

1 **Do small classes in higher education reduce performance gaps in STEM?**

2 Cissy J. Ballen^{1,2}, Stephanie M. Aguilon^{3,4}, Rebecca Brunelli⁵, Abby Grace Drake³, Deena
3 Wassenberg¹, Stacey L. Weiss⁶, Kelly R. Zamudio³, and Sehoya Cotner¹

4 ¹Department of Biology Teaching and Learning, University of Minnesota, Minneapolis, MN

5 ²balle027@umn.edu

6 ³Department of Ecology & Evolutionary Biology, Cornell University, Ithaca, NY

7 ⁴Cornell Lab of Ornithology, Ithaca, NY

8 ⁵Department of Biological Sciences, California State University, Chico, CA

9 ⁶Department of Biology, University of Puget Sound, Tacoma, WA

10

11

12

13

14

15

16

17 Manuscript published as an Article in *Bioscience*

18 Word count: 5,511

19

20 **Abstract**

21 Performance gaps in science are well-documented, and an examination of underlying
22 mechanisms that lead to underperformance and attrition of women and underrepresented
23 minorities (URM) may offer highly targeted means to promote such students. Determining
24 factors that influence academic performance may provide a basis for improved pedagogy and
25 policy development at the university level. We examined the impact of class size on students in
26 17 biology courses at four universities. While female students underperformed on high-stakes
27 exams compared to men as class size increased, women received higher scores than men on non-
28 exam assessments. URM students underperformed across grade measures compared to majority
29 students regardless of class size, suggesting that other characteristics of the education
30 environment affect learning. Student enrollment is expected to increase precipitously in the next
31 decade, underscoring the need to prioritize individual student potential rather than yield to
32 budget constraints when considering equitable pedagogy and caps on classroom sizes.

33

34 **Introduction**

35 Universities face the unique challenge of educating students from increasingly diverse
36 backgrounds who may excel in different educational contexts. Recent efforts to better serve
37 diverse classrooms include changes in instruction such as active learning (Ballen et al. 2017a;
38 Haak et al. 2011) and course-based undergraduate research experiences (Ballen et al. 2017b,
39 Lopatto 2007). To provide effective instructional practices for all, we must continue to identify
40 practical steps to promote the success of qualified students from historically underserved
41 demographics in STEM, such as women and underrepresented minority students (African
42 American, Hispanic, Native American, or Pacific Islander; hereafter URM).

43 If our goal is to achieve diversity in STEM, coursework should ideally nurture individual
44 potential rather than ‘weed out’ less prepared students at the start of an undergraduate degree
45 (Koester et al. 2016, Mervis 2011, Suresh 2006). Using 16 years of data from a liberal arts
46 college, Rask and Tiefenthaler (2008) demonstrated that students’ grades influenced their
47 decision to continue within their major. While lower grades led to lower persistence for all
48 students, female students with low grades were more likely to abandon the discipline and pursue
49 a different major than males. A second longitudinal study showed that negative experiences in
50 introductory science courses were cited as the primary reason for declining interests in obtaining
51 a science degree among women and URM students (Barr et al. 2008). Women and URM
52 students also face other well-documented challenges unrelated to academic competency, such as
53 discrimination (Grunspan et al. 2016, Milkman et al. 2015, Moss-Racusin et al. 2012, Steele
54 Jennifer et al. 2002), feelings of exclusion (Hall and Sandler 1982, Hurtado and Ruiz 2012),
55 imposter syndrome (Clance 1985), test anxiety (Ballen et al., 2017) and stereotype threat
56 (Schmader 2002, Steele 1997, Steele and Aronson 1995). All of these contribute to the well

57 documented higher attrition rates of women and URM students across STEM disciplines
58 (Alexander et al. 2009, Ballen and Mason 2017, Beede et al. 2011, Eddy et al. 2014, May and
59 Chubin 2003) and university campuses (Anderson and Kim 2006, Griffith 2010, Olson and
60 Riordan 2012, Smith 2000). Education research has also identified examples of learning contexts
61 that counteract the psychosocial barriers faced disproportionately by women and URM students,
62 including opportunities to interact with role models in and out of the classroom (Fried and
63 MacCleave 2009, Stout et al. 2011), interventions in social belonging (Walton et al. 2015), peer
64 mentoring (Snyder and Wiles 2015), and for females, schools with higher percentages of female
65 STEM graduate students (Griffith 2010). Thus, it is essential we identify obstacles that
66 specifically affect underrepresented students as a means of finding interventions that promote all
67 students' success in STEM.

68 Class size, an often overlooked variable, is worthy of careful consideration because
69 previous research suggests it influences student performance (Glass 1982, Ho and Kelman 2014,
70 Kokkelenberg et al. 2008) and, unlike other variables, is subject to legislative action. At least 24
71 states have mandated or incentivized class size reduction in American K-12 classrooms
72 (Whitehurst and Chingos 2011). At the undergraduate level, universities are constantly faced
73 with decisions on how to allocate faculty time to best serve their undergraduate population.
74 Recent changes in course content delivery – e.g., the rise of online classes (such as massive open
75 online courses or MOOCs) and hybrid online courses – are the direct result of an increased
76 demand for access to education (Kena 2016). The imminent growth in enrollment to degree-
77 granting institutions (Kena 2016) underscores the urgent need to quantify the effects of class
78 sizes on undergraduate students. Here, using data from 17 biology courses at four institutions, we

79 examine the extent that class size impacts achievement gaps for female and URM
80 undergraduates.

