

1 Measuring machine learning harms from  
2 stereotypes requires understanding who is being  
3 harmed by which errors in what ways

4 Angelina Wang\*, Xuechunzi Bai†, Solon Barocas‡§, Su Lin Blodgett§

5 **Abstract**

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

---

\*Department of Computer Science, Princeton University

†Department of Psychology, Princeton University

‡Department of Information Science, Cornell University

§Microsoft Research

## **21** Introduction

**22** Machine learning systems are increasingly playing a central role in everyday life.  
**23** They revolutionize how humans communicate information, generate ideas in arts  
**24** and science, and make decisions in hiring, education, medical diagnosis, and beyond.  
**25** Accompanying this rapid proliferation is an increasing attention on the potential harm  
**26** these systems may cause, and movements toward developing them in a fair, ethical, and  
**27** inclusive way [5, 6]. The first step in mitigating harm is to be precise about who exactly  
**28** experiences what kind of harm and why [5, 10, 11, 21, 43, 47, 87]. Despite the impor-  
**29** tance of human psychological experiences in thinking about harm, we know relatively  
**30** little about how humans react, evaluate, and reason about machine learning classi-  
**31** fication outputs and their potential harms during everyday interactions with these  
**32** systems (with exceptions in decision-making systems such as credit assignment, job  
**33** allocation, etc. [19]). Drawing on psychological theories of social stereotypes, this paper  
**34** presents empirical evidence that underscores the complexity of harmful experiences  
**35** when machine learning models inevitably make mistakes.

**36** Stereotypes are frequently invoked to explain why some machine learning classifica-  
**37** tions are more harmful than others [1, 4, 9, 88]; however, researchers rarely investigate  
**38** the concrete connection between stereotypes and harm. Without fully interrogating  
**39** this relationship, false assumptions slip through the cracks, both in terms of where  
**40** the stereotypes come from and how they relate to harm. Some studies overly rely  
**41** on researchers' own worldviews. For example, researchers identified an object recog-  
**42** nition model as harmful because it amplifies the degree to which labels for kitchen  
**43** items like "knife, fork, and spoon" are incorrectly assigned to photos featuring women,  
**44** and labels for technology-related items like "keyboard and mouse" are incorrectly  
**45** assigned to photos featuring men [96]. This rationale not only extends to incorrect  
**46** assignments of even more neutral objects like tables, but these remarked-upon errors  
**47** themselves may not even be genuinely harmful: While computer scientists working in

48 the male-dominated technology space tend to find technology-related items like key-  
49 boards highly male-stereotyped, broader audiences do not actually share this idea, as  
50 we find in our study. Other studies heavily rely on occupation data from the Ameri-  
51 can Bureau of Labor Statistics, e.g., WinoBias [97]. While more grounded than relying  
52 on researchers' assumptions, this approach has one large limitation (beyond only rep-  
53 resenting occupation data in America). It primarily captures descriptive stereotypes,  
54 meaning actual overrepresentations of groups in an occupation. However, it misses  
55 prescriptive stereotypes, which entail beliefs about what occupations people of dif-  
56 ferent groups should be in. Prescriptive and descriptive stereotypes often diverge in  
57 practice [14, 58]. While some prior work has annotated or drawn from more grounded  
58 stereotypes [12, 16, 17, 80], researchers commonly assign equal levels of harm to every  
59 stereotype, or even to every misclassification. Our study does the necessary work of  
60 explicitly connecting stereotypes to harm [8]. This connection is pivotal because, with-  
61 out a deeper understanding grounded in psychological experiences, bias mitigation  
62 may inadvertently *increase* the number of harmful errors in a well-intentioned but  
63 ultimately misguided attempt to reduce other kinds of errors.

64 Our first conceptual contribution is in differentiating between machine learning  
65 errors which are stereotype-reinforcing, stereotype-violating, or neutral. We posit that  
66 errors which reinforce social stereotypes can be more harmful than errors that do not.  
67 While stereotypes are cognitive beliefs in people's minds, they can have an influence  
68 on attitudes (i.e., prejudice) and behaviors (i.e., discrimination) [2, 42, 45, 51, 56]. For  
69 example, people may have *cognitive beliefs* that women are more warm but less compe-  
70 tent, and thus *emotionally* express protective attitudes and pity for women [38]. People  
71 then *behave* in ways that maintain women's warmth and discount their competence,  
72 such as being less likely to promote women to leadership positions [30, 33]. There-  
73 fore, stereotypes of certain social groups can prompt shifts in attitudes and behaviors  
74 that can ultimately harm the stereotyped group. Although researchers implicitly

75 acknowledge the mediating role of stereotypes in machine learning harm, we draw  
 76 on the psychological framework of stereotypes to provide a concrete and systematic  
 77 assessment (Fig. 1).

Prompt	What Errors	Which Harms	Why	By Whom
	<p>Stereotype-reinforcing</p>  <p>Stereotype-violating</p>  <p>Stereotype-neutral</p> 	<p>Pragmatic</p>  <p>Annotated image caption: "The woman is preparing the meal. The man is standing and watching"</p> <p>Experiential</p>  <p>"On a scale from 1 to 7, how irritated do you feel?"</p>	<p>"Women are considered the ones to do the cooking the most. I am a female so obviously I would think this is a bit offensive."</p> <p>"Because its mostly women that bake. Its not harmful because its true."</p> <p>"Some people believe women belong in the kitchen. It can be harmful because it can make men feel like they aren't allowed to cook."</p>	     

**Fig. 1 Summary of our work.** We distinguish false positive errors to be those that are: stereotype-reinforcing, -violating, or -neutral. This label comes from human annotators, and depends on what gender group the prediction target (e.g., *oven*) is marked to be associated with. Then, we measure two types of harms: *pragmatic* which are changes about a stereotyped group in cognitive beliefs, attitude, or behavior and *experiential* which are personal self-reports of harm. We find that participants find stereotypes to be harmful for a number of contrasting reasons, and also that this harm is different between different gender groups. While we rely on the social categories of men and women in this work due to the prevalence of stereotypes about both groups, we acknowledge this as a limitation and do not endorse the binarization of gender.

78 Our second conceptual contribution is in defining harm. Prior work in the machine  
 79 learning fairness space has rarely been concrete about what harm actually means [10].  
 80 We distinguish between two types of harm as the most likely to result from stereotype-  
 81 reinforcing errors: *pragmatic harms* involve measurable changes in someone's cognitive  
 82 beliefs, attitudes, or behaviors toward the group being stereotyped, while *experiential*  
 83 *harms* involve self-reports of negative affect (Fig. 1). Pragmatic harms are motivated  
 84 by prior research showing that, for example, people express envy and passively harm  
 85 groups that are stereotyped as competent but untrustworthy (e.g., lawyers), or express

86 contempt and actively attack groups that are stereotyped as incompetent and unreliable  
87 (e.g., homeless [20]). In the domain of machine learning, prior work has considered  
88 components of this framework and found that exposure to gender-biased image search  
89 results can lead to more biased estimations of gender representation of that occupation  
90 and decreased sense of belonging [52, 60]. To examine if people experience pragmatic  
91 harms, we measure cognitive, emotional, and behavioral changes between people who  
92 experience machine learning outputs containing stereotype-reinforcing errors com-  
93 pared to those who experience stereotype-neutral or stereotype-violating errors. We  
94 hypothesize that the former will result in pragmatic harm.

