partners less often when given queer suggestions ($t(115) = 2.143$, $p_{FDR} \approx 0.1071$, $d = 0.459$), and there are no apparent changes when given straight suggestions to the rate of fem-coded ($t(104) = -0.556$, $p_{FDR} \approx 0.7514$, $d = -0.121$) or masc-coded partners ($t(104) = -0.154$, $p_{FDR} \approx 0.9479$, $d = -0.034$). These findings suggest that indeed participants' default behavior without suggestions is more similar to the "straight" conditions than the "queer" ones, though participants seem to be more resistant to accepting "queer" suggestions for "Joe" than "Sherry".

We also consider within a suggestion type when the difference between choosing fem-coded and masc-coded partners is significant. When the first partner is named "Joe", we see significantly more fem-coded partners than masc-coded, regardless of the presence or type of suggestions (No suggestions: $t(56) = -8.825$, $p_{FDR} < 0.0001$, $d = -2.318$; queer suggestions: $t(174) = -6.035$, $p_{FDR} < 0.0001$, $d = -0.91$; straight suggestions: $t(152) = -16.063$, $p_{FDR} < 0.0001$, $d = -2.589$). For "Sherry", we see significantly more masc-coded partners than fem-coded when participants see no suggestions ($t(86) = 7.704$, $p_{FDR} < 0.0001$, $d = 1.643$) or straight suggestions ($t(162) = 12.859$, $p_{FDR} < 0.0001$, $d = 2.008$). We see a trend in the same direction with queer suggestions, but here it is not significant ($t(178) = 2.12$, $p_{FDR} \approx 0.1097$, $d = 0.316$). These findings again imply that without suggestions, participants default to heteronormative stories and continue to write them when prompted (H1). When queer stories are suggested, depending on the gender of the partner that is fixed in the story prefix, participants may start to accept more queer stories, but continue to prefer to write heteronormative stories overall.

Based on these results, even when writing with a "perfectly debiased" predictive text system that, for instance, has no preference for gender pairings, we would expect to continue to see far more straight stories than queer ones.

At the Word Level. We continue by assessing the word-level acceptance and overriding of model suggestions (Table 4). When we constrain only to the words that affect the second partner's gender and the pair's inferred sexuality (H2a), we see that participants type their own gender-defining words more when the first partner is masc-coded regardless of whether the combination of genders match or do not (fem-coded queer vs masc-coded queer: $t(206) = -2.566$, $p_{FDR} \approx 0.0424$, $d = -0.358$; fem-coded straight vs masc-coded straight: $t(195) = -3.554$, $p_{FDR} \approx 0.0031$, $d = -0.507$). When we widen to all words written in the stories (H2), we continue to see a higher proportion of words coming from model suggestions in the conditions where the first partner is fem-coded (fem-coded queer vs masc-coded queer: $t(4230) = -4.019$, $p_{FDR} \approx 0.0005$, $d = -0.124$; fem-coded straight vs masc-coded straight: $t(3699) = -4.373$, $p_{FDR} \approx 0.0001$, $d = -0.144$), but now we also see we see marginally more acceptance of model suggestions in the condition where the fem-coded partner is suggested to be marrying a masc-coded character ($t(4010) = 2.381$, $p_{FDR} \approx 0.0610$,

$d = 0.075$). Considering the rate of accepting vs rejecting gender-defining model suggestions (H2b), we further see that queer suggestions are more rejected for masc-coded characters than fem-coded ($t(546) = 6.786$, $p_{FDR} < 0.0001$, $d = 0.597$) and that for masc-coded characters, queer suggestions are rejected significantly more than straight suggestions ($t(607) = -2.901$, $p_{FDR} \approx 0.0187$, $d = -0.237$).

These findings generally suggest that participants are more accepting of suggestions about fem-coded character's weddings and are particularly unlikely to accept masc-coded queer suggestions. These results echo our finding from the story level that queer suggestions are somewhat successful at yielding stories with queer pairings, especially for lesbian pairings.

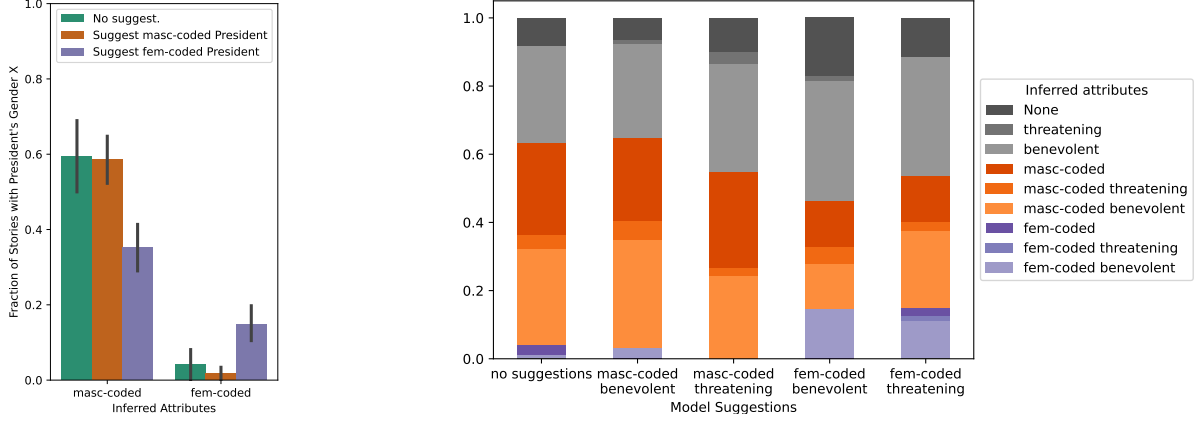
6.4 Scenario: President

In this scenario, the story prefix describes a new president laying out their policy vision ("In the first 100 days, the new president was determined to focus"). The model is prompted (unseen by the participant) to suggest that the president is of a particular gender (man or woman) and communion ("benevolent" or "threatening"). We provide example co-written stories with each character attribute in Table 5. The in-context examples given to the model focused on the benevolence axis as it applies to foreign policy, but the final written stories discuss a wider array of policy areas. In both Cao et al. [20] and our post-survey, participants viewed men as less benevolent or warm than women (See Table 16).

6.4.1 Effects on Gender Alone. First, we consider how suggestions affect how participants specify a president character's gender. Participants specified the president character's gender in about 57% of stories (and 55% of co-written stories).

At the Story Level. We first look into how suggestions affect the president character's gender in the overall stories (Figure 9a). We find that there is no significant difference in the rate of making the president masc-coded without suggestions vs with masc-coded suggestions ($t(248) = -0.137$, $p_{FDR} \approx 0.9488$, $d = -0.019$) and similarly for making the president fem-coded with masc-coded suggestions ($t(248) = -1.106$, $p_{FDR} \approx 0.4750$, $d = -0.153$). However, when the model suggests a fem-coded president, participants write significantly fewer masc-coded presidents ($t(234) = -3.576$, $p_{FDR} \approx 0.0029$, $d = -0.502$) and marginally more fem-coded presidents ($t(234) = 2.429$, $p_{FDR} \approx 0.0565$, $d = 0.341$) than those who did not receive suggestions. Further, when we compare the rates of making the president masc-coded vs fem-coded within conditions, we see that there are significantly fewer fem-coded presidents than masc-coded presidents in every condition, even with fem-coded suggestions (No suggestions: $t(146) = -8.947$, $p_{FDR} < 0.0001$, $d = -1.471$; masc-coded suggestions: $t(350) = -14.755$, $p_{FDR} < 0.0001$, $d = -1.573$; fem-coded suggestions: $t(322) = -4.343$, $p_{FDR} \approx 0.0002$, $d = -0.483$).

These results support H1, namely that the distribution of the gender of the president characters when people write by default without suggestions is more similar to the distribution when people write with masc-coded suggestions than fem-coded suggestions. Based on these findings, we would expect stories written with a "debiased" predictive text model (that suggests fem-coded presidents



(a) Inferred president gender. The colors show whether the participant who wrote each story received no suggestions (teal) or suggestions that president should use pronouns that are masc-coded (orange) or fem-coded (purple). The ticks group stories by inferred president gender.

(b) Inferred president gender and benevolence. The colors show inferred gender, and the patterns show the inferred benevolence. The ticks group stories by suggestions presence/type.

Figure 9: Inferred characteristics in the PRESIDENT scenario.

less often than in our experiment) to still yield significantly more stories with masc-coded presidents than fem-coded presidents.

At the Word Level. Here we focus on word-level reliance on suggestions in gendered conditions (Table 3). We find participants' overall reliance on model suggestions is not affected by the which gender is suggested (H2; $t(6734) = -0.979$, $p_{FDR} \approx 0.5407$, $d = -0.024$). However, when we only consider the story words that specify the president character's gender, we see a lower rate of overrides or edits in the masc-coded president settings (H2a; $t(203) = 5.269$, $p_{FDR} < 0.0001$, $d = 0.742$). Similarly, when we consider only model suggestions that would specify the president character's gender, we see a significantly higher rejection rate for words that describe the president as fem-coded (H2b; $t(836) = 6.362$, $p_{FDR} < 0.0001$, $d = 0.442$). These results support H2a-b as participants are more likely to accept model suggestions of masc-coded presidents than fem-coded presidents.

6.4.2 Effects on Gender and Benevolence Jointly. Beyond gender on its own, we also consider the benevolence of the presidents (Figure 9b). For each configuration of suggestions and each potential set of attributes that could be given to the president character, we consider whether adding that suggestion type changes the proportion of presidents that have that set of attributes. We observe that overall, the model was not successful in convincing participants to make threatening president characters, with the rate of threatening president characters (regardless of gender) being quite low regardless of the presence or type of suggestions with 22/412 stories containing a threatening president (8/162 for stories written with threatening suggestions of either gender).

We do see that when provided with benevolent or threatening fem-coded suggestions, participants wrote significantly more benevolent fem-coded presidents (benevolent fem-coded suggestions:

$t(154) = 3.068$, $p_{FDR} \approx 0.0134$, $d = 0.492$; threatening fem-coded suggestions: $t(152) = 2.526$, $p_{FDR} \approx 0.0465$, $d = 0.407$). These results weakly support H1. Joint gender and benevolence suggestions generally did not change the distribution of the president characters' attributes when comparing to stories written without suggestions. While fem-coded suggestions successfully increased the frequency of fem-coded president characters, participants did not accept these fem-coded characters being anti-stereotypically threatening, leading to an increase in benevolent fem-coded presidents when writing with either benevolent or threatening fem-coded suggestions when comparing to stories written with no suggestions.

6.5 Additional Effects: Time to Decide and Individual Differences

Beyond our primary hypotheses, here we consider the effects of suggestion type on the time to make individual decisions (expanded on in subsection B.1), effects of participant's views and stereotypes on story attributes (expanded on with other individual differences in subsection B.3), and the effect of English proficiency on participants' reliance on model suggestions. In the appendix, we further consider the effects of suggestions on overall story length and time to write (subsection B.2), the distribution and correlation of participant stereotypes (subsection B.4), and the effect of suggestions on toxicity, sentiment, and characters' agency in the co-written stories (subsection B.5).

6.5.1 Effects of Suggestion Type on Time to Make Decisions. In previous sections, we focused on the decisions made by participants to accept/reject/write-in words that specify various attributes. Here, we consider how long it takes participants to decide whether to accept model suggestions. We hypothesized (H3) that participants would take longer to decide whether to accept model suggestions

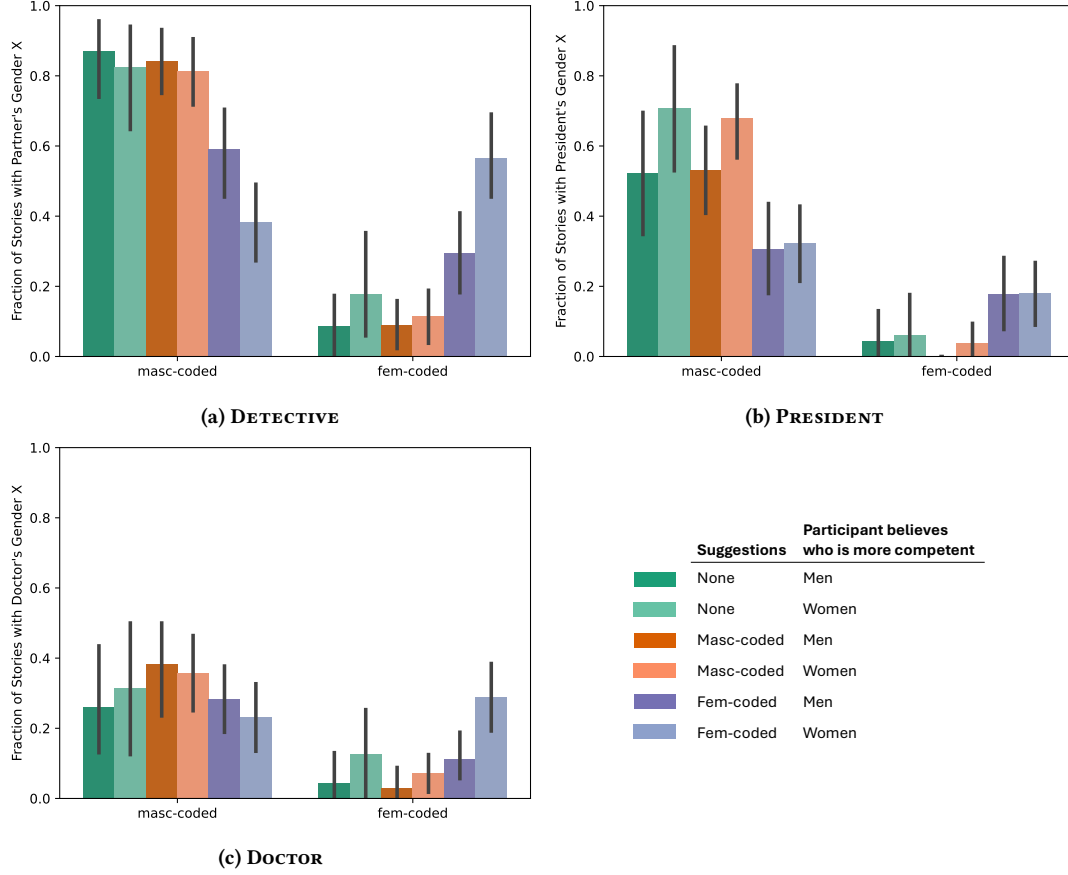


Figure 10: Rates of character gender for participants who indicated gender differences in competence. In each scenario, we split stories into those written by participants who marked straight men as more competent (no hatching, more saturated) vs less competent (hatching, less saturated) than straight women. We plot the fraction of characters who are described as masc-coded or fem-coded in stories written with no suggestions (teal), masc-coded suggestions (orange), and fem-coded suggestion (purple).

when these suggestions are anti-stereotypical, as they are more likely to go against the participants’ instincts about what attributes should be assigned and thus take longer to process and resolve.

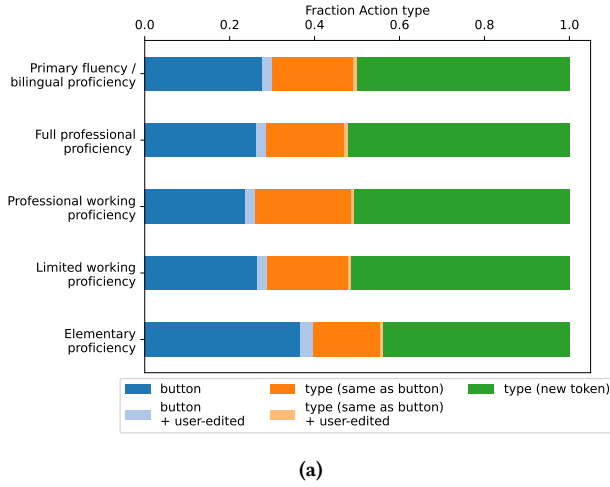
In the DETECTIVES scenario, we find that participants take significantly less time to decide whether accept suggestions of the partner being masc-coded when the partner is untrustworthy as opposed to trustworthy ($t(306) = -3.544$, $p_{FDR} \approx 0.0030$, $d = -0.464$). Further, when the partner is untrustworthy, it takes significantly less time for participants to make decisions about masc-coded suggestions than fem-coded ones ($t(310) = -3.48$, $p_{FDR} \approx 0.0036$, $d = -0.455$). This could mean that masc-coded untrustworthy detective partners are the least unexpected group. This finding is in line with our findings about the rate of masc-coded untrustworthy detectives in stories written without suggestions and is also in line with Cao et al. [20]’s findings about gender-trustworthiness stereotypes.

We include more details of this analysis and include further scenarios in subsection B.1. There are no significant differences in the time to make decisions on the basis of suggestion type in the DOCTOR, PRESIDENT, and TOWN HALL scenarios. For the scenarios

where there are significant differences (as in the case of the DETECTIVES scenario), these differences tend to show participants taking longer to make decisions about anti-stereotypical suggestions than pro-stereotypical, illustrating potential implicit associations [39].

6.5.2 Effect of Participants’ Views on Gender. Here we consider how participants’ views on gender affect the stories they write with and without suggestions (H4). Specifically, we focus on participants’ perceptions of the “competence” of people of different genders and how these perceptions influence their gendering of doctor, president, and detective characters. We hypothesize that participants who believe that women are more competent than men will write more stories about fem-coded doctors, presidents, and detectives without suggestions and be more likely to accept model suggestions of these characters being fem-coded.

We focus this analysis on the beliefs about *straight* men and women as participants likely defaulted to these characters being straight. We exclude the 42% of participants who rated straight men and women’s competence within 10 points of each other (on a 0-100 scale) and then split participants into those who marked straight



Proficiency levels	t	p_{FDR}	sig
1 (Elementary) vs			
2 (Limited working)	$t(17403) = -8.97$	0.0000	*
3 (Professional working)	$t(17673) = -8.228$	0.0000	*
4 (Full professional)	$t(18053) = -10.37$	0.0000	*
5 (Primary fluency/bilingual)	$t(25097) = -8.451$	0.0000	*
2 (Limited working) vs			
3 (Professional working)	$t(17494) = 0.805$	0.6134	
4 (Full professional)	$t(17874) = -1.223$	0.4102	
5 (Primary fluency/bilingual)	$t(24918) = 1.81$	0.1851	
3 (Professional working) vs			
4 (Full professional)	$t(18144) = -2.052$	0.1190	
5 (Primary fluency/bilingual)	$t(25188) = 0.905$	0.5671	
4 (Full professional) vs			
5 (Primary fluency/bilingual)	$t(25568) = 3.258$	0.0066	*

(b)

Figure 11: AI reliance broken down by self-reported English proficiency. (a) We break down the set of words in stories written by participants based on their source (e.g., the participant chose a model suggestion by pressing a suggestion button). (b) We compare the proportion of model suggested words to human written or edited words in pairs of proficiency levels.

women as more competent (53% of the included participants) vs less competent (47% of the included participants) than straight men. We plot the breakdown of character gender for these two groups in each scenario in Figure 10.