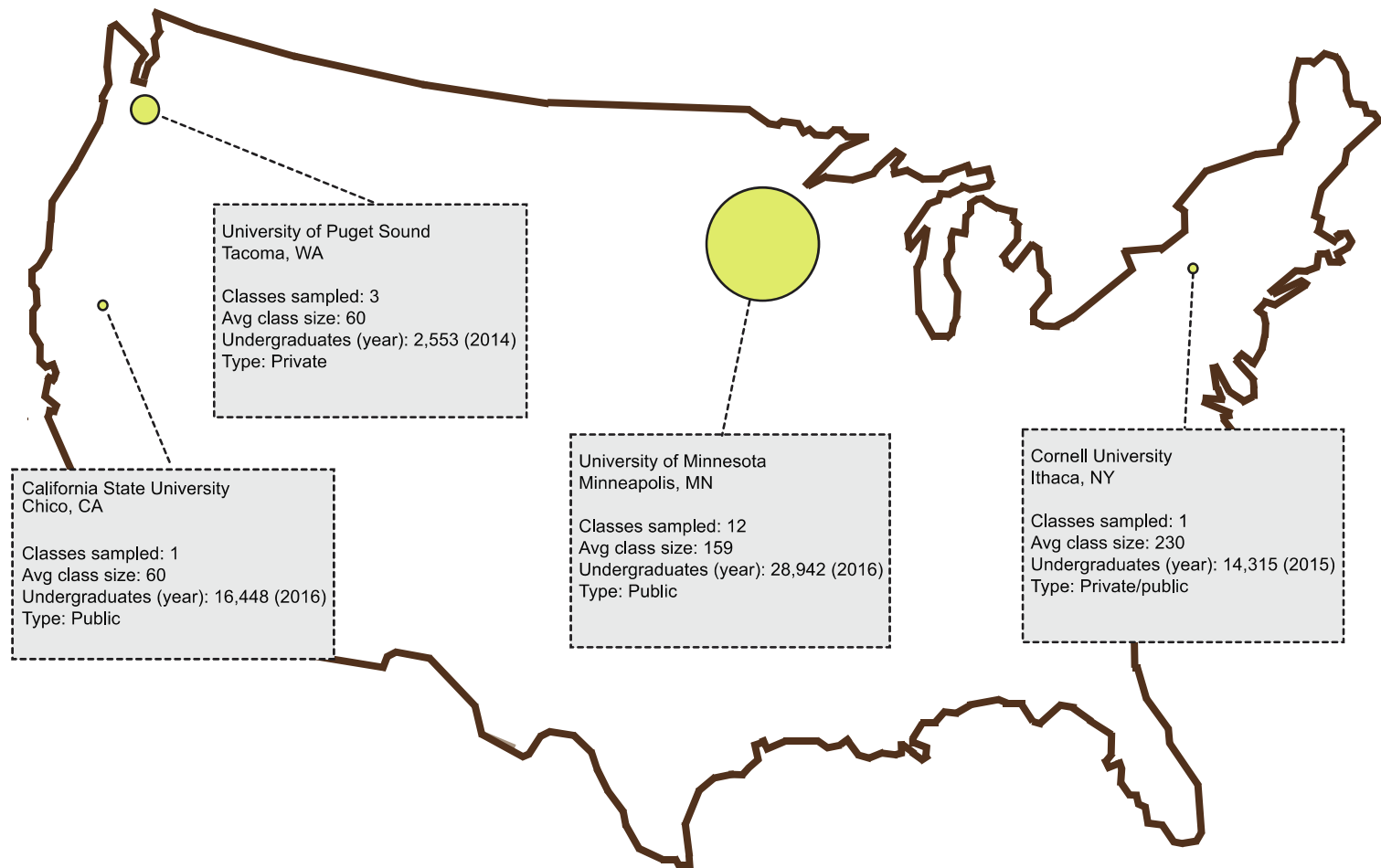
81

82 We address three questions by focusing on performance gaps between male and female students,
83 and URM and majority students: 1) Does class size influence performance on exams? 2) Does
84 class size influence performance on non-exam methods of assessment? 3) Does class size
85 influence final course grade?

86

87 **Data collection**

88 Administrative data were obtained from 17 lower division biology courses taken by 1836
89 students in fall 2016 (minimum class size $N = 40$, maximum $N = 239$; Figure 1). To establish a
90 collaborative research group, we solicited participation through an existing professional network
91 from biology instructors who teach majors or nonmajors from a diverse range of institutions, and
92 received data from California State University, Chico; Cornell University; University of
93 Minnesota, Twin Cities; and University of Puget Sound. We compared (1) pooled exam grades,
94 (2) pooled assessments of student knowledge other than exams (hereafter non-exam grades; e.g.,
95 discussion sections, laboratories, online activities, written assignments, low-stakes quizzes, as
96 well as active learning in-class activities), and (3) final course grades, which reflect cumulative
97 performance in all aspects of the course. We present analyses with transformed z-scores (a
98 measure of how many standard deviations a value is from the class section's mean score) for
99 ease of interpretation.



100

101

102 **Figure 1.** Four universities participated in the current study, representing diverse geographic locations across the US. Circle sizes are
 103 proportional to the number of classes sampled from each institution.

104

105 **Statistical Analyses**

106 *Linear mixed-effects model*

107 We used linear mixed-effects models to compare exam performance, performance on non-exam
108 assessments and total course performance, across the four universities. The data in this study are
109 hierarchically nested because a student's exam performance is likely to be more similar to a
110 classmate's performance than a student outside of their class, as students in the same class share
111 the same assessments (Kreft et al. 1998). Similarly, students in biology classes at one university
112 may perform or be assessed in the same way as students in biology classes at another university.
113 For this reason, we use multilevel modeling to account for the non-independence of data in
114 nested-data structures (Kreft et al. 1998, Paterson and Goldstein 1991).