95 In contrast to pragmatic harms which focus on external impositions towards a  
96 stereotyped group, experiential harms consider the subjective feeling of harm directly  
97 experienced by the stereotyped group member [91]. Subjective experiences of emo-  
98 tion have long been discounted as a legitimate source of knowledge, especially when  
99 expressed by social groups like women who are associated with emotion [48]. Addition-  
100 ally, these feelings can influence one’s own behaviors. For example, when women are  
101 given a math exam and told that the exam is diagnostic of their own intellectual abili-  
102 ties, stereotypes of women as less capable of math negatively impact their performance  
103 on the exam [81]. In conceptualizing the experiential harm of machine learning errors  
104 which may seem individually minor, we draw a parallel to the concept of microaggres-  
105 sions. Microaggressions are “small act[s] of insult or indignity, relating to a person’s  
106 membership in a socially oppressed group, which seems minor on its own but plays  
107 a part in significant systemic harm” [72]. Just like how a machine learning model’s  
108 classification error (e.g., of an oven on an image of a woman) may seem small on  
109 its own, and are “easily interpretable as inadvertent errors rather than as malevolent  
110 actions,” their negative effects on the target are real and should not be neglected [72].  
111 Important in this measure of harm is who the respondent is. Standpoint epistemology  
112 emphasizes the importance of the experiences of the individuals being stereotyped,

113 and the difficulty in establishing the legitimacy of this as a measure of harm thus far  
114 can be at least partially attributed to testimonial injustice [31, 35, 67, 93]. Hence, we  
115 hypothesize greater reports of experiential harm on stereotype-reinforcing errors for  
116 the stereotyped group.

117 Our third and more nuanced conceptual contribution is a call for an increased  
118 appreciation of the diversity of reasons that can lead to the same measured harm  
119 (Fig. 1). While prior work often uses human judgments, they do not always incor-  
120 porate the potential divergent reasons that individuals have which may lead to the  
121 same annotation. In our work, we find complexity in what people find to be stereo-  
122 typical and harmful. This complements prior work studying how human annotators  
123 bring different subjective experiences in their labeling of data [22, 23, 65, 89], intro-  
124 ducing strong associations between annotator identity and annotations [76]. In more  
125 subjective tasks such as labeling text as toxic or not, annotations are often divergent.  
126 While taking the majority vote is a common way to reconcile differences in annota-  
127 tions, there is a growing consensus to use a more representative system [26, 41, 50, 71].  
128 Understanding harm faces similar nuances when incorporating divergent perspectives  
129 on the same issue. However, simply incorporating representative annotations is not  
130 enough; it misses the personalized reasonings behind each response. For example, in  
131 a heterosexual gender normative society, some people think that men wearing skirts  
132 is harmful and should be regulated [15, 74]. Careless incorporation of this perspective  
133 could lead to a system that treats misclassifications of skirts on men as harmful errors  
134 on par with those which reinforce sexist stereotypes. If not carefully examined, naive  
135 additions of more voices may even exacerbate bias.

136 We conduct human studies to concretely measure the presence of harms when  
137 people experience machine learning errors. As a concrete application to ground our  
138 human studies in, we consider gender stereotypes in the popular machine learning  
139 task of object recognition as used in photo search engines. We use the COCO [55]

and OpenImages [54] datasets (Fig. 2), and design survey experiments with online American participants from Amazon Mechanical Turk through Cloud Research [57] (Methods). We first ask participants whether objects in COCO and OpenImages are stereotypically associated with different gender groups in order to distinguish which kinds of errors are stereotype-reinforcing, stereotype-violating, or neutral (Results - Study 1). Using those unveiled distinctions, we then expose participants to synthesized search result pages which contain different kinds of errors. We find little immediate evidence of pragmatic harms, but sizable evidence that stereotype-reinforcing errors are experientially harmful – a finding that is more pronounced among participants who identify as women compared to those who identify as men (Results - Study 2). In addition to stereotype-reinforcing errors (e.g., `oven` on women), we explore stereotype-violating errors (e.g., `oven` on men), which have received scarce attention in the machine learning fairness literature. We find that while the stereotyped group (e.g., women) generally finds it more harmful for the error to reinforce rather than violate stereotypes, this is not true when it comes to clothing-related items typically associated with women (e.g., `cosmetics`, `necklaces`) being misclassified on men. Here, we see a backlash towards violations of the norms around gender presentation where men tend to find these misclassifications of, e.g., `cosmetics`, more harmful on men rather than women, calling into question the idea that it is always normatively desirable to reduce errors perceived as more harmful due to their relationship to stereotypes (Results - Study 3). Finally, our qualitative analysis reveals the plurality of why participants think certain objects are stereotypes, and why those stereotypes may be harmful or not (Results - Study 4).

All studies are approved by our institution IRB, protocol number 14738. Studies 1 (<https://osf.io/cpyn4>), 2 (<https://osf.io/m9akd>, <https://osf.io/v2w4m>), and part of Study 3 (<https://osf.io/xpv5j>) are pre-registered on OSF, while Study 4 is more exploratory. By bringing greater clarity to different types of machine learning errors



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

167 based on their relationship to a stereotype and embracing the rich psychological expe-  
 168 riences behind them, we urge researchers and practitioners to more carefully consider  
 169 different kinds of classification errors, potential harms, and the relevant relationships  
 170 between them. We believe that identifying psychological experiences with machine  
 171 learning outputs is critical to understanding the potential harm of a system, and  
 172 in turn, mitigating it. Without doing so, we may inadvertently prioritize an overall  
 173 decrease of errors at the expense of increasing the number of harmful errors.

## 174 Results

175 We explore a popular task in machine learning known as object recognition (i.e.,  
 176 classifying the objects present in an image). To make it concrete for our human studies,  
 177 we use it in the context of a smart phone's photo search engine, and examine gender  
 178 stereotypes. Specifically, we consider one type of machine learning error called a false  
 179 positive: when an object is predicted to be present in an image when it is in fact

180 not there. This causes the image with a false positive to be wrongly surfaced on an  
181 image search results page.<sup>1</sup> In our work, we are only concerned with the effect of  
182 the misclassification, and not why the model may have made the mistake, or what  
183 the participant thinks is the reason the model made the mistake. Unlike prior work  
184 auditing search engines [52, 60, 68, 69, 86], our sole focus is on tracing the concrete  
185 effects that search results can have.

186 **Study 1: Distinguishing which machine learning errors reflect  
187 social stereotypes**

188 To understand the social stereotypes held by American society relevant to our machine  
189 learning task, we first elicit human judgments ( $N = 80$ ) on Common Objects in  
190 Context (COCO) [55]. COCO has 80 objects and perceived binary gender expression of  
191 pictured people annotated across the images [95]. In the survey, we ask the participants  
192 whether each object (e.g., `keyboard`, `zebra`) is stereotypically associated with men,  
193 women, or neither. As expected, not all objects reflect gender stereotypes. This is  
194 already in contrast to a somewhat common assumption in ML fairness research that  
195 any difference between groups is an amplification of a stereotype [11].

196 Among 80 objects, 13 objects are marked as stereotypes by more than half of the  
197 participants (Figs. 3). Some examples of stereotypically gendered objects are `handbag`  
198 with women, `wine glass` with women, `tie` with men, and `truck` with men. Among the  
199 remaining objects, 18 objects (e.g., `keyboard`, `carrot`, `traffic light`) are marked  
200 by zero participants as stereotypes with any gender group. If an object was marked  
201 to be a stereotype, we also asked participants whether they believed it was harmful  
202 in the abstract. Complete results are in the Supplementary Material, but we use  
203 these initial findings to select experimental stimuli in subsequent studies. In Study  
204 2a the stereotype-reinforcing condition includes women and `oven` (marked to be most

---

<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.

205 harmful), women and **hair dryer** (marked to be least harmful), and the associated  
 206 control conditions include women and **bowl**, women and **toothbrush**. In Study 2b  
 207 we also include in the stereotype-reinforcing conditions of men and **baseball glove**  
 208 (marked to be more harmful) and men and **necktie** (marked to be less harmful) with  
 209 the control conditions of men and **bench** and men and **cup**.