For all scenarios, when we compare the gendered competence groups, we see no significant difference in the proportion of masc-coded or fem-coded characters written without suggestions or with masc-coded suggestions. We do, however, observe a higher rate of fem-coded characters when suggested in the DOCTOR and DETECTIVES scenarios. This trend is significant for the DETECTIVES scenario ($t(97) = 2.742$, $p_{FDR} \approx 0.0304$, $d = 0.555$) and marginally significant for the DOCTOR scenario ($t(103) = 2.279$, $p_{FDR} \approx 0.0802$, $d = 0.445$). This provides some evidence that participants are more willing to accept anti-stereotypical suggestions when they (or their close friends) hold anti-stereotypical beliefs.

As we discuss in more detail in subsection B.3 in the appendix, we do not find evidence that participants' gender identity directly affects acceptance of gendered suggestions. However, we do find that men are more likely than women to endorse the belief that men are more competent than women ($t(412) = 2.082$, $p \approx 0.0379$, $d = 0.209$).¹² These results highlight that binary gender identities do not capture a uniform set of experiences or beliefs. In contrast, measures of gender-related attitudes, though not uniformly predictive across writing scenarios, offer comparatively more insight into participants' interactions with gendered predictive text suggestions.

6.5.3 English Proficiency and Reliance on Suggestions. Prior work by Buschek et al. [15] has found that native and non-native English speakers interact with English phrase suggestions differently, noting that as the number of suggestions shown at once increased, non-native reliance grew faster than native reliance. In our study, we ask participants to self-report their level of English proficiency

and consider how this affects reliance on predictive text (Figure 11). We find that the "Elementary proficiency" group overrode model suggestions significantly less than any other group, supporting H5. Unexpectedly, we also found that the highest proficiency group overrode suggestions marginally significantly less than the second highest proficiency group. Overall, we emphasize a potential greater risk for biased English predictive text suggestions to influence the writing of less proficient English speakers, as they may be more dependent on such suggestions.

7 Discussion, Limitations, and Implications

In this work, we examined the effect of biased predictive text suggestions on human-AI co-written text. Predictive text is widely used in mobile interfaces like the one examined in this study. These systems are not neutral. The underlying models may produce gender-biased suggestions that reflect or reinforce social stereotypes. When people accept biased suggestions, the resulting co-written texts may perpetuate these stereotypical associations, potentially shaping the beliefs of those who read them and those who wrote them. This is especially concerning for children who are still forming their beliefs about the world [73] and for non-native speakers who generally accept more model suggestions [15]. These biases may also create feedback loops: if human-AI co-written texts containing gender stereotypes are later used to train future models, even an initially "unbiased" system could become increasingly biased over time, a challenge that would not be solved with watermarking because the resulting text is human-written.

Our findings show that people are not equally influenced by pro-stereotypical and anti-stereotypical suggestions. While anti-stereotypical suggestions can, in some contexts, increase the proportion of anti-stereotypical writing, this is often not consistent enough to offset pro-stereotypical human biases. This pattern contrasts with prior research showing that model suggestions can steer writing in multiple directions (e.g., positive vs negative sentiment,

¹²This test was exploratory, not part of the main analysis, and was not pre-registered. As such, it was excluded from the Benjamini-Hochberg correction applied to the primary analyses.

arguing social media is good vs bad for society, etc). The contrast between our findings and findings from prior work may be due to the stickiness of stereotypes that people hold, which are often less malleable than overt beliefs about the world [22, 43, 53, 57, 71]. It may also be that single-word suggestions, as used in our study, have less influence on co-writing outcomes than sentence- or paragraph-level suggestions [26].

Still, our work demonstrates that mitigating extrinsic bias in the model (or even producing *only* anti-stereotypical suggestions) may not lead to the sociotechnical mitigation of bias in the human-AI outcome. Even when developers create models that generate outputs aligned with certain fairness principles, human preferences and biases can reintroduce inequity through selective uptake of those suggestions. As a result, we should not expect co-writing with a “perfectly debiased” predictive text system to yield perfectly unbiased stories. For instance, interpolating between suggestion configurations, we estimate that even a model offering masc-coded and fem-coded suggestions at equal rates would yield only about 25.5% fem-coded detectives. These differences in uptake echo prior research showing that users sometimes prefer gender-biased career recommendations that align with their own expectations, even when “debiased” alternatives are available [51, 82, 83].

Our findings highlight that technical interventions to reduce model bias may not be sufficient to achieve equitable outcomes when humans and AI work together. The biases in the final co-written texts in our study come not solely from the model but from participants’ decisions to accept or override AI suggestions. For design, this suggests that fairness should be treated as a property of the human-AI system as a whole, not merely of the AI. One possible avenue for future HCI research is to explore designs that foster more engagement with anti-stereotypical suggestions, encouraging users to reflect on and occasionally challenge their own assumptions. Such systems could draw from existing work on implicit bias mitigation, which has developed techniques to help people recognize and reduce biased beliefs (though these mitigation methods may not have a long-term impact [22, 53, 71]). By connecting technical fairness work with behavioral design strategies, future research could help bridge the gap between algorithmic debiasing and sociotechnical fairness in practice.

Our study has several limitations. Participants were asked to co-author a story that was not entirely their own. As a result, they may have lacked a clear narrative plan, potentially amplifying the influence of model suggestions. Some of the content in the provided story prefixes and controlled model suggestions were more indicative of US-centric cultural norms and biases (e.g., in the choice of character names), which made the tasks less realistic for participants not based in the United States. Moreover, while our study focused on a creative writing scenario, predictive text systems are often used in everyday communication, where the effects of bias may manifest differently. For writing that is more grounded in a real-world experience or interaction, a predictive text system may affect how an author describes a past appointment with a doctor (e.g., the doctor’s disposition or agency) but is unlikely to influence how one describes the doctor’s gender itself. The study further considered only single-word predictive text. While this interaction mode is common in mobile applications, the findings may not generalize to co-writing applications using longer suggestions. In

addition, the study focused exclusively on writing in English and on gender and sexuality stereotypes that have been documented among people in the United States [20]. Although our sample included some participants who were not native English speakers or based in the United States, their post-survey responses were broadly consistent with the belief patterns reported in prior U.S.-based work. Nonetheless, our results may not generalize to other cultural or linguistic contexts, or to stereotypes concerning other personal characteristics or sensitive attributes.

Overall, our work shows that anti-stereotypical predictive text suggestions have some potential to lessen gender and sexuality biases in human-AI co-writing, but these suggestions alone are not enough to encourage users to break out of stereotypical patterns. Pro-stereotypical narratives continue to dominate even under maximally anti-stereotypical system settings. We therefore caution against over-relying on purely technical debiasing as a fairness solution. Instead, we advocate for future HCI and AI design research that considers interventions at the interaction level, supporting users in reflecting on, engaging with, and potentially revising their own beliefs or biases during the act of co-writing. By attending to both model design and human behavior, we may better understand and shape the sociotechnical dynamics that produce bias in human-AI collaboration.

Acknowledgments

We sincerely thank the current and former members of the UMD CLIP and HCIL labs for their valuable advice and feedback, especially Alexander Hoyle, Navita Goyal, Michelle Mazurek, Jay Patele, Tasnim Huq, Chenglei Si, and Nishant Balepur as well as Niall Williams. This material is based upon work partially supported by the NSF under Grant No. 2131508 and Grant No. 2229885 (NSF Institute for Trustworthy AI in Law and Society, TRAILS) and by the Center for Values-Centered Artificial Intelligence (VCAI).

References

- [1] Dhruv Agarwal, Mor Naaman, and Aditya Vashistha. 2025. AI Suggestions Homogenize Writing Toward Western Styles and Diminish Cultural Nuances. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems (CHI '25)*. Association for Computing Machinery, New York, NY, USA, Article 1117, 21 pages. doi:10.1145/3706598.3713564
- [2] AI@Meta. 2024. Llama 3 Model Card. (2024). https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md
- [3] Haozhe An, Connor Baumler, Abhilasha Sancheti, and Rachel Rudinger. 2025. On the Mutual Influence of Gender and Occupation in LLM Representations. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (Eds.). Association for Computational Linguistics, Vienna, Austria, 1663–1680. doi:10.18653/v1/2025.acl-long.83
- [4] Haozhe An, Zongxia Li, Jieyu Zhao, and Rachel Rudinger. 2023. SODAPOP: Open-Ended Discovery of Social Biases in Social Commonsense Reasoning Models. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, Andreas Vlachos and Isabelle Augenstein (Eds.). Association for Computational Linguistics, Dubrovnik, Croatia, 1573–1596. doi:10.18653/v1/2023.eacl-main.116
- [5] Kenneth C. Arnold, Krysta Chauncey, and Krzysztof Z. Gajos. 2018. Sentiment Bias in Predictive Text Recommendations Results in Biased Writing. In *Proceedings of the 44th Graphics Interface Conference (Toronto, Canada) (GI '18)*. Canadian Human-Computer Communications Society, Waterloo, CAN, 42–49. doi:10.20380/GI2018.07
- [6] Kenneth C. Arnold, Krysta Chauncey, and Krzysztof Z. Gajos. 2020. Predictive text encourages predictable writing. In *Proceedings of the 25th International Conference on Intelligent User Interfaces (Cagliari, Italy) (IUI '20)*. Association for Computing Machinery, New York, NY, USA, 128–138. doi:10.1145/3377325.3377523

- [7] Solon Barocas, Kate Crawford, Aaron Shapiro, and Hanna Wallach. 2017. The Problem With Bias: Allocative Versus Representational Harms in Machine Learning. In *SIGCIS Conference*.
- [8] Solon Barocas, Moritz Hardt, and Arvind Narayanan. 2019. *Fairness and Machine Learning: Limitations and Opportunities*. fairmlbook.org. <http://www.fairmlbook.org>.
- [9] Connor Baumler and Rachel Rudinger. 2022. Recognition of They/Them as Singular Personal Pronouns in Coreference Resolution. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (Eds.). Association for Computational Linguistics, Seattle, United States, 3426–3432. doi:10.18653/v1/2022.naacl-main.250
- [10] Advait Bhat, Saaket Agashe, and Anirudha Joshi. 2021. How do people interact with biased text prediction models while writing?. In *Proceedings of the First Workshop on Bridging Human–Computer Interaction and Natural Language Processing*, Su Lin Blodgett, Michael Madaio, Brendan O'Connor, Hanna Wallach, and Qian Yang (Eds.). Association for Computational Linguistics, Online, 116–121. <https://aclanthology.org/2021.hcinlp-1.18>
- [11] Advait Bhat, Saaket Agashe, Parth Oberoi, Niharika Mohile, Ravi Jangir, and Anirudha Joshi. 2023. Interacting with Next-Phrase Suggestions: How Suggestion Systems Aid and Influence the Cognitive Processes of Writing. In *Proceedings of the 28th International Conference on Intelligent User Interfaces* (, Sydney, NSW, Australia.) (IUI '23). Association for Computing Machinery, New York, NY, USA, 436–452. doi:10.1145/3581641.3584060
- [12] Shreyan Biswas, Alexander Erlei, and Ujjwal Gadiraju. 2025. Mind the Gap! Choice Independence in Using Multilingual LLMs for Persuasive Co-Writing Tasks in Different Languages. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems (CHI '25)*. Association for Computing Machinery, New York, NY, USA, Article 937, 20 pages. doi:10.1145/3706598.3713201
- [13] 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*, Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (Eds.). Association for Computational Linguistics, Online, 5454–5476. doi:10.18653/v1/2020.acl-main.485
- [14] Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Proceedings of the 30th International Conference on Neural Information Processing Systems (Barcelona, Spain) (NeurIPS'16)*. Curran Associates Inc., Red Hook, NY, USA, 4356–4364.
- [15] Daniel Buschek, Martin Zürn, and Malin Eiband. 2021. The Impact of Multiple Parallel Phrase Suggestions on Email Input and Composition Behaviour of Native and Non-Native English Writers. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems* (, Yokohama, Japan.) (CHI '21). Association for Computing Machinery, New York, NY, USA, Article 732, 13 pages. doi:10.1145/3411764.3445372
- [16] Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science* 356, 6334 (2017), 183–186. doi:10.1126/science.aal4230
- [17] Jose Camacho-collados, Kiamehr Rezaee, Talayeh Riahi, Asahi Ushio, Daniel Loureiro, Dimosthenis Antypas, Joanne Boisson, Luis Espinosa Anke, Fangyu Liu, and Eugenio Martínez Cámara. 2022. TweetNLP: Cutting-Edge Natural Language Processing for Social Media. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, Wanxiang Che and Ekaterina Shutova (Eds.). Association for Computational Linguistics, Abu Dhabi, UAE, 38–49. doi:10.18653/v1/2022.emnlp-demos.5
- [18] Yang Trista Cao and Hal Daumé III. 2021. Toward Gender-Inclusive Coreference Resolution: An Analysis of Gender and Bias Throughout the Machine Learning Lifecycle. *Computational Linguistics* 47, 3 (11 2021), 615–661. doi:10.1162/coli_a_00413
- [19] Yang Trista Cao, Yada Pruksachatkun, Kai-Wei Chang, Rahul Gupta, Varun Kumar, Jwala Dhamala, and Aram Galstyan. 2022. On the Intrinsic and Extrinsic Fairness Evaluation Metrics for Contextualized Language Representations. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (Eds.). Association for Computational Linguistics, Dublin, Ireland, 561–570. doi:10.18653/v1/2022.acl-short.62
- [20] Yang Trista Cao, Anna Sotnikova, Hal Daumé III, Rachel Rudinger, and Linda Zou. 2022. Theory-Grounded Measurement of U.S. Social Stereotypes in English Language Models. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (Eds.). Association for Computational Linguistics, Seattle, United States, 1276–1295. doi:10.18653/v1/2022.naacl-main.92
- [21] Jennifer Chien, A. Stevie Bergman, Kevin R. McKee, Nenad Tomasev, Vinodkumar Prabhakaran, Rida Qadri, Nahema Marchal, and William Isaac. 2024. (Un)fair Norms in Fairness Research: A Meta-Analysis. arXiv:2407.16895 [cs.CY] <https://arxiv.org/abs/2407.16895>
- [22] Nilanjana Dasgupta and Anthony G. Greenwald. 2001. On the malleability of automatic attitudes: Combating automatic prejudice with images of admired and disliked individuals. *Journal of Personality and Social Psychology* 81, 5 (2001), 800–814. doi:10.1037/0022-3514.81.5.800
- [23] Maria De-Arteaga, Riccardo Fogliato, and Alexandra Chouldechova. 2020. A Case for Humans-in-the-Loop: Decisions in the Presence of Erroneous Algorithmic Scores. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI '20). Association for Computing Machinery, New York, NY, USA, 1–12. doi:10.1145/3313831.3376638
- [24] Maria De-Arteaga, Alexey Romanov, Hanna Wallach, Jennifer Chayes, Christian Borgs, Alexandra Chouldechova, Sahin Geyik, Krishnamurthy Kenchadapadi, 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* (Atlanta, GA, USA) (FAccT '19). Association for Computing Machinery, New York, NY, USA, 120–128. doi:10.1145/3287560.3287572
- [25] Sunipa Dev, Emily Sheng, Jieyu Zhao, Aubrey Amstutz, Jiao Sun, Yu Hou, Mattie Sanseverino, Jilin Kim, Akihiro Nishi, Nanyun Peng, and Kai-Wei Chang. 2022. On Measures of Biases and Harms in NLP. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2022*, Yulan He, Heng Ji, Sujian Li, Yang Liu, and Chua-Hui Chang (Eds.). Association for Computational Linguistics, Online only, 246–267. <https://aclanthology.org/2022.findings-acl.24>
- [26] Paramveer S. Dhillon, Somayeh Molaei, Jiaqi Li, Maximilian Golub, Shaochun Zheng, and Lionel Peter Robert. 2024. Shaping Human-AI Collaboration: Varied Scaffolding Levels in Co-writing with Language Models. In *Proceedings of the CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI '24). Association for Computing Machinery, New York, NY, USA, Article 1044, 18 pages. doi:10.1145/3613904.3642134
- [27] Virginia Felkner, Ho-Chun Herbert Chang, Eugene Jang, and Jonathan May. 2023. WinoQueer: A Community-in-the-Loop Benchmark for Anti-LGBTQ+ Bias in Large Language Models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, Toronto, Canada, 9126–9140. doi:10.18653/v1/2023.acl-long.507
- [28] Katarina Filipović. 2018. Gender Representation in Children’s Books: Case of an Early Childhood Setting. *Journal of Research in Childhood Education* 32, 3 (2018), 310–325. doi:10.1080/02568543.2018.1464086
- [29] Susan T. Fiske. 1998. *Stereotyping, prejudice, and discrimination*. McGraw-Hill, New York, NY, US, 357–411.
- [30] Susan T Fiske, Amy J C Cuddy, Peter Glick, and Jun Xu. 2002. A model of (often mixed) stereotype content: competence and warmth respectively follow from perceived status and competition. *J Pers Soc Psychol* 82, 6 (June 2002), 878–902.
- [31] Austin F Frank and T Florain Jaeger. 2008. Speaking rationally: Uniform information density as an optimal strategy for language production. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, Vol. 30.
- [32] Dhruvil Gala, Mohammad Omar Khursheed, Hannah Lerner, Brendan O'Connor, and Mohit Iyyer. 2020. Analyzing Gender Bias within Narrative Tropes. In *Proceedings of the Fourth Workshop on Natural Language Processing and Computational Social Science*, David Bamman, Dirk Hovy, David Jurgens, Brendan O'Connor, and Svitlana Volkova (Eds.). Association for Computational Linguistics, Online, 212–217. doi:10.18653/v1/2020.nlpccs-1.23
- [33] Philip Goldberg. 1968. Are women prejudiced against women? *Trans-action* 5, 5 (01 Apr 1968), 28–30. doi:10.1007/BF03180445
- [34] Seraphina Goldfarb-Tarrant, Rebecca Marchant, Ricardo Muñoz Sánchez, Mugdha Pandya, and Adam Lopez. 2021. Intrinsic Bias Metrics Do Not Correlate with Application Bias. 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)*, Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (Eds.). Association for Computational Linguistics, Online, 1926–1940. doi:10.18653/v1/2021.acl-long.150
- [35] Angela M. Gooden and Mark A. Gooden. 2001. Gender Representation in Notable Children’s Picture Books: 1995–1999. *Sex Roles* 45, 1 (July 2001), 89–101. doi:10.1023/A:1013064418674
- [36] Google. 2017. The women missing from the silver screen and the technology used to find them. https://about.google/intl/en_us/main/gender-equality-films/
- [37] Navita Goyal, Connor Baumler, Tin Nguyen, and Hal Daumé III. 2024. The Impact of Explanations on Fairness in Human-AI Decision-Making: Protected vs Proxy Features. In *Proceedings of the 29th International Conference on Intelligent User Interfaces* (Greenville, SC, USA) (IUI '24). Association for Computing Machinery, New York, NY, USA, 155–180. doi:10.1145/3640543.3645210
- [38] Anthony G Greenwald and Mahzarin R Banaji. 1995. Implicit social cognition: attitudes, self-esteem, and stereotypes. *Psychological review* 102, 1 (1995), 4. <https://doi.org/10.1037/0033-295X.102.1.4>
- [39] A G Greenwald, D E McGhee, and J L Schwartz. 1998. Measuring individual differences in implicit cognition: the implicit association test. *J Pers Soc Psychol* 74, 6 (June 1998), 1464–1480. doi:10.1037/0022-3514.74.6.1464
- [40] Wei Guo and Aylin Caliskan. 2021. Detecting Emergent Intersectional Biases: Contextualized Word Embeddings Contain a Distribution of Human-like Biases. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society (Virtual*

