

<sup>1</sup> Measuring machine learning harms from stereotypes:  
<sup>2</sup> Requires understanding who is being harmed by which  
<sup>3</sup> errors in what ways

<sup>4</sup> Anonymous

<sup>5</sup> **Abstract**

<sup>6</sup> As machine learning applications proliferate, we need an understanding of their potential  
<sup>7</sup> for harm. However, current fairness metrics are rarely grounded in human psychological  
<sup>8</sup> experiences of harm. Drawing on the social psychology of stereotypes, we use a case study  
<sup>9</sup> of gender stereotypes in image search to examine how people react to machine learning  
<sup>10</sup> errors. First, we use survey studies to show that not all machine learning errors reflect  
<sup>11</sup> stereotypes nor are equally harmful. Then, in experimental studies we randomly expose  
<sup>12</sup> participants to stereotype-reinforcing, -violating, and -neutral machine learning errors.  
<sup>13</sup> We find stereotype-reinforcing errors induce more experientially (i.e., subjectively) harm-  
<sup>14</sup> ful experiences, while having minimal changes to cognitive beliefs, attitudes, or behaviors.  
<sup>15</sup> This experiential harm impacts women more than men. However, certain stereotype-  
<sup>16</sup> violating errors are more experientially harmful for men, potentially due to perceived  
<sup>17</sup> threats to masculinity. We conclude that harm cannot be the sole guide in fairness miti-  
<sup>18</sup> gation, and propose a nuanced perspective depending on who is experiencing what harm  
<sup>19</sup> and why.

## 20 Introduction

21 Over the past decade, researchers have demonstrated that machine learning models run the risk of  
22 learning stereotypical associations. For example, natural language processing (NLP) models have  
23 been shown to associate women with homemakers and men with programmers [9], while computer  
24 vision models have been shown to associate women with shopping and men with driving [88]. These  
25 associations can cause machine learning models to make systematic errors, mistakenly reporting, for  
26 instance, that female doctors are nurses and that male nurses are doctors [71]. A rich literature has  
27 developed offering many more such examples, spurring calls to address the risk that machine learning  
28 models might make mistakes that perpetuate harmful stereotypes.

29 Unfortunately, portions of this literature have suffered from three main limitations that impede  
30 effective intervention. First, while some important prior work has asked crowd workers to annotate  
31 when errors invoke stereotypes or drawn on pre-existing inventories of stereotype from the psychology  
32 literature [7, 9, 12, 13, 73, 80], machine learning researchers often rely on their own moral intuitions  
33 to determine which categories of associations to investigate (e.g., associations between gender and  
34 occupation) and to judge which specific associations (e.g., the association between women and nurs-  
35 ing) are stereotypes. As a result, certain categories of associations that broader populations might  
36 perceive as stereotypical have not been investigated and many discovered associations within these  
37 categories are assumed to be stereotypical, even if broader populations would not perceive them to  
38 be so. For example, researchers identified an object recognition model as reproducing stereotypes  
39 because it amplifies the degree to which labels for kitchen items like “knife,” “fork,” and “spoon” are  
40 incorrectly assigned to photos featuring women, and labels for technology-related items like “key-  
41 board” and “mouse” are incorrectly assigned to photos featuring men [88]. However, these intuitions  
42 might not always be shared by the broader population; indeed, as we’ll show later in this paper,  
43 some of these associations are not commonly seen as invoking stereotypes among a sample of people  
44 in the United States.

45 Second, prior work has tended to ignore whether errors reinforce or violate stereotypes, treating  
46 each type of error as equally harmful. Researchers have even occasionally treated all spurious corre-  
47 lations between objects and specific social groups as harmful, even if such correlations may not be  
48 perceived as reinforcing or violating any stereotype. Blodgett et al [8] has documented many instances  
49 of this, for example how remarks following a statement about “Norwegian salmon” are nonsensically  
50 used to assess stereotypes. As a result, it remains unclear the degree to which the harmfulness of  
51 errors depends on whether and how these errors invoke stereotypes.

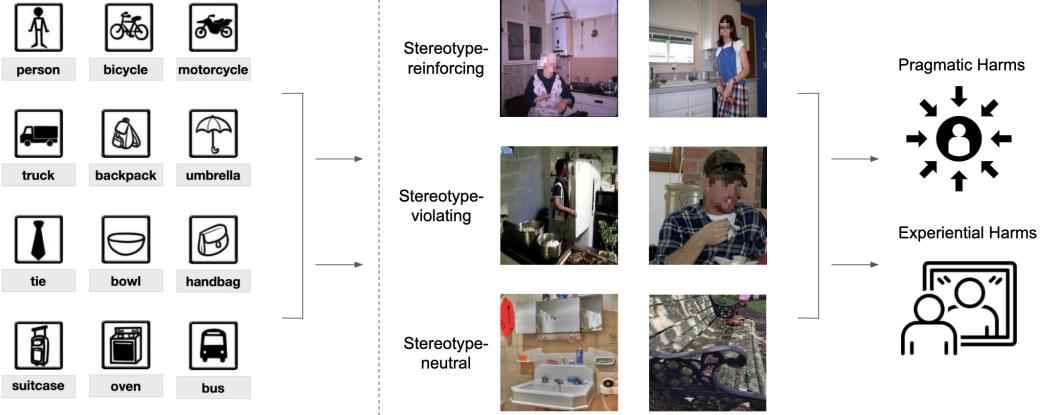
52 Third, while stereotypes are frequently invoked to explain why some machine learning errors are  
53 more harmful than others [1, 4, 6, 81], little has been done to establish the actual impact that these

54 stereotypes have on people in practice. Rather than measuring the harms that stereotypes bring  
55 about and identifying the particular mechanisms by which they do so, much of the literature simply  
56 presumes that stereotypes have negative impacts on the people so stereotyped or on society more  
57 generally. While some prior work has found that exposure to gender-biased image search results  
58 can lead both to more biased estimations of the representation of different gender groups in certain  
59 occupations and to a decreased sense of belonging [43, 50], there remains a paucity of empirical  
60 work that seeks to measure the psychological and practical effects of exposure to such errors and to  
61 characterize the nature of these harms.

62 In this paper, we draw on the psychological literature on social stereotypes to try to overcome each  
63 of these limitations, performing a series of empirical studies in which human subjects are exposed  
64 to machine learning errors. As a concrete application in which to ground these studies, we focus on  
65 gender stereotypes in the machine learning task of object recognition, which is now a common feature  
66 in photo management software like Apple Photos and Google Photos.

67 First, we investigate which associations people view as stereotypical. Recognizing from the psy-  
68 chological literature that the space of stereotypical associations is potentially much broader than  
69 what has been examined in the machine learning literature [24, 27, 37, 56], we make no a priori deci-  
70 sions about the categories of association worth investigating. We expand our scope of analysis beyond  
71 commonly focused-on categories—for example, occupations or activities—to include all objects that  
72 might appear in an image. We do so via two studies. In the first (Study 1), we present participants  
73 with images from popular computer vision datasets that include a range of objects (COCO [46] and  
74 OpenImages [45] datasets (Fig. 2)), and ask whether these objects are stereotypically associated  
75 with different gender groups. We further expand on this in a second study (Study 2), in which we  
76 conduct qualitative analysis on participants’ open-ended responses to the questions in Study 1 to  
77 better understand why certain objects are seen as stereotypes. We find that participants do not all  
78 agree on whether particular objects reflect stereotypical associations, and that even when objects  
79 are seen as reflecting stereotypes, the reasons that participants provide are varied. For example, the  
80 association between women and cats is variously explained by “cat lady,” “women are called *kitten*,”  
81 and “women are called *cougars*.” We additionally ask participants whether and why they find these  
82 stereotypical associations to be harmful. We find that while many participants describe stereotypes as  
83 self-evidently harmful, others differentiate reasons such as whether the error invokes a stereotypical  
84 association perceived as prescriptive (i.e., how members of certain gender groups should behave) or  
85 proscriptive (i.e., how they should not behave), providing a richer set of psychological reasons [63, 66].

86 Then, to overcome the second and third limitations of imprecision about which kind of stereo-  
87 typical errors may cause what kind of harm, we attempt to causally assess the effect that exposure  
88 to these stereotypes have on people in practice. To do so, we draw on psychological theories of



**Fig. 1 Summary of our studies.** The left side represents studies 1 and 2, where we ask human participants to mark which of the relevant objects in our application are stereotypically associated with which gender groups, as well as to qualitatively explain why that is and why it is harmful or not. The right side represents studies 3 and 4 where we randomly expose participants to machine learning errors which are stereotype-reinforcing, stereotype-violating, or stereotype-neutral—as determined by the annotations from study 1. Then, we measure two forms of harm we introduce: pragmatic (measurable changes in someone’s cognitive beliefs, attitudes, or behaviors toward the group being stereotyped) and experiential (self-reports of negative affect). The images shown are examples of misclassifications of *oven* where stereotype-reinforcing errors are when it is falsely predicted on a woman, stereotype-violating when on a man, and stereotype-neutral when no stereotypes are invoked.

