



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

## **Stereotype threat, gender and mathematics attainment: A conceptual replication of Stricker & Ward**

### **Citation for published version:**

Inglis, M & O'Hagan, S 2022, 'Stereotype threat, gender and mathematics attainment: A conceptual replication of Stricker & Ward', *PLoS ONE*, vol. 17, no. 5, e0267699.  
<https://doi.org/10.1371/journal.pone.0267699>

### **Digital Object Identifier (DOI):**

[10.1371/journal.pone.0267699](https://doi.org/10.1371/journal.pone.0267699)

### **Link:**

[Link to publication record in Edinburgh Research Explorer](#)

### **Document Version:**

Peer reviewed version

### **Published In:**

PLoS ONE

### **General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

### **Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



1

2

3

4       Stereotype threat, gender and mathematics attainment:

5                   A conceptual replication of Stricker & Ward

6

7

8   Matthew Inglis<sup>1</sup>, Steven O'Hagan<sup>2\*</sup>

9

10

11

12   <sup>1</sup> Centre for Mathematical Cognition, Loughborough University, United Kingdom

13   <sup>2</sup> School of Mathematics, University of Edinburgh, United Kingdom

14

15

16   \* Corresponding author

17   Email: [s.ohagan@ed.ac.uk](mailto:s.ohagan@ed.ac.uk) (SO'H)

## 18 **Abstract**

19           Stereotype threat has been proposed as one cause of gender differences in post-  
20 compulsory mathematics participation. Danaher and Crandall argued, based on a study  
21 conducted by Stricker and Ward, that enquiring about a student’s gender after they had  
22 finished a test, rather than before, would reduce stereotype threat and therefore increase the  
23 attainment of women students. Making such a change, they argued, could lead to nearly 5000  
24 more women receiving AP Calculus AB credit per year. We conducted a preregistered  
25 conceptual replication of Stricker and Ward’s study in the context of the UK Mathematics  
26 Trust’s Junior Mathematical Challenge, finding no evidence of this stereotype threat effect.  
27 We conclude that the ‘silver bullet’ intervention of relocating demographic questions on test  
28 answer sheets is unlikely to provide an effective solution to systemic gender inequalities in  
29 mathematics education.

30

## 31 **Introduction**

32           Mathematics education researchers have long been concerned that mathematics is  
33 experienced differently by men and women [1]. This concern is, in part, fueled by gender  
34 differences in post-compulsory participation rates in mathematical study and STEM  
35 careers [2].

36           One mechanism which some believe contributes to these observed gender differences  
37 in participation is *stereotype threat*. This account suggests that members of negatively  
38 stereotyped groups underperform when that stereotype is salient, perhaps because stereotype-  
39 related thoughts place an extra burden on stereotyped individuals’ cognitive resources [3].  
40 For example, Steele and Aronson [4] found that black participants underperformed on  
41 laboratory tests of verbal ability compared to white participants, but only when reminded of  
42 negative stereotypes concerning race and intelligence. Similarly, Spencer, Steele and

43 Quinn [5] found that women performed worse on a laboratory mathematics test than men, but  
44 only when they were told that the test usually revealed gender differences in achievement.  
45 Subsequently many similar lab-based studies have been conducted: a meta-analysis of 47  
46 such studies showed that women, on average, underperform on laboratory mathematics tests  
47 by 0.22 standard deviations when under stereotype threat conditions [6].

## 48 **Stereotype Threat in Real World Contexts**

49 Our goal in this paper is to discuss one particularly important context in which  
50 stereotype threat is hypothesized to negatively impact women's mathematics performance:  
51 authentic high-stakes tests (i.e., not in low-stakes laboratory tests, the context of the large  
52 majority of literature on stereotype threat). Stricker and Ward [7] investigated the extent to  
53 which stereotype threat plays a role in influencing women's real-world mathematics  
54 achievement by manipulating the location of demographic questions on two authentic high-  
55 stakes tests: the 1996 Advanced Placement Calculus AB examination and the Computerized  
56 Placement Test, both qualifications intended for students seeking college credit or  
57 placement [8]. Half the participants were asked to state their ethnicity and gender before  
58 answering any questions, and half gave their ethnicity and gender after answering the  
59 questions. Their hypothesis was that asking participants to state their gender before tackling  
60 the questions would increase the saliency of gender, and therefore provoke stereotype threat  
61 among the women in the sample. In contrast, they reasoned, moving these demographic  
62 questions to the end of the examination would reduce the chance of stereotype threat  
63 impacting women's performance.

64 Stricker and Ward [7] found the hypothesized significant manipulation-by-gender  
65 interaction on the AP Calculus AB exam (we focus on the 'formula score', which represents  
66 the number of correct answers, adjusted for guessing). The women who were asked about  
67 their gender in advance performed worse than those who were asked afterwards, with the

68 men showing a small trend in the opposite direction. However, Stricker and Ward also noted  
69 that the size of this effect was small (less than 10% of the variance in scores could be  
70 accounted for by this interaction effect), and therefore concluded that stereotype threat was  
71 not practically significant. A similar result was found for the reading comprehension  
72 component of the second standardized test Stricker and Ward studied, but not for the  
73 predicted mathematical components.