Stereotyped with Women							Stereotyped with Men								
handbag 21/23	hair dryer 17/24	wine glass 8/13	cat 11/19	potted plant 9/17	oven 12/23	cake 8/19	rose 7/18	ice 17/14	truck 14/18	motorcycle 8/11	baseball bat 15/21	baseball glove 11/18	sports ball 12/24	skateboard 10/19	fire hydrant 9/25
dining table 5/21	tennis racket 3/14	cow 3/25	sink 4/20	teddy bear 4/20	horse 5/26	bird 4/22	person 4/26	bear 6/17	snowboard 6/18	remote 7/22	car 6/23	suitcase 6/25	surfboard 5/23	sandwich 5/23	microwave 2/22
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	donut 3/19	person 4/26	boat 2/14	tennis racket 2/14	bicycle 2/20	bottle 3/20	tv 2/21	dog 2/21
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	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
dog 1/21	book 1/21	knife 1/22	remote 1/22	cup 1/22	mouse 1/22	toothbrush 1/24	chair 1/24	pizza 1/18	cake 1/19	hair dryer 1/20	teddy bear 1/20	train 1/20	elephant 1/20	toaster 1/21	stop sign 1/21
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	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
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	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
couch 0/22	train 0/20	kite 0/17	clock 0/24	giraffe 0/19	pizza 0/18	zebra 0/19	truck 0/18	cat 0/19	refrigerator 0/23	spoon 0/17	bus 0/16	oven 0/23	zebra 0/19	carrot 0/25	giraffe 0/19
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	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
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	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

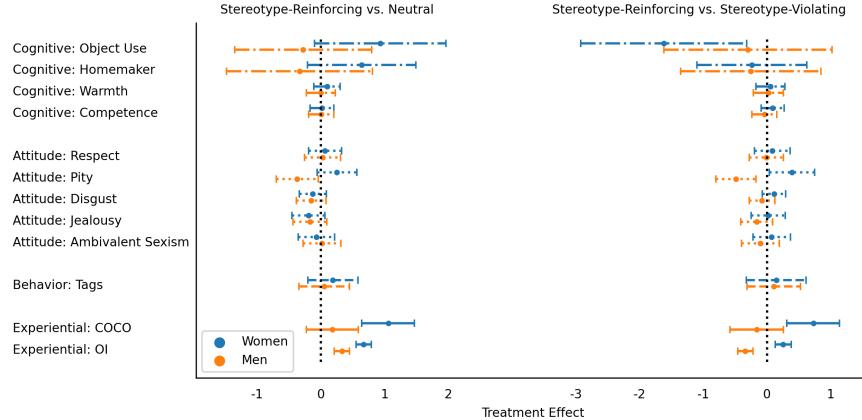
**Fig. 3 Study 1 Object Results.** Detailed participant responses for each of the 80 objects in COCO dataset. Fraction indicates number of participants asked about each object who marked it as stereotypically related to the gender group of women or men.

210 **Study 2a: Stereotype-reinforcing errors show no pragmatic  
 211 harm compared to both the stereotype-violating and neutral  
 212 conditions**

213 To test pragmatic harm in stereotype-reinforcing errors, we conduct a between-subject  
 214 survey experiment, using the stereotype-violating and neutral errors as control condi-  
 215 tions. The cover story instructs participants to look at our synthesized search result  
 216 page, imagining it is their personal phone photo album, and find a picture they had  
 217 taken of someone they saw with a particular object. The search result page looks dif-  
 218 ferent for each randomized condition. We randomly assign participants to one of the

219 three conditions ( $N = 600$ ): the stereotype-reinforcing condition exposes an image  
220 search result page with stereotype-reinforcing errors, e.g., false positive of **oven** on  
221 images of women; the stereotype-violating condition contains the same for stereotype-  
222 violating errors, e.g., false positive of **oven** on images of men; the stereotype-neutral  
223 condition contains neutral errors, e.g., false positive of **bowl** on images of women. We  
224 then measure participants' cognitive beliefs, attitudes, and behaviors to see if there  
225 are any changes because of such exposure (Methods). The behavioral measure is of  
226 particular interest, as we ask participants to undertake a realistic task they are liable  
227 to encounter by virtue of their jobs as online annotators: data labeling. We choose this  
228 measure because online participants are often the source of training labels in large-  
229 scale machine learning datasets. We ask participants to perform two common types of  
230 labeling on image data: tagging and captioning. If stereotype-reinforcing errors have  
231 an influence on participants' cognitive representations, attitudes, and tagging or cap-  
232 tioning behaviors, we should expect to see a statistically significant difference between  
233 participants who are exposed to search results with **oven**-women and those who are  
234 exposed to search results with **oven**-men or **bowl**-women.

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



**Fig. 4 Study 2, 3 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.

<sup>241</sup> **Study 2b: Stereotype-reinforcing errors show statistically  
242 significant experiential harm compared to both the  
243 stereotype-violating and neutral conditions**

<sup>244</sup> In terms of experiential harm, we design a within-subjects experiment ( $N = 100$ ).  
<sup>245</sup> We operationalize experiential harm by explicitly asking participants to rate how  
<sup>246</sup> personally harmful they find different kinds of errors (which are stereotype-reinforcing,  
<sup>247</sup> stereotype-violating, or neutral), on a scale from 0 (not at all) to 9 (extremely). This  
<sup>248</sup> experience of error is analogous to situations where one reads in the news about the  
<sup>249</sup> types of errors that artificial intelligence systems make [79], notices such a pattern of  
<sup>250</sup> errors themselves, or is informed by a friend.

<sup>251</sup> Comparing stereotype-reinforcing against neutral errors, an OLS regression shows  
<sup>252</sup> participants rate stereotype-reinforcing errors to be more harmful than neutral ones  
<sup>253</sup> ( $b = .62$ , 95% CI [.32, .91],  $p < .001$ ). However, when disaggregating by gender this  
<sup>254</sup> effect is only present among women participants (women:  $b = 1.06$ , 95% CI [.64, 1.47],

255  $p < .001$ ; men:  $b = .18$ , 95% CI [-.23, .59],  $p = .393$ ). When we use the stereotype-  
256 violating error as the control condition rather than the neutral error, we again find  
257 participants rate stereotype-reinforcing errors to be more harmful, though to a smaller  
258 degree, ( $b = .28$ , 95% CI [-.01, .58],  $p = .062$ ), with once again an effect only for  
259 women participants (women:  $b = .73$ , 95% CI [.31, 1.14],  $p = .001$ ; men:  $b = -.16$ ,  
260 95% CI [-.58, .26],  $p = .453$ ). Results are in Fig. 4.

261 In short, while we find little immediate evidence of pragmatic harms, we do find the  
262 existence of experiential harms resulting from stereotype-reinforcing errors, compared  
263 to both stereotype-violating and neutral errors. However, this pattern is present only  
264 among woman participants, and not men participants.

265 Prior work looking at a subset of what we call pragmatic harm has found very  
266 small effects in terms of cognitive belief changes about the representation of gendered  
267 occupations [52, 60], but we do not see the effects here, potentially because we have  
268 a coarser scale of measurement. Another line of work that finds a cognitive effect  
269 takes a different approach by studying occupations (e.g., peruker, lapidary) for which  
270 there are very few preconceived notions of stereotypes [86]. In our work, we focus on  
271 the activation of existing stereotypes, rather than the induction of novel stereotypes.  
272 Overall we find that the pragmatic harms are not measurable after exposure from  
273 repeated stereotypical errors in the current survey experiment, likely due to the fact  
274 that the effects of these harms are too diffuse and long-term, impacted by all of the  
275 facets of society we encounter in our lives [66]. Long-term observational studies are  
276 likely more well-suited to measure these kinds of impacts [34, 36, 49]. However, we  
277 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,  
279 consistent with the feelings of inclusivity in gender-biased occupations [60].

280 **Study 3: Stereotype-violating errors can be perceived as  
281 harmful too, but for system-justifying reasons**

282 In this study, we first test the generalizability of the previous findings by using a  
283 popular dataset in object recognition tasks which is much larger: OpenImages [54].  
284 We then explore a new hypothesis about gender presentation-aligned objects, e.g.,  
285 clothing, to dive deeper into our findings. OpenImages has 600 objects, annotated with  
286 perceived binary genders of people present in the image if applicable [77]. Following the  
287 same procedure as in the COCO dataset with new online participants ( $N = 120$ ), we  
288 find 249 of the 600 objects are marked as stereotypes by more than half participants,  
289 replicating the finding that not all objects are perceived as stereotypes (see more in  
290 Supplementary Materials). We then compile a list of 40 stereotypical objects (20 about  
291 men: e.g., **football**, **tool**; 20 about women: e.g., **doll**, **lipstick**), and 20 neutral  
292 objects (e.g., **balloon**, **goldfish**) for this study.