89 stereotypes to conceptualize harm as well as the relationship between stereotype and harm. We then  
90 investigate concretely the degree to which stereotypes cause such harm (Fig. 1).

91 The psychology literature describes stereotypes as cognitive beliefs in people’s minds, which can  
92 have an influence on attitudes (i.e., prejudice) and behaviors (i.e., discrimination) [2, 35, 37, 41, 47].  
93 For example, people may have *cognitive beliefs* that women are more warm but less competent than  
94 men, and thus *express* protective attitudes and pity for women [15, 32]. People then *behave* in ways  
95 that maintain women’s warmth and discount their competence, such as being less likely to promote  
96 women to leadership positions [24, 27]. Members of stereotyped social groups may also experience  
97 changes in their own beliefs and attitudes, causing them to behave differently; for example, when  
98 women are given a math exam and told that the exam is diagnostic of their own intellectual abilities,  
99 stereotypes of women as less capable of math negatively impact their performance on the exam [74].  
100 Therefore, stereotypes of certain social groups can prompt shifts in attitudes and behaviors that can  
101 ultimately harm the stereotyped group.

102 Beyond these changes in beliefs, attitudes, and behaviors, members of the stereotyped social group  
103 may also feel disrespected, demeaned, or otherwise discounted. Such experiences can be thought of as  
104 dignitary harms that bring about negative affect in members of the group so stereotyped [42]. While  
105 dignitary harms are frequently treated as less important than the concrete effects of changes in beliefs,  
106 attitudes, and behaviors, they may impose a substantial emotional toll, akin to that of microaggres-  
107 sions [64]. We therefore introduce a distinction between *pragmatic harms*, which involve measurable  
108 changes in someone’s cognitive beliefs, attitudes, or behaviors toward the group being stereotyped,  
109 and *experiential harms*, which involve self-reports of negative affect (Fig. 1). We additionally dif-  
110 ferentiate between errors that are stereotype-reinforcing, stereotype-violating, or stereotype-neutral.

111 Errors that invoke stereotypes may do so in different ways and may therefore have different effects.  
112 With these in place, we set out to measure the degree to which different kinds of errors bring about  
113 pragmatic and experiential harm. Concretely, we randomly assign participants to synthesized search  
114 result pages which contain different kinds of errors and we investigate three related hypotheses. First,  
115 we hypothesize that errors that reinforce social stereotypes will be perceived as more harmful than  
116 those that do not. Second, we hypothesize that stereotype-reinforcing errors will result in pragmatic  
117 harm, while stereotype-neutral or stereotype-violating errors will not. Third, we hypothesize greater  
118 reports of experiential harm on stereotype-reinforcing errors for the stereotyped group.

119 To examine if people experience pragmatic harms, we measure cognitive, attitudinal, and behav-  
120 ioral changes between people who experience machine learning outputs containing stereotypes,  
121 varying whether the errors reinforce or violate stereotypes. To measure experiential harms, we ask  
122 people to self-report negative affect when exposed to machine learning outputs containing stereotypes,  
123 again varying whether these errors reinforce or violate stereotypes.

124 We find little immediate causal effect of pragmatic harms, but sizable evidence that stereotype-  
125 reinforcing errors are experientially harmful—a finding that is more pronounced among participants  
126 who identify as women compared to those who identify as men (Study 3). We find that while the  
127 stereotyped group (e.g., women) generally finds it more experientially harmful for the error to rein-  
128 force rather than violate stereotypes, this is not true when it comes to clothing-related items typically  
129 associated with women (e.g., **cosmetics**, **necklaces**) being misclassified on men. Here, we see a  
130 backlash towards violations of the norms around gender presentation where men tend to find these  
131 misclassifications of, e.g., **cosmetics**, more harmful on men rather than women. This last observa-  
132 tion calls into question the idea that it is always normatively desirable to reduce errors perceived as  
133 more harmful due to their relationship to stereotypes (Study 4).

134 All studies are approved by our institution IRB. Studies 1 ([https://osf.io/cpyn4/?view\\_only=3ff6c9625f0c4fce864960ee47b0433a](https://osf.io/cpyn4/?view_only=3ff6c9625f0c4fce864960ee47b0433a)), 3 ([https://osf.io/m9akd/?view\\_only=b61bc54308a5481a96e43db2dac23498](https://osf.io/m9akd/?view_only=b61bc54308a5481a96e43db2dac23498), [https://osf.io/v2w4m/?view\\_only=5f68252ff9864ddfa198097fdd78e803](https://osf.io/v2w4m/?view_only=5f68252ff9864ddfa198097fdd78e803)), and part of Study 4 ([https://osf.io/xpv5j/?view\\_only=55aac464d7694e81ae69eccf86cd004f](https://osf.io/xpv5j/?view_only=55aac464d7694e81ae69eccf86cd004f)) are pre-registered on OSF, while Study 2 is more exploratory.  
135 By bringing greater clarity to different types of machine learning errors based on their relationship  
136 to a stereotype and embracing the rich psychological experiences behind them, we urge researchers  
137 and practitioners to more carefully consider different kinds of machine learning errors, potential  
138 harms, and the relevant relationships between them. Investigating the psychological experiences that  
139 people have when encountering machine learning errors is critical to understanding the potential  
140 harm of a system, and in turn, mitigating it.



**Fig. 2** COCO and Open Images object recognition datasets. We use two commonly used image recognition datasets to represent the application of a photo search engine. Both datasets contain annotations for perceived binary gender expression of the people in the images and the objects present in each image. The left panel shows one example figure from COCO annotated with objects such as `oven` and `bowl`. The right panel shows one example figure from Open Images annotated with objects such as `person` and `skirt`.

## 145 Results

146 We explore a popular task in machine learning known as object recognition (i.e., classifying the  
 147 objects present in an image). To make it concrete for our human studies, we use it in the context  
 148 of a smart phone’s photo search engine, and examine gender stereotypes. Specifically, we consider  
 149 one type of machine learning error called a false positive: when an object is predicted to be present  
 150 in an image when it is in fact not there. This causes the image with a false positive to be wrongly  
 151 surfaced on an image search results page.<sup>1</sup> In our work, we are only concerned with the *effect* of the  
 152 misclassification, and not why the model may have made the mistake, or what the participant thinks is  
 153 the reason the model made the mistake. Unlike prior work auditing search engines [43, 50, 60, 61, 79],  
 154 our sole focus is on tracing the concrete effects that search results can have.

155 **Study 1: Distinguishing which machine learning errors reflect social  
 156 stereotypes**

157 To understand the social stereotypes held by American society relevant to our machine learning task,  
 158 we first elicit human judgments ( $N = 80$ ) on Common Objects in Context (COCO) [46]. COCO has 80  
 159 objects and perceived binary gender expression of pictured people annotated across the images [87].  
 160 In the study, we ask the participants whether each object (e.g., `keyboard`, `zebra`) is stereotypically  
 161 associated with men, women, or neither. As expected, not all objects reflect gender stereotypes. This

<sup>1</sup>We note that false negatives are subsumed in this setting because enough false positives will crowd out the results page and ultimately have a similar effect as false negatives on images of the gender that does not have false positives.

### Stereotyped with Women

handbag 21/23	hair drier 17/20	wine glass 8/13	cat 11/19	potted plant 9/17	oven 12/23	cake 8/19	vase 7/18
dining table 5/21	tennis racket 3/14	cow 3/15	sink 4/20	teddy bear 4/20	horse 5/26	bird 4/22	person 4/26
sandwich 3/23	umbrella 3/25	parking meter 2/19	toaster 2/21	refrigerator 2/23	sheep 2/24	suitcase 2/25	backpack 2/26
bench 1/14	apple 1/17	snowboard 1/18	frisbee 1/18	scissors 1/18	baseball glove 1/18	bicycle 1/20	bottle 1/20
dog 1/21	book 1/21	knife 1/22	remote 1/22	cup 1/22	mouse 1/22	toothbrush 1/24	chair 1/24
orange 1/25	fork 1/26	cell phone 1/30	boat 0/14	bowl 0/18	bus 0/16	skis 0/25	toilet 0/15
stop sign 0/21	tie 0/14	keyboard 0/15	sports ball 0/21	baseball bat 0/21	spoon 0/17	carrot 0/25	donut 0/19
couch 0/22	train 0/20	kite 0/17	clock 0/24	giraffe 0/19	pizza 0/18	zebra 0/19	truck 0/18
traffic light 0/17	motorcycle 0/11	skateboard 0/19	microwave 0/12	car 0/23	bed 0/15	laptop 0/24	elephant 0/20
broccoli 0/14	bear 0/17	banana 0/26	hot dog 0/14	surfboard 0/21	fire hydrant 0/25	airplane 0/18	tv 0/21