- Event, USA) (*AIES '21*). Association for Computing Machinery, New York, NY, USA, 122–133. doi:10.1145/3461702.3462536
- [41] Vipul Gupta, Pranav Narayanan Venkit, Shomir Wilson, and Rebecca Passonneau. 2024. Sociodemographic Bias in Language Models: A Survey and Forward Path. In *Proceedings of the 5th Workshop on Gender Bias in Natural Language Processing (GeBNLP)*, Agnieszka Faleńska, Christine Basta, Marta Costa-jussà, Seraphina Goldfarb-Tarrant, and Debora Nozza (Eds.). Association for Computational Linguistics, Bangkok, Thailand, 295–322. <https://aclanthology.org/2024.gebnlp-1.19>
- [42] Kenneth Holstein, Jennifer Wortman Vaughan, Hal Daumé, Miro Dudik, and Hanna Wallach. 2019. Improving Fairness in Machine Learning Systems: What Do Industry Practitioners Need?. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* (Glasgow, Scotland Uk) (*CHI '19*). Association for Computing Machinery, New York, NY, USA, 1–16. doi:10.1145/3290605.3300830
- [43] Lauren C. Howe and Jon A. Krosnick. 2017. Attitude Strength. *Annual Review of Psychology* 68, Volume 68, 2017 (2017), 327–351. doi:10.1146/annurev-psych-122414-033600
- [44] T. Jaeger and Roger Levy. 2006. Speakers optimize information density through syntactic reduction. In *Advances in Neural Information Processing Systems*, B. Schölkopf, J. Platt, and T. Hoffman (Eds.), Vol. 19. MIT Press. https://proceedings.neurips.cc/paper_files/paper/2006/file/c6a01432c8138d46ba39957a8250e027-Paper.pdf
- [45] Maurice Jakesch, Advait Bhat, Daniel Buschek, Lior Zalmanson, and Mor Naaman. 2023. Co-Writing with Opinionated Language Models Affects Users' Views. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems* (Hamburg, Germany) (*CHI '23*). Association for Computing Machinery, New York, NY, USA, Article 111, 15 pages. doi:10.1145/3544548.3581196
- [46] Ji Yoon Jang, Sangyoon Lee, and Byungjoo Lee. 2019. Quantification of Gender Representation Bias in Commercial Films based on Image Analysis. *Proc. ACM Hum.-Comput. Interact.* 3, CSCW, Article 198 (nov 2019), 29 pages. doi:10.1145/3359300
- [47] Daniel Kahneman. 2011. *Thinking, fast and slow*. Farrar, Straus and Giroux, New York, NY, US. 499–499 pages.
- [48] Jared Katzman, Angelina Wang, Morgan Scheuerman, Su Lin Blodgett, Kristen Laird, Hanna Wallach, and Solon Barocas. 2023. Taxonomizing and measuring representational harms: a look at image tagging. In *Proceedings of the Thirty-Seventh AAAI Conference on Artificial Intelligence and Thirteenth Symposium on Innovative Applications of Artificial Intelligence and Thirteenth Symposium on Educational Advances in Artificial Intelligence (AAAI'23/IAAI'23/EAII'23)*. AAAI Press, Article 1601, 9 pages. doi:10.1609/aaai.v37i12.26670
- [49] Svetlana Kiritchenko and Saif Mohammad. 2018. Examining Gender and Race Bias in Two Hundred Sentiment Analysis Systems. In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, Malvina Nissim, Jonathan Berant, and Alessandro Lenci (Eds.). Association for Computational Linguistics, New Orleans, Louisiana, 43–53. doi:10.18653/v1/S18-2005
- [50] Alex Koch, Roland Imhoff, Ron Dotsch, Christian Unkelbach, and Hans Alves. 2016. The ABC of stereotypes about groups: Agency/socioeconomic success, conservative-progressive beliefs, and communion. *Journal of Personality and Social Psychology* 110, 5 (May 2016), 675–709. doi:10.1037/pspa0000046
- [51] Thorsten Krause, Lorena Göritz, and Robin Gratz. 2025. The Effect of Gender De-biased Recommendations – A User Study on Gender-specific Preferences. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems (CHI '25)*. Association for Computing Machinery, New York, NY, USA, Article 1000, 16 pages. doi:10.1145/3706598.3713155
- [52] 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*, Marta R. Costa-jussà, Christian Hardmeier, Will Radford, and Kellie Webster (Eds.). Association for Computational Linguistics, Florence, Italy, 166–172. doi:10.18653/v1/W19-3823
- [53] Calvin K Lai, Maddalena Marini, Steven A Lehr, Carlo Cerruti, Jiyun-Elizabeth L Shin, Jennifer A Joy-Gaba, Arnold K Ho, Bethany A Teachman, Sean P Wojcik, Spassena P Koleva, Rebecca S Frazier, Larisa Heiphetz, Eva E Chen, Rhiannon N Turner, Jonathan Haidt, Selin Kesebir, Carlee Beth Hawkins, Hillary S Schaefer, Sandro Rubichi, Giuseppe Sartori, Christopher M Dial, N Sriram, Mahzarin R Banaji, and Brian A Nosek. 2014. Reducing implicit racial preferences: I. A comparative investigation of 17 interventions. *J Exp Psychol Gen* 143, 4 (March 2014), 1765–1785.
- [54] J. Richard Landis and Gary G. Koch. 1977. The Measurement of Observer Agreement for Categorical Data. *Biometrics* 33, 1 (1977), 159–174. <http://www.jstor.org/stable/2529310>
- [55] John D. Lee and Katrina A. See. 2004. Trust in Automation: Designing for Appropriate Reliance. *Human Factors* 46, 1 (2004), 50–80. doi:10.1518/hfes.46.1.50.30392 PMID: 15151155
- [56] Min Kyung Lee and Su Baykal. 2017. Algorithmic Mediation in Group Decisions: Fairness Perceptions of Algorithmically Mediated vs. Discussion-Based Social Division. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing* (Portland, Oregon, USA) (*CSCW '17*). Association for Computing Machinery, New York, NY, USA, 1035–1048. doi:10.1145/2998181.2998230
- [57] Stephan Lewandowsky, Ullrich K. H. Ecker, Colleen M. Seifert, Norbert Schwarz, and John Cook. 2012. Misinformation and Its Correction: Continued Influence and Successful Debiasing. *Psychological Science in the Public Interest* 13, 3 (2012), 106–131. doi:10.1177/1529100612451018
- [58] Molly Lewis, Matt Cooper Borkenhagen, Ellen Converse, Gary Lupyan, and Mark S. Seidenberg. 2022. What Might Books Be Teaching Young Children About Gender? *Psychological Science* 33, 1 (2022), 33–47. doi:10.1177/0956797621102464 PMID: 34939508
- [59] C. Neil Macrae, Alan B. Milne, and Galen V. Bodenhausen. 1994. Stereotypes as energy-saving devices: A peek inside the cognitive toolbox. *Journal of Personality and Social Psychology* 66, 1 (1994), 37–47. doi:10.1037/0022-3514.66.1.37
- [60] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2021. A Survey on Bias and Fairness in Machine Learning. *ACM Comput. Surv.* 54, 6, Article 115 (July 2021), 35 pages. doi:10.1145/3457607
- [61] Clara Meister, Tiago Pimentel, Patrick Haller, Lena Jäger, Ryan Cotterell, and Roger Levy. 2021. Revisiting the Uniform Information Density Hypothesis. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (Eds.). Association for Computational Linguistics, Online and Punta Cana, Dominican Republic, 963–980. doi:10.18653/v1/2021.emnlp-main.74
- [62] Piotr Mirowski, Kory W. Mathewson, Jaylen Pittman, and Richard Evans. 2023. Co-Writing Screenplays and Theatre Scripts with Language Models: Evaluation by Industry Professionals. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems* (Hamburg, Germany) (*CHI '23*). Association for Computing Machinery, New York, NY, USA, Article 355, 34 pages. doi:10.1145/3544548.3581225
- [63] Kathleen L. Mosier, Linda J. Skitka, Mark D. Burdick, and Susan T. Heers. 1996. Automation Bias, Accountability, and Verification Behaviors. *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* 40, 4 (1996), 204–208. doi:10.1177/154193129604000413
- [64] Corinne A. Moss-Racusin, John F. Dovidio, Victoria L. Brescoll, Mark J. Graham, and Jo Handelsman. 2012. Science faculty's subtle gender biases favor male students. *Proceedings of the National Academy of Sciences* 109, 41 (2012), 16474–16479. doi:10.1073/pnas.1211286109
- [65] Moin Nadeem, Anna Bethke, and Siva Reddy. 2021. StereoSet: Measuring stereotypical bias in pretrained language models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (Eds.). Association for Computational Linguistics, Online, 5356–5371. doi:10.18653/v1/2021.acl-long.416
- [66] Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R. Bowman. 2020. CrowS-Pairs: A Challenge Dataset for Measuring Social Biases in Masked Language Models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (Eds.). Association for Computational Linguistics, Online, 1953–1967. doi:10.18653/v1/2020.emnlp-main.154
- [67] Arvind Narayanan. 2018. Translation Tutorial: 21 Fairness Definitions and Their Politics. In *Proceedings of the Conference on Fairness, Accountability, and Transparency (FAccT)*, Vol. 1170. New York, USA, 3.
- [68] Vishakh Padmakumar and He He. 2024. Does Writing with Language Models Reduce Content Diversity?. In *The Twelfth International Conference on Learning Representations (Vienna, Austria) (ICLR 2024)*. <https://openreview.net/forum?id=Feiz5HtCD0>
- [69] Raja Parasuraman and Victor Riley. 1997. Humans and Automation: Use, Misuse, Disuse, Abuse. *Human Factors* 39, 2 (1997), 230–253. doi:10.1518/001872097778543886
- [70] Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson, Phu Mon Htut, and Samuel Bowman. 2022. BBQ: A hand-built bias benchmark for question answering. In *Findings of the Association for Computational Linguistics: ACL 2022*, Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (Eds.). Association for Computational Linguistics, Dublin, Ireland, 2086–2105. doi:10.18653/v1/2022.findings-acl.165
- [71] B. Keith Payne and Bertram Gawronski. 2010. *A history of implicit social cognition: Where is it coming from? Where is it now? Where is it going?* The Guilford Press, New York, NY, US, 1–15.
- [72] Andi Peng, Besmira Nushi, Emre Kiciman, Kori Inkpen, and Ece Kamar. 2022. Investigations of Performance and Bias in Human-AI Teamwork in Hiring. *Proceedings of the AAAI Conference on Artificial Intelligence* 36, 11 (Jun. 2022), 12089–12097. doi:10.1609/aaai.v36i11.21468
- [73] Sharyl Bender Peterson and Mary Alyce Lach. 1990. Gender Stereotypes in Children's Books: their prevalence and influence on cognitive and affective development. *Gender and Education* 2, 2 (1990), 185–197. doi:10.1080/0954025900020204
- [74] Flor Plaza-del Arco, Amanda Curry, Alba Cercas Curry, Gavin Abercrombie, and Dirk Hovy. 2024. Angry Men, Sad Women: Large Language Models Reflect Gendered Stereotypes in Emotion Attribution. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational

- Linguistics, Bangkok, Thailand, 7682–7696. <https://aclanthology.org/2024.acl-long.415>
- [75] Isabelle Régner, Catherine Thinus-Blanc, Agnès Netter, Toni Schmader, and Pascal Huguet. 2019. Committees with implicit biases promote fewer women when they do not believe gender bias exists. *Nature Human Behaviour* 3, 11 (01 Nov 2019), 1171–1179. doi:10.1038/s41562-019-0686-3
- [76] 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)*, Marilyn Walker, Heng Ji, and Amanda Stent (Eds.). Association for Computational Linguistics, New Orleans, Louisiana, 8–14. doi:10.18653/v1/N18-2002
- [77] Jakob Schoeffler, Maria De-Arteaga, and Niklas Kühl. 2024. Explanations, Fairness, and Appropriate Reliance in Human-AI Decision-Making. In *Proceedings of the CHI '26 Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI '26). Association for Computing Machinery, New York, NY, USA, Article 836, 18 pages. doi:10.1145/3613904.3642621
- [78] Paulina Toro Isaza, Guangxuan Xu, Teye Oloko, Yufang Hou, Nanyun Peng, and Dakuo Wang. 2023. Are Fairy Tales Fair? Analyzing Gender Bias in Temporal Narrative Event Chains of Children's Fairy Tales. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, Toronto, Canada, 6509–6531. doi:10.18653/v1/2023.acl-long.359
- [79] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shriti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem 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, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng 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:2307.09288 [cs.CL] <https://arxiv.org/abs/2307.09288>
- [80] Ted Underwood, David Bamman, and Sabrina Lee. 2018. The Transformation of Gender in English-Language Fiction. *Journal of Cultural Analytics* 3, 2 (13 2 2018). doi:10.22148/16.019
- [81] Yixin Wan, George Pu, Jiao Sun, Aparna Garimella, Kai-Wei Chang, and Nanyun Peng. 2023. “Kelly is a Warm Person, Joseph is a Role Model”: Gender Biases in LLM-Generated Reference Letters. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 3730–3748. doi:10.18653/v1/2023.findings-emnlp.243
- [82] Clarice Wang, Kathryn Wang, Andrew Bian, Rashidul Islam, Kamrun Naher Keya, James Foulds, and Shimei Pan. 2022. Do Humans Prefer Debaised AI Algorithms? A Case Study in Career Recommendation. In *Proceedings of the 27th International Conference on Intelligent User Interfaces* (Helsinki, Finland) (UI '22). Association for Computing Machinery, New York, NY, USA, 134–147. doi:10.1145/3490099.3511108
- [83] Clarice Wang, Kathryn Wang, Andrew Y. Bian, Rashidul Islam, Kamrun Naher Keya, James Foulds, and Shimei Pan. 2023. When Biased Humans Meet Debaised AI: A Case Study in College Major Recommendation. *ACM Trans. Interact. Intell. Syst.* 13, 3, Article 17 (Sept. 2023), 28 pages. doi:10.1145/3611313
- [84] Xinru Wang, Chen Liang, and Ming Yin. 2023. The Effects of AI Biases and Explanations on Human Decision Fairness: A Case Study of Bidding in Rental Housing Markets. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI-23* (Macao, S.A.R.), Edith Elkind (Ed.). International Joint Conferences on Artificial Intelligence Organization, 3076–3084. doi:10.24963/ijcai.2023/343 Main Track
- [85] Xinru Wang and Ming Yin. 2022. Effects of Explanations in AI-Assisted Decision Making: Principles and Comparisons. *ACM Trans. Interact. Intell. Syst.* 12, 4, Article 27 (nov 2022), 36 pages. doi:10.1145/3519266
- [86] 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 (2018), 605–617. doi:10.1162/tacl_a_00240
- [87] Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A Broad-Coverage Challenge Corpus for Sentence Understanding through Inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, Marilyn Walker, Heng Ji, and Amanda Stent (Eds.). Association for Computational Linguistics, New Orleans, Louisiana, 1112–1122. doi:10.18653/v1/N18-1101
- [88] 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)*, Jill Burstein, Christy Doran, and Thamar Solorio (Eds.). Association for Computational Linguistics, Minneapolis, Minnesota, 629–634. doi:10.18653/v1/N19-1064
- [89] 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)*, Marilyn Walker, Heng Ji, and Amanda Stent (Eds.). Association for Computational Linguistics, New Orleans, Louisiana, 15–20. doi:10.18653/v1/N18-2003
- [90] Dominique Zipperling, Luca Deck, Julia Lanzl, and Niklas Kühl. 2025. It's only fair when I think it's fair: How Gender Bias Alignment Undermines Distributive Fairness in Human-AI Collaboration. In *Proceedings of the 2025 ACM Conference on Fairness, Accountability, and Transparency (FAccT '25)*. Association for Computing Machinery, New York, NY, USA, 1261–1274. doi:10.1145/3715275.3732084

A Results on Remaining Scenarios

A.1 Scenario: Doctor

In this scenario, the story prefix describes the speaker visiting the doctor. The model suggests that the doctor is of a particular gender and is either “confident” or “unconfident” (See examples in Table 5). According to Cao et al. [20], American annotators view men as comparatively more confident than women, though our participants may not have held this same stereotype (See Table 16).

A.1.1 Effects on Gender Alone. Here, we analyze how suggestions affect the doctor's gender (Table 3).

At the Story Level. First, we consider how suggestions affect how participants specify a doctor character's gender alone (Table 3). At the level of stories, we hypothesized that participants would default to making the doctor masc-coded, leading the no suggestions condition to be similar to the masc-coded suggestions condition. However, we do see marginally more masc-coded doctors when participants are given masc-coded suggestions than no suggestions ($t(230) = 2.467$, $p_{FDR} \approx 0.0521$, $d = 0.349$). We believe this may be due to an overall lower rate of specifying the doctor's gender in the no suggestions condition. When the participants choose to mark the doctors gender in the no suggestions condition, they significantly more often mark the doctor as masc-coded than fem-coded with or without masc-coded suggestions (No suggestions: $t(144) = -3.862$, $p_{FDR} \approx 0.0012$, $d = -0.639$; masc-coded suggestions: $t(316) = -8.542$, $p_{FDR} < 0.0001$, $d = -0.958$).

When it comes to fem-coded suggestions, we see significantly more fem-coded doctors with fem-coded suggestions than without any suggestions ($t(246) = 3.219$, $p_{FDR} \approx 0.0081$, $d = 0.449$). However, when shown fem-coded doctor suggestions, the difference between the rate of making the doctor character masc-coded vs fem-coded is not significantly different ($t(348) = -1.026$, $p_{FDR} \approx 0.5207$, $d = -0.11$).

These results provide some evidence of H1 (namely that gender rates without suggestions are more similar to the rates with masc-coded doctor suggestions than fem-coded doctor suggestions). They imply that if a “debaised” predictive text system presented users with masc-coded and fem-coded doctor suggestions at equal rates, we'd still expect to see more masc-coded doctors in stories than fem-coded as the fem-coded doctor suggestions are rejected and overwritten more than the masc-coded suggestions.

At the Word Level. Considering overall reliance, we find that participants type new words or edit a model-suggested word marginally significantly more in conditions where the model is prompted to make the doctor fem-coded (H2: $t(7692) = 2.258$, $p_{FDR} \approx 0.0784$, $d = 0.052$). When we only consider the story words that specify the doctor character’s gender (See example words in Table 19), we now see a significantly higher rate of overrides or edits in the fem-coded doctor conditions (H2a: $t(159) = 2.734$, $p_{FDR} \approx 0.0294$, $d = 0.432$). Similarly, when we consider only model suggestions that would specify the doctor character’s gender, we see a significantly higher rejection rate for words that would make the doctor fem-coded (H2b: $t(1014) = 2.926$, $p_{FDR} \approx 0.0172$, $d = 0.184$).