74 Danaher and Crandall [9] reanalyzed Stricker and Ward's [7] data, and argued that  
75 their criterion for 'practical significance' was too conservative. They pointed out that  
76 changing the location of demographic questions on examination scripts would be inexpensive  
77 and therefore that any evidence of a non-zero effect, regardless of its size, had policy  
78 implications. They calculated that changing the location of demographic questions would  
79 increase the number of U.S. women receiving AP Calculus AB credit by more than 4700 per  
80 year, and wrote that this "would be the single most cost-effective action our country could  
81 take to increase girls' performance on AP Calculus exams" (p. 1652). In response, Stricker  
82 and Ward [10] rejected Danaher and Crandall's argument, suggesting that they had  
83 selectively focused on women and mathematics (the original paper had also investigated  
84 ethnicity and verbal ability) and that, in view of traditional effect size guidelines, their  
85 decision to interpret a statistically significant but small effect as not practically significant  
86 was justified.

87 In our view, Danaher and Crandall's [9] argument is persuasive. Effects in education  
88 research found where standardized tests are the dependent measures are typically small  
89 [10, 11]. Because of this, using guidelines for effect size interpretation developed in the  
90 context of psychology may lead to effective interventions being dismissed [13]. A citation  
91 analysis suggests that Danaher and Crandall's interpretation is the more widely accepted:  
92 Scopus reports that Danaher and Crandall's [9] reanalysis has received more citations than

93 Stricker and Ward’s [10] original report (as of 21st May 2021). Furthermore, in the education  
94 literature, Danaher and Crandall’s [9] interpretation is often cited without mention of Stricker  
95 and Ward [e.g., 14, 15, 16] and, when both authors are cited, the fact that Stricker and Ward  
96 disagreed with Danaher and Crandall is not always highlighted [e.g., 17].

97         Our goal in this paper is to report a conceptual replication of Stricker and Ward’s [10]  
98 investigation of stereotype threat in an authentic high-stakes setting, using the analysis  
99 approach favored by Danaher and Crandall [9]. Such a replication is timely, as since Stricker  
100 and Ward’s [10,11] debate with Danaher and Crandall [9], several researchers have  
101 questioned the reliability of lab-based stereotype threat research. One reason is that attempts  
102 to replicate Spencer et al.’s [5] original lab study have not always been successful [e.g., 18].  
103 Stoet and Geary [19] reviewed 23 replication attempts, finding that only 55% had results  
104 consistent with Spencer et al.’s, and that half of these only had so when the researchers  
105 controlled for participants’ pre-existing mathematics achievement (an analytic choice not  
106 made by Spencer et al.).

107         Flore and Wicherts [6] pointed out that the lab-based literature on stereotype threat  
108 and mathematics has an excess of significant findings (more significant results than one  
109 would expect given the average statistical power of published studies). They investigated two  
110 possible reasons. First, earlier researchers may have engaged in *p*-hacking, by using  
111 questionable research practices (such as selectively including covariates) to obtain significant  
112 effects [20]. Second, the literature may be subject to publication bias, a phenomenon where  
113 articles which report significant results are more likely to be accepted for publication than  
114 those which do not [18]. Flore and Wicherts’s [6] meta-analysis of 47 lab studies that  
115 investigated stereotype threat and mathematics achievement found that publication bias might  
116 have “seriously distorted” the meta-analytic effect size estimate they derived from the  
117 literature. However, they found that questionable research practices such as *p*-hacking were

118 not, on their own, sufficient to have created the effect. They left open the possibility that a  
119 combination of publication bias and questionable research practices may be present in the  
120 literature.

121 In sum, there is now some doubt about the reliability of the lab-based literature on  
122 stereotype threat. While lab studies, on average, report small effects in the same direction as  
123 Spencer et al.'s [5] original experiment, it is unclear whether this effect is robust or an  
124 artefact of publication bias. If stereotype threat effects cannot be robustly found in well-  
125 controlled lab studies, it seems unlikely that they could be found in authentic contexts such as  
126 real-world high-stakes tests.

127 The particular history of Danaher and Crandall's [9] analysis means that the two  
128 factors identified by Flore and Wicherts [6], publication bias and *p*-hacking, are unlikely to  
129 apply. Stricker and Ward's [7] original article was published despite the authors claiming that  
130 they had found no effect. This means that Stricker and Ward did not have the usual  
131 motivation for *p*-hacking (finding significant effects where none exist), as although they  
132 reported a statistically significant *p* value, they interpreted it as indicating an effect so small  
133 as to be practically unimportant. Similarly, if publication bias had played a role then, given  
134 its negative conclusions, Stricker and Ward's original paper would not have made it into  
135 print. For these reasons, it seems unlikely that Flore and Wicherts's concerns about the  
136 general literature apply to this particular article.