115 Akaike's information criterion (AIC) was used to determine model fit in a multimodel
116 inference technique. AIC estimates the goodness of fit of each model given our sample (Akaike
117 1974), and allows us to rank models based on this estimation using AIC differences ($\Delta_i =$
118 $AIC_{\text{model } i} - \text{minAIC}$, where minAIC is the model with the smallest AIC value). Models with a Δ_i
119 > 10 are considered poor predictors compared to the best model, and so we only present results
120 with small Δ_i values for brevity (Table S1). We were interested in the interaction of class size
121 with gender (SGender, a factor with two levels) and with URM status (a factor with two levels).
122 Therefore, our model initially included those three main effects (SGender, URM status, and class
123 size) and two interaction effects (SGender*class size, URM status*class size).

124 In addition, we tested whether the following variables improved the fit of the model for
125 the given set of data: (1) an interaction between student gender identity and URM status
126 (SGender*URM status); (2) instructor gender identity (IGender, a factor with three levels

127 including female, male, or multiple instructor genders: in other words, more than one instructor
 128 for the course in question who did not identify as the same gender); (3) an interaction between
 129 student gender identity and instructor gender identity (SGender*IGender); (4) an interaction
 130 between student gender identity, URM status, and class size (SGender*URM status*class size);
 131 (5) age. Only students with a complete set of these variables were included in these analyses. All
 132 models included random effects for university, class ID (nested within university), and instructor
 133 ID (nested within classes and university). Random effects were tested for significance by
 134 removing one random factor at a time and taking the difference between the -2 log likelihoods.
 135 This was tested against a chi-square distribution with one degree of freedom (per removed
 136 random factor). Instructor ID was removed from the analysis as a random effect.

137 We explored all possible models and chose the most parsimonious model that best fit the
 138 data in accordance to AIC model-selection statistics (Table 1). The AIC estimates indicated that
 139 the elimination of the URM*class size interaction resulted in better fit models, and so the
 140 interaction was backwards eliminated from the final models ($P > 0.25$; see results). We used
 141 Bonferroni corrected post-hoc pairwise comparisons to clarify performance outcomes of students
 142 based on gender and URM status. We performed all statistical analyses using SPSS software
 143 version 24 (SPSS Inc., Chicago, IL, USA).

144

Rank	Model: Combined exam grades	AIC	Δi	Relative likelihoods	w_i
1	URM status + class size + SGender + class size*Sgender	4885.468	0.000	1.000	0.935
2	URM status + class size + Sgender + class size*Sgender + age	4891.961	6.493	0.039	0.036
3	URM status + class size + Sgender + class size*Sgender + Sgender*URM status + age	4892.347	6.879	0.032	0.030

145

Rank	Model: Non-exam grades	AIC	Δi	Relative	w_i
------	------------------------	-----	------------	----------	-------

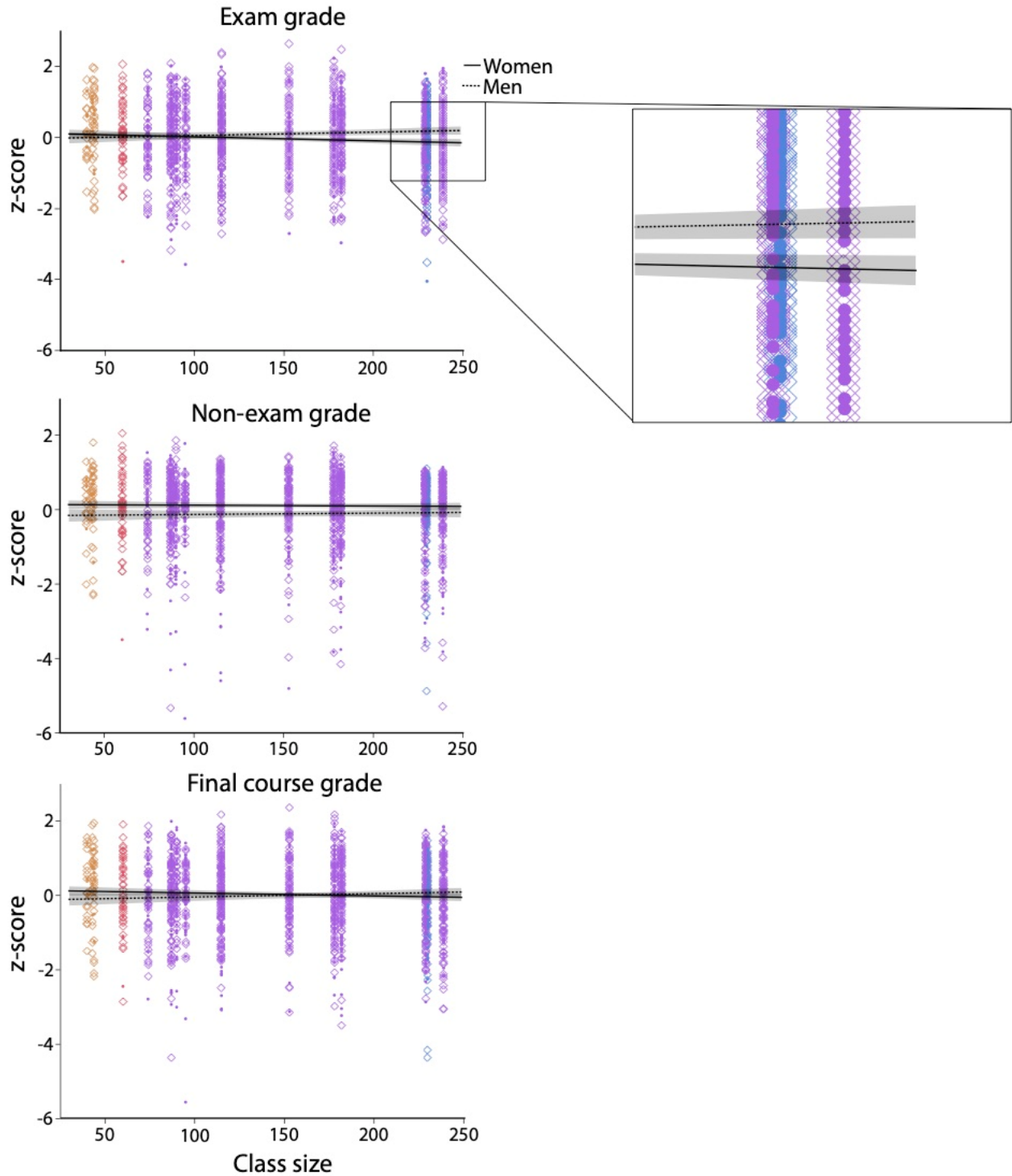
				likelihoods	
1	URM status + class size + SGender + class size*Sgender	4835.231	0.000	1.000	0.885
2	URM status + class size + SGender + class size*Sgender age	4839.840	4.609	0.100	0.088
3	URM status + class size + SGender + class size*Sgender + class size* URM status + age	4842.562	7.331	0.026	0.023