### Stereotyped with Men

tie 11/14	truck 14/18	motorcycle 8/11	baseball bat 15/21	baseball glove 11/18	sports ball 12/21	skateboard 10/19	fire hydrant 9/25
bear 6/17	snowboard 6/18	remote 7/22	car 6/23	suitcase 6/25	surfboard 5/21	sandwich 5/23	microwave 2/12
donut 3/19	person 4/26	boat 2/14	tennis racket 2/14	bicycle 2/20	bottle 2/20	tv 2/21	dog 2/21
couch 2/22	knife 2/22	sheep 2/24	backpack 2/26	bench 1/14	hot dog 1/14	vase 1/18	frisbee 1/18
pizza 1/18	cake 1/19	hair drier 1/20	teddy bear 1/20	train 1/20	elephant 1/20	toaster 1/21	stop sign 1/21
chair 1/24	toothbrush 1/24	umbrella 1/25	skis 1/25	horse 1/26	cell phone 1/30	potted plant 0/17	cup 0/22
scissors 0/18	traffic light 0/17	dining table 0/21	fork 0/26	book 0/21	orange 0/25	toilet 0/15	mouse 0/22
cat 0/19	refrigerator 0/23	spoon 0/17	bus 0/16	oven 0/23	zebra 0/19	carrot 0/25	giraffe 0/19
bed 0/15	laptop 0/24	sink 0/20	broccoli 0/14	banana 0/26	clock 0/24	bowl 0/18	wine glass 0/13
parking meter 0/19	airplane 0/18	kite 0/17	handbag 0/23	cow 0/15	keyboard 0/15	bird 0/22	apple 0/17

**Fig. 3 Study 1 Object Results.** Detailed participant responses for each of the 80 objects in COCO dataset. Fraction indicates number of participants asked about each object who marked it as stereotypically related to the gender group of women or men.

is already in contrast to a somewhat common assumption in ML fairness research that *any* difference between groups is an amplification of a stereotype [8].

Among 80 objects, 13 objects are marked as stereotypes by more than half of the participants (Figs. 3). Some examples of stereotypically gendered objects are **handbag** with women, **wine glass** with women, **tie** with men, and **truck** with men. Among the remaining objects, 18 objects (e.g., **keyboard**, **carrot**, **traffic light**) are marked by zero participants as stereotypes with any gender group, challenging prior assumptions on what is seen as a stereotype [88]. If an object was marked to be a stereotype, we also asked participants whether they believed it was harmful in the abstract. Complete results are in the Supplementary Material, but we use these initial findings to select experimental stimuli in subsequent studies. In Study 3a the stereotype-reinforcing condition includes women and **oven** (marked to be most harmful), women and **hair dryer** (marked to be least harmful), and the associated control conditions include women and **bowl**, women and **toothbrush**. In Study 3b we also include in the stereotype-reinforcing conditions of men and **baseball glove** (marked to be more harmful) and men and **necktie** (marked to be less harmful) with the control conditions of men and **bench** and men and **cup**.

### Study 2: Plurality of stereotypes and harms with image recognition objects

Next, we report qualitative analyses on open-ended responses from participants' annotations, where they explain why certain objects are seen as stereotypes and harmful or not. While prior work in gender stereotypes has often focused on social roles and traits [22, 32], our data provides insights as to how objects (e.g., **oven**, **hair dryer**) can also be associated with stereotypes. This is an important departure because it expands the scope of machine learning tasks for which stereotypes are relevant

183 beyond its current more narrow framing. Specifically, when a participant from Study 1 responds that  
184 an object is a stereotype, we follow up and ask: “Please describe in 1-2 sentences a) why you marked  
185 the above as a stereotype, and b) why you found it to be harmful or not.”

186 One of the authors coded the responses for why an object is a stereotype into roughly six cate-  
187 gories. The most prevalent reasons were: descriptive (45%), e.g., for **handbag** and women: “women are  
188 often seen wearing handbags and buying them”; occupation/role (22%), e.g., for **oven** and women:  
189 “women are stereotyped to always be in the kitchen cooking while the men go out and work”; trait  
190 (11%), e.g., for **chair** and men: “sometimes men would be seen as coming home and just being lazy  
191 and lounging in their chair.” The full analysis is in the Supplementary Material. It is interesting to  
192 note that an object’s association to a stereotype is frequently mediated by its connection to a role  
193 or trait, which are the more common sites of inquiry when it comes to stereotypes [22, 24, 27]. We  
194 also found that associations between a group and an object can exist through a number of paths.  
195 For example, explanations for stereotypical associations between cats and women include: “cat lady,”  
196 “women are called *kitten*,” “women like cats more than dogs,” “cats are a feminine animal,” and  
197 “women are called *cougars*.”

198 When asked why a stereotype was harmful or not, many respondents simply reiterated that the  
199 object was a stereotype. Dropping those responses, one of the authors coded the free responses  
200 of why a stereotype was marked to be harmful into seven categories, with the top three being:  
201 proscriptive (40%), e.g., for **dining table** and women: “it makes it looked down upon if a man  
202 cooks dinner”; prescriptive (26%), e.g., for **dining table** and women: “I think it puts women in a  
203 box that says they must prepare dinner”; negative trait (13%), e.g., for **handbag** and women: “it is  
204 harmful because it implies that women cares more about looks and their appearance.” The remaining  
205 response categories are in the Supplementary Material. There seems to be a disparity in responses  
206 based on the participant’s gender regarding whom they believe is harmed. When women specify  
207 which of the men group or women group are harmed, they say it is the women group 79% (95% CI  
208 [.67, .88]) of the time, while men say it is the women group only 67% (95% CI [.51, .80]) of the time.

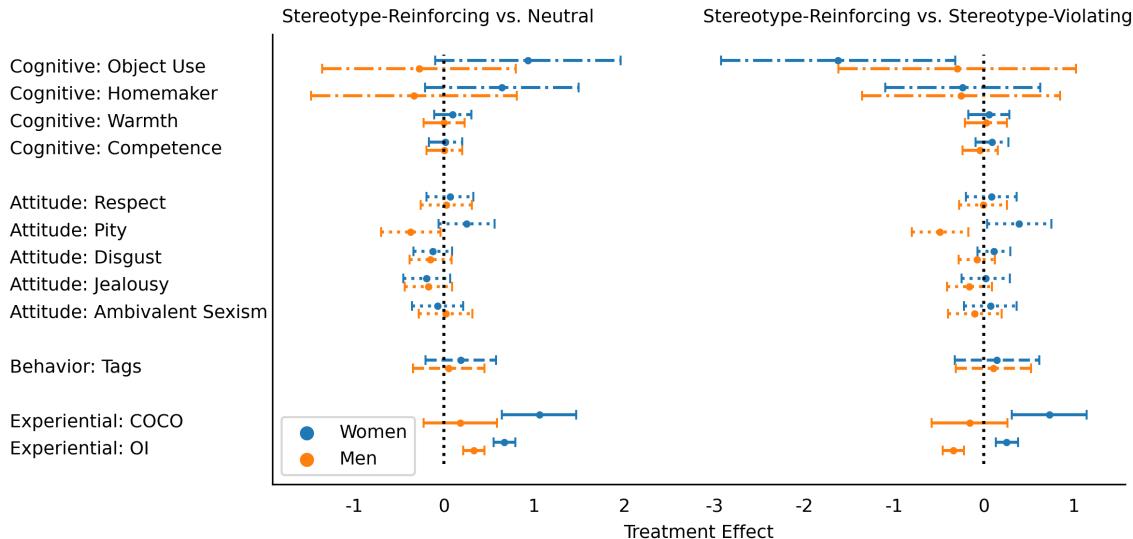
209 Building on Study 1’s finding that participants do not even all agree on whether an object is a  
210 stereotype or not (and if it is, whether it is harmful), this analysis further shows that even when  
211 participants are in agreement that an object is a stereotype, they are not necessarily in agreement  
212 about why. The same holds true for whether a stereotype is harmful. One potential implication of this  
213 is considering whether different reasonings should lead to different bias mitigation. For example, if  
214 the reason an object is a stereotype is descriptive, then mitigation should aim to change the cognitive  
215 representations of people. To change these descriptive statistics, while we can work to alter the model  
216 outputs, we should also work to change society, the burden of which falls on a much larger group than  
217 just machine learning practitioners, e.g., policymakers. On the other hand, if particular stereotypes

218 are deemed harmful because they are prescriptive and seem to restrict people from various avenues,  
219 we can consider ways to break free of gender norms.