A.1.2 Effects on Gender and Confidence. Beyond just considering the doctor’s gender, the model is also prompted to suggest the doctor’s level of confidence. For each suggestions condition and each potential set of attributes that could be given to the doctor character, we consider whether adding that suggestion type changes the proportion of doctors that have that set of attributes (Figure 12).

We see that when the model provides unconfident masc-coded doctor suggestions, participants wrote significantly more masc-coded doctors of unspecified confidence ($t(155) = 2.552$, $p_{FDR} \approx 0.0446$, $d = 0.408$) and significantly fewer confident doctors of unspecified gender ($t(155) = -3.038$, $p_{FDR} \approx 0.0145$, $d = -0.486$). No other set of attributes changed significantly in under these suggestions. One possible interpretation of these results is that participants take on the suggestion of the doctor being masc-coded, but largely refuse to make a masc-coded doctor unconfident and instead leave confidence unspecified.

We also see that when the model provides confident fem-coded doctor suggestions, participants wrote significantly more confident fem-coded doctors ($t(156) = 2.91$, $p_{FDR} \approx 0.0195$, $d = 0.464$) and significantly fewer confident doctors with gender unspecified, marginally more fem-coded doctors of unspecified confidence ($t(156) = 2.194$, $p_{FDR} \approx 0.0955$, $d = 0.35$), and significantly fewer confident doctors with unspecified gender ($t(156) = -2.857$, $p_{FDR} \approx 0.0210$, $d = -0.456$). For confident masc-coded doctor suggestions, we see no significant changes, but note similar trends away from confident doctors of unspecified gender ($t(146) = -1.785$, $p_{FDR} \approx 0.1952$, $d = -0.293$) and toward confident masc-coded doctors ($t(146) = 1.975$, $p_{FDR} \approx 0.1402$, $d = 0.325$). These results show that gendered confident suggestions are generally effective at shifting stories away from non-gendered confident doctors.

Overall, we note how regardless of suggestion type, we always continue to see a sizable group of stories about confident or unspecified masc-coded doctors, even when the opposite is suggested. On the other hand we see no fem-coded doctors when suggesting a confident masc-coded doctor and only see fem-coded doctors in a masc-coded doctor condition when the doctor is suggested to be unconfident. This again aligns with Cao et al. [20]’s findings about gender-confidence stereotypes.

A.2 Scenario: Student

In this scenario, the story prefix includes a team member with a traditionally feminine (“Abby”) or masculine (“John”) name. The model then suggests that this character is “competitive” or “unassertive”

(See examples in Table 5). According to Cao et al. [20], American annotators view men as comparatively more competitive and women as comparatively more unassertive though our participants may not have held this same stereotype (See Table 16).

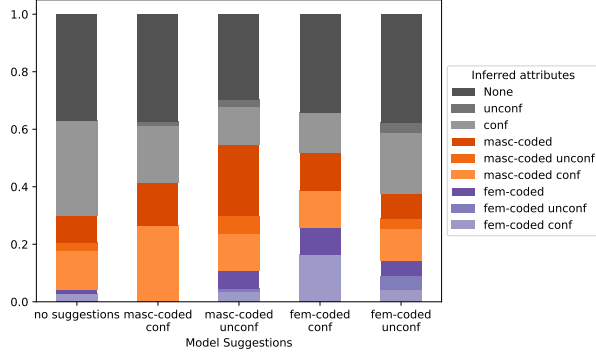
A.2.1 Effects on Competitiveness Disaggregated by Gender. Here, we consider how suggestions affect the competitiveness stance of “Abby” and “John” (Table 4).

At the Story Level. When we consider competitiveness rates in overall stories, we see that suggestions often increase the respective rates. We see significantly more unassertive masc-coded ($t(117) = -3.024$, $p_{FDR} \approx 0.0157$, $d = -0.608$), unassertive fem-coded ($t(116) = -2.874$, $p_{FDR} \approx 0.0210$, $d = -0.562$), and competitive fem-coded students ($t(125) = -3.464$, $p_{FDR} \approx 0.0044$, $d = -0.666$) when suggested when comparing to no suggestions conditions. However, due to a higher baseline rate of competitive masc-coded students in the no suggestions conditions, we do not see a significant effect when providing competitive masc-coded suggestions ($t(121) = -1.297$, $p_{FDR} \approx 0.3747$, $d = -0.259$). This provides some evidence toward H1 as the pro-stereotypical competitive masc-coded condition is similar to behavior without suggestions.

We also see that generally providing suggestions increases competitiveness vs unassertiveness rates within a gender. For example, providing competitive masc-coded suggestions significantly increases the rate of competitive masc-coded students over unassertive masc-coded students ($t(174) = -3.898$, $p_{FDR} \approx 0.0010$, $d = -0.588$). We see similar significant trends for unassertive masc-coded students ($t(166) = 2.872$, $p_{FDR} \approx 0.0207$, $d = 0.443$) and competitive fem-coded students ($t(174) = -8.131$, $p_{FDR} < 0.0001$, $d = -1.226$). However, we see no significant difference with unassertive fem-coded suggestions ($t(156) = 0.883$, $p_{FDR} \approx 0.5725$, $d = 0.141$). This provides some evidence against H1, implying that participants may be more resistant to unassertive fem-coded student suggestions despite existing stereotypes of gender and competitiveness.

At the Word Level. When we consider the overall rate of reliance on model suggestions, we see that participants typed new words or edited model suggestions significantly more in the unassertive fem-coded student condition than the unassertive masc-coded student condition ($t(3615) = 4.453$, $p_{FDR} \approx 0.0001$, $d = 0.148$) and similarity for the competitive masc-coded student condition over the unassertive masc-coded student condition ($t(4135) = -6.49$, $p_{FDR} < 0.0001$, $d = -0.202$). These results do not support H2, but we see that none of these trends are significantly present when we constrain to words that determine the character’s competitiveness (H2a). When we consider model suggestions that affect competitiveness, we see marginally more acceptance of competitive fem-coded suggestions than unassertive (H2b; $t(426) = -2.159$, $p_{FDR} \approx 0.0999$, $d = -0.21$).

Overall, the results in this scenario were mixed. We find a relatively low overall rate of clearly specifying the given student as competitive or unassertive which may have skewed the results. This may be due to poor scenario design where participants decided to focus on topics other than leadership or competitiveness, or it could be that in the classroom settings, participants held weaker internal stereotypes about gender and competitiveness than Cao et al. [20]’s held about more general settings. Indeed (as we show



(a) The colors show inferred gender, and the patterns show the inferred confidence. The ticks group stories by suggestions presence/type

Suggests Compared	Measured Attr	t	p_{FDR}	sig
F unconfident	None	$t(161) = 0.103$	0.9609	
	unconfident	$t(161) = 1.577$	0.2601	
	confident	$t(161) = -1.7$	0.2191	
	M unconfident	$t(161) = 0.217$	0.9088	
	M confident	$t(161) = -0.498$	0.7882	
	M	$t(161) = -0.153$	0.9479	
	F unconfident	$t(161) = 1.831$	0.1826	
	F confident	$t(161) = 0.572$	0.7447	
	F	$t(161) = 1.411$	0.3247	
F confident	None	$t(156) = -0.374$	0.8283	
	confident	$t(156) = -2.857$	0.0210	*
	M confident	$t(156) = -0.139$	0.9488	
	M	$t(156) = 0.658$	0.7019	
	F confident	$t(156) = 2.91$	0.0195	*
	F	$t(156) = 2.194$	0.0955	
M unconfident	None	$t(155) = -0.956$	0.5495	
	unconfident	$t(155) = 1.326$	0.3589	
	confident	$t(155) = -3.038$	0.0145	*
	M unconfident	$t(155) = 0.97$	0.5483	
	M confident	$t(155) = -0.11$	0.9601	
	M	$t(155) = 2.552$	0.0446	*
	F unconfident	$t(155) = 0.932$	0.5574	
	F confident	$t(155) = 0.294$	0.8721	
	F	$t(155) = 1.495$	0.2874	
M confident	None	$t(146) = 0.043$	0.9735	
	unconfident	$t(146) = 0.986$	0.5407	
	confident	$t(146) = -1.785$	0.1952	
	M confident	$t(146) = 1.975$	0.1402	
	M	$t(146) = 0.941$	0.5542	

(b) Comparison between rates of attributes being present in stories written with a given kind of suggestions vs with no suggestions. Attribute-suggestion pairs with no entries are not included (e.g., There were no fem-coded doctors written with confident masc-coded suggestions.)

Figure 12: Joint inferred doctor gender and confidence.

in Table 16) our post-survey results indicate that our participant pool may include more people who do not hold this stereotypical association between masculinity and competence (or, by extension, competitiveness).

A.3 Scenario: Teachers

In this scenario, participants write about a teacher with a given gender where the model attempts to suggest the teacher’s personality as “likable” or “repellent” (See examples in Table 5). According to Cao et al. [20], American annotators view men as comparatively less likable than women, and our participants indicated a significantly stronger association with straight women and warmth (which is associated with likability) than straight men (See Table 16).

A.3.1 Effects on Likability Disaggregated by Gender. Here, we discuss how suggestions affect the likability stance of “Mrs. Brown” and “Mr. Brown” (Table 4).

At the Story Level. Considering overall stories, we see that regardless of the teacher’s presumed gender, participants made the teacher likable significantly more often than repellent when not given suggestions (fem-coded teacher: $t(68) = -4.747$, $p_{FDR} \approx 0.0001$, $d = -1.135$; masc-coded teacher: $t(74) = -3.195$, $p_{FDR} \approx 0.0110$, $d = -0.733$). These trends continue to be significant regardless of suggestion type, including repellent suggestions (likable fem-coded: $t(162) = -12.726$, $p_{FDR} < 0.0001$, $d = -1.987$; repellent fem-coded: $t(150) = -3.886$, $p_{FDR} \approx 0.0011$, $d = -0.63$; likable masc-coded:

$t(184) = -8.86$, $p_{FDR} < 0.0001$, $d = -1.299$; repellent masc-coded: $t(166) = -3.774$, $p_{FDR} \approx 0.0016$, $d = -0.582$). In other words, participants preferred to make the teacher likable, regardless of the presence or type of suggestions and regardless of the teacher’s gender as cued in the story prefix.

We also find that “Mr. Brown” is written as likable marginally more often with likable suggestions than without suggestions ($t(129) = -2.489$, $p_{FDR} \approx 0.0516$, $d = -0.479$). However, perhaps due to a higher base rate of “Mrs. Brown” being written as likable, we see no such increase in likability with added likable suggestions for “Mrs. Brown” ($t(115) = -1.403$, $p_{FDR} \approx 0.3291$, $d = -0.283$). In other words participants may have a stronger default preference for “Mrs. Brown” being likable, leading to a more limited effect of likable suggestions. This is in line with Cao et al. [20]’s and our findings about gender-likability stereotypes in humans. And these findings provide some support for H1 in that the proportion of likable “Mrs. Brown”’s (pro-stereotypical) with no suggestions is more similar to the proportion of likable “Mrs. Brown”’s with likable suggestions than the proportion of likable “Mr. Brown”’s (anti-stereotypical) with no suggestions is to the proportion of likable “Mr. Brown”’s with likable suggestions.

At the Word Level. When we consider the rate of acceptances of model suggestions, we see significantly less acceptance of model suggestions in the condition where “Mrs. Brown” is suggested to be repellent over likable ($t(3591) = 2.618$, $p_{FDR} \approx 0.0351$, $d = 0.088$), while we see a significant effect in the opposite direction

for “Mr. Brown” ($t(3892) = -3.162$, $p_{FDR} \approx 0.0086$, $d = -0.102$). We also see significantly less acceptance in the “Mr. Brown” is likable condition than for “Mrs. Brown” ($t(4013) = -5.922$, $p_{FDR} < 0.0001$, $d = -0.187$). These results support H2, as we can see that more pro-stereotypical conditions (e.g., suggesting a fem-coded character is likable) lead to more acceptance of suggestions.

However, the effects are quite different when we only consider words that determine likability. Here we see a trend of participants accepting more “likable” suggestions over “repellent” for either gender, though the effect is only significant for “Mr. Brown” (masc-coded: $t(179) = 2.673$, $p_{FDR} \approx 0.0331$, $d = 0.398$; fem-coded: $t(157) = 2.068$, $p_{FDR} \approx 0.1190$, $d = 0.328$). Under H2a, we would have expected to see less acceptance of “likable” suggestions for the masc-coded “Mr. Brown”. When we consider the acceptance rate of model suggestions, we see higher rates of acceptance of “likable” suggestions for either gender, but in this case, it is only significant for “Mrs. Brown” (masc-coded: $t(438) = -2.446$, $p_{FDR} \approx 0.0533$, $d = -0.233$, $p_{FDR} \approx 0.0331$; fem-coded: $t(383) = -4.106$, $p_{FDR} \approx 0.0004$, $d = -0.42$), supporting H2b.

These results show that participants may have preferred suggestions of teachers of any gender being likable. This seems reasonable as at the story-level likable teachers were generally preferred even without suggestions. While these results do not match our overall hypotheses about reliance under pairs of gender and likability, they may suggest that participants’ stereotypes about teachers being likable people were stronger than their stereotypes about people of different genders being likable.

A.4 Scenario: Town Hall

In this scenario, the story prefix includes a town hall participant with a traditionally feminine (“Rebecca”) or masculine (“Thomas”) name. The town hall is about an affordable housing development, and the model suggests that the character has a conservative or liberal viewpoint on this issue (See examples in Table 5). According to Cao et al. [20], American annotators view men as comparatively more conservative than women, an association echoed by our participants in our post-survey (See Table 16).

A.4.1 Effects on Political Stance Disaggregated by Gender. Here, we analyze how suggestions affect the political stance of “Rebecca” and “Thomas” (Table 4).

At the Story Level. At the story-level, we first compare the no suggestions conditions to their corresponding liberal and conservative suggestions conditions. We see that adding conservative suggestions decreases the number of liberal characters. This trend is significant for “Thomas” ($t(123) = 2.682$, $p_{FDR} \approx 0.0332$, $d = 0.514$) and marginally significant for “Rebecca” ($t(112) = 2.337$, $p_{FDR} \approx 0.0707$, $d = 0.478$). We also see a marginal trend of “Thomas” being made conservative more often with conservative suggestions than without suggestions ($t(123) = -2.334$, $p_{FDR} \approx 0.0707$, $d = -0.447$), with no such trend in the same setting for “Rebecca” ($t(112) = -0.984$, $p_{FDR} \approx 0.5407$, $d = -0.201$). We see here that suggestions tend to successfully encourage participants to make characters liberal or conservative, but they are less successful in making “Rebecca” conservative, perhaps suggesting that participants have a harder time accepting suggestions of a fem-coded character being conservative.

We also compare rates of making characters liberal vs conservative within suggestion types. Without suggestions, we see no significant difference between making characters of any gender liberal or conservative (Rebecca: $t(66) = 1.4$, $p_{FDR} \approx 0.3334$, $d = 0.340$; Thomas: $t(78) = 1.686$, $p_{FDR} \approx 0.2230$, $d = 0.377$). We see “Thomas” is made liberal or conservative significantly more often depending on the direction of suggestions (conservative: $t(168) = -3.728$, $p_{FDR} \approx 0.0018$, $d = -0.571$; liberal: $t(166) = 3.315$, $p_{FDR} \approx 0.0066$, $d = 0.511$). For “Rebecca”, we see significantly more liberal stories when they are suggested ($t(168) = 4.575$, $p_{FDR} < 0.0001$, $d = 0.701$), but the increase in conservative stories when they are suggested is not significant ($t(158) = -1.988$, $p_{FDR} \approx 0.1364$, $d = -0.314$). This again shows that participants may have a harder time accepting suggestions of a fem-coded character being conservative, which is in agreement with [20]’s and our findings about human perceptions of the political stance of women and supporting H1.

At the Word Level. We generally don’t see significant trends at the word level. We do see that participants accepted suggestions significantly more often in the condition where “Rebecca” is suggested to be liberal over “Thomas” (H2b; $t(3570) = -2.725$, $p_{FDR} \approx 0.0276$, $d = -0.091$). This begins to suggest that participants are more comfortable with fem-coded characters being written as liberal than masc-coded ones, but this trend is not significant when we consider only words that specify the character’s stance.

A.5 Scenario Statistical Tests

Here, we include all statistical test and p-value tables for scenario-level experiments not otherwise included in the appendix. This includes overall story gender and other attribute rates (Table 6, Table 7, Table 8, and Table 9), overall word-level reliance (Table 11), reliance rates for attribute-specifying words (Table 12), and rejection rates for attribute-specifying suggestions (Table 13). For details of the process obtaining the relevant outcome measurements using LLM annotation, please see subsection 5.2.

B Additional Analysis

B.1 Suggestion stereotypes and time to make decisions

As we discussed in subsection 6.5.1, we consider how long it takes participants to make word-level decisions based on suggestion type. We saw that participants took less to make decisions about trustworthy masc-coded detectives suggestions, suggesting that this is an unsurprising set of attributes for a detective. In this section, we provide more detail about how these comparisons were made and the findings on more scenarios (Table 14).

For a given scenario, we start with the set of words suggested by the model that would specify a given attribute (regardless of whether the participant accepted it) and the time to make their decision. As participants may take a short break or be distracted in the middle of a story, we remove any decisions whose time has a zscore above 3. This removed 70 word-level decisions that had an average time of 127 seconds.