146

Rank	Model: Final course grade	AIC	Δi	Relative likelihoods	w_i
1	URM status + SGender	4826.220	0.000	1.000	0.926
2	URM status + class size + SGender	4831.668	5.448	0.066	0.061
3	URM status + class size + Sgender + age	4836.220	10.000	0.007	0.006
4	Sgender	4837.260	11.040	0.004	0.004
5	URM status + class size + SGender + class size*Sgender + class size* URM status	4837.736	11.516	0.003	0.003

147 **Table 1.** Best models for predicting performance metrics across four universities using AIC
148 model selection. Compared to the first model, models with an $\Delta i > 10$ are considered poor
149 predictors and so we do not report them here. The Akaike weights, w_i , represent probabilities
150 that a given model is the best model under repeated sampling.

151



152

153 **Figure 2.** The effects of class size on exam grade z-scores, non-exam grade z-scores, and final
 154 course grade z-scores for women (solid line) and men (dashed line). Colors represent different
 155 universities: University of Puget Sound (yellow), California State University, Chico (red),
 156 University of Minnesota, Twin Cities (purple), and Cornell University (blue).

157

158

159 **Results**

160 We used mixed model analyses to compare students' combined exam grade, non-exam
161 grade, and total course grade in the fall 2016 semester (Figure 2, Table S1-S3). First, we
162 observed a nonsignificant interaction effect of URM status and class size on metrics of
163 performance.

164 When we removed the interaction from the models, URM status became a significant
165 predictor of performance (Combined exam grade $B = 0.417$, $t(1377) = 6.01$, $P < 0.001$, $SE =$
166 0.069 ; Non-exam grade $B = 0.262$, $t(1533) = 3.83$, $P < 0.001$, $SE = 0.069$; Final course score B
167 $= 0.407$, $t(1522) = 5.87$, $P < 0.001$, $SE = 0.069$). These results suggest that URM students'
168 exam scores on average was 0.42 standard deviation lower than non-URM students, and
169 their non-exam scores were on average 0.26 standard deviation lower than non-URM
170 students. Bonferroni corrected post-hoc pairwise comparisons, presented from the final models,
171 show URM students underperforming on all performance metrics compared to non-URM
172 students (Table 2; hereafter 'underperform' is used to describe raw gaps, and not those for which
173 some measure of student academic ability/preparation is controlled). Second, we observed a
174 significant interaction between gender and class size, such that as class size increased, women
175 underperformed on exams ($SGender * class\ size\ B = -0.145$, $t(1599) = -2.89$, $P = 0.004$, $SE =$
176 0.050 ; Figure 2 inset) and in the course overall ($B = -0.108$, $t(1649) = -2.16$, $P = 0.031$, $SE =$
177 0.050) compared to men. We also found that women obtained higher non-exam grades ($B =$
178 0.217 , $t(1731) = 4.60$, $P < 0.001$, $SE = 0.047$) compared to men, regardless of class size.

179 Next, we explored whether women are underperforming on exams because they are
180 higher stakes in larger classes; i.e., they account for a larger proportion of the grade. To
181 investigate this, we examined the correlation between class size and the percentage of students'

182 final course grades that are from their performance on exams. We did not find a strong
 183 correlation (Pearson correlation = -0.386; $P = 0.126$). This result runs counter to what one would
 184 expect due to the courses included in this sample, and is probably not representative of most
 185 lower division lecture courses, in which exams generally account for a larger proportion of final
 186 course grade (Koester et al. 2016). Finally, to test whether our results are the same within one
 187 institution, we isolated twelve lower division classes from the University of Minnesota that
 188 varied in class size. In these classes, all exams had identical multiple choice format. We found
 189 the same main results across assessment types within one institution as we observed across all
 190 institutions (Tables S4-S6). Thus, as was the case across universities, increasing class size was
 191 negatively correlated with female performance, and URM status significantly predicted
 192 performance outcomes within our most sampled university.