220 **Study 3a: Stereotype-reinforcing errors show no pragmatic harm  
221 compared to both the stereotype-violating and neutral conditions**

222 To test pragmatic harm in stereotype-reinforcing errors, we conduct a between-subject survey experi-  
223 ment, using the stereotype-violating and neutral errors as control conditions. The cover story instructs  
224 participants to look at our synthesized search result page, imagining it is their personal phone photo  
225 album, and find a picture they had taken of someone they saw with a particular object. The search  
226 result page looks different for each randomized condition. We randomly assign participants to one  
227 of the three conditions ( $N = 600$ ): the stereotype-reinforcing condition exposes an image search  
228 result page with stereotype-reinforcing errors, e.g., false positive of **oven** on images of women; the  
229 stereotype-violating condition contains the same for stereotype-violating errors, e.g., false positive  
230 of **oven** on images of men; the stereotype-neutral condition contains neutral errors, e.g., false pos-  
231 itive of **bowl** on images of women. We then measure participants' cognitive beliefs, attitudes, and  
232 behaviors [27] to see if there are any changes because of such exposure (Methods). The behavioral  
233 measure is of particular interest, as we ask participants to undertake a realistic task they are liable to  
234 encounter by virtue of their jobs as online annotators: data labeling. We choose this measure because  
235 online participants are often the source of training labels in large-scale machine learning datasets. We  
236 ask participants to perform two common types of labeling on image data: tagging and captioning. If  
237 stereotype-reinforcing errors have an influence on participants' cognitive representations, attitudes,  
238 and tagging or captioning behaviors, we should expect to see a statistically significant difference  
239 between participants who are exposed to search results with **oven**-women and those who are exposed  
240 to search results with **oven**-men or **bowl**-women.

241 Contrary to what we had expected, after adjusting for multiple comparisons we do not find  
242 hypothesized statistically significant differences. We run an Ordinary-Least-Square (OLS) regression  
243 with the control condition coded as 0 and the experimental condition coded as 1, composite scores for  
244 beliefs, attitudes, and behaviors respectively as the dependent variables. Results are shown in Fig. 4  
245 with further details of the descriptive analysis of the captioning task in the Supplementary Material.



**Fig. 4 Study 3, 4 Results** The effect sizes and 95% confidence intervals are reported for 10 of our 11 measures of pragmatic harm (for the behavior measure of captioning, we provide a descriptive analysis), experiential harm on COCO, and experiential harm on our larger dataset of OpenImages. Deviations from zero indicate that exposure to the stereotype-reinforcing stimulus resulted in our measured harm compared to exposure to the control condition.

246 **Study 3b: Stereotype-reinforcing errors show statistically significant  
247 experiential harm compared to both the stereotype-violating and neutral  
248 conditions**

249 In terms of experiential harm, we design a within-subjects survey experiment ( $N = 100$ ). We oper-  
250 ationalize experiential harm by explicitly asking participants to rate how personally harmful they  
251 find different kinds of errors (which are stereotype-reinforcing, stereotype-violating, or neutral), on  
252 a scale from 0 (not at all) to 9 (extremely). This experience of error is analogous to situations where  
253 one reads in the news about the types of errors that artificial intelligence systems make [72], notices  
254 such a pattern of errors themselves, or is informed by a friend.

255 Comparing stereotype-reinforcing against neutral errors, an OLS regression shows participants  
256 rate stereotype-reinforcing errors to be more harmful than neutral ones ( $b = .62$ , 95% CI [.32,  
257 .91],  $p < .001$ ). However, when disaggregating by gender this effect is only present among women  
258 participants (women:  $b = 1.06$ , 95% CI [.64, 1.47],  $p < .001$ ; men:  $b = .18$ , 95% CI [-.23, .59],  
259  $p = .393$ ). When we use the stereotype-violating error as the control condition rather than the neutral  
260 error, we again find participants rate stereotype-reinforcing errors to be more harmful, though to a  
261 smaller degree, ( $b = .28$ , 95% CI [-.01, .58],  $p = .062$ ), with once again an effect only for women  
262 participants (women:  $b = .73$ , 95% CI [.31, 1.14],  $p = .001$ ; men:  $b = -.16$ , 95% CI [-.58, .26],  
263  $p = .453$ ). Results are in Fig. 4.

264 In short, while we find little immediate evidence of pragmatic harms, we do find the existence of  
265 experiential harms resulting from stereotype-reinforcing errors, compared to both stereotype-violating

266 and neutral errors. However, this pattern is present only among woman participants, and not men  
267 participants.

268 Prior work looking at a subset of what we call pragmatic harm has found very small effects in terms  
269 of cognitive belief changes about the representation of gendered occupations [43, 50]. Another line of  
270 work that finds a cognitive effect takes a different approach by studying occupations (e.g., peruker,  
271 lapidary) for which there are very few preconceived notions of stereotypes [79]. In our work, we focus  
272 on the activation of existing stereotypes, rather than the induction of novel stereotypes. Overall  
273 we find that the pragmatic harms are not measurable after exposure from repeated stereotypical  
274 errors in the current survey experiment, likely due to the fact that the effects of these harms are too  
275 diffuse and long-term, impacted by all of the facets of society we encounter in our lives [58]. Long-  
276 term observational studies are likely more well-suited to measure these kinds of impacts [28, 30, 39].  
277 However, we do find consistent evidence that members of the oppressed group report a significant  
278 experiential harm in the form of negative affect on stereotypical errors made on them, consistent  
279 with the feelings of inclusivity in gender-biased occupations [50].

#### 280 Study 4: Stereotype-violating errors can be perceived as harmful too

281 In this study, we first test the generalizability of the previous findings by using a popular dataset in  
282 object recognition tasks which is much larger: OpenImages [45]. We then explore a new hypothesis  
283 about gender presentation-aligned objects, e.g., clothing, to dive deeper into our findings. OpenImages  
284 has 600 objects, annotated with perceived binary genders of people present in the image if applica-  
285 ble [69]. Following the same procedure as in the COCO dataset with new online participants ( $N =$   
286 120), we find 249 of the 600 objects are marked as stereotypes by more than half of the participants,  
287 replicating the finding that not all objects are perceived as stereotypes (see more in Supplementary  
288 Materials). We then compile a list of 40 stereotypical objects (20 about men: e.g., **football**, **tool**; 20  
289 about women: e.g., **doll**, **lipstick**), and 20 neutral objects (e.g., **balloon**, **goldfish**) for this study.

290 To test whether participants experience more experiential harm when they are exposed to  
291 stereotype-reinforcing (e.g., **skirt** on women), stereotype-violating (e.g., **skirt** on men), and neu-  
292 tral (e.g., **toothbrush** on women) errors, we use a similar procedure as in Study 3b. Rather than  
293 asking simply about “personal harm” as we did in Study 3b, here we draw from the Positive and  
294 Negative Affect Schedule (PANAS; [14, 83]) and provide more details by asking about if they expe-  
295 rience harm such as feeling upset, irritated, ashamed, or distressed. We conduct a within-subjects  
296 study and ask participants ( $N = 300$ ) to report their subjective experiences on a Likert scale from 0  
297 to 9 for a variety of errors (see more in Methods). The analysis uses a mixed-effects regression with  
298 experimental conditions as the independent variable, a composite score of experiential harm as the

299 dependent variable, participants' gender as the covariate variable, and error terms clustered at the  
300 individual level.

301 Replicating Study 3b, we find that participants experience stereotype-reinforcing errors to be  
302 more harmful than neutral ones ( $b = .50$ , 95% CI [.42, .59],  $p < .001$ ). Again, this pattern is more  
303 pronounced among women participants ( $b = .67$ , 95% CI [.55, .79],  $p < .001$ ), with now a small  
304 effect among men participants ( $b = .33$ , 95% CI [.21, .45],  $p < .001$ ). Unlike Study 3b, we do not see  
305 differences in experiential harm between stereotype-reinforcing and stereotype-violating conditions  
306 ( $b = -.04$ , 95% CI [-.13, .05],  $p = .338$ ). The effect is canceled out by the opposite effects for women  
307 ( $b = .25$ , 95% CI [.13, .38],  $p < .001$ ) and men ( $b = -.34$ , 95% CI [-.46, -.22],  $p < .001$ ) participants.  
308 In other words, while women participants feel upset, irritated, ashamed, and distressed when they  
309 see stereotype-reinforcing errors (e.g., skirt on women), men participants feel that way when they  
310 see stereotype-violating errors (e.g., skirt on men). Results are in Fig. 4.