137 In sum, we believe that Stricker and Ward's [7] study merits replication. It is an  
138 influential study about a topic of societal importance which is often cited in policy debates  
139 and the education literature [e.g., 14, 15, 16, 17, 21]; the wider literature has cast a degree of  
140 doubt about the reliability of the theoretical mechanism proposed to underly the effect; and  
141 the dispute between Stricker and Ward [7,10] and Danaher and Crandall [9] about how to

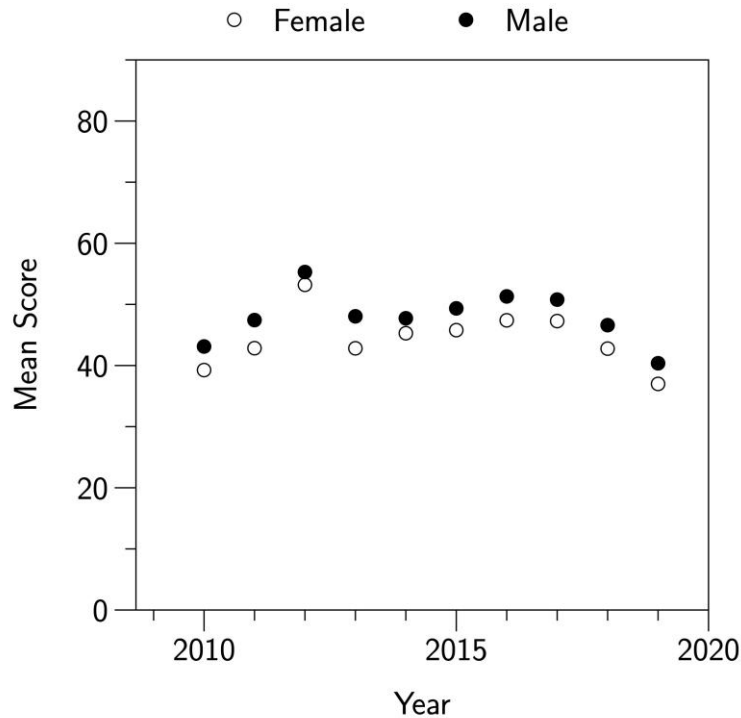
142 legitimately interpret the study's findings, means that if the observed effect were a false  
143 positive, it is unlikely to be due to publication bias or questionable research practices.

## 144 **The Context of the Study**

145 Our study took place in the context of the UKMT Junior Mathematical Challenge  
146 (JMC), a 60-minute, multiple-choice competition consisting of 25 mathematical problems.  
147 The UKMT describes the challenge as encouraging “mathematical reasoning, precision of  
148 thought, and fluency in using basic mathematical techniques to solve interesting problems.” It  
149 is aimed at students across the UK, particularly those aged 11-13. Exactly which young  
150 people participate is decided by their schools; some enter whole year groups and others enter  
151 only certain classes. The questions on the JMC are intended to be accessible to students  
152 studying the English National Curriculum (and the other national curricula of the UK). The  
153 top-scoring 40% of participants are awarded Bronze, Silver and Gold certificates in the ratio  
154 3:2:1. The JMC is also the first round of a suite of problem solving competitions for students  
155 of this age. The roughly 1200 highest scorers are invited to participate in the Junior  
156 Mathematical Olympiad and the next highest 8000 or so are invited to participate in the  
157 Junior Kangaroo.

158 The JMC has a history of gender differences in outcomes, which have not been  
159 satisfactorily explained. Fig 1 shows the average scores for male and female participants  
160 since 2010. Across this period male participants scored between 0.10 and 0.24 standard  
161 deviations higher than female participants (mean 0.19). Because participants in the JMC are  
162 asked to state their demographic information, including gender, prior to answering questions  
163 we hypothesized that the stereotype threat effect discussed by Stricker and Ward [7,10] and  
164 Danaher and Crandall [9] might be contributing to these disparities.





165

166 **Fig 1. The mean scores of male and female participants in the JMC between 2010 and 2019.**

167 Because of the large samples (mean  $N/year = 247,566$ ), error bars are not visible.

168 Clearly the JMC is a different context to the AP Calculus AB examination discussed  
 169 by Stricker and Ward [7,10] and Danaher and Crandall [9]. We return to this issue, and  
 170 outline more precisely the differences between the respective contexts, in the discussion.

## 171 **Participants and Procedure**

172 Of the 2642 schools scheduled to take part in the 2019 JMC as of 14th December  
 173 2018, the largest 13 state schools (by expected number of entries) were approached by the  
 174 UKMT to take part in the study and 6 agreed. We stopped inviting schools to participate after  
 175 6 agreed, as we felt this would provide adequate power to detect the hypothesized interaction  
 176 (see sensitivity analysis below). The final sample consisted of students from 5 English state  
 177 schools, because one school administered the test on the wrong day (and so their students  
 178 were not eligible for JMC certificates, as they may have had prior access to the questions).

179 Four of the five remaining schools were coeducational, one only taught girls. Three of the  
180 schools were selective, two were not. Students at the schools were entered in the JMC in the  
181 normal way, and took part on the same day as participants at other schools. Three schools  
182 (two selective) entered all their students in Years 7 and 8, two entered a selection of higher-  
183 attaining students to participate, based on their own assessment data.

184 The JMC is administered and invigilated by teachers (or other members of staff) in  
185 their own schools, typically in a school hall. The UKMT sends detailed instructions to each  
186 participating school and requires that these are followed. Our instructions asked invigilators  
187 to explain to students that they had the opportunity to contribute to a research study designed  
188 to help the UKMT improve its competitions in the future, and that their answer sheet may be  
189 different to those around them. No information was given to participants about the purpose of  
190 the study in advance.

191 Two different versions of the (optical mark recognition) answer sheet were distributed  
192 to schools. The versions were supplied to schools in a random order, and teachers were asked  
193 to ensure that they were handed out in order to students according to where they were sitting  
194 (which followed the school's normal policy for examinations). This ensured that our two  
195 experimental groups were formed randomly.