293 To test whether participants experience more experiential harm when they are  
294 exposed to stereotype-reinforcing (e.g., **skirt** on women), stereotype-violating (e.g.,  
295 **skirt** on men), and neutral (e.g., **toothbrush** on women) errors, we use a similar  
296 procedure as in Study 2b. Rather than asking simply about “personal harm” as we did  
297 in Study 2b, here we draw from the Positive and Negative Affect Schedule (PANAS;  
298 [18, 90]) and provide more details by asking about if they experience harm such as  
299 feeling upset, irritated, ashamed, or distressed. We conduct a within-subjects study  
300 and ask participants ( $N = 300$ ) to report their subjective experiences on a Likert  
301 scale from 0 to 9 for a variety of errors (see more in Methods). The analysis uses a  
302 mixed-effects regression with experimental conditions as the independent variable, a  
303 composite score of experiential harm as the dependent variable, participants’ gender  
304 as the covariate variable, and error terms clustered at the individual level.

305 Replicating Study 2b, we find that participants experience stereotype-reinforcing  
306 errors to be more harmful than neutral ones ( $b = .50$ , 95%  $CI[.42, .59]$ ,  $p < .001$ ).

307 Again, this pattern is more pronounced among women participants ( $b = .67$ , 95%  
308  $CI [.55, .79]$ ,  $p < .001$ ), with now a small effect among men participants ( $b = .33$ ,  
309 95%  $CI [.21, .45]$ ,  $p < .001$ ). Different from Study 2b, we do not see differences in  
310 experiential harm between stereotype-reinforcing and stereotype-violating conditions  
311 ( $b = -.04$ , 95%  $CI [-.13, .05]$ ,  $p = .338$ ). The effect is canceled out by the opposite  
312 effects for women ( $b = .25$ , 95%  $CI [.13, .38]$ ,  $p < .001$ ) and men ( $b = -.34$ , 95%  $CI [-$   
313  $.46, -.22]$ ,  $p < .001$ ) participants. In other words, while women participants feel upset,  
314 irritated, ashamed, and distressed when they see stereotype-reinforcing errors (e.g.,  
315 skirt on women), men participants feel that way when they see stereotype-violating  
316 errors (e.g., skirt on men). Results are in Fig. 4.

317 To better understand this finding, we conduct an exploratory analysis that digs  
318 deeper into the 40 stereotypical objects to understand why stereotype-violating errors  
319 are sometimes perceived to be more experientially harmful than stereotype-reinforcing  
320 ones. According to the gender trouble framework, costume (i.e., body and appearance)  
321 and script (i.e., behavior, traits, and preferences) are two aspects of gender perfor-  
322 mance, and reactions to androgynous or conventionally contradictory components can  
323 differ depending on which of the two it manifests in [15, 40, 63, 84]. We thus hypothe-  
324 size that in our study, conventionally contradictory costume objects may be evoking a  
325 more negative reaction compared to conventionally contradictory script objects [74].  
326 So, we add an additional independent variable we call “wearable.” We determined  
327 the value of this variable by manually marking 13 of the 40 stereotypical objects to  
328 be conventionally wearable by a person. These include objects like **football helmet**  
329 and **lipstick**, and exclude those like **truck** or **wine glass**. After introducing this  
330 independent variable, we find that overall participants do rate stereotype-reinforcing  
331 errors to be more harmful than stereotype-violating ones ( $b = .23$  95%  $CI [.12, .34]$ ,  
332  $p < .001$ ), though again this is true of women participants ( $b = .49$ , 95%  $CI [.34,$   
333  $.64]$ ,  $p < .001$ ) rather than men participants ( $b = -.03$ , 95%  $CI [-.18, .12]$ ,  $p = .726$ ).

334 Very interestingly, for the interaction effect of a “wearable” object with the condi-  
335 tion type, we find that wearable stereotype-violating errors have higher experiential  
336 harm than wearable stereotype-reinforcing errors ( $b=.80$ , 95% CI [.62, .99],  $p < .001$ ),  
337 which is higher for men participants ( $b=.94$ , 95% CI [.67, 1.12],  $p < .001$ ) than women  
338 participants ( $b=.69$ , 95% CI [.43, .94],  $p < .001$ ).

339 In addition to this result being a consequence of backlash effects [75], we raise  
340 two more possible mechanisms. First, it could be seen as an expression of precari-  
341 ous manhood; a concept that suggests manhood is precarious and needs continuous  
342 social validation such that threats to traditional masculinity can provoke anxiety in  
343 men [85]. Second, these results may reflect elements of transphobia, which involves  
344 a negative reaction to the apparent incongruity between a person’s perceived gender  
345 and a wearable gender presentation item [15, 63]. The divergent effect between men  
346 and women participants aligns with research indicating that transphobia is higher  
347 amongst cisgender men when judging transgender women due to the perceived threat  
348 to masculinity [59, 64]. This analysis pushes us to reevaluate how we should think  
349 about reducing experiential harm, as it may encompass intolerances we do not wish  
350 to support.

### 351 **Study 4: Plurality of stereotypes and harms with image 352 recognition objects**

353 Finally, we report qualitative analyses on open-ended responses from participants’  
354 annotations, where they explain why certain objects are seen as stereotypes and harm-  
355 ful or not. While prior work in gender stereotypes has often focused on social roles and  
356 traits [28, 38], our data provides insights as to how objects (e.g., oven, hair dryer) can  
357 also be associated with stereotypes. This is an important departure because it expands  
358 the scope of machine learning tasks for which stereotypes are relevant beyond its cur-  
359 rent more narrow framing. Specifically, when a participant from Study 1 responds that

360 an object is a stereotype, we follow up and ask: “Please describe in 1-2 sentences a)  
361 why you marked the above as a stereotype, and b) why you found it to be harmful or  
362 not.”

363 One of the authors coded the responses for why an object is a stereotype into  
364 roughly six categories. The most prevalent reasons were: descriptive (45%), e.g., for  
365 **handbag** and women: “women are often seen wearing handbags and buying them”;  
366 occupation/role (22%), e.g., for **oven** and women: “women are stereotyped to always  
367 be in the kitchen cooking while the men go out and work”; trait (11%), e.g., for  
368 **chair** and men: “sometimes men would be seen as coming home and just being lazy  
369 and lounging in their chair.” The full analysis is in the Supplementary Material. It is  
370 interesting to note that an object’s association to a stereotype is frequently mediated  
371 by its connection to a role or trait, which are the more common sites of inquiry when it  
372 comes to stereotypes. We also found that associations between a group and an object  
373 can exist through a number of paths. For example, explanations for stereotypical  
374 associations between cats and women include: “cat lady,” “women are called *kitten*,”  
375 “women like cats more than dogs,” “cats are a feminine animal,” and “women are  
376 called *cougars*.”

377 When asked why a stereotype was harmful or not, many respondents simply reit-  
378 erated that the object was a stereotype. Dropping those responses, one of the authors  
379 coded the free responses of why a stereotype was marked to be harmful into seven  
380 categories, with the top three being: proscriptive (40%), e.g., for **dining table** and  
381 women: “it makes it looked down upon if a man cooks dinner”; prescriptive (26%),  
382 e.g., for **dining table** and women: “I think it puts women in a box that says they  
383 must prepare dinner”; negative trait (13%), e.g., for **handbag** and women: “it is harm-  
384 ful because it implies that women cares more about looks and their appearance.” The  
385 remaining response categories are in the Supplementary Material. There seems to be a  
386 disparity in responses based on the participant’s gender regarding whom they believe

387 is harmed. When women specify which of the men group or women group are harmed,  
388 they say it is the women group 79% (95% CI [.67, .88]) of the time, while men say it  
389 is the women group only 67% (95% CI [.51, .80]) of the time.