311 To better understand this finding, we conduct an exploratory analysis that digs deeper into the  
312 40 stereotypical objects to understand why stereotype-violating errors are sometimes perceived to  
313 be more experientially harmful than stereotype-reinforcing ones. According to the gender trouble  
314 framework, costume (i.e., body and appearance) and script (i.e., behavior, traits, and preferences)  
315 are two aspects of gender performance, and reactions to androgynous or conventionally contradictory  
316 components can differ depending on which of the two it manifests in [11, 34, 53, 54, 77]. We thus  
317 hypothesize that conventionally contradictory costume objects may evoke more negative reactions  
318 compared to conventionally contradictory script objects [67]. To test this hypothesis, we explore an  
319 additional independent variable we call “wearable.” We determined the value of this variable by  
320 manually marking 13 of the 40 stereotypical objects to be conventionally wearable by a person. These  
321 include objects like `football helmet` and `lipstick`, and exclude those like `truck` or `wine glass`.  
322 With this “wearable” distinction, we find that participants do rate stereotype-reinforcing errors to be  
323 more harmful than stereotype-violating ones ( $b = .23$  95% CI [.12, .34],  $p < .001$ ), though again this  
324 effect exists in women participants ( $b = .49$ , 95% CI [.34, .64],  $p < .001$ ) rather than men participants  
325 ( $b = -.03$ , 95% CI [-.18, .12],  $p = .726$ ). Notably, for the interaction effect of a “wearable” object with  
326 the condition type, we find that wearable stereotype-violating errors have higher experiential harm  
327 than wearable stereotype-reinforcing errors ( $b=.80$ , 95% CI [.62, .99],  $p < .001$ ), which is higher for  
328 men participants ( $b=.94$ , 95% CI [.67, 1.12],  $p < .001$ ) than women participants ( $b=.69$ , 95% CI  
329 [.43, .94],  $p < .001$ ). In other words, men participants tend to find it more harmful than women  
330 participants do when lipstick is misclassified on a man than on a woman.

331 Stereotype-violating errors seem to cause harm too, possibly through different mechanisms. In  
332 addition to this result being a consequence of backlash effects [68], we raise two more possible mech-  
333 anisms. First, it could be seen as an expression of precarious manhood; a concept that suggests

manhood is precarious and needs continuous social validation such that threats to traditional masculinity can provoke anxiety in men, thus resulting in higher reports of harm [78]. Second, these results may reflect elements of transphobia, which involves a negative reaction to the apparent incongruity between a person's perceived gender and a wearable gender presentation item [11, 54]. The divergent effect between men and women participants aligns with research indicating that transphobia is higher amongst cisgender men when judging transgender women due to the perceived threat to masculinity [49, 55]. This analysis pushes us to reevaluate how we should think about reducing experiential harm, as it may encompass intolerances we do not wish to support.

## Discussion

In summary, our studies have three key contributions: we investigate the kinds of associations people believe to be stereotypical; we distinguish between machine learning errors that are stereotype-reinforcing, stereotype-violating, or stereotype-neutral; we formulate harm as pragmatic or experiential to empirically study the effect of stereotypes. Overall, while stereotype-reinforcing errors do not lead to more pragmatic harm in the lab setting we use, we do find that stereotype-reinforcing errors are consistently found to be more experientially harmful. Such experiential harm is unequally distributed, impacting more participants who are women than who are men. Formulating concrete notions of harm as we have done has implications beyond just machine learning: legal documents like the European AI Act is beginning to incorporate notions of psychological harm but lacking definitions to ground regulation in [5, 62]. We also find stereotype-violating errors to be experientially harmful, especially when these errors pertain to wearable items associated with gender presentation. This effect is stronger for participants who identify as men compared to those who identify as women. This final point warrants an especially nuanced discussion, as we find ourselves qualifying a prior claim that we should take people's words at face value when they indicate something is personally harmful. To navigate this complexity, we turn to the notions of epistemic injustice [29] and standpoint epistemology [25, 59, 85]. If we interpret the negative reactions to misclassifications of stereotypically feminine clothing items on men as a manifestation of precarious manhood [78] or transphobia [11], then we should down weight these concerns in designing mitigation algorithms. Respecting people's experiential harms may not be as simple as accepting them at face value for direct measurement, but rather involves understanding which groups are likely to be harmed by what kinds of errors and why.

Our findings call for reconsidering fairness measurement in supervised machine learning tasks. This involves considering how we can leverage human-driven insights to inform model training and evaluation [10]. Traditionally, fairness evaluations tend to focus on stereotypes only in relation to occupations or traits. However our work expands this idea by showing that labels such as objects can also give rise to such harms. Additionally, most prior work has only considered the implications of

368 errors that reinforce stereotypes, which is relatively more intuitive to think of as harmful. However,  
369 both practically and normatively, it is important to understand the implications of stereotype-  
370 violating errors. Practically, strategies aimed at mitigating stereotype-reinforcing errors which act  
371 upon the target label will inevitably impact the occurrence of stereotype-violating errors as well.  
372 And normatively, there are also questions about whether stereotype-violating errors may even play  
373 a role in reducing stereotypical associations by counteracting them. This finding that not only are  
374 certain labels more liable to cause harm than others, but that it matters for *which* demographic  
375 group that label is misclassified, suggests that generic approaches like having a higher threshold for  
376 the classification of certain labels are insufficient. Instead, more nuanced fairness-through-awareness  
377 approaches [21] will need to be taken. While adopting simply a cost-sensitive framework [44] (e.g.,  
378 different costs are associated with false positives and false negatives) is a simplified interpretation of  
379 our findings, it could be a starting point as one grapples with the questions of whose levels of harms  
380 we would prioritize reducing in a bias mitigation framework.

381 Understanding whose levels of harms we should prioritize, and why, will come from stronger  
382 understandings of the psychological basis and reasoning of different harms. Our finding from Study  
383 2 that stereotypical associations between a single group and object can emerge from many paths  
384 (e.g., the many reasonings behind the association between cat and women), each with different  
385 normative valences, illustrates what an oversimplification it is to only label an association as “good” or  
386 “bad,” and the limitations of mitigations simply aiming to sever the associations deemed “bad.” This  
387 underscores the importance of work about diversity in annotators’ perspectives [16, 17, 20, 40, 57, 82],  
388 and how much complexity is reduced by the use of discrete labels. Qualitative follow-up questions  
389 supplemented our annotations, where a lack of consensus is not a weakness or artifact to be averaged  
390 out, but rather a point for deeper inquiry on how to prioritize differential experiences of harm. This  
391 also indicates that even if the growing power of large language models enables us to predict with  
392 higher accuracy which objects are stereotypes, we likely still may want to ensure these annotations  
393 come from people themselves [3, 38, 86], thus allowing room for positionality, explanation, and critical  
394 reflection.

395 Our findings are limited by the methodological choices we made: First, we focused on gender  
396 stereotypes as a case study. We do not know to what extent this finding generalizes to other groups  
397 such as race and age. Second, we recruited online participants who identify as men and women  
398 and who speak English without an extensive inclusion of non-binary participants or who come from  
399 a different cultural background. Given that stereotypes are culture-sensitive, and our work also  
400 shows that the harm perception is identity-sensitive, future work needs to study the interaction  
401 between participants’ identity, culture, and harm perceptions. Third, by setting a threshold of 50%  
402 for respondents indicating an object is a stereotype, we are in some senses privileging the majority

403 opinion which may further reify marked stereotypes to be those for the majority subset [31, 51].  
404 Fourth, the survey experiment does not capture harms beyond the two we measure (e.g., stereotype-  
405 threat [75, 76]), nor the longitudinal effects of machine learning effects. Future work needs to capture  
406 not only the plurality in harm of machine learning errors but also how its' effect emerges and endures  
407 over time.

408 Overall, our work offers a rigorous empirical study connecting machine learning outputs to con-  
409 crete harms by understanding the impact of stereotypical misclassifications. Rather than gesturing  
410 at harm as a justification for fairness measurement, we are very concrete in our analysis of the  
411 effects on people. Our finding that stereotype-reinforcing errors are experientially harmful for women  
412 underscores the importance for machine learning fairness interventions to be more rooted in social  
413 contexts, moving beyond objectives like just achieving equal prediction performance across groups.  
414 The diversity of responses we've presented, each influenced by participants' unique rationales, sug-  
415 gests the need for greater exploration of human psychological experiences in understanding how  
416 machine learning can cause harm.

## 417 Methods

### 418 Analysis

419 We use a mixture of qualitative and regression analyses to report our findings. For our within-subjects  
420 surveys, we regress with a mixed-effect model whose parameter estimations are adjusted by the group  
421 random effects for each individual. We report the coefficients from our regression analyses, which  
422 represent the effect size of that independent variable.

### 423 Participants

424 While men and women generally tend to hold the same gender stereotypes [23, 36, 48, 84], we still  
425 collect equal numbers of participants who identify as men and women, and use this variable as a  
426 covariate throughout. Due to limitations in the survey platform which only allow us to specify gender  
427 as "male" or "female," this formulation excludes people who identify as non-binary, which is a harmful  
428 limitation. Because we do not control for race in the recruitment of participants, our sample diverges  
429 from a nationally representative sample. For the gender stereotype scope of our current work, we  
430 find this to be an acceptable limitation, especially given that one defining feature of stereotypes is  
431 they are largely shared through a cultural consensus [41].