193

Class size	URM			non-URM		
	50	150	250	50	150	250
Combined exam grade	-0.285 (0.08)	-0.295 (0.07)	-0.305 (0.08)	0.140 (0.06)	0.130 (0.05)	0.120 (0.06)
Non-exam grade	-0.165 (0.08)	-0.165 (0.08)	-0.166 (0.09)	0.086 (0.07)	0.086 (0.06)	0.085 (0.07)
Total course grade	-0.262 (0.09)	-0.264 (0.08)	-0.266 (0.09)	0.132 (0.07)	0.130 (0.06)	0.128 (0.07)
<i>N</i>		261			1575	

194

195 **Table 2.** Least-squares means comparison of relative performance of students who differ based
 196 on their racial minority status (URM or non-URM) in different class sizes (50 students, 150
 197 students, or 250 students). Measures are standardized, and reflect performance relative to the
 198 mean of the class; positive scores are students who overperformed in standard deviations from
 199 the mean, and negative scores represent those who underperformed relative to the mean. Our
 200 data indicate that URM students underperform across all metrics, compared to non-URM
 201 students, but unlike female students, their performance is not affected by class size, suggesting
 202 factors other than class size negatively influence URM student performance. Standard errors are
 203 shown in parentheses.

204

205 One possibility is that the positive effects we observe from students in small classes is
206 due to increased active learning and student interactions with the instructor in smaller classes,
207 which may influence student performance (e.g., Ballen et al. 2017, Haak et al. 2014). Using data
208 collected for nine of the seventeen courses (Table S7), we used a linear regression to examine the
209 relationship between class size and total number of student-instructor interactions per class
210 period. Results from the linear regression were not conclusive. First, when we included all of the
211 schools in our analysis we found a significant relationship between the two variables (Figure S1;
212 Pearson correlation = -0.72; $P = 0.028$), such that students interacted more with their instructors
213 in smaller classes. However, when we isolated classes within the University of Minnesota, the
214 correlation was no longer significant (Pearson correlation = 0.24; $P = 0.645$). Class size likely
215 influences the frequency in which students interact with their instructor, and this may be why
216 small class sizes appear to disproportionately benefit women in our sample. Future work will
217 profit from a thorough examination of the relationship between class size, active learning, and
218 performance gaps.

219

220 **Discussion**

221 We compared female and male exam performance, non-exam performance, and total
222 course performance across four universities and found that as class sizes increased, women
223 underperformed on exams and final course grades compared to men in their classes. However,
224 female students outperformed males regardless of class size on non-exam scores that contributed
225 to total course grade. We did not find a similar effect of class size on students based on minority
226 status. Across class size and assessment type, URM students underperformed relative to non-
227 URM students (Table 2).

228 Reasons for the pervasive disparity between URM and non-URM students are likely
229 complex and multifaceted, but may include differences in incoming academic preparation
230 (Ballen and Mason 2017), economic hardship (Cabrera et al. 1992), university campus social
231 climate (Gloria et al. 1999), and low representation in the classroom or discipline (Braxton et al.
232 2011). The underrepresentation of URM individuals in the STEM workforce (Landivar 2013)
233 underscores the urgent need for effective approaches that promote students who are racial or
234 ethnic minorities (Brewer and Smith 2011).

235 While our findings do not suggest tractable solutions to racial disparities in STEM, they
236 do suggest strategies for mitigating gender biases. Specifically, to increase female retention in
237 STEM, we recommend offering smaller classes and emphasizing non-exam points—especially in
238 lower division classes that serve as gateway courses to students’ major field of study. In these
239 gateway courses students are often ‘weeded out’ because students’ perceived or actual academic
240 performance suffers in those environments (Baker et al. 2016).

241 A review by Cuseo (2007) identified five reasons that large classes have adverse effects
242 on some students: (1) fewer opportunities for students to interact with course material, (2) fewer
243 opportunities for students to interact with the instructor, (3) reduced opportunities for instructors
244 to challenge students, (4) lower overall student satisfaction with the learning experience, and (5)
245 lower satisfaction with the instructor according to student evaluations (Cuseo 2007). Future
246 research will benefit from a close examination of the consequences of these factors, and whether
247 they respond to experimental class-size manipulations. We do recognize the reality of budgetary
248 constraints, and the fact that larger classes are often the simplest solution to fiscal crises.
249 However, when large classes are a “necessary evil,” instructors can minimize the negative
250 consequences of large classes via evidence-based interventions. For example, in large-lecture

251 settings students can have more opportunities to interact with lecture material and the instructor
252 via numerous instant-feedback strategies (e.g. the Immediate Feedback Assessment Technique
253 [Cotner et al 2008a], classroom response systems [Cotner et al 2008b, Lewin et al 2016, Knight
254 et al 2016], plicker cards [Howell et al 2017], etc.) and low-stakes—or no-stakes—formative
255 assessments (e.g., one-minute papers, worksheets, and concept maps; Angelo and Cross 1993).

256 Because in our dataset female students excelled at non-exam assessments of the course
257 material regardless of class size, an alternative strategy to promote women in STEM may be to
258 make non-exam scores a larger component of the final course grade (Koester et al. 2016). Recent
259 work shows traditional exams do not accurately capture student mastery of the cognitive skills
260 required to do science and exacerbate existing gaps in performance (Moneta-Koehler et al. 2017,
261 Stanger-Hall 2012). Further, women are adversely affected by test anxiety, which in itself is
262 higher in women than in their male counterparts (Ballen et al 2017). Thus, if our aim is to reward
263 ongoing preparation and cooperative group work rather than performance on a few, high-stakes
264 exams, these assignments will nurture those qualities and work habits in developing scientists.
265 For instructors who teach large classes, the challenge will be to develop scalable assignments
266 that can effectively evaluate students' learning. Despite these challenges, our data show that an
267 effective way for instructors to reduce gender gaps in their classrooms is to experiment with
268 strategies to tailor the learning environment to their student population.