390 Building on Study 1's finding that participants do not even all agree on whether an  
391 object is a stereotype or not (and if it is, whether it is harmful), this analysis further  
392 shows that even when participants are in agreement that an object is a stereotype,  
393 they are not necessarily in agreement about why. The same holds true for whether a  
394 stereotype is harmful. One potential implication of this is considering whether different  
395 reasonings should lead to different bias mitigation. For example, if the reason an object  
396 is a stereotype is descriptive, then mitigation should aim to change the cognitive  
397 representations of people. To change these descriptive statistics, while we can work  
398 to alter the model outputs, we should also work to change society, the burden of  
399 which falls on a much larger group than just machine learning practitioners, e.g.,  
400 policymakers. On the other hand, if particular stereotypes are deemed harmful because  
401 they are prescriptive and seem to restrict people from various avenues, we can consider  
402 ways to break free of gender norms.

## 403 Discussion

404 In summary, our studies have three key findings regarding our three conceptual contribu-  
405 tions: a meaningful distinction between machine learning errors is whether they are  
406 stereotype-reinforcing, stereotype-violating, or neutral; harm formulated as pragmatic  
407 or experiential; and showcasing how harm annotations can stem from a diversity of rea-  
408 sons that require critical engagement. First, we find that harm is different depending on  
409 a machine learning error's relation to a stereotype. Second, while stereotype-reinforcing  
410 errors do not lead to more pragmatic harm in the lab setting we use, we do find that  
411 stereotype-reinforcing errors are consistently found to be more experientially harm-  
412 ful. Such experiential harm is unequally distributed, impacting more participants who

413 are women than who are men. Formulating concrete notions of harm as we have done  
414 has implications beyond just machine learning: legal documents like the European  
415 AI Act is beginning to incorporate notions of psychological harm but lacking defini-  
416 tions to ground regulation in [7, 70]. Third, we find stereotype-violating errors are also  
417 experientially harmful, especially when these errors pertain to wearable items asso-  
418 ciated with gender presentation. This effect is stronger for participants who identify  
419 as men compared to those who identify as women. This final point warrants an espe-  
420 cially nuanced discussion, as we find ourselves qualifying a prior claim that we should  
421 take people’s words at face value when they indicate something is personally harm-  
422 ful. To navigate this complexity, we turn to the notions of epistemic injustice [35] and  
423 standpoint epistemology [31, 67, 93]. If we interpret the negative reactions to mis-  
424 classifications of stereotypically feminine clothing items on men as a manifestation of  
425 precarious manhood [85] or transphobia [15], then we should down weight these con-  
426cerns in designing mitigation algorithms. Respecting people’s experiential harms may  
427 not be as simple as accepting them at face value for direct measurement, but rather  
428 involves understanding which groups are likely to be harmed by what kinds of errors  
429 and why.

430 Our findings call for reconsidering fairness measurement in supervised machine  
431 learning tasks. This involves considering how we can leverage human-driven insights  
432 to inform model training and evaluation [13]. Traditionally, fairness evaluations tend  
433 to focus on stereotypes only in relation to occupations or traits. However our work  
434 expands this idea by showing that labels such as objects can also give rise to such  
435 harms. Additionally, most prior work has only considered the implications of errors  
436 that reinforce stereotypes, which is relatively more intuitive to think of as harm-  
437 ful. However, both practically and normatively, it is important to understand the  
438 implications of stereotype-violating errors. Practically, strategies aimed at mitigating  
439 stereotype-reinforcing errors which act upon the target label will inevitably impact

440 the occurrence of stereotype-violating errors as well. And normatively, there are also  
441 questions about whether stereotype-violating errors may even play a role in reduc-  
442 ing stereotypical associations by counteracting them. This finding that not only are  
443 certain labels more liable to cause harm than others, but that it matters for *which*  
444 demographic group that label is misclassified, suggests that generic approaches like  
445 having a higher threshold for the classification of certain labels are insufficient. Instead,  
446 more nuanced fairness-through-awareness approaches [27] will need to be taken. While  
447 adopting simply a cost-sensitive framework [53] (e.g., different costs are associated  
448 with false positives and false negatives) is a simplified interpretation of our findings, it  
449 could be a starting point as one grapples with the questions of whose levels of harms  
450 we would prioritize reducing in a bias mitigation framework.

451 Understanding whose levels of harms we should prioritize, and why, will come from  
452 stronger understandings of the psychological basis and reasoning of different harms.  
453 Our finding from Study 4 that stereotypical associations between a single group and  
454 object can emerge from many paths (e.g., the many reasonings behind the associa-  
455 tion between cat and women), each with different normative valences, illustrates  
456 what an oversimplification it is to only label an association as “good” or “bad,”  
457 and the limitations of mitigations simply aiming to sever the associations deemed  
458 “bad.” This underscores the importance of work about diversity in annotators’ per-  
459 spectives [22, 23, 26, 50, 65, 89], and how much complexity is reduced by the use of  
460 discrete labels. Qualitative follow-up questions supplemented our annotations, where  
461 a lack of consensus is not a weakness or artifact to be averaged out, but rather a point  
462 for deeper inquiry on how to prioritize differential experiences of harm. This also indi-  
463 cates that even if the growing power of large language models enables us to predict  
464 with higher accuracy which objects are stereotypes, we likely still may want to ensure  
465 these annotations come from people themselves [3, 46, 94], thus allowing room for  
466 positionality, explanation, and critical reflection.

467 Our findings are limited by the methodological choices we made: First, we focused  
468 on gender stereotypes as a case study. We do not know to what extent this finding  
469 generalizes to other groups such as race and age. Second, we recruited online partic-  
470 ipants who identify as men and women and who speak English without an extensive  
471 inclusion of non-binary participants or who come from a different cultural background.  
472 Given that stereotypes are culture-sensitive, and our work also shows that the harm  
473 perception is identity-sensitive, future work needs to study the interaction between  
474 participants' identity, culture, and harm perceptions. Third, by setting a threshold of  
475 50% for respondents indicating an object is a stereotype, we are in some senses privi-  
476 leging the majority opinion which may further reify marked stereotypes to be those for  
477 the majority subset [37, 61]. Fourth, the survey experiment does not capture harms  
478 beyond the two we measure (e.g., stereotype-threat [82, 83]), nor the longitudinal  
479 effects of machine learning effects. Future work needs to capture not only the plurality  
480 in harm of machine learning errors but also how its' effect emerges and endures over  
481 time.

482 Overall, our work offers a rigorous empirical study connecting machine learning  
483 outputs to concrete harms by understanding the impact of stereotypical misclassifica-  
484 tions. Rather than gesturing at harm as a justification for fairness measurement, we  
485 are very concrete in our analysis of the effects on people. Our finding that stereotype-  
486 reinforcing errors are experientially harmful for women underscores the importance for  
487 machine learning fairness interventions to be more rooted in social contexts, moving  
488 beyond objectives like just achieving equal prediction performance across groups. The  
489 diversity of responses we've presented, each influenced by participants' unique ratio-  
490 nales, suggests the need for greater exploration of human psychological experiences in  
491 understanding how machine learning can cause harm.

492 **Methods**

493 **Analysis**

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

499 **Participants**

500 While men and women generally tend to hold the same gender stereotypes [29, 44,  
501 58, 92], we still collect equal numbers of participants who identify as men and women,  
502 and use this variable as a covariate throughout. Due to limitations in the survey  
503 platform which only allow us to specify gender as “male” or “female,” this formulation  
504 excludes people who identify as non-binary, which is a harmful limitation. Because we  
505 do not control for race in the recruitment of participants, our sample diverges from a  
506 nationally representative sample. For the gender stereotype scope of our current work,  
507 we find this to be an acceptable limitation, especially given that one defining feature  
508 of stereotypes is they are largely shared through a cultural consensus [51].

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

513 **Studies 1, 3: Distinguishing Errors by Stereotype**

514 When asking about which machine learning errors are stereotypes, we make sure to  
515 ask participants about their perception of stereotypes held by Americans, rather than  
516 for their personal beliefs [24].