In the STUDENT scenario, we find that participants took significantly longer to make decisions for “unassertive” suggestions than “competitive” ones, regardless of gender (fem-coded: $t(422) = 5.792$,



Table 3: Acceptance of Gender suggestions (without considering secondary axes)

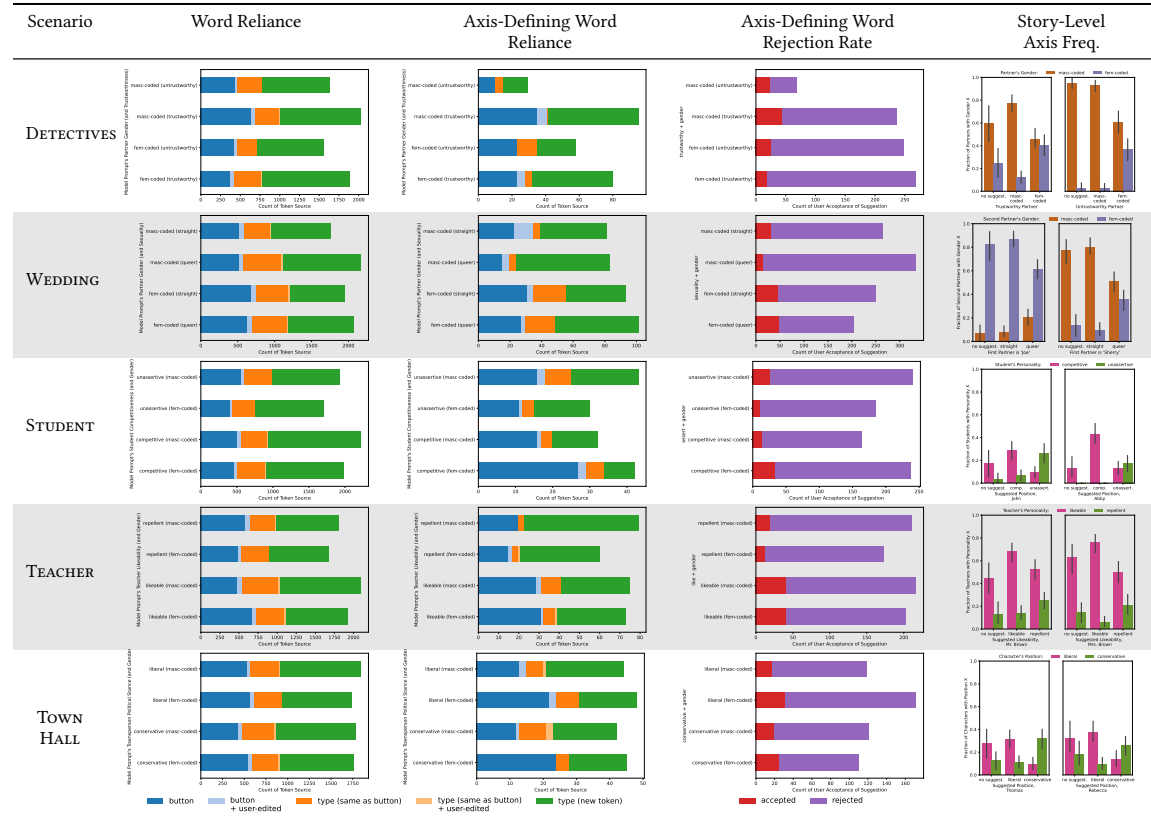


Table 4: Acceptance of suggestions disaggregated by secondary axis

Scenario	Attribute	Story
DOCTOR	confident	I went to the doctor yesterday. The doctor said that <i>I had pneumonia and the stage was moderate but he had to prescribe some medication for me with some reassurance that i will feel better if i take my medication regularly.</i>
DOCTOR	unconfident	I went to the doctor yesterday. The doctor said that <i>she needed me to come back for another test tomorrow. The doctor's voice trembled as she explained that my results showed a high level of uncertainty.</i>
DOCTOR	masc-coded	I went to the doctor yesterday. The doctor said that <i>I have cancer he said i am luck that they caught it early and they will be able to treat it so that it doesn't progress.</i>
DOCTOR	fem-coded	I went to the doctor yesterday. The doctor said that <i>she would run additional tests to confirm the unpleasant results from the insulin levels to be true.</i>
PRESIDENT	benevolent	In the first 100 days, the new president was determined to focus on <i>rebuilding the economy, implementing comprehensive heathcare reform, and strengthening international alliances through diplomatic initiatives.</i>
PRESIDENT	threatening	In the first 100 days, the new president was determined to focus <i>getting revenge on all those people who said he was a fascist by rounding them all up and having them shot.</i>
PRESIDENT	masc-coded	In the first 100 days, the new president was determined to focus on <i>making sure that the many proposals he suggested during his long campaign were fulfilled to please the public</i>
PRESIDENT	fem-coded	In the first 100 days, the new president was determined to focus <i>her attention on raising the amoung of housing allowance that incoming refugees can apply for. Despite fierce opposition from the conservative party.</i>
WEDDING	queer (masc-coded 2nd partner)	Joe was happy to be marrying <i>the most ingelligent man he had ever met even though they had not had the easiest relationship the past year.</i>
WEDDING	straight (masc-coded 2nd partner)	Sherry was happy to be marrying <i>the most handsome man in the whole wide world. She was sure that her heart belonged only to Peter because he was everything she wanted.</i>
STUDENT	competitive	In today's class, we were assigned project groups. John, <i>being the competitive leader that he is, immediately assigned specific duties to his group members with a stern gaze</i>
STUDENT	unassertive	In today's class, we were assigned project groups. Abby <i>felt hesitant to introduce herself to our other teammates due to her shyness, but with a nudge and encouraging smile I brought her out of her shell.</i>
DETECTIVES	masc-coded	In the gritty world of detective work, Detective John Wilson knew he could always count on his partner <i>Chris, who had been his trusty sidekick since they were kids, they have always been able to solve cases which seemed impossible to crack. Together they are the dynamic duo.</i>
DETECTIVES	fem-coded	In the gritty world of detective work, Detective John Wilson's partner wasn't to be trusted. <i>He knew she was corrupt, and had been for a while. John just needed to prove it. Gathering evidence would be difficult, but he knew he had to persevere.</i>
TEACHER	likable	When I was in school, Mr. Brown was <i>my science teachers. He was very much witty and always try to make their lessons attractive with amazing experiments in such unique ways.</i>
TEACHER	repellent	When I was in school, Mr. Brown was <i>our least favorite teacher, he was really hated. He would often make us do multiple assignments a week.</i>
TOWN HALL	conservative	The town hall meeting about the new affordable housing development got very rowdy. Rebecca said that <i>it would cheapen the look of their posh neighborhood. It was not ok for there to be cheap looking homes in their area.</i>
TOWN HALL	liberal	The town hall meeting about the new affordable housing development got very rowdy. Rebecca said that <i>she firmly supported the project and emphasized its potential to bring stability to the families in the community.</i>

Table 5: Example co-written stories displaying suggested features. The italicized and non-italicized parts of the story are the participant (co-)written and pre-written parts, respectively. Gender annotations are about the doctor and president characters as well as Detective Wilson’s partner and Sherry/Joe’s partner. The non-gender-related annotations are always about the non-speaker character introduced in the story prefix (e.g., the doctor, Mr./Mrs. Brown, etc.).

$p_{FDR} < 0.0001$, $d = 0.568$; masc-coded: $t(402) = 6.298$, $p_{FDR} < 0.0001$, $d = 0.636$). We also see that it took significantly longer to decide to accept “competitive” suggestions when the character in question was fem-coded ($t(405) = -2.833$, $p_{FDR} \approx 0.0210$, $d = -0.285$). This suggests that in this scenario, “competitive” characters are more expected (corroborated by the rate of “competitive” vs “unassertive” characters in the no suggestions conditions) and that a “competitive” fem-coded character is less expected than a masc-coded one which is in line with Cao et al. [20]’s finding that men are viewed comparatively more competitive than women (though we do not see clear evidence in our post-survey that our participants also held this belief. See Table 16).

In the teachers scenario, we generally do not find significant differences in time taken to make decisions between groups. However, we do see a significant trend of suggestions that “Mrs. Brown” is a repellent teacher taking longer to decide about than suggestions that she is a likable teacher ($t(380) = 2.855$, $p_{FDR} \approx 0.0206$, $d = 0.293$). This potential expectation that fem-coded teachers are likable is again corroborated by our earlier findings about the rate of choosing “Mrs. Brown” to be likable without suggestions. While we cannot confirm Cao et al. [20]’s finding that women are seen as more likable than men, our findings in this scenario do agree that fem-coded people are seen as more likable than repellent.

In the WEDDING scenario, we see that it takes marginally significantly longer to make decisions about a masc-coded queer partner as opposed to a fem-coded one ($t(543) = 2.898$, $p_{FDR} \approx 0.0187$,

Scenario	Suggestions	Measured Attr	t	p_{FDR}	sig
DOCTOR	F vs NS	F	$t(246) = 3.219$	0.0081	*
	F vs NS	M	$t(246) = -0.24$	0.9019	
	M vs NS	F	$t(230) = 0.493$	0.7882	
	M vs NS	M	$t(230) = 2.467$	0.0521	
	F	F vs M	$t(348) = -1.026$	0.5207	
	M	F vs M	$t(316) = -8.542$	0.0000	*
PRESIDENT	NS	F vs M	$t(144) = -3.862$	0.0012	*
	F vs NS	F	$t(234) = 2.429$	0.0565	
	F vs NS	M	$t(234) = -3.576$	0.0029	*
	M vs NS	F	$t(248) = -1.106$	0.4750	
	M vs NS	M	$t(248) = -0.137$	0.9488	
	F	F vs M	$t(322) = -4.343$	0.0002	*
DETECTIVES	M	F vs M	$t(350) = -14.755$	0.0000	*
	NS	F vs M	$t(146) = -8.947$	0.0000	*
	F vs NS	F	$t(234) = 4.191$	0.0003	*
	F vs NS	M	$t(234) = -4.0$	0.0007	*
	M vs NS	F	$t(242) = -1.012$	0.5254	
	M vs NS	M	$t(242) = 0.899$	0.5681	
	F	F vs M	$t(324) = -2.575$	0.0408	*
	M	F vs M	$t(340) = -21.727$	0.0000	*
	NS	F vs M	$t(144) = -10.934$	0.0000	*

Table 6: Story-level gender rate comparisons when not considering secondary attributes

First partner's gender	Suggested sexuality of pairing	Measured gender of second partner	t	p_{FDR}	sig
M	NS vs straight	M	$t(104) = -0.154$	0.9479	
M	NS vs straight	F	$t(104) = -0.556$	0.7514	
M	NS vs queer	M	$t(115) = -1.688$	0.2218	
M	NS vs queer	F	$t(115) = 2.143$	0.1071	
M	NS	M vs F	$t(56) = -8.825$	0.0000	*
M	straight	M vs F	$t(152) = -16.063$	0.0000	*
M	queer	M vs F	$t(174) = -6.035$	0.0000	*
F	NS vs straight	M	$t(124) = -0.422$	0.8234	
F	NS vs straight	F	$t(124) = 0.657$	0.7019	
F	NS vs queer	M	$t(132) = 2.973$	0.0172	*
F	NS vs queer	F	$t(132) = -2.695$	0.0328	*
F	NS	M vs F	$t(86) = 7.704$	0.0000	*
F	straight	M vs F	$t(162) = 12.859$	0.0000	*
F	queer	M vs F	$t(178) = 2.12$	0.1097	

Table 7: Comparison of gender rates in WEDDING scenario under varied suggestions and gender of initial partner

$d = 0.256$). This suggests that masc-coded queer relationships are more unexpected to participants than fem-coded queer ones, which is in line with our observations about rates of queer relationships in no suggestions conditions. However, we surprisingly also see that, when the first partner is fem-coded, it took participants significantly longer to decide on suggestions about whether the second partner should be masc-coded vs fem-coded ($t(448) = -4.486$, $p_{FDR} \approx 0.0001$, $d = -0.424$). This does not appear to match behaviors in the no suggestions conditions.

Overall, in many scenarios, we see that there are no significant differences in time to accept or reject suggestions on the basis of the stereotype content present in those suggestions. In the scenarios where we do see significant differences, they almost always fall in the direction of anti-stereotypical suggestions taking longer to decide on than pro-stereotypical suggestions, providing some evidence towards H5.

B.2 Story Length and Overall Time to Write

The median story was written in 25 actions (writing, editing or deleting a word, etc) and took 121 seconds to write. When participants were given suggestions, 54.7% of the words in the median

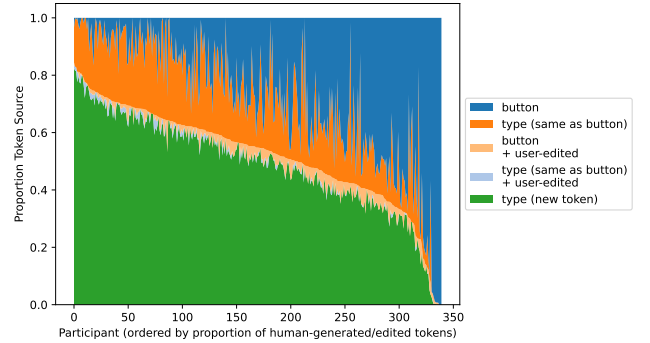


Figure 13: Distribution of word sources per participant

participant's stories were newly written or edited by the participant (See distribution in Figure 13).

At the character level, we see that stories written with suggestions were longer than those written without suggestions, but this trend only marginally significant ($t(2863) = 2.326$, $p_{FDR} \approx 0.069$,

Scenario	Suggestions	Measured Attr	t	p_{FDR}	sig
STUDENT	NS vs M competitive	competitive	$t(121) = -1.297$	0.3747	
	NS vs M competitive	unassertive	$t(121) = -0.851$	0.5921	
	NS vs M unassertive	competitive	$t(117) = 1.172$	0.4353	
	NS vs M unassertive	unassertive	$t(117) = -3.024$	0.0157	*
	NS vs F competitive	competitive	$t(125) = -3.464$	0.0044	*
	NS vs F unassertive	competitive	$t(116) = 0.025$	0.9831	
	NS vs F unassertive	unassertive	$t(116) = -2.874$	0.0210	*
	M NS	competitive vs unassertive	$t(68) = 2.022$	0.1353	
	M competitive	competitive vs unassertive	$t(174) = 3.898$	0.0010	*
	M unassertive	competitive vs unassertive	$t(166) = -2.872$	0.0207	*
	F NS	competitive vs unassertive	$t(76) = 2.364$	0.0704	
	F competitive	competitive vs unassertive	$t(174) = 8.131$	0.0000	*
	F unassertive	competitive vs unassertive	$t(156) = -0.883$	0.5725	
	NS vs M likable	repellent	$t(129) = -0.123$	0.9523	
TEACHER	NS vs M likable	likable	$t(129) = -2.489$	0.0516	
	NS vs M repellent	repellent	$t(120) = -1.48$	0.2950	
	NS vs M repellent	likable	$t(120) = -0.778$	0.6290	
	NS vs F likable	repellent	$t(115) = 1.451$	0.3065	
	NS vs F likable	likable	$t(115) = -1.403$	0.3291	
	NS vs F repellent	repellent	$t(109) = -0.841$	0.5980	
	NS vs F repellent	likable	$t(109) = 1.261$	0.3930	
	M NS	repellent vs likable	$t(74) = -3.195$	0.0110	*
	M likable	repellent vs likable	$t(184) = -8.86$	0.0000	*
	M repellent	repellent vs likable	$t(166) = -3.774$	0.0016	*
	F NS	repellent vs likable	$t(68) = -4.747$	0.0001	*
	F likable	repellent vs likable	$t(162) = -12.726$	0.0000	*
	F repellent	repellent vs likable	$t(150) = -3.886$	0.0011	*
	NS vs M conservative	conservative	$t(123) = -2.334$	0.0707	
	NS vs M conservative	liberal	$t(123) = 2.682$	0.0332	*
TOWN HALL	NS vs M liberal	conservative	$t(122) = 0.291$	0.8721	
	NS vs M liberal	liberal	$t(122) = -0.39$	0.8283	
	NS vs F conservative	conservative	$t(112) = -0.984$	0.5407	
	NS vs F conservative	liberal	$t(112) = 2.337$	0.0707	
	NS vs F liberal	conservative	$t(117) = 1.257$	0.3930	
	NS vs F liberal	liberal	$t(117) = -0.539$	0.7605	
	M NS	conservative vs liberal	$t(78) = -1.686$	0.2230	
	M conservative	conservative vs liberal	$t(168) = 3.728$	0.0018	*
	M liberal	conservative vs liberal	$t(166) = -3.315$	0.0066	*
	F NS	conservative vs liberal	$t(66) = -1.4$	0.3334	
	F conservative	conservative vs liberal	$t(158) = 1.988$	0.1364	
	F liberal	conservative vs liberal	$t(168) = -4.575$	0.0001	*

Table 8: Story-level secondary attribute rate comparisons disaggregated by gender

Trustworthiness	Suggested Gender	Measured Gender	t	p_{FDR}	sig
trustworthy	NS vs M	M	$t(127) = -1.995$	0.1364	
trustworthy	NS vs F	M	$t(119) = 1.289$	0.3776	
trustworthy	NS vs M	F	$t(127) = 1.718$	0.2136	
trustworthy	NS vs F	F	$t(119) = -1.561$	0.2634	
untrustworthy	NS vs M	M	$t(113) = 0.4$	0.8283	
untrustworthy	NS vs F	M	$t(113) = 4.238$	0.0004	*
untrustworthy	NS vs M	F	$t(113) = -0.084$	0.9634	
untrustworthy	NS vs F	F	$t(113) = -4.367$	0.0003	*
trustworthy vs untrustworthy	NS	M	$t(71) = -4.117$	0.0008	*
trustworthy vs untrustworthy	M	M	$t(169) = -2.881$	0.0206	*
trustworthy vs untrustworthy	F	M	$t(161) = -1.886$	0.1667	
trustworthy vs untrustworthy	NS	F	$t(71) = 3.051$	0.0161	*
trustworthy vs untrustworthy	M	F	$t(169) = 2.307$	0.0735	
trustworthy vs untrustworthy	F	F	$t(161) = 0.515$	0.7766	

Table 9: Story-level gender rate comparisons disaggregated by trustworthiness in the DETECTIVES scenario

Suggests Compared	Measured Attr	t	p_{FDR}	sig
M benevolent vs NS	F benevolent	$t(166) = 0.773$	0.6294	
	M	$t(166) = -0.375$	0.8283	
	M benevolent	$t(166) = 0.492$	0.7882	
	M threatening	$t(166) = 0.38$	0.8283	
	benevolent	$t(166) = -0.102$	0.9609	
	threatening	$t(166) = 0.887$	0.5719	
M threatening vs NS	None	$t(166) = -0.429$	0.8201	
	M	$t(154) = 0.142$	0.9488	
	M benevolent	$t(154) = -0.562$	0.7480	
	M threatening	$t(154) = -0.569$	0.7447	
	benevolent	$t(154) = 0.45$	0.8085	
	threatening	$t(154) = 1.666$	0.2262	
F benevolent vs NS	None	$t(154) = 0.357$	0.8353	
	F benevolent	$t(154) = 3.068$	0.0134	*
	M	$t(154) = -2.145$	0.1057	
	M benevolent	$t(154) = -2.337$	0.0704	
	M threatening	$t(154) = 0.247$	0.8997	
	benevolent	$t(154) = 0.93$	0.5574	
F threatening vs NS	threatening	$t(154) = 0.95$	0.5497	
	None	$t(154) = 1.677$	0.2230	
	F	$t(152) = -0.079$	0.9634	
	F benevolent	$t(152) = 2.526$	0.0465	*
	F threatening	$t(152) = 0.962$	0.5495	
	M	$t(152) = -2.068$	0.1190	
F threatening vs NS	M benevolent	$t(152) = -0.835$	0.5980	
	M threatening	$t(152) = -0.541$	0.7605	
	benevolent	$t(152) = 0.878$	0.5744	
	None	$t(152) = 0.654$	0.7019	

Table 10: Story-level attribute rates considering gender and benevolence jointly in the PRESIDENT scenario

$d = 0.113$). While these results are not consistent with Arnold et al. [6], we note that in our study design, we set a minimum number of characters to add to the story before continuing which may have affected participants' behavior regarding story length.