269 Research demonstrating the negative impacts of large classes on students reinforce
270 conceptual arguments against these classes (Achilles 2012, Baker et al. 2016, Glass 1982, Glass
271 and Smith 1979, Ho and Kelman 2014, Schanzenbach 2014), and can inform policy related to
272 education. The state of Minnesota, in which the majority of classes were sampled, has
273 historically taken innovative approaches to improving its schools (Mazzoni 1993). In fact, the

274 state's former Governor, Jesse Ventura, campaigned on an education platform that declared "the
275 best way to solve most of our educational problems is to reduce class size" (Ventura 2000).
276 Nationally, schools aim to keep class sizes low, but according to the National Center for
277 Education Statistics, total enrollment at public and private degree-granting post-secondary
278 institutions is expected to increase 15 percent between 2014 and 2025 (Kena 2016). While it may
279 be tempting to increase the number of students per class section in order to decrease costs, the
280 consequences on student learning and performance must be carefully considered. Note that our
281 classes range in size from 40 to over 200 students. Thus, a class of 50-100 students is associated,
282 in our model, with more equitable performance than is one with 200 or more students; in other
283 words, a "smaller" class is likely still cost-effective. Future work will conduct similar
284 investigations into the effects of class size on students of low-socioeconomic status and first-
285 generation college students.

286 This work has limitations that warrant consideration. First, we were unable to control for
287 incoming student preparation (e.g., pre-course measures such as SAT or cumulative GPA) for all
288 students across universities. Previous work finds that incoming preparation predicts performance
289 and retention across institutions (Ballen and Mason 2017, Ballen et al. 2017, Bonous-Hammarth
290 2000, Easton et al. 2017). However, by normalizing performance across cohorts—we show the
291 achievement gaps in course grades as they are corrected in magnitude. Second, to test the
292 generality of these results it will be important to test a wider range of universities nationally and
293 internationally. While our dataset is subject to some biases, these collaborative efforts among
294 universities allow for much larger datasets--across a broad sample of university types--that
295 would not be possible within one institution. Thus, multi-institution efforts allow for meaningful
296 comparisons, and have considerable potential to illuminate the nature of persistent demographic

297 gaps within classrooms, as well as gaps in institutional representation in the STEM workforce.
298 Finally, many other variables may contribute to student performance that we did not include in
299 our analysis, including teaching strategy (e.g., active or traditional lecturing; Haak et al. 2011),
300 classroom social climate (Crawford and MacLeod 1990, Grunspan et al. 2016), campus social
301 climate (Hall and Sandler 1984), and opportunity for academic support outside of the classroom
302 (e.g., tutorials or peer mentoring; Snyder et al. 2016). Future work will also benefit from a focus
303 on the underlying mechanisms that explain the observed gender gaps in large classes at the
304 undergraduate level.

305 Despite these limitations, we detect an interaction effect between gender and class size,
306 such that women are negatively affected by large class sizes in ways that men are not. These
307 findings add an equity dimension to previous work citing the benefits of smaller classes. This
308 aspect of smaller-class impacts may be especially compelling to administrators, curriculum
309 committees, or legislators who are motivated to eliminate gender gaps in performance that
310 plague higher education.

311

312 **Acknowledgements.** We thank Daniel Baltz for help with data organization and interpretation;
313 J.D. Walker and Lauren Sullivan for statistical support; Gregor Siegmund, Paula Soneral, Brian
314 Wisenden, Daniel Stovall, Michelle Mabry, Steven Karafit, Denise Monti, Danielle Grunzke,
315 Leslie Saucedo, Gregory Johnson, Jorge Tomasevic, and Heather Patterson for help with student
316 data. We received IRB exemption to work with student data at all universities (IRB protocol
317 UMN: 1405E50826; CU: 1410005010; UPS: 1617-006; CSU: 461).