**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 4	7	1.75	Women	0	3	5	0	0	25	6	1	40
			Men	1	4	2	2	0	30	1	0	40
2a	10	2.50	Women	1	11	32	8	0	229	19	0	300
			Men	0	19	35	10	1	211	22	2	300
2b	5	1.25	Women	0	4	7	3	1	35	5	0	50
			Men	0	4	2	3	1	35	5	0	50
3	4	1	Women	0	5	8	0	0	42	4	1	60
			Men	0	2	6	5	1	44	2	0	60
(Labeling)	5	1.25	Women	0	5	15	1	0	120	7	2	150
			Men	1	9	17	6	1	107	9	0	150
(Harms)												

## 517 Study 2a: Measuring Pragmatic Harm

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

529 The stimuli take the form of an image search result and are pictured in Fig. 5  
 530 with teal and orange colored boxes around the component of the image that changes  
 531 between conditions. The search bar contains the search query, and then eight images

---

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

532 that may or may not be correctly retrieved are shown. Each of the eight images is  
 533 annotated with either “In image” or “Not in image” to make it clear to the partici-  
 534 pant which images are correct or not. The stereotype-reinforcing condition on the left  
 535 contains the search query of “oven” with five correctly identified ovens, and three false  
 536 positive images that all contain women. In other words, this classifier erroneously (and  
 537 stereotypically) assumes there are ovens in images of women. The stereotype-violating  
 538 condition contains the same search query, but the mistakes are replaced with false pos-  
 539 itive images that all contain men. The neutral condition contains all of the exact same  
 540 images as the stereotype-reinforcing condition, with the only change being that the  
 541 search query is now “bowl” instead of “oven.” This is because the five correct images  
 542 were deliberately chosen to contain both bowls and ovens, which allows us to control  
 543 for the variance between the different search conditions. All false positive images were  
 544 selected from the actual errors of a Vision Transformer (ViT) model [25] trained on  
 545 COCO so that they are as realistic as possible to a computer vision model’s errors,  
 546 and not completely egregious, e.g., a picture of a woman in a sports field as a false  
 547 positive for “oven” or “bowl.”



**Fig. 5 Study 2 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.

548 For *cognitive beliefs*, we ask three sets of questions which span the spectrum of  
 549 stereotype-specific to more generically about gendered beliefs. Concretely, we ask

about: estimations of who uses ovens and bowls more between men and women; estimations of who tends to be the homemaker more between men and women; and perceived levels of warmth and competence [33] of women. To assess *attitude*, we ask two sets of questions. The first is about how participants feel about women in terms of four emotional components that are believed to mediate interactions between cognitive beliefs and behaviors: a) respect or admiration, b) pity or sympathy, c) disgust or sicken-  
ing, and d) jealousy or envy [20, 32, 78]. The second asks about sexist attitudes via a shortened scale focused on benevolent sexism [38, 39, 73].<sup>3</sup> Finally, for *behavioral* measures, we ask participants to undertake a realistic task they are liable to encounter which can cause harm: data labeling [62]. We chose this behavior measure because online participants are often the source of training labels in large-scale machine learning datasets. We ask participants to perform two common types of labeling on image data: tagging and captioning (Fig. 6). In the tagging task, we ask participants to label the top three most relevant tags in an image which contains both the stereotype object (e.g., oven) and neutral object (e.g., bowl). We 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 swap whether there is a bowl or oven present in the image. This is to understand if the annotators will differently describe who is interacting with the object depending on whether it is stereotypically associated with women or not. All images are generated and/or manipulated by DALL-E 2.

## Dependent Variables

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

---

<sup>3</sup>We ask questions from the Ambivalent Sexism Inventory [38] about benevolent sexism, as opposed to hostile sexism, because the latter is believed to suffer heavily from social desirability bias.



**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.

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

585        **Behavior - Captions.** We offer some descriptive statistics about the captions in  
 586        the Supplementary Material. This analysis was mostly exploratory, and we do not find  
 587        any statistically significant differences. We first ran Study 2a looking at pragmatic  
 588        harms on the stereotype of women and **oven** (with **bowl** as the control). In this iter-  
 589        ation, we asked that respondents please describe each person in the image in separate  
 590        sentences. However, there was too much noise in how respondents interpreted this set  
 591        of instructions, such that the data became hard to interpret. Thus, in our second iter-  
 592        ation of this study using the stereotype of women and hair dryer (with **toothbrush**  
 593        as the control), we have two separate text entry boxes to caption each person in the

594 image. We only present the results of this iteration in the table, as we were unable to  
595 parse anything differentiating in the first iteration.

596 **Cognitive - Object Use.** In this measurement, we have a value from -10 (mostly  
597 men) to 10 (mostly women) for both the stereotypical and neutral object. The  
598 dependent variable is the summation of both values. Again, this is intended to cap-  
599 ture whether the stereotype-reinforcing condition is able to change the value of its  
600 associated object more than the control condition is able to.

## 601 Study 2b, 3: Measuring Experimental Harm

602 In Study 2b, in addition to personal discomfort, we also ask about societal harm.  
603 This way, even if the participant does not personally feel harmed, they may feel it  
604 on behalf of the stereotyped group. However, we find that participants' responses to  
605 both personal and societal harm are extremely correlated, and leave the results for  
606 the latter in the Supplementary Material.

607 **Acknowledgments.** This material is based upon work supported by the National  
608 Science Foundation Graduate Research Fellowship to Angelina Wang. We are grateful  
609 to funding from the Data-Driven Social Science Initiative at Princeton University. We  
610 thank Orly Bareket, Molly Crockett, Sunnie S. Y. Kim, Anne Kohlbrenner, Danaë  
611 Metaxa, Vikram V. Ramaswamy, Olga Russakovsky, Hanna Wallach, and members of  
612 the Visual AI Lab at Princeton, Fiske Lab at Princeton, and Perception and Judgment  
613 Lab at the University of Chicago for feedback.

## 614 References

- 615 [1] Abbasi M, Friedler SA, Scheidegger C, et al (2019) Fairness in representa-  
616 tion: quantifying stereotyping as a representational harm. Siam International  
617 Conference on Data Mining

- 618 [2] Allport GW, Clark K, Pettigrew T (1954) The nature of prejudice
- 619 [3] Argyle LP, Busby EC, Fulda N, et al (2023) Out of one, many: Using language  
620 models to simulate human samples. Political Analysis
- 621 [4] Barlas P, Kyriakou K, Guest O, et al (2021) To "see" is to stereotype: Image  
622 tagging algorithms, gender recognition, and the accuracy-fairness trade-off.  
623 Proceedings of the ACM on Human-Computer Interaction (CSCW)
- 624 [5] Barocas S, Crawford K, Shapiro A, et al (2017) The Problem With Bias: Alloca-  
625 tive Versus Representational Harms in Machine Learning. In: Proceedings of  
626 SIGCIS, Philadelphia, PA
- 627 [6] Barocas S, Hardt M, Narayanan A (2019) Fairness and Machine Learning:  
628 Limitations and Opportunities. fairmlbook.org, <http://www.fairmlbook.org>
- 629 [7] Bayefsky R (2016) Psychological harm and constitutional standing. Brooklyn Law  
630 Review
- 631 [8] Beeghly EM (2014) Seeing difference: The ethics and epistemology of stereotyp-  
632 ing. UC Berkeley PhD Thesis
- 633 [9] Bhaskaran J, Bhallamudi I (2019) Good secretaries, bad truck drivers? occupa-  
634 tional gender stereotypes in sentiment analysis. Proceedings of the First Workshop  
635 on Gender Bias in Natural Language Processing
- 636 [10] Blodgett SL, Barocas S, III HD, et al (2020) Language (technology) is power: A  
637 critical survey of "bias" in nlp. Association for Computational Linguistics (ACL)
- 638 [11] Blodgett SL, Lopez G, Olteanu A, et al (2021) Stereotyping norwegian salmon:  
639 An inventory of pitfalls in fairness benchmark datasets. Proceedings of the 59th  
640 Annual Meeting of the Association for Computational Linguistics and the 11th