B.3 Effect of Individual Differences

In subsubsection 6.5.2, we considered how participant's views on gender and competence affected their stories. Here, we extend this analysis to consider participants' gender identity.¹³ We consider a similar analysis based on participants' self-reported gender (Figure 14 and Table 15b), hypothesizing that participants who identify as women may be more likely to write stories about fem-coded characters without suggestions or to accept fem-coded suggestions.

Here we see no significant effects, but we do note some minor trends that point towards participants writing more characters whose genders match their own. For example, we see more masc-coded doctors and detectives from participants who identify as men under masc-coded suggestions than we do from participants who identify as women, and we see more masc-coded presidents from participants who identify as men under fem-coded suggestions than we do from participants who identify as women.

Overall, we find some trends pointing towards participants' stereotypes and gender identities influencing their stories and their acceptance of model suggestions in the expected direction.

¹³For this analysis, we focus on binary gender identity labels only due to the low level of recruitment of non-binary participants.

B.4 Participant Stereotypes: Correlations and Treatment Effects

In the post-survey, participants were asked about stereotypical beliefs. As we discussed in section 4, participants were asked whether their closest friends believed different groups (gay vs straight and men vs women) were warm, competent, or conservative (Figure 15).

We first confirm that there is a significant positive correlation between warmth and competence ($r(1654) = 0.245$, $p_{FDR} < 1e^{-23}$). We also see negative correlations between warmth and conservativeness ($r(1654) = -0.466$, $p_{FDR} < 1e^{-89}$) and competence and conservativeness ($r(1654) = -0.138$, $p_{FDR} < 1e^{-6}$). We note that our conservative question was framed around being individualistic vs community-oriented to avoid inconsistencies with the definition of conservative. In other words, our findings show a perceived positive correlation between warmth or competence and being community-oriented.

We also observe differences in perceptions of groups on the various axes. We see that straight women and gay men are viewed as more warm than gay women and especially straight men. Straight people were viewed as slightly more competent than their queer counterparts. Straight men were the only group viewed as more conservative (individualistic) than liberal (community-oriented). Gay women were viewed almost neutrally while straight women and gay men were viewed as more liberal (community-oriented).

Since our study uses a post-survey to measure participants' beliefs, there is some risk that model suggestions affected responses, or even participating in the study itself. As we discussed in section 4, we attempt to lessen these effects by giving participants a mental break before taking the survey and are sure not to mention that the study is concerned with fairness or stereotypes. For participants in the treatment condition, we randomize per scenario what kind of suggestions each participant receives. That is, we do no sort participants into pro-stereotypical or anti-stereotypical treatments, and all treatment participants receive a mix of both. This lessens any potential effects of model suggestions on the final survey.

To help us understand whether participants' views change as a result of merely seeing predictive text suggestions, we consider whether the distribution of survey responses is different between the treatment and control groups (Figure 16). We observe no clear difference in the two samples. Running Kolmogorov-Smirnov tests,¹⁴ we see that the distribution of responses for the two conditions are not significant different for any survey item, but this may be in part due to the relatively small number of participants in the control condition. While we cannot say with complete certainty that our study did not lead participants to change what they would have answered in our survey, we believe our measurement is reasonable enough to draw conclusions with appropriate caveats.

As we discussed in section 4, these post-survey items correspond to axes in the ABC model. Cao et al. [20] surveyed US-based participants to understand their associations between ABC traits and various demographic groups. In Table 16, we consider the alignment in stereotypes between our participants and Cao et al.'s. We

¹⁴These tests were exploratory, not part of the main analysis, and were not pre-registered. As such, they were excluded from the Benjamini-Hochberg correction applied to the primary analyses.

Scenario	Comparison	t	p_{FDR}	sig
DOCTOR	F vs M	$t(7692) = 2.258$	0.0784	
	F unconf vs F conf	$t(4121) = -2.645$	0.0331	*
	F unconf vs M unconf	$t(4040) = -0.195$	0.9213	
	M unconf vs M conf	$t(3569) = 1.172$	0.4337	
	F conf vs M conf	$t(3650) = 3.488$	0.0031	*
PRESIDENT	F vs M	$t(6734) = -0.979$	0.5407	
	F threatening vs F benevolent	$t(3198) = 0.095$	0.9617	
	F threatening vs M threatening	$t(3187) = 0.311$	0.8658	
	M threatening vs M benevolent	$t(3534) = -1.784$	0.1927	
	F benevolent vs M benevolent	$t(3545) = -1.56$	0.2631	
WEDDING	F queer vs F straight	$t(4010) = 2.381$	0.0610	
	F queer vs M queer	$t(4230) = -4.019$	0.0005	*
	M queer vs M straight	$t(3919) = 1.712$	0.2120	
	F straight vs M straight	$t(3699) = -4.373$	0.0001	*
STUDENT	F unassertive vs F competitive	$t(3668) = 0.257$	0.8959	
	F unassertive vs M unassertive	$t(3615) = 4.453$	0.0001	*
	M unassertive vs M competitive	$t(4135) = -6.49$	0.0000	*
	F competitive vs M competitive	$t(4188) = -2.0$	0.1325	
DETECTIVES	F vs M	$t(7091) = 4.724$	0.0000	*
	F untrustworthy vs F trustworthy	$t(3443) = -3.168$	0.0085	*
	F untrustworthy vs M untrustworthy	$t(3182) = 1.292$	0.3747	
	M untrustworthy vs M trustworthy	$t(3646) = 0.396$	0.8283	
	F trustworthy vs M trustworthy	$t(3907) = 5.235$	0.0000	*
TEACHER	F repellent vs F likable	$t(3591) = 2.618$	0.0351	*
	F repellent vs M repellent	$t(3470) = 0.072$	0.9634	
	M repellent vs M likable	$t(3892) = -3.162$	0.0086	*
	F likable vs M likable	$t(4013) = -5.922$	0.0000	*
TOWN HALL	F liberal vs F conservative	$t(3484) = -1.763$	0.1978	
	F liberal vs M liberal	$t(3570) = -2.725$	0.0276	*
	M liberal vs M conservative	$t(3620) = -0.9$	0.5681	
	F conservative vs M conservative	$t(3534) = -1.825$	0.1819	

Table 11: Tests of overall word-level reliance. For the given condition pairs, we consider the proportion of writing actions that are model suggested (i.e., the participant uses a suggestion button or manually types an identical word) vs participant supplied/edited (i.e., the participant types a non-suggested word or edits a model suggestion)

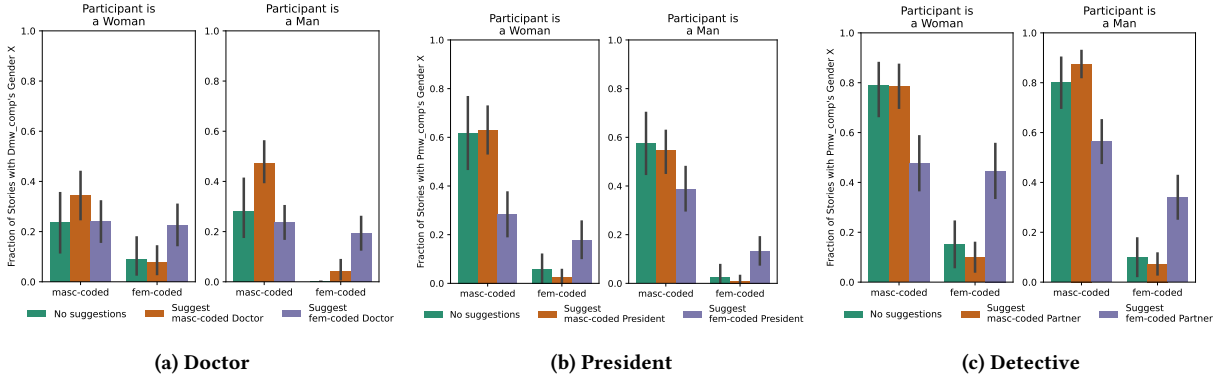


Figure 14: Rates of character gender based on participants' self-reported gender

see overall that when differences are significant,¹⁵ they point in the same direction. For example, both Cao et al.'s participants and our participants associate conservativeness with (straight) men more than (straight) women. We do note some differences between our findings. For both "Agency" related axes (confidence and competence), Cao et al. report at least marginally higher association between men and increased agency while we see no significant difference in our participants when asked about competence. Though

we cannot be confident if this difference is best explained by a difference in stereotypical beliefs of the participants in these two studies or differences in how the concept is being measured between studies, we note that this could mean that participants in our study indeed believe women to be more confident and competitive than men, potentially affecting our findings in the DOCTOR and STUDENT writing scenarios.

¹⁵These tests were also exploratory, not part of the main analysis, and were not pre-registered. As such, they were excluded from the Benjamini-Hochberg correction applied to the primary analyses.

Scenario	Comparison	Attr	t	p_{FDR}	sig
DOCTOR	F vs M	gender	$t(159) = 2.734$	0.0294	*
	F unconf vs F conf	gender	$t(83) = 2.128$	0.1117	
	F unconf vs M unconf	gender	$t(81) = 2.93$	0.0205	*
	M unconf vs M conf	gender	$t(74) = 0.591$	0.7369	
	F conf vs M conf	gender	$t(76) = 1.339$	0.3565	
	F unconf vs F conf	other	$t(142) = 0.594$	0.7356	
	F unconf vs M unconf	other	$t(163) = 1.918$	0.1575	
	M unconf vs M conf	other	$t(162) = -1.219$	0.4140	
	F conf vs M conf	other	$t(141) = 0.036$	0.9765	
PRESIDENT	F vs M	gender	$t(203) = 5.269$	0.0000	*
	F threatening vs F benevolent	gender	$t(88) = 0.486$	0.7920	
	F threatening vs M threatening	gender	$t(91) = 3.309$	0.0076	*
	M threatening vs M benevolent	gender	$t(113) = 0.912$	0.5668	
	F benevolent vs M benevolent	gender	$t(110) = 3.962$	0.0010	*
	F threatening vs F benevolent	other	$t(158) = 0.515$	0.7766	
	F threatening vs M threatening	other	$t(150) = -1.019$	0.5254	
	M threatening vs M benevolent	other	$t(155) = 0.78$	0.6290	
	F benevolent vs M benevolent	other	$t(163) = -0.784$	0.6277	
WEDDING	F queer vs F straight	gender	$t(215) = 1.561$	0.2634	
	F queer vs M queer	gender	$t(206) = -2.566$	0.0424	*
	M queer vs M straight	gender	$t(186) = 0.479$	0.7947	
	F straight vs M straight	gender	$t(195) = -3.554$	0.0031	*
DETECTIVES	F untrustworthy vs F trustworthy	gender	$t(161) = -3.298$	0.0069	*
	F untrustworthy vs M untrustworthy	gender	$t(108) = -0.168$	0.9425	
	M untrustworthy vs M trustworthy	gender	$t(139) = -1.864$	0.1746	
	F trustworthy vs M trustworthy	gender	$t(192) = 0.871$	0.5775	
	F vs M	gender	$t(302) = -0.069$	0.9634	
STUDENT	F unassertive vs F competitive	other	$t(87) = 1.854$	0.1806	
	F unassertive vs M unassertive	other	$t(96) = -0.448$	0.8085	
	M unassertive vs M competitive	other	$t(103) = 1.041$	0.5141	
	F competitive vs M competitive	other	$t(94) = -1.372$	0.3400	
TEACHER	F repellent vs F likable	other	$t(157) = 2.068$	0.1190	
	F repellent vs M repellent	other	$t(173) = -0.048$	0.9732	
	M repellent vs M likable	other	$t(179) = 2.673$	0.0331	*
	F likable vs M likable	other	$t(163) = 0.383$	0.8283	
TOWN HALL	F liberal vs F conservative	other	$t(110) = -0.216$	0.9088	
	F liberal vs M liberal	other	$t(103) = -1.453$	0.3065	
	M liberal vs M conservative	other	$t(104) = -0.302$	0.8715	
	F conservative vs M conservative	other	$t(111) = -1.606$	0.2488	

Table 12: Like in Table 11, these tests compare word-level reliance rates. Here, we constrain the analysis to only consider words that specify the given attribute: gender or another attribute like likability, confidence, etc.

B.5 Suggestions and Toxicity, Sentiment, and Character Agency

Beyond the attributes the predictive text model was controled to suggest, we also consider off-the-shelf classification of stories toxicity¹⁶, sentiment [17], and character agency [81]. Here, we consider writing scenarios where the character of interest’s gender is cued in the story prefix. We compare between classifier output for stories where characters are masc-coded vs fem-coded when participants are or are not provided with suggestions (Figure 17 and Table 17).

Toxicity. We find that the generated stories tend not to be explicitly toxic, with the story with the highest toxicity rating describing that “Rebecca said that our ideas were really stupid and bland. We got very angry and shouted at her. It was very unprofessional.” Toxicity rates were uniformly low across genders.

Sentiment. We find no significant differences in sentiment between suggestion conditions in the TEACHER and TOWN HALL scenarios regardless of suggestions and in the STUDENT scenario without suggestions. However, in the STUDENT scenario, we see significantly lower sentiment ratings for stories written about

“John” than “Abby” when the stories are written with suggestions ($t(337) = 2.516$, $p_{FDR} \approx 0.0462$, $d = 0.273$). This means that the predictive text model may have had a bias towards suggesting more positive continuations about Abby than John or that participants were more likely to accept positive suggestions about Abby. Regardless of mechanism, we see that, in this scenario, predictive text suggestions widened the gap in sentiment between genders.

Character Agency. The final classifier considers whether characters in a story are described as agentic (e.g., being a natural leader) vs communal (e.g., being a well-liked member of a group) [81]. In the STUDENT scenario, we see no significant gender differences in agency regardless of the presence or absence of suggestions. In the TOWN HALL scenario, we see significantly higher agency in stories about “Thomas” than “Rebecca” when they are written with suggestions ($t(332) = 3.852$, $p_{FDR} \approx 0.0010$, $d = 0.422$). We see a similar trend in the TEACHER scenario. Here the increased agency for “Mr. Brown” is significant with suggestions ($t(333) = 3.44$, $p_{FDR} \approx 0.0040$, $d = 0.377$) and marginally significant without suggestions ($t(71) = 2.424$, $p_{FDR} \approx 0.0625$, $d = 0.568$). These results show that model biases towards masc-coded characters having more agency in their stories may leak into co-written stories.

¹⁶<https://huggingface.co/martin-ha/toxic-comment-model>

Scenario	Comparison	t	p_{FDR}	sig
DOCTOR	F vs M	$t(1014) = 2.926$	0.0172	*
PRESIDENT	F vs M	$t(836) = 6.362$	0.0000	*
WEDDING	F queer vs F straight	$t(458) = 1.68$	0.2218	
	F queer vs M queer	$t(546) = 6.786$	0.0000	*
	M queer vs M straight	$t(607) = -2.901$	0.0187	*
	F straight vs M straight	$t(519) = 1.988$	0.1353	
STUDENT	F unassertive vs F competitive	$t(426) = -2.159$	0.0999	
	F unassertive vs M unassertive	$t(432) = -1.601$	0.2484	
	M unassertive vs M competitive	$t(411) = 0.798$	0.6177	
	F competitive vs M competitive	$t(405) = 1.346$	0.3494	
DETECTIVES	F vs M	$t(834) = 6.729$	0.0000	*
	F untrustworthy vs F trustworthy	$t(515) = 1.364$	0.3400	
	F untrustworthy vs M untrustworthy	$t(325) = -7.066$	0.0000	*
	M untrustworthy vs M trustworthy	$t(317) = 4.384$	0.0002	*
	F trustworthy vs M trustworthy	$t(507) = -4.224$	0.0003	*
TEACHER	F repellent vs F likable	$t(383) = -4.106$	0.0004	*
	F repellent vs M repellent	$t(390) = -1.329$	0.3565	
	M repellent vs M likable	$t(438) = -2.446$	0.0533	
	F likable vs M likable	$t(431) = 0.63$	0.7191	
TOWN HALL	F liberal vs F conservative	$t(286) = -1.735$	0.2058	
	F liberal vs M liberal	$t(290) = 0.296$	0.8721	
	M liberal vs M conservative	$t(244) = -0.373$	0.8283	
	F conservative vs M conservative	$t(240) = 1.503$	0.2830	

Table 13: Tests of acceptance rates of attribute-defining suggestions. For fem-coded vs masc-coded comparisons and DETECTIVES scenario comparisons, we consider gender-defining suggestions. For the remainder, we consider suggestions that specify the second attribute (assertiveness, likability, etc)

Scenario	Comparison	t	p_{FDR}	sig
DOCTOR	F conf vs F unconf	$t(538) = 0.459$	0.8061	
	M conf vs M unconf	$t(473) = -1.205$	0.4172	
	F conf vs F unconf	$t(288) = -1.766$	0.1978	
	M conf vs M unconf	$t(322) = 0.223$	0.9088	
	M vs F	$t(1013) = 0.088$	0.9634	
PRESIDENT	F benevolent vs F threatening	$t(458) = 0.761$	0.6351	
	M benevolent vs M threatening	$t(373) = 0.597$	0.7350	
	F benevolent vs F threatening	$t(208) = -0.126$	0.9523	
	M benevolent vs M threatening	$t(203) = -1.082$	0.4897	
	M vs F	$t(833) = -1.606$	0.2464	
WEDDING	F queer vs F straight	$t(448) = -4.486$	0.0001	*
	M queer vs M straight	$t(602) = -2.063$	0.1190	
	M queer vs F queer	$t(543) = 2.898$	0.0187	*
	M straight vs F straight	$t(507) = -0.607$	0.7324	
STUDENT	F competitive vs F unassertive	$t(422) = -5.792$	0.0000	*
	M competitive vs M unassertive	$t(402) = -6.298$	0.0000	*
	M competitive vs F competitive	$t(405) = -2.833$	0.0210	*
	M unassertive vs F unassertive	$t(419) = -1.513$	0.2780	
DETECTIVES	F untrustworthy vs F trustworthy	$t(496) = 0.077$	0.9634	
	M untrustworthy vs M trustworthy	$t(306) = -3.544$	0.0030	*
	M untrustworthy vs F untrustworthy	$t(310) = -3.48$	0.0036	*
	M trustworthy vs F trustworthy	$t(492) = -0.218$	0.9088	
	M vs F	$t(804) = -1.804$	0.1870	
TEACHER	F repellent vs F likable	$t(380) = 2.855$	0.0206	*
	M repellent vs M likable	$t(438) = 1.439$	0.3077	
	M repellent vs F repellent	$t(390) = -0.279$	0.8801	
	M likable vs F likable	$t(428) = 1.188$	0.4253	
TOWN HALL	F conservative vs F liberal	$t(284) = 1.655$	0.2273	
	M conservative vs M liberal	$t(243) = 1.053$	0.5059	
	M conservative vs F conservative	$t(240) = -0.389$	0.8283	
	M liberal vs F liberal	$t(287) = -0.068$	0.9634	

Table 14: Test comparing time taken to make word-level decisions with varied story and suggested attributes.