318

319

320

321

322 **References**

- 323 Achilles CM. 2012. Class-Size Policy: The STAR Experiment and Related Class-Size
324 Studies. NCPEA Policy Brief. 1(2): NCPEA Publications.
- 325 Akaike H. 1974. A new look at the statistical model identification. *IEEE transactions on*
326 *automatic control* 19:716-723.
- 327 Alexander C, Chen E, Grumbach K. 2009. How leaky is the health career pipeline? Minority
328 student achievement in college gateway courses. *Academic Medicine* 84:797-802.
- 329 Anderson E, Kim D. 2006. Increasing the Success of Minority Students in Science and
330 Technology, American Council on Education. Washington, DC: American Council on
331 Education.
- 332 Angelo TA, Cross KP. 1993. *Classroom assessment techniques: A handbook for faculty.*
333 Ann Arbor, MI: National Center for Research to Improve Postsecondary Teaching and
334 Learning.
- 335 Baker BD, Farrie D, Sciarra DG. 2016. Mind the Gap: 20 Years of Progress and
336 Retrenchment in School Funding and Achievement Gaps. ETS Research Report Series
337 2016:1-37.
- 338 Ballen CJ, Blum JE, Brownell S, Hebert S, Hewlett J, Klein JR, McDonald EA, Monti DL,
339 Nold SC, Slemmons KE. 2017a. A Call to Develop Course-Based Undergraduate Research
340 Experiences (CUREs) for Nonmajors Courses. *CBE-Life Sciences Education* 16(2): mr2.
- 341 Ballen CJ, Mason NA. 2017. Longitudinal Analysis of a Diversity Support Program in
342 Biology: A National Call for Further Assessment. *Bioscience* 67:367-373.
- 343 Ballen CJ*, Salehi S*, Cotner S. 2017. Exams disadvantage women in introductory biology.
344 *PLoS ONE* 12(10):e0186419.
- 345 Ballen CJ, Wieman C, Salehi S, Searle JB, Zamudio KR. 2017b. Enhancing diversity in
346 undergraduate science: Self-efficacy drives performance gains with active learning. *CBE-*
347 *Life Sciences Education* 16(4):ar56.
- 348 Barr DA, Gonzalez ME, Wanat SF. 2008. The leaky pipeline: Factors associated with early
349 decline in interest in premedical studies among underrepresented minority undergraduate
350 students. *Academic Medicine* 83:503-511.
- 351 Beede D, Julian T, Langdon D, McKittrick G, Khan B, Doms M. 2011. Women in STEM: A
352 Gender Gap to Innovation. *Economics and Statistics Administration Issue Brief* 4:1-11.
- 353 Bonous-Hammarth M. 2000. Pathways to success: Affirming opportunities for science,
354 mathematics, and engineering majors. *Journal of Negro Education* 69:92-111.
- 355 Braxton JM, Hirschy AS, McClendon SA. 2011. *Understanding and Reducing College*
356 *Student Departure: ASHE-ERIC Higher Education Report, Volume 30, Number 3.* John
357 Wiley & Sons.
- 358 Brewer CA, Smith D. 2011. Vision and change in undergraduate biology education: a call to
359 action. American Association for the Advancement of Science, Washington, DC.
- 360 Cabrera AF, Nora A, Castaneda MB. 1992. The role of finances in the persistence process: A
361 structural model. *Research in higher education* 33:571-593.
- 362 Clance PR. 1985. *The impostor phenomenon: Overcoming the fear that haunts your success.*
363 Peachtree Pub Ltd.
- 364 Cotner S, Baepler P, Kellerman A. 2008. Scratch This! The IF-AT as a Technique for
365 Stimulating Group Discussion and Exposing Misconceptions. *Journal of College Science*
366 *Teaching* 37: 48–53.

367 Cotner SH, Fall BA, Wick SM, Walker JD, Baeppler PM. 2008. Rapid feedback assessment
368 methods: Can we improve engagement and preparation for exams in large-enrollment
369 courses? *Journal of Science Education and Technology* 17: 437–443. doi:10.1007/s10956-
370 008-9112-8

371 Crawford M, MacLeod M. 1990. Gender in the college classroom: An assessment of the
372 “chilly climate” for women. *Sex Roles* 23:101-122.

373 Cuseo J. 2007. The empirical case against large class size: adverse effects on the teaching,
374 learning, and retention of first-year students. *The Journal of Faculty Development* 21:5-21.

375 Easton JQ, Johnson E, and Sartain L. 2017. The predictive power of ninth-grade GPA.
376 University of Chicago Consortium on School Research.

377 Eddy SL, Brownell SE, Wenderoth MP. 2014. Gender gaps in achievement and participation
378 in multiple introductory biology classrooms. *CBE-Life Sciences Education* 13:478-492.

379 Fried T, MacCleave A. 2009. Influence of role models and mentors on female graduate
380 students' choice of science as a career. *Alberta Journal of Educational Research* 55:482.

381 Glass GV. 1982. School class size: Research and policy. *School class size: Research and*
382 *policy*. Beverly Hills, CA: Sage Publications.

383 Glass GV, Smith ML. 1979. Meta-analysis of research on class size and achievement.
384 *Educational Evaluation and Policy Analysis* 1:2-16.

385 Gloria AM, Kurpius SER, Hamilton KD, Willson MS. 1999. African American students'
386 persistence at a predominantly White university: Influences of social support, university
387 comfort, and self-beliefs. *Journal of College Student Development* 40:257.

388 Griffith AL. 2010. Persistence of women and minorities in STEM field majors: Is it the
389 school that matters? *Economics of Education Review* 29:911-922.

390 Grunspan DZ, Eddy SL, Brownell SE, Wiggins BL, Crowe AJ, Goodreau SM. 2016. Males
391 under-estimate academic performance of their female peers in undergraduate biology
392 classrooms. *PLoS ONE* 11:1-16.

393 Haak DC, HilleRisLambers J, Pitre E, Freeman S. 2011. Increased structure and active
394 learning reduce the achievement gap in introductory biology. *Science* 332:1213-1216.

395 Hall RM, Sandler BR. 1982. *The Classroom Climate: A Chilly One for Women?* Association
396 of American Colleges, Washington, DC. Project on the Status and Education of Women.

397 Ho DE, Kelman MG. 2014. Does class size affect the gender gap? a natural experiment in
398 law. *The Journal of Legal Studies* 43:291-321.

399 Hurtado S, Ruiz A. 2012. The climate for underrepresented groups and diversity on campus.
400 *American Academy of Political and Social Science* 634:190-206.

401 Kena G, Hussar, W, McFarland, J, de Brey, C, Musu-Gillette, L, Wang, X, Zhang, J,
402 Rathbun, A, Wilkinson-Flicker, S, Diliberti, M, Barmer, A, Bullock Mann, F, and Dunlop
403 Velez, E. 2016. *The Condition of Education 2016*. National Center for Educational
404 Statistics.