432 We did not use quality check questions in any of our surveys, because our pilot studies showed  
433 high quality responses. Instead, we used filters on Cloud Research to only recruit participants who  
434 have had at least 50 HITs approved, and have a HIT approval rate of 98%.

**Table 1** The time, pay, and reported races of the participants for each of our five studies. The full column names of races from left to right are: American Indian or Alaska Native, Asian, Black or African American, Hispanic or Latinx, Native Hawaiian or Other Pacific Islander, White, Multi-Racial / Other, and Prefer not to say.

Study	Time (min)	Pay (\$)	Gender	AI/AN	Asian	Black	H/L	NHOPI	White	MR/O	PNTS	Total
1 and 2	7	1.75	Women	0	3	5	0	0	25	6	1	40
			Men	1	4	2	2	0	30	1	0	40
3a	10	2.50	Women	1	11	32	8	0	229	19	0	300
			Men	0	19	35	10	1	211	22	2	300
3b	5	1.25	Women	0	4	7	3	1	35	5	0	50
			Men	0	4	2	3	1	35	5	0	50
4	4	1	Women	0	5	8	0	0	42	4	1	60
			Men	0	2	6	5	1	44	2	0	60
(Labeling)	4	5	Women	0	5	15	1	0	120	7	2	150
			Men	1	9	17	6	1	107	9	0	150

## 435 Studies 1, 4: Distinguishing Errors by Stereotype

436 When asking about which machine learning errors are stereotypes, we make sure to ask participants  
 437 about their perception of stereotypes held by Americans, rather than for their personal beliefs [18].

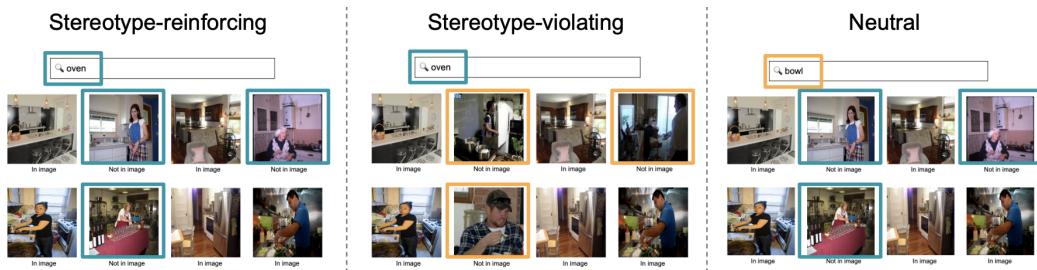
## 438 Study 3a: Measuring Pragmatic Harm

439 We conduct a between-subjects survey experiment on participants who are exposed to an image search  
 440 result page that contain one of three types of errors: stereotype-reinforcing, stereotype-violating, or  
 441 neutral (Fig. 5).<sup>2</sup> To have the participants engage with these results we ask them to describe it in  
 442 3-4 sentences. Next, we ask them the behavior questions, then re-expose them to the stimulus before  
 443 asking them the cognitive belief and attitude questions. We analyze changes in cognitive beliefs,  
 444 attitudes, and behaviors as pragmatic harms resulting from stereotype-reinforcing errors compared  
 445 to the two other conditions as controls. In this section when describing our method, we will use as  
 446 examples oven and women for the stereotype-reinforcing error, oven and men for the stereotype-  
 447 violating error, and bowl and women for the neutral one. Each question we ask is carefully grounded  
 448 in the social psychology literature.

449 The stimuli take the form of an image search result and are pictured in Fig. 5 with teal and orange  
 450 colored boxes around the component of the image that changes between conditions. The search bar  
 451 contains the search query, and then eight images that may or may not be correctly retrieved are  
 452 shown. Each of the eight images is annotated with either “In image” or “Not in image” to make it clear  
 453 to the participant which images are correct or not. The stereotype-reinforcing condition on the left  
 454 contains the search query of “oven” with five correctly identified ovens, and three false positive images  
 455 that all contain women. In other words, this classifier erroneously (and stereotypically) assumes there  
 456 are ovens in images of women. The stereotype-violating condition contains the same search query,

<sup>2</sup>The people pictured in our search results pages are predominantly White, which is the majority group in the dataset we employ.

457 but the mistakes are replaced with false positive images that all contain men. The neutral condition  
 458 contains all of the exact same images as the stereotype-reinforcing condition, with the only change  
 459 being that the search query is now “bowl” instead of “oven.” This is because the five correct images  
 460 were deliberately chosen to contain both bowls and ovens, which allows us to control for the variance  
 461 between the different search conditions. All false positive images were selected from the actual errors  
 462 of a Vision Transformer (ViT) model [19] trained on COCO so that they are as realistic as possible  
 463 to a computer vision model’s errors, and not completely egregious, e.g., a picture of a woman in a  
 464 sports field as a false positive for “oven” or “bowl.”



**Fig. 5 Study 3 Stimuli.** Our three different stimuli are shown for the conditions: stereotype-reinforcing, stereotype-violating, and neutral. They are all image search results containing minimal changes from each other, each of which indicates whether the search query is pictured in the image, i.e., if the image search retrieval was correct or not. The teal and orange squares indicate that the only difference between the stimuli, as all images which contain an oven also contain a bowl, and all which do not contain an oven also do not contain a bowl. This was a deliberate choice to control for all potential confounding factors from the images in the study.

465 For *cognitive beliefs*, we ask three sets of questions which span the spectrum of stereotype-specific  
 466 to more generically about gendered beliefs. Concretely, we ask about: estimations of who uses ovens  
 467 and bowls more between men and women; estimations of who tends to be the homemaker more  
 468 between men and women; and perceived levels of warmth and competence [27] of women. To assess  
 469 *attitude*, we ask two sets of questions. The first is about how participants feel about women in terms  
 470 of four emotional components that are believed to mediate interactions between cognitive beliefs and  
 471 behaviors: a) respect or admiration, b) pity or sympathy, c) disgust or sickening, and d) jealousy or  
 472 envy [15, 26, 70]. The second asks about sexist attitudes via a shortened scale focused on benevolent  
 473 sexism [32, 33, 65].<sup>3</sup> Finally, for *behavioral* measures, we ask participants to undertake a realistic task  
 474 they are liable to encounter which can cause harm: data labeling [52]. We chose this behavior measure  
 475 because online participants are often the source of training labels in large-scale machine learning  
 476 datasets. We ask participants to perform two common types of labeling on image data: tagging and  
 477 captioning (Fig. 6). In the tagging task, we ask participants to label the top three most relevant tags  
 478 in an image which contains both the stereotype object (e.g., **oven**) and neutral object (e.g., **bowl**). We  
 479 alter the perceived gender of the person to assess whether this changes what is tagged in the image.  
 For the captioning task we show two people, one who looks masculine and another feminine, and

<sup>3</sup>We ask questions from the Ambivalent Sexism Inventory [32] about benevolent sexism, as opposed to hostile sexism, because the latter is believed to suffer heavily from social desirability bias.

481 swap whether there is a bowl or oven present in the image. This is to understand if the annotators  
 482 will differently describe who is interacting with the object depending on whether it is stereotypically  
 483 associated with women or not. All images are generated and/or manipulated by DALL-E 2.



**Fig. 6** To measure behavioral tendencies, we ask participants to complete a realistic data annotation task on images which are created and manipulated by DALL-E2. The left pair is for the annotation of image tags, and the right pair is for image captions. Each participant is shown one image from each pair, and then we perform a between-subjects analysis to understand whether perceived gender expression affects the tags, and whether object shown influences how people of different perceived genders are described.

#### 484 **Dependent Variables**

485 For most of our measurements, we simply use the measure directly (e.g., the value for competence of  
 486 women) as the dependent variable to regress on. For the measurements that we do something more  
 487 complicated, we describe below.

488 **Behavior - Tags.** Each participant produces a set of three ordered tags associated with an image  
 489 of a feminine-presenting person and a set associated with a counterfactual image of a masculine-  
 490 presenting person. We convert this set of tags by scoring the presence of the object in question,  
 491 e.g., “hair dryer” (along with common misspellings such as “hair drier”) based on its position in  
 492 the ordered list of tags. When the word is present in the first spot it is given 3 points, second  
 493 spot 2 points, third spot 1 point, otherwise no points. The dependent variable is the score of both  
 494 the stereotypical and neutral object on the feminine-presenting person. This is intended to capture  
 495 whether the stereotype-reinforcing condition is able to increase the presence of the stereotype tag  
 496 more than just the priming effect captured by the neutral object.