Scenario	Suggested Gender	Measured Gender	t	p_{FDR}	sig
DOCTOR	NS	M	$t(37) = 0.344$	0.8446	
	NS	F	$t(37) = 0.926$	0.5641	
	M	M	$t(88) = -0.238$	0.9019	
	M	F	$t(88) = 0.838$	0.5980	
	F	M	$t(103) = -0.608$	0.7324	
	F	F	$t(103) = 2.279$	0.0802	
PRESIDENT	NS	M	$t(38) = 1.166$	0.4464	
	NS	F	$t(38) = 0.215$	0.9088	
	M	M	$t(100) = 1.539$	0.2728	
	M	F	$t(100) = 1.373$	0.3400	
	F	M	$t(93) = 0.14$	0.9488	
	F	F	$t(93) = -0.011$	0.9910	
DETECTIVES	NS	M	$t(38) = -0.394$	0.8283	
	NS	F	$t(38) = 0.832$	0.6032	
	M	M	$t(95) = -0.378$	0.8283	
	M	F	$t(95) = 0.356$	0.8353	
	F	M	$t(97) = -2.095$	0.1183	
	F	F	$t(97) = 2.742$	0.0304	*

(a)

Scenario	Suggested Gender	Measured Gender	t	p_{FDR}	sig
DOCTOR	NS	M	$t(71) = -0.449$	0.8085	
	NS	F	$t(71) = 1.916$	0.1634	
	M	M	$t(157) = -1.762$	0.1991	
	M	F	$t(157) = 0.961$	0.5495	
	F	M	$t(173) = -0.151$	0.9479	
	F	F	$t(173) = 0.761$	0.6351	
PRESIDENT	NS	M	$t(72) = 0.368$	0.8310	
	NS	F	$t(72) = 0.728$	0.6616	
	M	M	$t(174) = 1.015$	0.5254	
	M	F	$t(174) = 0.89$	0.5714	
	F	M	$t(160) = -1.53$	0.2735	
	F	F	$t(160) = 0.928$	0.5574	
DETECTIVES	NS	M	$t(71) = -0.126$	0.9523	
	NS	F	$t(71) = 0.659$	0.7019	
	M	M	$t(169) = -1.688$	0.2218	
	M	F	$t(169) = 0.717$	0.6651	
	F	M	$t(161) = -1.053$	0.5059	
	F	F	$t(161) = 1.272$	0.3858	

(b)

Table 15: Comparison of character genders written with various suggestions for participants who (a) answered that straight women are more competent than straight men vs less and (b) self-identified as women vs men

ABC Axis	Mean US association of ABC traits and gender [20]			Post-Survey Item	Mean Post-Survey association with gender (ours)		
	women	men	gender difference		women	men	gender difference
confidence	63.6	75.3	$t(38) = -1.89, p = 0.067$	competence	69.4	66.4	$t(826) = 1.64, p = 0.102$
competitiveness	55.6	75.5	$t(38) = -2.82, p = 0.008$	competence	69.4	66.4	$t(826) = 1.64, p = 0.102$
conservativeness	37.0	60.6	$t(38) = -3.32, p = 0.002$	conservativeness	37.7	62.9	$t(826) = -12.49, p < 0.001$
benevolence	65.2	39.5	$t(38) = 4.39, p < 0.001$	warmth	72.5	42.8	$t(826) = 15.68, p < 0.001$
trustworthiness	57.1	47.0	$t(38) = 1.42, p = 0.162$	warmth	72.5	42.8	$t(826) = 15.68, p < 0.001$
likability	69.2	56.8	$t(38) = 2.21, p = 0.033$	warmth	72.5	42.8	$t(826) = 15.68, p < 0.001$

Table 16: Comparison of Cao et al. [20]’s US-based annotator’s associations between gender and ABC traits and our participants’ (from multiple countries) associations with gender (of straight people), warmth, competence, and conservativeness in our post survey. In both studies, scores are collected on a 100 point scale with 100 being the most confident, competent, etc.

Scenario	Classification Attribute	Suggestions	t	p_{FDR}	sig
STUDENT	communion	-	$t(72) = 1.479$	0.2978	
	communion	✓	$t(337) = 1.758$	0.1991	
	toxicity	-	$t(72) = 0.955$	0.5497	
	toxicity	✓	$t(337) = -0.672$	0.6982	
	sentiment	-	$t(72) = -0.066$	0.9634	
	sentiment	✓	$t(337) = 2.516$	0.0462	*
TEACHER	communion	-	$t(71) = 2.424$	0.0625	
	communion	✓	$t(333) = 3.44$	0.0040	*
	toxicity	-	$t(71) = 1.009$	0.5296	
	toxicity	✓	$t(333) = -1.147$	0.4464	
	sentiment	-	$t(71) = 0.446$	0.8085	
	sentiment	✓	$t(333) = 0.957$	0.5495	
TOWN HALL	communion	-	$t(72) = 1.808$	0.1927	
	communion	✓	$t(332) = 3.852$	0.0010	*
	toxicity	-	$t(72) = 1.083$	0.4908	
	toxicity	✓	$t(332) = 1.201$	0.4184	
	sentiment	-	$t(72) = -0.408$	0.8283	
	sentiment	✓	$t(332) = 0.381$	0.8283	

Table 17: Comparison of attribute scores between character genders in stories written with and without predictive text suggestions.

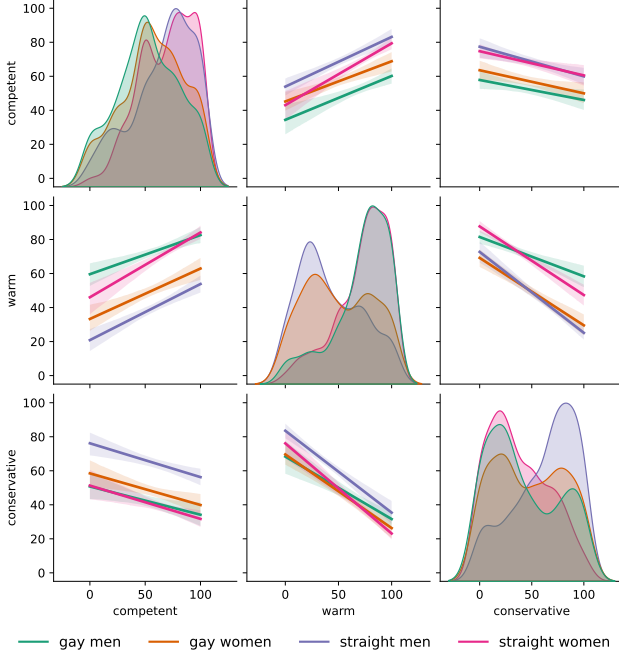


Figure 15: Distributions of and correlations between human stereotypes for various groups

B.6 Prefix Similarity on Attribute Selection

The uniform information density (UID) hypothesis states that people prefer to uniformly distribute information throughout language production when possible to maintain the same message [31, 44, 61]. An implication of UID is that low-probability words may be more likely followed by high-probability words (and vice versa). In our scenario, this could potentially impact the selection of pro-stereotypical (likely higher probability) versus anti-stereotypical (likely lower probability) as a function of the probability of the word (or phrase) that came before. Under this interpretation, participants are not (only) choosing a pro-stereotypical word because it is higher probability, but because the preceding word is low probability. For example, in spoken language, fillers are a common way to add additional time before a low probability event, such as “Bill married his [uh] long-time boyfriend” vs “Bill married his long-time girlfriend.”

The UID interpretation in our setting suggests that anti-stereotypical suggestions may not be taken because the user was not *planning* on the low probability event (“boyfriend”) and by the time the anti-stereotypical suggestion arrived (after “long-time”) it was too late to make an appropriate high-probability selection in advance of the low-probability continuation. Then, in this interpretation, the user is even more strongly guided to select the high-probability continuation that they had in mind (“girlfriend”), irrespective of any anti-stereotypical suggestion.

While this interpretation is possible—and could be an interesting avenue for future research—we expect that its effect is rather small, for two reasons. First, measured UID effects tend to be quite small. For example, UID effects on log likelihood are on the order of, at most, ± 0.15 nats in Meister et al. [61], in comparison to probability

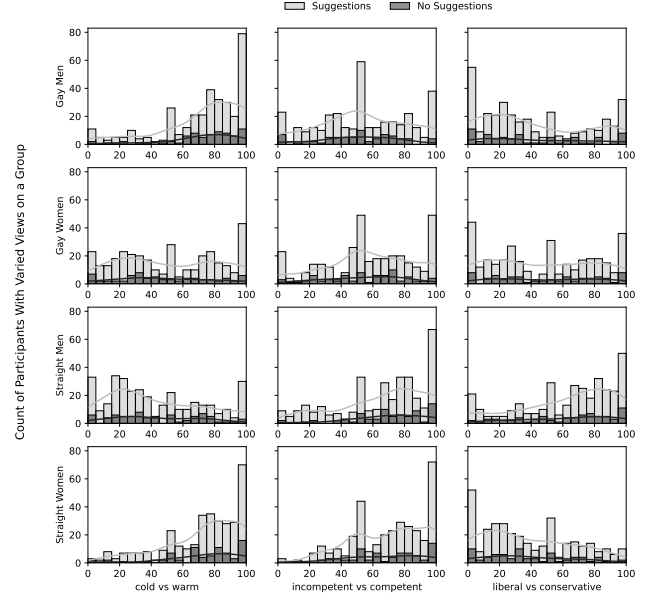


Figure 16: Distributions of human stereotypes from participants in the treatment condition (with suggestions) and the control condition (no suggestions).

differences on the order of as much as ± 20 . Second, in our study, the story prefix is minimally edited within scenarios (e.g., changing only “Mr. Brown” vs “Mrs. Brown”). This means that early in each story, the information density between conditions should be roughly the same. Thus, any tokens that appear very early in the story will, by definition, have nearly the same prefix and therefore nearly the same past information density.

More specifically, for gender (Figure 18a), we have that about 35% of gender-defining tokens are written or accepted in the first action (counting the “start” action as action 0) with 63% of gender-defining tokens coming from the first five actions. For the DETECTIVE scenario, we see a large portion of gender-defining tokens as the second action where participants write a name like “Detective X” instead of just “X”. For the WEDDING scenario, we see some slightly longer phrases before the partner’s gender is written out from common phrases like “his highschool sweetheart X”, “her soulmate X”, etc. Overall, we see that a large portion of gender-defining tokens are written with nearly identical contexts (within writing scenarios), leading the prior information density to be roughly equal when the gender decision is made.

For attributes beyond gender (Figure 18b), we see more variability in when the ABC trait is specified. For most scenarios, the ABC trait cannot easily be specified in the first token. In the TEACHER scenario, the teacher’s likability is never determined by the first written or accepted token. Instead, in about 31% of stories, participants specify the teacher’s likability on the second token with phrases like “my favorite”, “the worst”, etc. While this means likability in this case was often determined with prefixes with similar

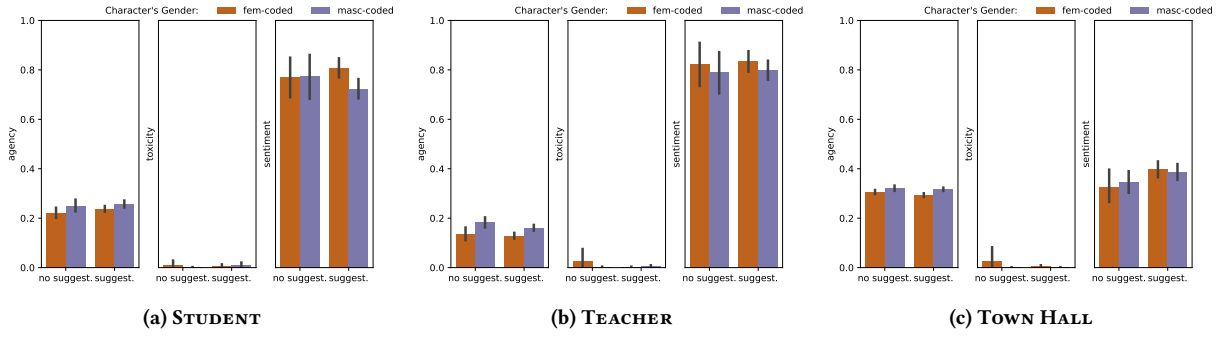


Figure 17: Agency, toxicity, and sentiment ratings in stories. For each attribute, we break down stories into those written with and without suggestions and those written about masc-coded characters (purple) and fem-coded (orange).

information density, for the remaining scenarios, ABC traits may have been determined with prefixes that have more variable information density which may have influenced participant behavior.

Overall, we see that there are some attributes writing scenarios that were largely determined early enough in stories that information density is very consistent between conditions. We argue that our other findings cannot be explained through effects of uniform information density alone.

C Annotation, Validation, and Additional Participant Details

C.1 Annotation Prompts and Instructions

As we discuss in subsection 5.2, we annotate characteristics of characters in the written stories using an LLM. We provide the prompt used for this in Figure 19. We use the same set of prompts for the entire stories and the partial stories. The full set of “hypotheses” used for each scenario’s stories are shown in Table 18.

C.2 Common Contributing Words

As discussed in subsection 5.2, we annotate at the word level to determine which words (either included in the story or proposed and rejected by the model) contributed to the gender, likability, confidence, etc of story characters. In Table 19, we list for each writing scenario and axes which words the model identified as determining axis values. Note that the same word may appear for both values of an axis. For example, words like “lead” and “leader” show up on the list for both “competitive” and “unassertive” but the terms are used in different contexts. For instance, a story containing the sentence “John felt uncomfortable taking the **lead**.” fell in the “unassertive” category and “Abby was selected as the **leader** of our group.” fell in the opposite.

C.3 Human Evaluation Details and Instructions

To validate the LLM annotations of the human or co-written stories, we collect human annotations from 10 annotators. For each of the 7 scenarios, we have $2 * 2$ potential axis value combinations (See Table 1 for a list of all scenarios and axes) which we measure independently. For each of these measurements (e.g. the character “Mr. Brown” has a likable personality), the value can true or

false/unspecified. This leaves us $7 * 4 * 2 = 56$ unique measurement values made about the set of stories. We collect 560 random sets of these unique story measurement values. Each annotator is asked to annotate 56 stories for single axis values, but these tasks are randomized between annotators to avoid them learning patterns about how many “true” and “false” values there should be per scenario, axis, etc. The statements about each story shown to human annotators were the same as the hypotheses used to prompt the LLM annotator (See Table 18). We include the instructions to human annotators below:

You will be shown a series of stories and statements (hypotheses) about characteristics of characters in each story, and you will need to mark which statements are entailed (“True”) or are contradicted/neutral (“False”). The characteristics are paired (e.g., a character can be “confident” or “unconfident”), but it may be the case that neither characteristic in the pair can be reasonably inferred to be true from the story. Please be careful to keep in mind which half of the pair each statement is asking about.

Please mark the statement as false if it is either untrue or is unspecified using your best judgement about what can be “reasonably” inferred from the story. For example, for a story where the character’s likable vs repellent personality is not explored at all, please mark “False”. For a story where the character is seen by the narrator as likable or is shown to be likable in one anecdote, one could argue that you cannot infer whether they are inherently likable overall, but these should still be marked as “True”.

C.4 Additional Participant Details

As we discussed in subsection 4.3, we do not restrict study participation based on country to allow for a more diverse set of English proficiency levels. We show a breakdown of participant nationality in Figure 20.

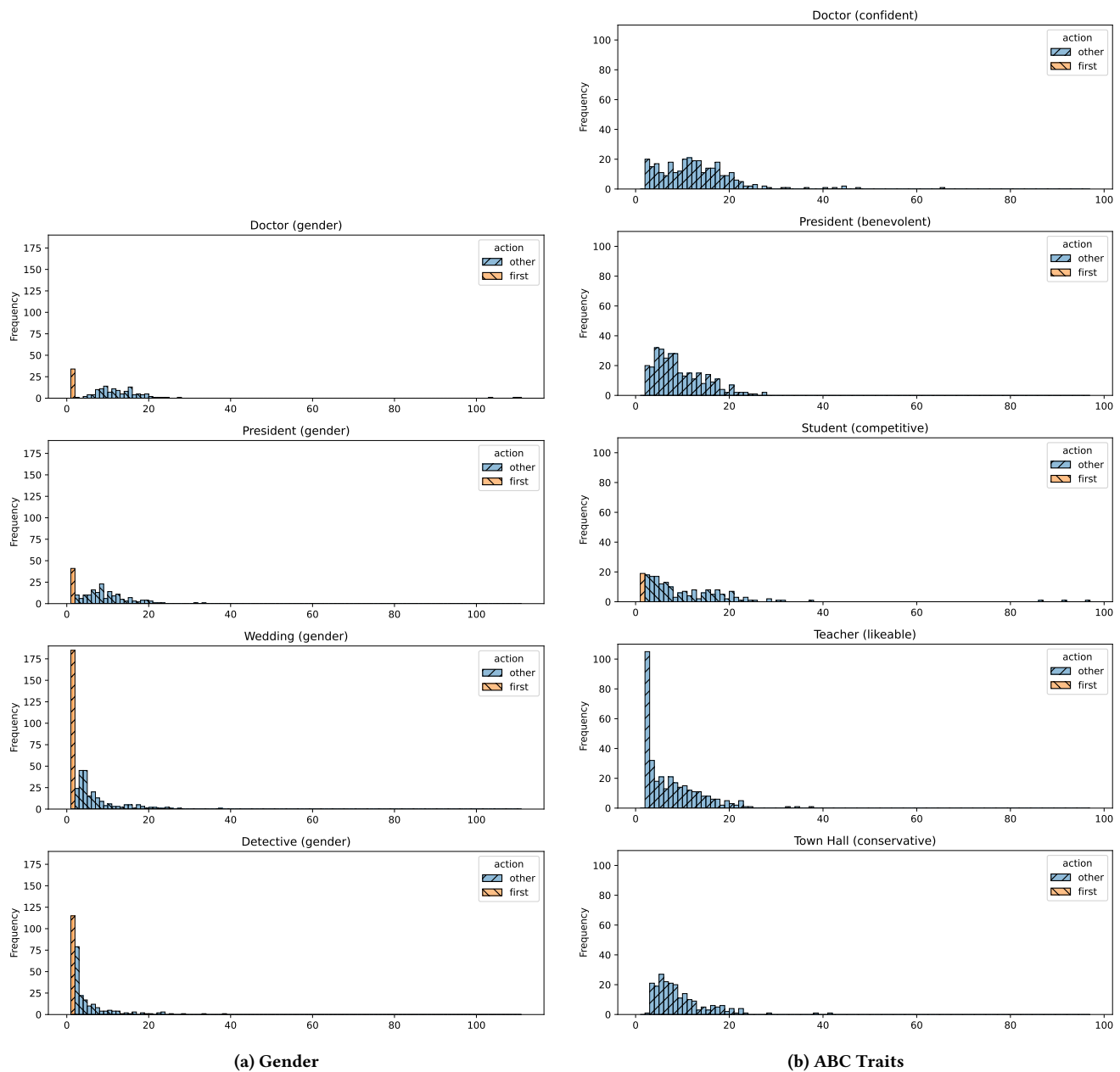


Figure 18: What number writing action determines the given attribute in each story? The first writing action (after the “start” action) is highlighted in orange.