196 Each answer sheet had four sections. Section 1 of the gender-first version asked  
197 participants to state their first name, surname and gender; Section 2 asked them to give their  
198 answers to the JMC questions (all multiple choice questions supplied on a separate sheet);  
199 Section 3 asked them to state their school's name and their year group; Section 4 asked  
200 participants to state whether or not they gave consent for their data to be included in the  
201 analysis. The gender-last version of the answer sheet switched the order of Sections 1 and 3.  
202 Sections 1 and 3 were always on different sides of the answer sheet, and participants were  
203 instructed not to turn over their answer sheets. We used the UKMT's standard phrasing for

204 all questions, which in the case of gender was “Please indicate your gender by putting a solid  
205 line through one of the options”, with the options being “female”, “male” and “unspecified”  
206 in that order. (The wording of the questions used to gather demographic information on the  
207 JMC has changed over the years. In the early years of the competition participants were  
208 asked for their “sex” with three options: female, male, unspecified. In the past 20 years the  
209 wording was changed to ask participants for their “gender” but with the same three options.  
210 Since we conducted the study the UKMT have changed the text of the question used to ask  
211 for students’ gender. It now reads: “GENDER (optional): [ ] Boy or young man, [ ] Girl or  
212 young woman”, and a third option with a blank box to allow participants to self-describe their  
213 gender. For clarity, in the remainder of the paper, we use the terms “male” and “female” to  
214 refer to participants who identified themselves as such on their answer sheet.) Participants  
215 were given five minutes to complete Section 1, one hour for Section 2, and five minutes for  
216 Sections 3 and 4.

217 As noted above, in Section 4 test takers were asked whether or not they gave consent  
218 for their data to be used in the study. All those who explicitly refused consent ( $N = 265$ ), or  
219 who failed to answer this question ( $N = 94$ ) were omitted from the analysis, but were  
220 nevertheless eligible for achievement certificates from the UKMT. Importantly, there was no  
221 significant association between whether or not participants consented for their data to be  
222 analyzed, and which answer sheet they had received (gender-first or gender-last),  
223  $\chi^2(1) = 1.351, p = .245$ , suggesting that different consent levels between conditions was not a  
224 threat to the validity of our results.

225 As preregistered, after each school had returned their students’ answer sheets, they  
226 were asked to confirm that they followed our instructions to the letter. Two schools provided  
227 a list of 44 scripts from students who did not follow the instructions (e.g., they turned over  
228 their answer sheet during the Section 1 phase). Although our preregistration had only

229 anticipated excluding entire schools where the invigilation had not proceeded in line with the  
230 instructions, we felt it appropriate to exclude these 44 participants. We did not preregister any  
231 other exclusions, but it was necessary to exclude 23 participants who did not report their  
232 gender (20 did not answer the question, 3 chose “unspecified”). This left 1169 participants:  
233 719 females and 450 males. A sensitivity analysis indicated that this sample size gave us 80%  
234 power to detect a gender by answer-sheet version (gender-first/gender-last) of  $\eta^2 = 0.0065$   
235 [22].

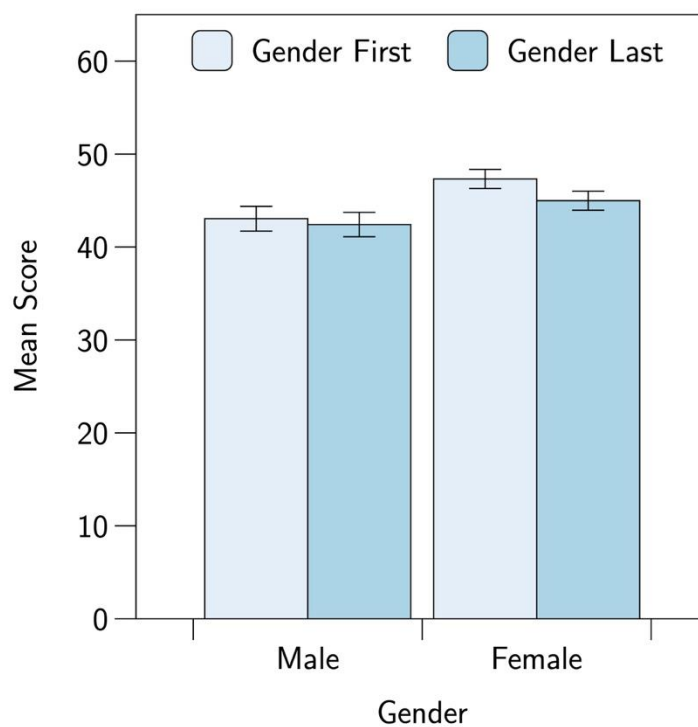
236 The study was approved by the Loughborough University Ethics (Human  
237 Participants) Subcommittee (reference R18-P140). Our analysis plan was preregistered at  
238 AsPredicted.org (#21450) prior to data collection, and can be inspected at  
239 <https://aspredicted.org/hf3jc.pdf>. The 2019 JMC examination paper, the two experimental  
240 response forms, the response form used by challenge participants not involved in the study,  
241 the invigilator’s instructions, and the raw data and analyses scripts are available at  
242 <https://doi.org/10.17605/OSF.IO/UMJ4H>.