641 International Joint Conference on Natural Language Processing

642 [12] Bolukbasi T, Chang KW, Zou J, et al (2016) Man is to computer programmer  
643 as woman is to homemaker? debiasing word embeddings. Conference on Neural  
644 Information Processing Systems (NeurIPS)

645 [13] Boykin CM, Dasch ST, Jr. VR, et al (2021) Opportunities for a more interdisci-  
646 plinary approach to measuring perceptions of fairness in machine learning. Equity  
647 and Access in Algorithms, Mechanisms, and Optimization (EAAMO)

648 [14] Burgess D, Borgida E (1999) Who women are, who women should be: descriptive  
649 and prescriptive gender stereotyping in sex discrimination. Psychology, Public  
650 Policy, and Law 5

651 [15] Butler J (1990) Gender trouble: Feminism and the subversion of identity.  
652 Routledge

653 [16] Caliskan A, Bryson JJ, Narayanan A (2017) Semantics derived automatically  
654 from language corpora contain human-like biases. Science

655 [17] Cao Y, Sotnikova A, Hal Daumé III RR, et al (2022) Theory-grounded mea-  
656 surement of u.s. social stereotypes in english language models. Conference of  
657 the North American Chapter of the Association for Computational Linguistics:  
658 Human Language Technologies

659 [18] Crawford JR, Henry JD (2004) The positive and negative affect schedule (panas):  
660 construct validity, measurement properties and normative data in a large non-  
661 clinical sample. British Journal of Clinical Psychology

662 [19] Crawford K (2017) The trouble with bias. Conference on Neural Information  
663 Processing Systems Keynote

- 664 [20] Cuddy AJC, Fiske ST, Glick P (2007) The BIAS map: behaviors from intergroup  
665 affect and stereotypes. *Journal of Personality and Social Psychology* 92
- 666 [21] Danks D, London AJ (2017) Algorithmic bias in autonomous systems. Interna-  
667 tional Joint Conference on Artificial Intelligence (IJCAI)
- 668 [22] Davani AM, Díaz M, Prabhakaran V (2022) Dealing with disagreements: Look-  
669 ing beyond the majority vote in subjective annotations. *Transactions of the*  
670 *Association for Computational Linguistics*
- 671 [23] Denton E, Díaz M, Kivlichan I, et al (2021) Whose ground truth? accounting for  
672 individual and collective identities underlying dataset annotation. *NeurIPS 2021*  
673 Workshop on Data-Centric AI
- 674 [24] Devine PG, Elliot AJ (1995) Are racial stereotypes really fading? the princeton  
675 trilogy revisited. *Personality and Social Psychology Bulletin* 21
- 676 [25] Dosovitskiy A, Beyer L, Kolesnikov A, et al (2021) An image is worth 16x16  
677 words: Transformers for image recognition at scale. International Conference on  
678 Learning Representations (ICLR)
- 679 [26] Dumitrache A, Aroyo L, Welty C (2018) Capturing ambiguity in crowdsourc-  
680 ing frame disambiguation. AAAI Conference on Human Computation and  
681 Crowdsourcing (HCOMP)
- 682 [27] Dwork C, Hardt M, Pitassi T, et al (2012) Fairness through awareness. Proceed-  
683 ings of the 3rd Innovations in Theoretical Computer Science Conference
- 684 [28] Eagly AH (1987) Sex differences in social behavior: A social-role interpretation.  
685 Lawrence Erlbaum Associates, Inc

- 686 [29] Eagly AH, Nater C, Miller DI, et al (2020) Gender stereotypes have changed:  
687 A cross-temporal meta-analysis of u.s. public opinion polls from 1946 to 2018.  
688 American Psychologist
- 689 [30] Ellemers N, et al (2018) Gender stereotypes. Annual review of psychology 69:275–  
690 298
- 691 [31] Fatima S (2020) I know what happened to me: The epistemic harms of microag-  
692 gression. Microaggressions and Philosophy
- 693 [32] Fiske ST, Cuddy AJC, Glick P (2002) Emotions up and down: Intergroup emo-  
694 tions result from status and competition. Prejudice to Intergroup Emotions:  
695 Differentiated Reactions to Social Groups
- 696 [33] Fiske ST, Cuddy AJC, Glick P, et al (2002) A model of (often mixed) stereotype  
697 content: Competence and warmth respectively follow from perceived status and  
698 competition. Journal of Personality and Social Psychology 82
- 699 [34] Ford TE (1997) Effects of stereotypical television portrayals of african-americans  
700 on person perception. Social Psychology Quarterly
- 701 [35] Fricker M (2009) Epistemic injustice: Power and the ethics of knowing. Oxford  
702 University Press
- 703 [36] Fujioka Y (1999) Television portrayals and african-american stereotypes: Exam-  
704 ination of television effects when direct contact is lacking. Journalism and Mass  
705 Communication Quarterly
- 706 [37] Ghavami N, Peplau LA (2012) An intersectional analysis of gender and ethnic  
707 stereotypes: Testing three hypotheses. Psychology of Women Quarterly 37

- 708 [38] Glick P, Fiske ST (1996) The ambivalent sexism inventory: Differentiating hostile  
709 and benevolent sexism. *Journal of Personality and Social Psychology* 70
- 710 [39] Glick P, Whitehead J (2010) Hostility toward men and the perceived stability of  
711 male dominance. *Social Psychology* 41
- 712 [40] Goffman E (1959) *The presentation of self in everyday life*. Doubleday
- 713 [41] Gordon ML, Lam MS, Park JS, et al (2022) Jury learning: Integrating dissenting  
714 voices into machine learning models. Conference on Human Factors in Computing  
715 Systems (CHI)
- 716 [42] Hamilton DL, Sherman JW (2014) Stereotypes. In: *Handbook of social cognition*.  
717 Psychology Press, p 17–84
- 718 [43] Hellman D (2011) When is discrimination wrong? Harvard University Press
- 719 [44] Hentschel T, Heilman ME, Peus CV (2019) The multiple dimensions of gender  
720 stereotypes: A current look at men's and women's characterizations of others and  
721 themselves. *Frontiers in Psychology*
- 722 [45] Hilton JL, Von Hippel W (1996) Stereotypes. *Annual review of psychology*  
723 47(1):237–271
- 724 [46] Hääläinen P, Tavast M, Kunnari A (2023) Evaluating large language models  
725 in generating synthetic hci research data: a case study. Conference on Human  
726 Factors in Computing Systems (CHI)
- 727 [47] Jacobs AZ, Wallach H (2021) Measurement and fairness. Conference on Fairness,  
728 Accountability and Transparency (FAccT)
- 729 [48] Jaggar AM (1989) Love and knowledge: Emotion in feminist epistemology. *Inquiry*

- 730 [49] Jennings-Walstedt J, Geis FL, Brown V (1980) Influence of television commercials  
731 on women's self-confidence and independent judgment. *Journal of Personality and*  
732 *Social Psychology*
- 733 [50] Kairam S, Heer J (2016) Parting crowds: Characterizing divergent interpretations  
734 in crowdsourced annotation tasks. *ACM Conference On Computer-Supported*  
735 *Cooperative Work And Social Computing (CSCW)*
- 736 [51] Katz D, Braly K (1933) Racial stereotypes of one hundred college students. *The*  
737 *Journal of Abnormal and Social Psychology* 28
- 738 [52] Kay M, Matuszek C, Munson SA (2015) Unequal representation and gender  
739 stereotypes in image search results for occupations. *Conference on Human Factors*  
740 *in Computing Systems (CHI)*
- 741 [53] Kukar M, Kononenko I (1998) Cost-sensitive learning with neural networks.  
742 European Conference on Artificial Intelligence
- 743 [54] Kuznetsova A, Rom H, Alldrin N, et al (2020) The open images dataset v4:  
744 Unified image classification, object detection, and visual relationship detection at  
745 scale. *International Journal of Computer Vision (IJCV)*
- 746 [55] Lin TY, Maire M, Belongie S, et al (2014) Microsoft COCO: Common objects in  
747 context. European Conference on Computer Vision (ECCV)
- 748 [56] Lippmann W (1922) Public opinion.
- 749 [57] Litman L, Robinson J, Abberbock T (2017) Turkprime.com: A versatile crowd-  
750 sourcing data acquisition platform for the behavioral sciences. *Behavior Research*  
751 *Methods* 42