405 Knight JK, Wise SB, Sieke S. 2016. Group random call can positively affect student in-class
406 clicker discussions. *CBE-Life Sciences Education* 15(4). doi:10.1187/cbe.16-02-0109

407 Koester BP, Grom G, and McKay TA. (2016) Patterns of gendered performance difference in
408 introductory STEM courses, *arXiv preprint arXiv:1608.07565*.

409 Kokkelenberg EC, Dillon M, Christy SM. 2008. The effects of class size on student grades at
410 a public university. *Economics of Education Review* 27:221-233.

411 Kreft IG, Kreft I, de Leeuw J. 1998. *Introducing multilevel modeling*. Sage.

412 Landivar LC. 2013. Disparities in STEM employment by sex, race, and Hispanic origin.
413 Education Review 29:911-922.

414 Lewin JD, Vinson EL, Stetzer MR, Smith MK. 2016. A Campus-Wide Investigation of
415 Clicker Implementation: The Status of Peer Discussion in STEM Classes 15: 1–12.
416 doi:10.1187/cbe.15-10-0224

417 Lopatto D. 2007. Undergraduate research experiences support science career decisions and
418 active learning. CBE-Life Sciences Education 6:297-306.

419 May GS, Chubin DE. 2003. A retrospective on undergraduate engineering success for
420 underrepresented minority students. Journal of Engineering Education 92:27-39.

421 Mazzone TL. 1993. The changing politics of state education policy making: A 20-year
422 Minnesota perspective. Educational evaluation and policy analysis 15:357-379.

423 Mervis J. 2011. Weed-out courses hamper diversity. Science 334:1333-1333.

424 Milkman KL, Akinola M, Chugh D. 2015. What happens before? A field experiment
425 exploring how pay and representation differentially shape bias on the pathway into
426 organizations. Journal of Applied Psychology, 100(6):1678.

427 Moneta-Koehler L, Brown AM, Petrie KA, Evans BJ, Chalkley R. 2017. The limitations of
428 the GRE in predicting success in biomedical graduate school. PLOS One, 12(1):e0166742.

429 Moss-Racusin CA, Dovidio JF, Brescoll VL, Graham MJ, Handelsman J. 2012. Science
430 faculty's subtle gender biases favor male students. Proceedings of the National Academy of
431 Sciences 109:16474-16479.

432 Olson S, Riordan DG. 2012. Engage to Excel: Producing One Million Additional College
433 Graduates with Degrees in Science, Technology, Engineering, and Mathematics. Report to
434 the President. Executive Office of the President.

435 Paterson L, Goldstein H. 1991. New statistical methods for analysing social structures: an
436 introduction to multilevel models. British educational research journal 17:387-393.

437 Rask K, Tiefenthaler J. 2008. The role of grade sensitivity in explaining the gender
438 imbalance in undergraduate economics. Economics of Education Review 27:676-687.

439 Schanzenbach DW. 2014. Does Class Size Matter? National Education Policy Center.

440 Schmader T. 2002. Gender identification moderates stereotype threat effects on women's
441 math performance. Journal of Experimental Social Psychology 38:194-201.

442 Smith TY. 2000. 1999–2000 SMET Retention Report: The Retention and Graduation Rates
443 of 1992–98 Entering Science, Mathematics, Engineering and Technology Majors in 119
444 Colleges and Universities. University of Oklahoma Center for Institutional Data Exchange
445 and Analysis.

446 Snyder JJ, Sloane JD, Dunk RD, Wiles JR. 2016. Peer-led team learning helps minority
447 students succeed. PLoS Biology 14:e1002398.

448 Snyder JJ, Wiles JR. 2015. Peer led team learning in introductory biology: Effects on peer
449 leader critical thinking skills. PLoS ONE 10:e0115084.

450 Stanger-Hall KF. 2012. Multiple-choice exams: an obstacle for higher-level thinking in
451 introductory science classes. CBE-Life Sciences Education 11(3):294-306.

452 Steele CM. 1997. A threat in the air. How stereotypes shape intellectual identity and
453 performance. American Psychology 52:613-629.

454 Steele CM, Aronson J. 1995. Stereotype threat and the intellectual test performance of
455 African Americans. Journal of Personality and Social Psychology 69:797.

456 Steele J, James JB, Barnett RC. 2002. Learning in a man's world: Examining the perceptions
457 of undergraduate women in male-dominated academic areas. *Psychology of Women*
458 *Quarterly* 26:46-50.

459 Stout JG, Dasgupta N, Hunsinger M, McManus MA. 2011. STEMing the tide: using ingroup
460 experts to inoculate women's self-concept in science, technology, engineering, and
461 mathematics (STEM). *Journal of Personality and Social Psychology* 100:255.

462 Suresh R. 2006. The relationship between barrier courses and persistence in engineering.
463 *Journal of College Student Retention: Research, Theory & Practice* 8:215-239.

464 Ventura J. 2000. I ain't got time to bleed: Reworking the body politic from the bottom up.
465 Villard.

466 Walton GM, Logel C, Peach JM, Spencer SJ, Zanna MP. 2015. Two brief interventions to
467 mitigate a “chilly climate” transform women’s experience, relationships, and achievement in
468 engineering. *Journal of Educational Psychology* 107:468.

469 Whitehurst GJ, Chingos MM. 2011. *Class size: What research says and what it means for*
470 *state policy*. Brookings Institution.

471

472

473

474