## 243 **Results: Preregistered Analyses**

244 To create our dependent variable we used the standard, and longstanding, JMC  
245 scoring system, which awarded 5 points for each correct answer to the first 15 questions, 6  
246 points for each correct answer to the last 10 questions. One point and 2 points were deducted  
247 for incorrect answers to questions 15-20 and questions 21-25 respectively subject to a  
248 minimum total score of zero. Participants’ scores varied from 0 to 110, and their mean,  $M =$   
249 44.9,  $SD = 19.7$ , was slightly higher than the overall average,  $N = 251,064$ ,  $M = 38.71$ ,  $SD =$   
250 19.2,  $t(1168) = 10.7$ ,  $p < .001$ ,  $d = 0.312$ .

251 Participants’ mean scores, split by answer-sheet version and gender, are shown in  
252 Fig 2. As stated in our preregistration, these scores were subjected to a 2 (version) by 2  
253 (gender) between-subjects Analysis of Variance (ANOVA). This revealed a significant main

254 effect of gender,  $F(1,1165) = 8.410, p = .004, \eta^2 = .007$ , which reflected that female  
255 participants had a higher mean score than male participants, 46.2 versus 42.7,  $d = 0.177$ .  
256 There was no significant main effect of version,  $F(1,1165) = 1.586, p = .208, \eta^2 = .001$ ,  
257 (means 45.7, 44.0,  $d = 0.91$ . Crucially, we did not find the hypothesized version-by-gender  
258 interaction effect,  $F(1, 1165) = 0.525, p = .469, \eta^2 = .000$ . Indeed, contrary to the prediction  
259 of the stereotype threat account, female participants in the gender-first condition had slightly  
260 (but non-significantly) higher scores than those in the gender-last condition, 47.3 v 45.0,  
261  $t(717) = 1.61, p = .108, d = 0.120$ .



262  
263 **Fig 2. A plot showing participants' mean scores, split by gender and answer sheet version.** Error  
264 bars show  $\pm 1$  SE of the mean.

265 To be consistent with Stricker and Ward [7], our primary preregistered analysis  
266 involved performing an ANOVA. However, we also ran a multilevel analysis to take account  
267 of between-school variation. We compared models with (i) random intercepts (where  
268 intercepts were able to vary across schools) and (ii) random intercepts and slopes (where both

269 intercepts and slopes were able to vary across schools). Allowing intercepts to vary across  
270 schools yielded a significantly better fit ( $BIC = 10037.65$ ) than a model where intercepts  
271 were identical across schools ( $BIC = 10298.66$ ),  $\chi^2(2) = 268.1$ ,  $p < .001$ . However, allowing  
272 slopes to vary across schools did not significantly improve the fit of the model that included  
273 version, gender and the version-by-gender interaction (both  $BICs = 10047.86$ ),  $\chi^2(2) = 0.00$ ,  
274  $p = 1$ . In this model (in which gender was coded 0 for males, 1 for females; and version was  
275 coded 0 for gender-first and 1 for gender-last), neither version,  $b = -0.902$ ,  $t(1161) = -0.551$ ,  $p$   
276  $= .582$ , nor gender,  $b = -2.99$ ,  $t(1161) = -1.863$ ,  $p = .063$ , nor the version-by-gender  
277 interaction effect,  $b = -1.04$ ,  $t(1161) = -0.499$ ,  $p = .618$ , were significant predictors of  
278 participants' scores (intercept  $b = 46.95$ ,  $t(1161) = 8.96$ ,  $p < .001$ ). In sum, analyzing the data  
279 in this fashion again provided no evidence of the hypothesized version-by-gender interaction.

280         Finally, we conducted a preregistered Bayesian version of our main ANOVA [24].  
281 This required us to specify a model for the null hypothesis. As specified in our  
282 preregistration, we ran two analyses, with Cauchy prior widths of 0.2 and 0.5. Both analyses  
283 provided strong support for the model that only included gender as a predictor over the model  
284 which captured the predicted stereotype threat effect (i.e. the model which included gender,  
285 version and the version-by-gender interaction effect),  $BF_{01s} = 8.123$ , 46.052 respectively.

## 286 **Results: Exploratory Analyses**

287         To explore whether our inclusion of a single-gender school in the sample effected the  
288 results (perhaps, for example, students at single-gender schools are not as affected by societal  
289 stereotypes as those at coeducational schools [cf. 23]), we conducted an exploratory analysis  
290 with the 329 participants from this school omitted. This resulted in an essentially identical  
291 pattern of results. In particular we again found no significant version-by-gender interaction  
292 effect,  $F(1,836) = 0.059$ ,  $p = .809$ ,  $\eta^2 = .000$ .

293 To explore whether or not our decision to use the standard method of scoring the JMC  
294 affected the results, we conducted the primary ANOVA analysis again using (i) number of  
295 problems answered correctly and (ii) percentage accuracy (number of problems answered  
296 correctly as a percentage of problems attempted) as dependent variables. Our primary  
297 conclusions remained for both these dependent variables. Specifically, neither version-by-  
298 gender interaction effects with these two dependent variables was significant: number  
299 correct,  $F(1,1165) = 0.253$ ,  $p = .615$ ,  $\eta_p^2 = .000$ ; percentage accuracy,  $F(1,1165) = 0.326$ ,  $p =$   
300  $.568$ ,  $\eta_p^2 = .000$ .

301 In sum, we found no evidence in favor of the hypothesis that female participants  
302 scored lower when they received the gender-first version of the answer sheet in any of our  
303 analyses, and a Bayesian analysis provided strong evidence against this hypothesis.