497 **Behavior - Captions.** We offer some descriptive statistics about the captions in the Supplemen-  
 498 tary Material. This analysis was mostly exploratory, and we do not find any statistically significant  
 499 differences. We first ran Study 3a looking at pragmatic harms on the stereotype of women and oven  
 500 (with bowl as the control). In this iteration, we asked that respondents please describe each person in  
 501 the image in separate sentences. However, there was too much noise in how respondents interpreted  
 502 this set of instructions, such that the data became hard to interpret. Thus, in our second iteration of  
 503 this study using the stereotype of women and hair dryer (with toothbrush as the control), we have  
 504 two separate text entry boxes to caption each person in the image. We only present the results of  
 505 this iteration in the table, as we were unable to parse anything differentiating in the first iteration.

506       **Cognitive - Object Use.** In this measurement, we have a value from -10 (mostly men) to  
507       10 (mostly women) for both the stereotypical and neutral object. The dependent variable is the  
508       summation of both values. Again, this is intended to capture whether the stereotype-reinforcing  
509       condition is able to change the value of its associated object more than the control condition is able  
510       to.

511       **Study 3b, 4: Measuring Experimental Harm**

512       In Study 3b, in addition to personal discomfort, we also ask about societal harm. This way, even if  
513       the participant does not personally feel harmed, they may feel it on behalf of the stereotyped group.  
514       However, we find that participants' responses to both personal and societal harm are extremely  
515       correlated, and leave the results for the latter in the Supplementary Material.

516       **References**

- 517       [1] Abbasi M, Friedler SA, Scheidegger C, et al (2019) Fairness in representation: quantifying  
518       stereotyping as a representational harm. Siam International Conference on Data Mining
- 519       [2] Allport GW, Clark K, Pettigrew T (1954) The nature of prejudice
- 520       [3] Argyle LP, Busby EC, Fulda N, et al (2023) Out of one, many: Using language models to simulate  
521       human samples. Political Analysis
- 522       [4] Barlas P, Kyriakou K, Guest O, et al (2021) To "see" is to stereotype: Image tagging algorithms,  
523       gender recognition, and the accuracy-fairness trade-off. Proceedings of the ACM on Human-  
524       Computer Interaction (CSCW)
- 525       [5] Bayefsky R (2016) Psychological harm and constitutional standing. Brooklyn Law Review
- 526       [6] Bhaskaran J, Bhallamudi I (2019) Good secretaries, bad truck drivers? occupational gender  
527       stereotypes in sentiment analysis. Proceedings of the First Workshop on Gender Bias in Natural  
528       Language Processing
- 529       [7] Bianchi F, Kalluri P, Durmus E, et al (2023) Easily accessible text-to-image generation ampli-  
530       fies demographic stereotypes at large scale. ACM Conference on Fairness, Accountability, and  
531       Transparency (FAccT)
- 532       [8] Blodgett SL, Lopez G, Olteanu A, et al (2021) Stereotyping norwegian salmon: An inventory  
533       of pitfalls in fairness benchmark datasets. Proceedings of the 59th Annual Meeting of the Asso-  
534       ciation for Computational Linguistics and the 11th International Joint Conference on Natural

535 Language Processing

- 536 [9] Bolukbasi T, Chang KW, Zou J, et al (2016) Man is to computer programmer as woman is to  
537 homemaker? debiasing word embeddings. Conference on Neural Information Processing Systems  
538 (NeurIPS)
- 539 [10] Boykin CM, Dasch ST, Jr. VR, et al (2021) Opportunities for a more interdisciplinary approach  
540 to measuring perceptions of fairness in machine learning. Equity and Access in Algorithms,  
541 Mechanisms, and Optimization (EAAMO)
- 542 [11] Butler J (1990) Gender trouble: Feminism and the subversion of identity. Routledge
- 543 [12] Caliskan A, Bryson JJ, Narayanan A (2017) Semantics derived automatically from language  
544 corpora contain human-like biases. Science
- 545 [13] Cao Y, Sotnikova A, III HD, et al (2022) Theory-grounded measurement of u.s. social stereotypes  
546 in english language models. Conference of the North American Chapter of the Association for  
547 Computational Linguistics: Human Language Technologies
- 548 [14] Crawford JR, Henry JD (2004) The positive and negative affect schedule (panas): construct  
549 validity, measurement properties and normative data in a large non-clinical sample. British  
550 Journal of Clinical Psychology
- 551 [15] Cuddy AJC, Fiske ST, Glick P (2007) The BIAS map: behaviors from intergroup affect and  
552 stereotypes. Journal of Personality and Social Psychology 92
- 553 [16] Davani AM, Díaz M, Prabhakaran V (2022) Dealing with disagreements: Looking beyond the  
554 majority vote in subjective annotations. Transactions of the Association for Computational  
555 Linguistics
- 556 [17] Denton E, Díaz M, Kivlichan I, et al (2021) Whose ground truth? accounting for individual and  
557 collective identities underlying dataset annotation. NeurIPS 2021 Workshop on Data-Centric AI
- 558 [18] Devine PG, Elliot AJ (1995) Are racial stereotypes really fading? the princeton trilogy revisited.  
559 Personality and Social Psychology Bulletin 21
- 560 [19] Dosovitskiy A, Beyer L, Kolesnikov A, et al (2021) An image is worth 16x16 words: Transformers  
561 for image recognition at scale. International Conference on Learning Representations (ICLR)
- 562 [20] Dumitracă A, Aroyo L, Welty C (2018) Capturing ambiguity in crowdsourcing frame  
563 disambiguation. AAAI Conference on Human Computation and Crowdsourcing (HCOMP)

- 564 [21] Dwork C, Hardt M, Pitassi T, et al (2012) Fairness through awareness. Proceedings of the 3rd  
565 Innovations in Theoretical Computer Science Conference
- 566 [22] Eagly AH (1987) Sex differences in social behavior: A social-role interpretation. Lawrence  
567 Erlbaum Associates, Inc
- 568 [23] Eagly AH, Nater C, Miller DI, et al (2020) Gender stereotypes have changed: A cross-temporal  
569 meta-analysis of u.s. public opinion polls from 1946 to 2018. American Psychologist
- 570 [24] Ellemers N, et al (2018) Gender stereotypes. Annual review of psychology 69:275–298
- 571 [25] Fatima S (2020) I know what happened to me: The epistemic harms of microaggression.  
572 Microaggressions and Philosophy
- 573 [26] Fiske ST, Cuddy AJC, Glick P (2002) Emotions up and down: Intergroup emotions result from  
574 status and competition. Prejudice to Intergroup Emotions: Differentiated Reactions to Social  
575 Groups
- 576 [27] Fiske ST, Cuddy AJC, Glick P, et al (2002) A model of (often mixed) stereotype content:  
577 Competence and warmth respectively follow from perceived status and competition. Journal of  
578 Personality and Social Psychology 82
- 579 [28] Ford TE (1997) Effects of stereotypical television portrayals of african-americans on person  
580 perception. Social Psychology Quarterly
- 581 [29] Fricker M (2009) Epistemic injustice: Power and the ethics of knowing. Oxford University Press
- 582 [30] Fujioka Y (1999) Television portrayals and african-american stereotypes: Examination of  
583 television effects when direct contact is lacking. Journalism and Mass Communication Quarterly
- 584 [31] Ghavami N, Peplau LA (2012) An intersectional analysis of gender and ethnic stereotypes:  
585 Testing three hypotheses. Psychology of Women Quarterly 37
- 586 [32] Glick P, Fiske ST (1996) The ambivalent sexism inventory: Differentiating hostile and benevolent  
587 sexism. Journal of Personality and Social Psychology 70
- 588 [33] Glick P, Whitehead J (2010) Hostility toward men and the perceived stability of male dominance.  
589 Social Psychology 41
- 590 [34] Goffman E (1959) The presentation of self in everyday life. Doubleday

- 591 [35] Hamilton DL, Sherman JW (2014) Stereotypes. In: Handbook of social cognition. Psychology  
592 Press, p 17–84
- 593 [36] Hentschel T, Heilman ME, Peus CV (2019) The multiple dimensions of gender stereotypes:  
594 A current look at men's and women's characterizations of others and themselves. Frontiers in  
595 Psychology
- 596 [37] Hilton JL, Von Hippel W (1996) Stereotypes. Annual review of psychology 47(1):237–271
- 597 [38] Hämläinen P, Tavast M, Kunnari A (2023) Evaluating large language models in generating  
598 synthetic hci research data: a case study. Conference on Human Factors in Computing Systems  
599 (CHI)
- 600 [39] Jennings-Walstedt J, Geis FL, Brown V (1980) Influence of television commercials on women's  
601 self-confidence and independent judgment. Journal of Personality and Social Psychology
- 602 [40] Kairam S, Heer J (2016) Parting crowds: Characterizing divergent interpretations in crowd-  
603 sourced annotation tasks. ACM Conference On Computer-Supported Cooperative Work And  
604 Social Computing (CSCW)
- 605 [41] Katz D, Braly K (1933) Racial stereotypes of one hundred college students. The Journal of  
606 Abnormal and Social Psychology 28
- 607 [42] Katzman J, Wang A, Scheuerman M, et al (2023) Taxonomizing and measuring representational  
608 harms: A look at image tagging. AAAI Conference on Artificial Intelligence
- 609 [43] Kay M, Matuszek C, Munson SA (2015) Unequal representation and gender stereotypes in image  
610 search results for occupations. Conference on Human Factors in Computing Systems (CHI)
- 611 [44] Kukar M, Kononenko I (1998) Cost-sensitive learning with neural networks. European Confer-  
612 ence on Artificial Intelligence
- 613 [45] Kuznetsova A, Rom H, Alldrin N, et al (2020) The open images dataset v4: Unified image  
614 classification, object detection, and visual relationship detection at scale. International Journal  
615 of Computer Vision (IJCV)
- 616 [46] Lin TY, Maire M, Belongie S, et al (2014) Microsoft COCO: Common objects in context.  
617 European Conference on Computer Vision (ECCV)
- 618 [47] Lippmann W (1922) Public opinion.