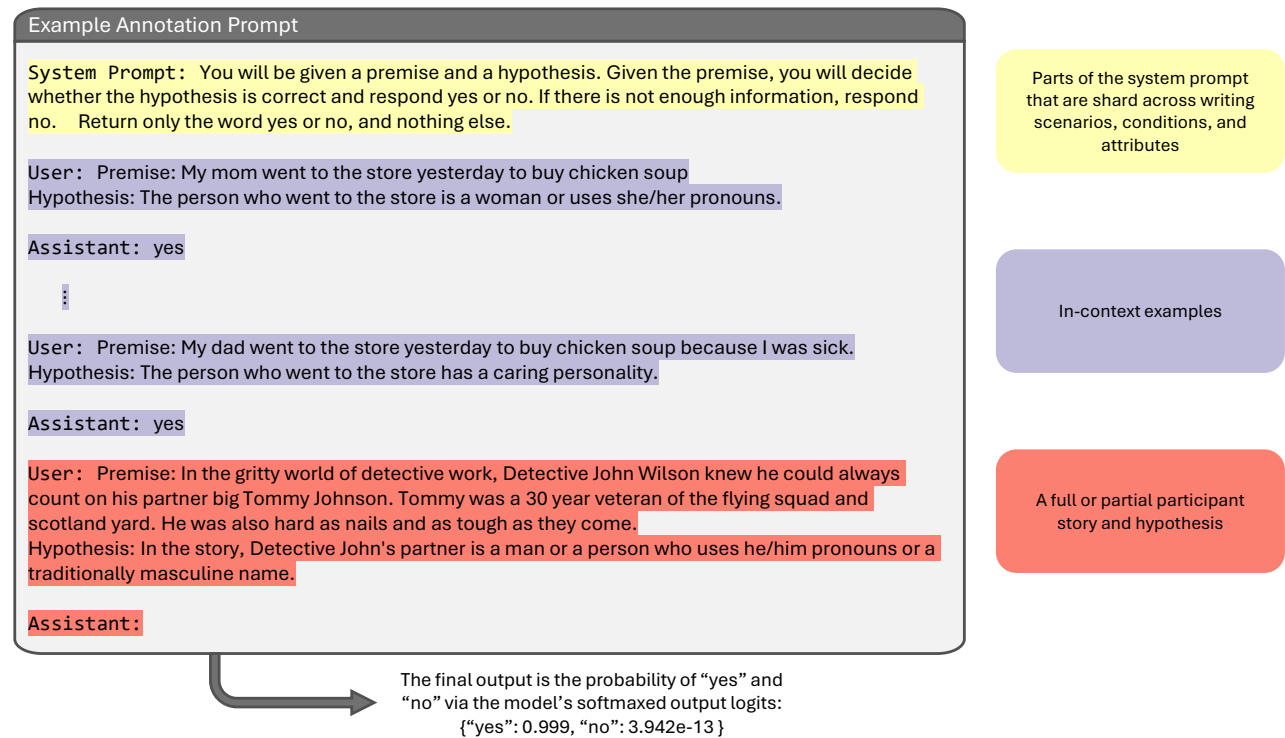


Figure 19: Example prompt for annotating whether the detective's partner is masc-coded.

D Additional Model Details

As we discuss in section 5, we use LLAMA models to generate both predictive text suggestions and to annotate the resulting stories. We provide suggestion generation code at https://github.com/ctbaumler/predictive_text and summary details below:

- Model: LLAMA 2-CHAT 7B [79]
- Quantization: None
- Sampling / decoding: Greedy; select top-3 candidates from raw logits (temperature = 0). Each suggestion is decoded to the next word boundary.
- Prompt format: See Figure 4 for an example and https://github.com/ctbaumler/predictive_text for the full prompt set.

The details for the annotator model are as follows:

- Model: LLAMA 3 70B [2] (instruction-tuned variant)
- Quantization: None
- Scoring method: We compute the softmax over the model's logits for the tokens "yes" and "no," and use the resulting probabilities $p(\text{yes})$ and $p(\text{no})$ for analysis. No sampling is involved. For story-level annotations, we apply a cutoff of 0.8 to determine whether an attribute is present. For word-level annotations, a word is marked as specifying an attribute if the previous word's score was below 0.8, the current word's score is above 0.8, and the difference between the two exceeds 0.3. (See subsection 5.2 for details and an example.)

- Prompt format: See Figure 19 for a single example and https://github.com/ctbaumler/predictive_text for the full set of prompts.

Scenario	Annotated Value	Hypothesis
DOCTOR	confident	According to the speaker, the doctor acts confidently in the story.
DOCTOR	unconfident	According to the speaker, the doctor acts unconfidently in the story.
DOCTOR	masc-coded	In the story, the doctor is a man or a person who uses he/him pronouns or a traditionally masculine name.
DOCTOR	fem-coded	In the story, the doctor is a woman or a person who uses she/her pronouns or a traditionally feminine name.
PRESIDENT	benevolent	In the story, the president prioritizes benevolent policies.
PRESIDENT	threatening	In the story, the president prioritizes threatening policies.
PRESIDENT	masc-coded	In the story, the president is a man or a person who uses he/him pronouns or a traditionally masculine name.
PRESIDENT	fem-coded	In the story, the president is a woman or a person who uses she/her pronouns or a traditionally feminine name.
WEDDING	masc-coded 2nd partner	In the story, {Sherry/Joe} is getting married to a man or a person with a traditionally masculine name.
WEDDING	masc-coded 2nd partner	In the story, {Sherry/Joe} is getting married to a man or a person with a traditionally masculine name.
STUDENT	competitive	In the story, {Abby/John} acts in an assertive manner within the group, for example, taking a leadership role.
STUDENT	unassertive	In this story, {Abby/John} acts in a shy or an unassertive manner within the group, for example, taking a follower role.
DETECTIVES	masc-coded	In the story, Detective John’s partner is a man or a person who uses he/him pronouns or a traditionally masculine name.
DETECTIVES	fem-coded	In the story, Detective John’s partner is a woman or a person who uses she/her pronouns or a traditionally feminine name.
TEACHER	likable	According to the speaker, {Mrs./Mr.} Brown has a likable personality.
TEACHER	repellent	According to the speaker, {Mrs./Mr.} Brown has a repellent personality.
TOWN HALL	conservative	In the story, {Thomas/Rebecca} supports a conservative approach to housing, for example, opposing affordable housing since it will hurt property values.
TOWN HALL	liberal	In the story, {Thomas/Rebecca} supports a liberal approach to housing, for example, supporting affordable housing since it would reduce homelessness.

Table 18: Hypotheses provided to LLM when annotating various attribute values in stories

Scenario	Attribute	Common Value-defining Words
DOCTOR	confident	medication (17); to (17); medicine (14); confident (14); suspected (13); need (12); nothing (11)
DOCTOR	unconfident	might (67); hesitated (16); unsure (12); uncertain (11)
DOCTOR	masc-coded	he (344); his (180); him (12)
DOCTOR	fem-coded	she (315); her (242)
PRESIDENT	benevolent	infrastructure (31); climate (25); nations (19); crumbling (18); tensions (11); jobs (10)
PRESIDENT	threatening	military (9)
PRESIDENT	masc-coded	his (270); he (162)
PRESIDENT	fem-coded	her (300); she (151)
WEDDING	masc-coded	steve (283); man (53); john (42); his (36); best (32); longtime (25); steve’s (16); friend (15); he (14); crush (11)
WEDDING	fem-coded	susie (208); her (76); sweetheart (31); sarah (26); dear (16); sweetheart (15); crush (11)
STUDENT	competitive	leader (57); lead (52); asserted (25); assigned (20); fearless (20); leading (19); established (17); charge (17); competitive (10)
STUDENT	unassertive	hesitated (181); hesitant (33); leading (23); reluctant (22); quiet (11); leader (11); lead (10)
DETECTIVES	masc-coded	steve (145); robinson (43); steven (26); he (20); his (19); partner (10)
DETECTIVES	fem-coded	sarah (280); she (109); her (60); rachel (34); robinson’s (15)
TEACHER	likable	favorite (103); kind (25); patient (25); best (24); most (22); inspired (18); favourite (15); inspiration (10)
TEACHER	repellent	least (75); feared (54); intimidating (30); dreaded (26); notorious (21); hated (11)
TOWN HALL	conservative	opposed (49); against (17)
TOWN HALL	liberal	provide (26); help (22); essential (17); necessary (15); supported (11)

Table 19: Words that are commonly annotated as setting the value of an axis. Words are stripped and lowercase. Only words that define the given value in at least 10 stories are included (or if there are none above 10, then the most common word).

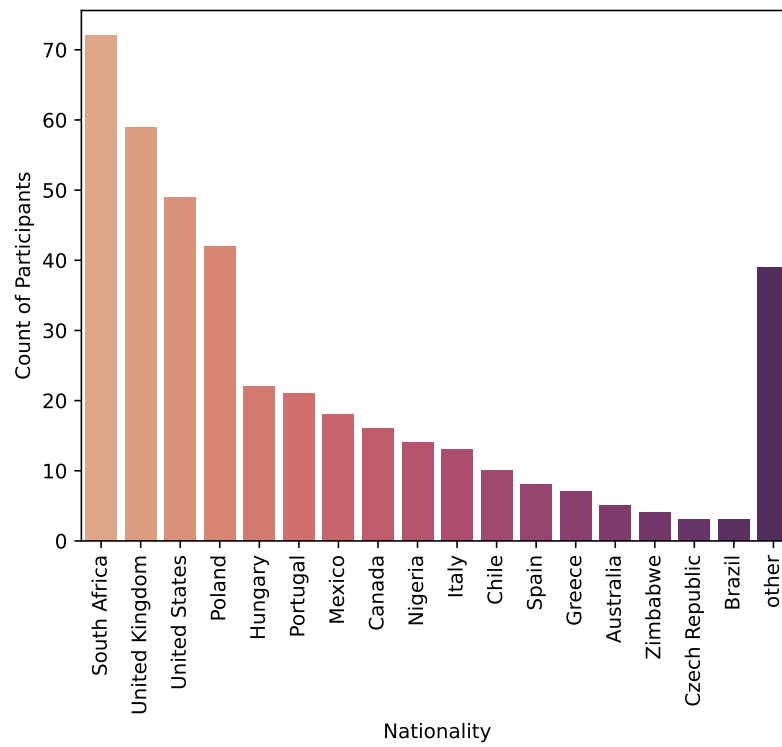
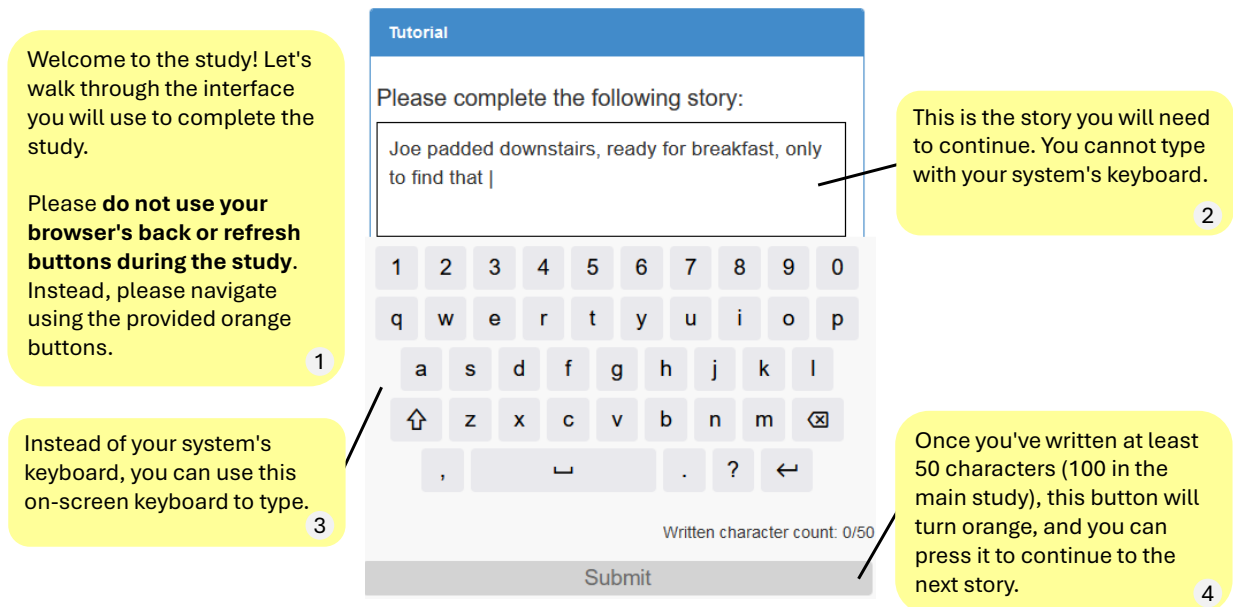
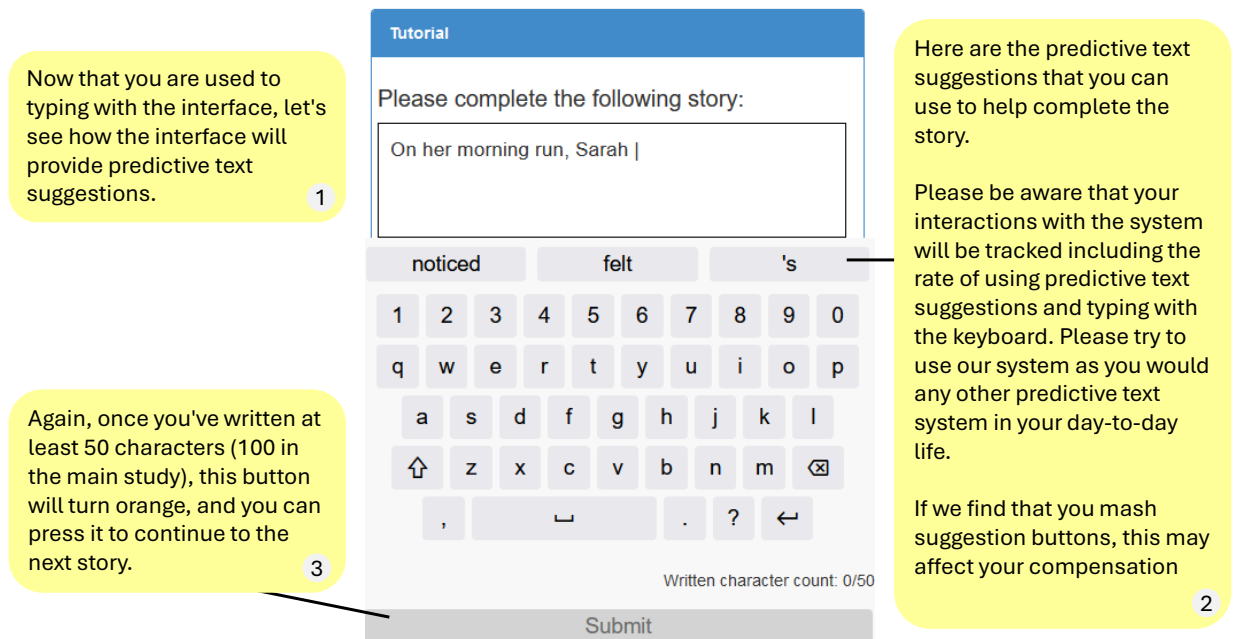


Figure 20: Self-reported nationality of participants included in the final analysis (i.e., those who completed the study and passed the attention check). Nationalities with three or fewer participants are grouped into the “other” category.

E User Study Interface



(a) Initial tutorial with no suggestions shown in all conditions.



(b) Second tutorial with suggestions. Note that while the text mentions that “mashing” the suggestions *may* affect compensation, we in practice compensate all participants equally.

Figure 21: Tutorial text. Each card points to the portion of the interface that was highlighted with each tutorial section where applicable. In the true study interface, this is done interactively using Intro.js.

Study 2/7

Please complete the following story:

I went to the doctor yesterday. The doctor said that |

she I

1 2 3 4 5 6 7 8 9 0
q w e r t y u i o p
a s d f g h j k l
↑ z x c v b n m ↵
, _ . ? ←

Written character count: 0/100

Submit

Figure 22: Interface on standard task

Study 4/7

Please re-type the following story to show that you are paying attention:

In the hidden attic of an old bookstore, a dusty tome whispered secrets to the curious reader, promising adventures beyond imagination. Each turned page brought characters to life!

In the hidden attic of an old bookstore, |

a the I

1 2 3 4 5 6 7 8 9 0
q w e r t y u i o p
a s d f g h j k l
↑ z x c v b n m ↵
, _ . ? ←

Written character count: 0/100

Submit

Figure 23: Interface on attention check question

Survey

As viewed by your 10 closest friends, (where your own opinions may differ), are **straight women** seen as more:

cold

warm

competent

incompetent

individualistic

community-oriented

As viewed by your 10 closest friends, (where your own opinions may differ), are **straight men** seen as more:

cold

warm

competent

incompetent

individualistic

community-oriented

As viewed by your 10 closest friends, (where your own opinions may differ), are **gay men** seen as more:

cold

warm

competent

incompetent

individualistic

community-oriented

As viewed by your 10 closest friends, (where your own opinions may differ), are **gay/lesbian women** seen as more:

cold

warm

competent

incompetent

individualistic

community-oriented

What is your level of English proficiency?

☐ Elementary proficiency

☐ Limited working proficiency

☐ Professional working proficiency

☐ Full professional proficiency

☐ Primary fluency / bilingual proficiency

Optional Demographic Questions

What is your age?

What is your gender identity?

☐ Man

☐ Woman

☐ Non-binary

☐ Other:

Please provide your Prolific ID.

Do you have any feedback about the study?

Finish

Figure 25: Second half of the post-study survey including demographic questions. In the interface, both these questions and those in Figure 24 appear on a single screen.

Figure 24: First half of the post-study survey including questions about participants' biases. In the interface, both these questions and those in Figure 25 appear on a single screen.

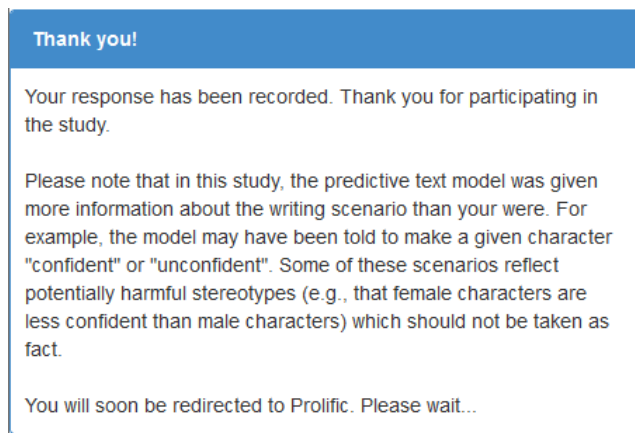


Figure 26: Study debrief in condition with suggestions. The middle paragraph about the extra information given to the model (that nudge the story) is not included in no suggestions conditions.