## 304 **Discussion**

305 Danaher and Crandall's [9] reanalysis of Stricker and Ward's [7] study of stereotype  
306 threat in real-world conditions suggested that moving the location of demographic questions  
307 in test answer sheets could contribute to reducing the impact of stereotype threat on women's  
308 performance. We conducted a conceptual replication of Stricker and Ward's study in the  
309 context of the UKMT JMC and found no evidence for this effect.

310 Are there any reasons to doubt the validity of our findings? There are at least two  
311 ways in which our sample could be said to be unrepresentative. First, the students in our  
312 sample had significantly higher scores than the overall average in the 2019 JMC. Second,  
313 although there was a significant difference in overall male and female performance in the  
314 2019 JMC (shown in Fig 1), in our sample we found the reverse pattern: our female  
315 participants significantly outperformed our male participants. We discuss each of these  
316 factors in turn.

317           Could the relatively high performance, in comparison with the wider population, of  
318 students in our sample account for us not having found the hypothesized stereotype threat  
319 effect? We think not. In fact, some earlier findings suggest that this factor should *increase* the  
320 size of the effect. For example, Pronin, Steele and Ross [25] found that female students with  
321 higher calculus GPAs and who identified more strongly with mathematics, were more  
322 impacted by a stereotype threat manipulation than those with lower calculus GPAs, or lower  
323 levels of mathematics identification. In other words, because our sample consisted of  
324 relatively successful mathematics students, we might expect the stereotype threat to be  
325 magnified.

326           Might the small female advantage found in our sample, compared to the small male  
327 advantage found nationally, account for the lack of a stereotype effect in our data? Again, we  
328 doubt this. This difference was driven by the inclusion of a high achieving girls-only school  
329 in our sample. This school had the highest mean score of any which participated. When this  
330 school was excluded from our analysis, we found a small male advantage consistent with the  
331 overall picture,  $t(838) = 2.391, p = .017, d = 0.165$ . As noted above, our substantive  
332 conclusions remain if the analysis is conducted on this restricted sample ( $N = 840$ ).

333           So what accounts for the difference between our results and those reported by  
334 Danaher and Crandall [9]? Does this failure to replicate indicate that their finding was a false  
335 positive? Not necessarily. Certain factors are theorized to increase the effect of stereotype  
336 threat [26]. Specifically, stereotype threat is thought to be more prominent (i) when the test is  
337 more difficult [27, 28]; (ii) when stereotyped groups show higher levels of domain  
338 identification with the content of the test [28]; (iii) when (in the context of  
339 mathematics/gender stereotype threat) participants have higher levels of mathematics anxiety  
340 [17, 29]; and (iv) when participants consider membership of the stereotyped groups to be an



341 important part of their identity [30]. What can be said about these factors in the context of our  
342 study?

343 *Test difficulty.* As noted above, our sample scored slightly above the overall average  
344 for the 2019 JMC. However, we should not conclude from this that the test was not  
345 challenging for our sample. The mean number of correct answers offered by participants was  
346 9.5 (from a total of 25 questions), and the mean accuracy (expressed as a percentage of  
347 attempted questions) was 49%. These figures suggest that it is likely that our sample found  
348 the 2019 JMC paper to be challenging.

349 *Domain identification.* We do not have direct evidence about the extent to which our  
350 participants identified with mathematics as a domain, or the extent to which they were  
351 anxious about mathematics. However, the structure of our dataset does allow us to approach  
352 these questions, albeit indirectly. Specifically, as noted in the introduction, some schools  
353 enter their entire cohort into the JMC while others restrict entry to pupils in classes for higher  
354 attaining students (the large majority of English schools teach mathematics in attainment  
355 groups). In our sample, two of the five participating schools restricted entry in this fashion.  
356 On the assumption that female participants from high-prior-attainment classes were more  
357 likely to have higher domain identification than female participants from low-prior-  
358 attainment classes, we repeated our primary analysis, restricted to participants from just these  
359 two schools ( $N = 660$ ). There was again no significant gender by version interaction,  $F(1,$   
360  $656) = 0.00, p = .996, \eta^2 = .000$ . Critically, the female participants from these schools had  
361 very similar mean scores for the two version conditions: gender-first mean 44.4, gender-last  
362 mean 43.1,  $t(313) = 0.661, p = .509, d = 0.07$ . In sum, when our analysis was restricted to a  
363 subsample that we might expect to show higher levels of domain identification we again  
364 found no evidence of the hypothesized stereotype threat effect. Of course, prior attainment is

365 not a perfect proxy for domain identification, so this analysis should not be considered  
366 definitive.

367 *Mathematics anxiety and gender identification.* Unfortunately we do not have any  
368 evidence concerning the extent to which the female participants in our study were anxious  
369 about mathematics, or whether they considered gender to be an important part of their self-  
370 identity, so cannot speak to these moderators.

371 There are also differences between the context of our study and the context of Stricker  
372 and Ward's [7] that must be highlighted. Indeed, Schoenfeld [31] pointed out that direct and  
373 conceptual replications can serve different purposes. He argued that while direct replications  
374 help us guard against 'statistical accidents' (Type I errors where significant results are  
375 obtained in the sample despite there being no true effect in the population), a main goal of  
376 conceptual replications is to help guard against results that do not generalize beyond their  
377 initial contexts. Our study was a conceptual replication, not direct replication, and our context  
378 differs from the original in significant ways. We note five: (i) JMC participants are younger  
379 than AP Calculus AB participants (early high school rather than late); (ii) the JMC takes  
380 place in the UK rather than the US; (iii) our study took place in 2019 rather than 1996; (iv)  
381 the JMC is likely a lower stakes examination than the AP Calculus AB examination; and (v)  
382 the ethnic makeup of our sample may have been different to Stricker and Ward's.