- 619 [48] López-Sáez M, Lisbona A (2014) Descriptive and prescriptive features of gender stereotyping.  
620 relationships among its components. International Journal of Social Psychology 24
- 621 [49] Makwana AP, Dhont K, keersmaecker JD, et al (2018) The motivated cognitive basis of  
622 transphobia: The roles of right-wing ideologies and gender role beliefs. Sex Roles 79
- 623 [50] Metaxa D, Gan MA, Goh S, et al (2021) An image of society: Gender and racial representation  
624 and impact in image search results for occupations. ACM Conference on Human-Computer  
625 Interaction (CSCW)
- 626 [51] Mill JS (1859) On liberty. Longman, Roberts, Green Co
- 627 [52] van Miltenburg E (2016) Stereotyping and bias in the flickr30k dataset. Proceedings of the  
628 Workshop on Multimodal Corpora
- 629 [53] Morgenroth T, Ryan MK (2018) Gender trouble in social psychology: How can butler's work  
630 inform experimental social psychologists' conceptualization of gender? Frontiers in Psychology
- 631 [54] Morgenroth T, Ryan MK (2020) The effects of gender trouble: An integrative theoretical frame-  
632 work of the perpetuation and disruption of the gender/sex binary. Perspectives on Psychological  
633 Science 16
- 634 [55] Nagoshi CT, Cloud JR, Lindley LM, et al (2019) A test of the three-component model of gender-  
635 based prejudices: Homophobia and transphobia are affected by raters' and targets' assigned sex  
636 at birth. Sex Roles 80
- 637 [56] Nicolas G, Bai X, Fiske ST (2022) A spontaneous stereotype content model: Taxonomy,  
638 properties, and prediction. Journal of personality and social psychology
- 639 [57] Noble JA (2012) Minority voices of crowdsourcing: why we should pay attention to every mem-  
640 ber of the crowd. ACM Conference On Computer-Supported Cooperative Work And Social  
641 Computing (CSCW)
- 642 [58] Noble SU (2018) Algorithms of oppression: How search engines reinforce racism. NYU Press
- 643 [59] O'Dowd O (2018) Microaggressions: A kantian account. Ethical Theory and Moral Practice 21
- 644 [60] Otterbacher J, Bates J, Clough P (2017) Competent men and warm women: Gender stereotypes  
645 and backlash in image search results. Conference on Human Factors in Computing Systems  
646 (CHI)

- 647 [61] Otterbacher J, Checco A, Demartini G, et al (2018) Investigating user perception of gender bias  
648 in image search: The role of sexism. ACM SIGIR Conference on Research and Development in  
649 Information Retrieval (SIGIR)
- 650 [62] Palka P (2023) Ai, consumers & psychological harm. Cambridge University Press
- 651 [63] Prentice DA, Carranza E (2002) What women and men should be, shouldn't be, are allowed to  
652 be, and don't have to be: The contents of prescriptive gender stereotypes. Psychology of women  
653 quarterly 26(4):269–281
- 654 [64] Rini R (2020) The ethics of microaggression. Routledge Taylr & Francis Group
- 655 [65] Rollero C, Glick P, Tartaglia S (2014) Psychometric properties of short versions of the ambiva-  
656 lent sexism inventory and ambivalence toward men inventory. TPM-Testing, Psychometrics,  
657 Methodology in Applied Psychology 21
- 658 [66] Rudman LA, Glick P (2001) Prescriptive gender stereotypes and backlash toward agentic women.  
659 Journal of social issues 57(4):743–762
- 660 [67] Rudman LA, Moss-Racusin CA, Glick P, et al (2012) Reactions to vanguards: Advances in  
661 backlash theory. Advances in experimental social psychology 45
- 662 [68] Rudman LA, Moss-Racusin CA, Phelan JE, et al (2012) Status incongruity and backlash  
663 effects: Defending the gender hierarchy motivates prejudice against female leaders. Journal of  
664 Experimental Social Psychology 48
- 665 [69] Schumann C, Ricco S, Prabhu U, et al (2021) A step toward more inclusive people annotations  
666 for fairness. ACM Conference on Artificial Intelligence, Ethics, and Society (AIES)
- 667 [70] Seger CR, Banerji I, Park SH, et al (2017) Specific emotions as mediators of the effect of inter-  
668 group contact on prejudice: findings across multiple participant and target groups. Cognition  
669 and Emotion 31
- 670 [71] Selvaraju RR, Cogswell M, Das A, et al (2017) Grad-cam: Visual explanations from deep net-  
671 works via gradient-based localization. In: Proceedings of the IEEE international conference on  
672 computer vision, pp 618–626
- 673 [72] Simonite T (2018) When It Comes to Gorillas, Google Photos Remains Blind. Wired, January
- 674 [73] Sotnikova A, Cao YT, III HD, et al (2021) Analyzing stereotypes in generative text inference  
675 tasks. Findings of the Association for Computational Linguistics: ACL-IJCNLP

- 676 [74] Spencer SJ, Steele CM, Quinn DM (1999) Stereotype threat and women's math performance.  
677 Journal of Experimental Social Psychology 35
- 678 [75] Spencer SJ, Logel C, Davies PG (2015) Stereotype threat. Annual Review of Psychology 67
- 679 [76] Steele CM, Aronson J (1995) Stereotype threat and the intellectual test performance of african  
680 americans. Journal of Personality and Social Psychology 69
- 681 [77] Stern C, Rule NO (2017) Physical androgyny and categorization difficulty shape political  
682 conservatives' attitudes toward transgender people. Social Psychological and Personality Science
- 683 [78] Vandello JA, Bosson JK, Cohen D, et al (2008) Precarious manhood. Journal of Personality and  
684 Social Psychology
- 685 [79] Vlasceanu M, Amodio DM (2022) Propagation of societal gender inequality by internet search  
686 algorithms. Proceedings of the National Academy of Sciences of the United States of America  
687 (PNAS) 119
- 688 [80] Wan Y, Pu G, Sun J, et al (2023) "kelly is a warm person, joseph is a role model": Gender  
689 biases in llm-generated reference letters. Conference on Empirical Methods in Natural Language  
690 Processing (EMNLP Findings)
- 691 [81] Wang A, Liu A, Zhang R, et al (2022) REVISE: A tool for measuring and mitigating bias in  
692 visual datasets. International Journal of Computer Vision (IJCV)
- 693 [82] Waseem Z (2016) Are you a racist or am i seeing things? annotator influence on hate speech  
694 detection on twitter. Proceedings of the First Workshop on NLP and Computational Social  
695 Science
- 696 [83] Watson D, Clark LA, Tellegen A (1988) Development and validation of brief measures of positive  
697 and negative affect: the panas scales. Journal of Personality and Social Psychology 54
- 698 [84] Williams JE, Best DL (1977) Sex stereotypes and trait favorability on the adjective check list.  
699 Educational and Psychological Measurement 37
- 700 [85] Wyylie A (2003) Why standpoint matters. Science and Other Cultures: Issues in Philosophies of  
701 Science and Technology
- 702 [86] Yaghini M, Krause A, Heidari H (2021) A human-in-the-loop framework to construct  
703 context-aware mathematical notions of outcome fairness. AAAI/ACM Conference on Artificial  
704 Intelligence, Ethics, and Society

- 705 [87] Zhao D, Wang A, Russakovsky O (2021) Understanding and evaluating racial biases in image  
706 captioning. International Conference on Computer Vision (ICCV)
- 707 [88] Zhao J, Wang T, Yatskar M, et al (2017) Men also like shopping: Reducing gender bias ampli-  
708 fication using corpus-level constraints. Proceedings of the Conference on Empirical Methods in  
709 Natural Language Processing (EMNLP)