- 752 [58] López-Sáez M, Lisbona A (2014) Descriptive and prescriptive features of gender  
753 stereotyping. relationships among its components. International Journal of Social  
754 Psychology 24
- 755 [59] Makwana AP, Dhont K, keersmaecker JD, et al (2018) The motivated cognitive  
756 basis of transphobia: The roles of right-wing ideologies and gender role beliefs.  
757 Sex Roles 79
- 758 [60] Metaxa D, Gan MA, Goh S, et al (2021) An image of society: Gender and  
759 racial representation and impact in image search results for occupations. ACM  
760 Conference on Human-Computer Interaction (CSCW)
- 761 [61] Mill JS (1859) On liberty. Longman, Roberts, Green Co
- 762 [62] van Miltenburg E (2016) Stereotyping and bias in the flickr30k dataset. Proceed-  
763 ings of the Workshop on Multimodal Corpora
- 764 [63] Morgenroth T, Ryan MK (2020) The effects of gender trouble: An integrative the-  
765 oretical framework of the perpetuation and disruption of the gender/sex binary.  
766 Perspectives on Psychological Science 16
- 767 [64] Nagoshi CT, Cloud JR, Lindley LM, et al (2019) A test of the three-component  
768 model of gender-based prejudices: Homophobia and transphobia are affected by  
769 raters' and targets' assigned sex at birth. Sex Roles 80
- 770 [65] Noble JA (2012) Minority voices of crowdsourcing: why we should pay atten-  
771 tion to every member of the crowd. ACM Conference On Computer-Supported  
772 Cooperative Work And Social Computing (CSCW)
- 773 [66] Noble SU (2018) Algorithms of oppression: How search engines reinforce racism.  
774 NYU Press

- 775 [67] O'Dowd O (2018) Microaggressions: A kantian account. Ethical Theory and Moral  
776 Practice 21
- 777 [68] Otterbacher J, Bates J, Clough P (2017) Competent men and warm women:  
778 Gender stereotypes and backlash in image search results. Conference on Human  
779 Factors in Computing Systems (CHI)
- 780 [69] Otterbacher J, Checco A, Demartini G, et al (2018) Investigating user perception  
781 of gender bias in image search: The role of sexism. ACM SIGIR Conference on  
782 Research and Development in Information Retrieval (SIGIR)
- 783 [70] Pałka P (2023) Ai, consumers & psychological harm. Cambridge University Press
- 784 [71] Peterson JC, Battleday RM, Griffiths TL, et al (2019) Human uncertainty makes  
785 classification more robust. International Conference on Computer Vision (ICCV)
- 786 [72] Rini R (2020) The ethics of microaggression. Routledge Taylr & Francis Group
- 787 [73] Rollero C, Glick P, Tartaglia S (2014) Psychometric properties of short versions  
788 of the ambivalent sexism inventory and ambivalence toward men inventory. TPM-  
789 Testing, Psychometrics, Methodology in Applied Psychology 21
- 790 [74] Rudman LA, Moss-Racusin CA, Glick P, et al (2012) Reactions to vanguards:  
791 Advances in backlash theory. Advances in experimental social psychology 45
- 792 [75] Rudman LA, Moss-Racusin CA, Phelan JE, et al (2012) Status incongruity  
793 and backlash effects: Defending the gender hierarchy motivates prejudice against  
794 female leaders. Journal of Experimental Social Psychology 48
- 795 [76] Sap M, Swayamdipta S, Vianna L, et al (2022) Annotators with attitudes: How  
796 annotator beliefs and identities bias toxic language detection. Conference of  
797 the North American Chapter of the Association for Computational Linguistics:

- 799 [77] Schumann C, Ricco S, Prabhu U, et al (2021) A step toward more inclusive people  
800 annotations for fairness. ACM Conference on Artificial Intelligence, Ethics, and  
801 Society (AIES)
- 802 [78] Seger CR, Banerji I, Park SH, et al (2017) Specific emotions as mediators of the  
803 effect of intergroup contact on prejudice: findings across multiple participant and  
804 target groups. *Cognition and Emotion* 31
- 805 [79] Simonite T (2018) When It Comes to Gorillas, Google Photos Remains Blind.  
806 *Wired*, January
- 807 [80] Sotnikova A, Cao YT, III HD, et al (2021) Analyzing stereotypes in generative  
808 text inference tasks. Findings of the Association for Computational Linguistics:  
809 ACL-IJCNLP
- 810 [81] Spencer SJ, Steele CM, Quinn DM (1999) Stereotype threat and women's math  
811 performance. *Journal of Experimental Social Psychology* 35
- 812 [82] Spencer SJ, Logel C, Davies PG (2015) Stereotype threat. *Annual Review of*  
813 *Psychology* 67
- 814 [83] Steele CM, Aronson J (1995) Stereotype threat and the intellectual test per-  
815 formance of african americans. *Journal of Personality and Social Psychology*  
816 69
- 817 [84] Stern C, Rule NO (2017) Physical androgyny and categorization difficulty shape  
818 political conservatives' attitudes toward transgender people. *Social Psychological*  
819 *and Personality Science*

- 820 [85] Vandello JA, Bosson JK, Cohen D, et al (2008) Precarious manhood. *Journal of*  
821 *Personality and Social Psychology*
- 822 [86] Vlasceanu M, Amadio DM (2022) Propagation of societal gender inequality by  
823 internet search algorithms. *Proceedings of the National Academy of Sciences of*  
824 *the United States of America (PNAS)* 119
- 825 [87] Wang A, Barocas S, Laird K, et al (2022) Measuring representational harms  
826 in image captioning. *ACM Conference on Fairness, Accountability, and Trans-*  
827 *parency (FAccT)*
- 828 [88] Wang A, Liu A, Zhang R, et al (2022) REVISE: A tool for measuring and  
829 mitigating bias in visual datasets. *International Journal of Computer Vision*  
830 (*IJCV*)
- 831 [89] Waseem Z (2016) Are you a racist or am i seeing things? annotator influence on  
832 hate speech detection on twitter. *Proceedings of the First Workshop on NLP and*  
833 *Computational Social Science*
- 834 [90] Watson D, Clark LA, Tellegen A (1988) Development and validation of brief  
835 measures of positive and negative affect: the panas scales. *Journal of Personality*  
836 *and Social Psychology* 54
- 837 [91] Wenzel K, Devireddy N, Davidson C, et al (2023) Can voice assistants be microag-  
838 gressors? cross-race psychological responses to failures of automatic speech  
839 recognition. *Conference on Human Factors in Computing Systems (CHI)*
- 840 [92] Williams JE, Best DL (1977) Sex stereotypes and trait favorability on the  
841 adjective check list. *Educational and Psychological Measurement* 37

- 842 [93] Wylie A (2003) Why standpoint matters. *Science and Other Cultures: Issues in*  
843       *Philosophies of Science and Technology*
- 844 [94] Yaghini M, Krause A, Heidari H (2021) A human-in-the-loop framework to  
845       construct context-aware mathematical notions of outcome fairness. *AAAI/ACM*  
846       *Conference on Artificial Intelligence, Ethics, and Society*
- 847 [95] Zhao D, Wang A, Russakovsky O (2021) Understanding and evaluating racial  
848       biases in image captioning. *International Conference on Computer Vision (ICCV)*
- 849 [96] Zhao J, Wang T, Yatskar M, et al (2017) Men also like shopping: Reducing gender  
850       bias amplification using corpus-level constraints. *Proceedings of the Conference*  
851       *on Empirical Methods in Natural Language Processing (EMNLP)*
- 852 [97] Zhao J, Wang T, Yatskar M, et al (2018) Gender bias in coreference resolution:  
853       Evaluation and debiasing methods. *North American Chapter of the Association*  
854       *for Computational Linguistics*