383 Might one of these factors have caused our failure to replicate? We doubt either of the  
384 first two are responsible. First, the stereotype threat effect has been observed in young  
385 children as well as college students [e.g., 32], so there is no reason to suppose that early high  
386 school students would not be affected. Second, analyses of cultural differences in implicit  
387 associations between gender and science have found that the UK and US have similar  
388 profiles [33], suggesting that both countries have similar societal stereotypes with respect to  
389 gender and mathematics.

390 Drawing on evidence from draw-a-scientist studies, Miller, Nolla, Eagly and Uttal  
391 [34] reported evidence that the children's gender-science stereotypes may have reduced over  
392 time, at least in the US. This result perhaps lends credence to the possibility that the fact that  
393 our study took place two decades after Stricker and Ward's [7] might account for the  
394 different results. Perhaps gender-mathematics stereotypes are simply not as strong as they  
395 were in the 1990s. If this account were correct we would indeed expect to see reduced  
396 stereotype threat effects.

397 The fourth difference noted above concerns the extent to which the JMC can be  
398 considered a high-stakes setting. There are several reasons to suppose that participants do  
399 regard the JMC as being an important test. First, the test takes place in formal examination  
400 conditions (i.e. in an examination hall outside of normal classes). Second, the JMC is a  
401 national competition that is used to award certificates and select participants for subsequent  
402 competitions. Third, schools often enthusiastically highlight strong performances on the JMC  
403 by their students to parents, and experience suggests that students commonly discuss their  
404 participation in the JMC on the 'personal statement' section of their applications to study at  
405 university. Nevertheless, it seems plausible to suppose that the stakes involved in the JMC  
406 are not as high as the AP Calculus AB examination, which is used by some students to gain  
407 college credit.

408 Finally, it is worth noting that – in line with the UKMT's normal practice – we did  
409 not enquire about the ethnicity of the participants in our sample. In Stricker and Ward's [7]  
410 study there was a significant version-by-gender-by-ethnicity interaction (for the formula  
411 score dependent variable on the AP Calculus AB examination), with some suggestion that the  
412 reduction in performance in the gender-first condition compared to the gender-last condition  
413 was greater for black female participants than white female participants. In the absence of  
414 ethnicity data for our sample we are unable to explore this potential factor further.

415           Despite these considerations, the other possibility highlighted by Schoenfeld’s [31]  
416 discussion – that the original result could be a ‘statistical accident’ – should not be dismissed.  
417 The disagreement between Stricker and Ward [7,10] and Danaher and Crandall [9] about the  
418 correct interpretation of the original findings, coupled with subsequent growing doubts about  
419 how robust the stereotype threat effect is more generally, suggests that we should take  
420 seriously the possibility that Danaher and Crandall’s conclusions were based on a false  
421 positive.

422           While we cannot definitively determine whether our results differ from Danaher and  
423 Crandall’s [9] because of a ‘statistical accident’, or because our context was different, we can  
424 conclude, in line with Stricker and Ward [10], that the ‘silver bullet’ intervention of  
425 relocating demographic questions on test answer sheets is unlikely to be effective in all  
426 contexts. Instead, we suggest that much more deep-seated interventions are required if we are  
427 to successfully address systemic gender inequities in mathematics.

## 428 **References**

- 429 1. Fenema E. Mathematics learning and the sexes: A review. *Journal for Research in*  
430 *Mathematics Education*, 1974;5:126-139.
- 431 2. Ellis J, Fosdick BK, Rasmussen C. Women 1.5 times more likely to leave STEM  
432 pipeline after calculus compared to men: Lack of mathematical confidence a potential  
433 culprit. *PLOS ONE*. 2016;11:e0157447.
- 434 3. Schmader T, Johns M. Converging evidence that stereotype threat reduces working  
435 memory capacity. *Journal of Personality and Social Psychology*. 2003;85:440-452.
- 436 4. Steele CM, Aronson J. Stereotype threat and the intellectual test performance of  
437 African Americans. *Journal of Personality and Social Psychology*. 1995;69:797–811.
- 438 5. Spencer SJ, Steele CM, Quinn DM. Stereotype threat and women’s math performance.  
439 *Journal of Experimental Social Psychology*. 1999;35:4–28.

- 440 6. Flore PC., Wicherts JM. Does stereotype threat influence performance of girls in  
441 stereotyped domains? A meta-analysis. *Journal of School Psychology*. 2015;53:25-44.
- 442 7. Stricker LJ, Ward WC. Stereotype threat, inquiring about test takers' ethnicity and  
443 gender, and standardized test performance. *Journal of Applied Social Psychology*.  
444 2004;34:665–693.
- 445 8. College Board. Advanced Placement course description, Mathematics, Calculus AB,  
446 Calculus BC – May 1995, May 1996. New York, NY: College Board. 1994.
- 447 9. Danaher K, Crandall CS. Stereotype threat in applied settings re- examined. *Journal of*  
448 *Applied Social Psychology*. 2008;38:1639-1655.
- 449 10. Stricker LJ, Ward WC. Stereotype threat in applied settings re- examined: A reply.  
450 *Journal of Applied Social Psychology*. 2008;38:1656-1663.
- 451 11. Cheung AC., Slavin RE. How methodological features affect effect sizes in education.  
452 *Educational Researcher*. 2016;45:283-292.
- 453 12. Lortie-Forgues H, Inglis M. Rigorous large-scale educational RCTs are often  
454 uninformative: Should we be concerned? *Educational Researcher*. 2019;48:158-166.
- 455 13. Bakker A, Cai J, English L., Kaiser G, Mesa V, Van Dooren W. Beyond small,  
456 medium, or large: Points of consideration when interpreting effect sizes. *Educational*  
457 *Studies in Mathematics*. 2019;102:1–8.
- 458 14. Kricheli-Katz T, Regev T. The effect of language on performance: do gendered  
459 languages fail women in maths? *NPJ Science of Learning*. 2021;9.
- 460 15. Browne RK, Allen PJ, Noam GG. The double-dip: quality discrepancies in out-of-  
461 school time STEM programs. *International Journal of Science Education B*.  
462 2021;11:35-54.

- 463 16. Seo E, Lee YK. Stereotype threat in high school classrooms: how it links to teacher  
464 mindset climate, mathematics anxiety and achievement. *Journal of Youth and*  
465 *Adolescence*. 2021;50:1410-1423.
- 466 17. Maloney EA, Schaeffer MW, Beilock SL. Mathematics anxiety and stereotype threat:  
467 shared mechanisms, negative consequences, and promising interventions. *Research in*  
468 *Mathematics Education*. 2013;15:115-128
- 469 18. Ganley CM, Mingle LA, Ryan AM, Ryan K, Vasilyeva M, Perry M. An examination of  
470 stereotype threat effects on girls' mathematics performance. *Developmental*  
471 *Psychology*. 2013;49:1886–1897.
- 472 19. Stoet G, Geary DC. Can stereotype threat explain the gender gap in mathematics  
473 performance and achievement?. *Review of General Psychology*. 2012;16:93-102.
- 474 20. Simonsohn U, Nelson LD, Simmons JP. P-curve: A key to the file drawer. *Journal of*  
475 *Experimental Psychology: General*. 2013;143:534-547.
- 476 21. Walton, GM, Spencer SJ, Erman S. Affirmative meritocracy. *Social Issues and Policy*  
477 *Review*. 2013;7:1-35.
- 478 22. Lakens D, Caldwell, AR. Simulation-based power-analysis for factorial ANOVA  
479 designs. 2019. <https://doi.org/10.31234/osf.io/baxsf>
- 480 23. Inzlicht M, Ben-zeev T. A threatening intellectual environment: Why females are  
481 susceptible to experiencing problem-solving deficits in the presence of males.  
482 *Psychological Science*. 2000;11:365–371.
- 483 24. Wagenmakers EJ, Love J, Marsman M, Jamil T, Ly A, Verhagen J, Selker R, Gronau  
484 QF, Dropmann D, Boutin B, Meerhoff F. Bayesian inference for psychology. Part II:  
485 Example applications with JASP. *Psychonomic Bulletin & Review*. 201;25:58-76.

- 486 25. Pronin E, Steele CM, Ross L. Identity bifurcation in response to stereotype threat:  
487 Women and mathematics. *Journal of Experimental Social Psychology*. 2004;40:152–  
488 168.
- 489 26. Flore PC, Mulder J, Wicherts JM. The influence of gender stereotype threat on  
490 mathematics test scores of Dutch high school students: A registered report.  
491 *Comprehensive Results in Social Psychology*. 2018;3:140-174.
- 492 27. Neuville E, Croizet JC. Can salience of gender identity impair math performance  
493 among 7-8 years old girls? The moderating role of task difficulty. *European Journal of*  
494 *Psychology of Education*. 2007;22:307-316.
- 495 28. Nguyen HHD, Ryan AM. Does stereotype threat affect test performance of minorities  
496 and women? A meta-analysis of experimental evidence. *Journal of Applied*  
497 *Psychology*. 2008;93:1314-1334.
- 498 29. Delgado AR, Prieto G. Stereotype threat as a validity threat: The anxiety-sex-threat  
499 interaction. *Intelligence*. 2008;36:635-640.
- 500 30. Schmader T. Gender identification moderates stereotype threat effects on women’s  
501 math performance. *Journal of Experimental Social Psychology*. 2002;38:194-201.
- 502 31. Schoenfeld AH. On replications. *Journal for Research in Mathematics Education*.  
503 2018;49:91-97.
- 504 32. Jordan AH, Lovett BJ. Stereotype threat and test performance: A primer for school  
505 psychologists. *Journal of School Psychology*. 2007;45:45-59.
- 506 33. Nosek BA, Smyth FL, Sriram N, Lindner NM, Devos T, Ayala A, Bar-Anan Y, Bergh  
507 R, Cai H, Gonsalkorale K, Kesebir S. National differences in gender–science  
508 stereotypes predict national sex differences in science and math achievement.  
509 *Proceedings of the National Academy of Sciences*. 2009;106:10593-7.

- 510 34. Miller DI, Nolla KM, Eagly AH, Uttal DH. The development of children's gender-  
511 science stereotypes: a meta- analysis of 5 decades of US draw- a- scientist studies.  
512 Child Development. 2018;89:1943-1955.