

Bethel University

Spark

---

Biological Sciences Faculty Publications

Biological Sciences Department

---

8-1-2019

## Smaller classes promote equitable student participation in STEM

Cissy J. Ballen  
*Auburn University*

Stephanie M. Aguilon  
*Cornell University*

Azza Awwad  
*American University in Cairo*

Anne E. Bjune  
*Universitetet i Bergen*

Daniel Challou  
*University of Minnesota Twin Cities*

*See next page for additional authors*

Follow this and additional works at: <https://spark.bethel.edu/biology-faculty>



Part of the [Biology Commons](#), [Scholarship of Teaching and Learning Commons](#), and the [Science and Mathematics Education Commons](#)

---

### Recommended Citation

Ballen, Cissy J.; Aguilon, Stephanie M.; Awwad, Azza; Bjune, Anne E.; Challou, Daniel; Drake, Abby Grace; Driessen, Michelle; Ellozy, Aziza; Ferry, Vivian E.; Goldberg, Emma E.; Harcombe, William; Jensen, Steve; Jørgensen, Christian; Koth, Zoe; McGaugh, Suzanne; Mitry, Caroline; Mosher, Bryan; Mostafa, Hoda; Petipas, Renee H.; Soneral, Paula A.G.; Watters, Shana; Wassenberg, Deena; Weiss, Stacey L.; Yonas, Azariah; Zamudio, Kelly R.; and Cotner, Sehoya, "Smaller classes promote equitable student participation in STEM" (2019). *Biological Sciences Faculty Publications*. 1.  
<https://spark.bethel.edu/biology-faculty/1>

This Article is brought to you for free and open access by the Biological Sciences Department at Spark. It has been accepted for inclusion in Biological Sciences Faculty Publications by an authorized administrator of Spark. For more information, please contact [kent-gerber@bethel.edu](mailto:kent-gerber@bethel.edu).

---

## Authors

Cissy J. Ballen, Stephanie M. Aguilon, Azza Awwad, Anne E. Bjune, Daniel Challou, Abby Grace Drake, Michelle Driessen, Aziza Ellozy, Vivian E. Ferry, Emma E. Goldberg, William Harcombe, Steve Jensen, Christian Jørgensen, Zoe Koth, Suzanne McGaugh, Caroline Mitry, Bryan Mosher, Hoda Mostafa, Renee H. Petipas, Paula A.G. Soneral, Shana Watters, Deena Wassenberg, Stacey L. Weiss, Azariah Yonas, Kelly R. Zamudio, and Sehoya Cotner

1 **Smaller classes promote equitable student participation in STEM**

2  
3 Cissy J. Ballen<sup>1,2,\*</sup>, Stepfanie M. Aguillon<sup>3,4</sup>, Azza Awwad<sup>5</sup>, Anne E. Bjune<sup>6</sup>, Daniel Challou<sup>7</sup>, Abby  
4 Grace Drake<sup>3</sup>, Michelle Driessen<sup>8</sup>, Aziza Ellozy<sup>5</sup>, Vivian E. Ferry<sup>9</sup>, Emma E. Goldberg<sup>10</sup>, William  
5 Harcombe<sup>10</sup>, Steve Jensen<sup>7</sup>, Christian Jørgensen<sup>6</sup>, Zoe Koth<sup>2</sup>, Suzanne McGaugh<sup>10</sup>, Caroline  
6 Mitry<sup>5</sup>, Bryan Mosher<sup>11</sup>, Hoda Mostafa<sup>5</sup>, Renee H. Petipas<sup>12</sup>, Paula A.G. Soneral<sup>13</sup>, Shana  
7 Watters<sup>7</sup>, Deena Wassenberg<sup>2</sup>, Stacey L. Weiss<sup>14</sup>, Azariah Yonas<sup>2</sup>, Kelly R. Zamudio<sup>3</sup>, Sehoya  
8 Cotner<sup>2</sup>

9  
10  
11 <sup>1</sup>Department of Biological Sciences, Auburn University, Auburn, AL, USA

12 <sup>2</sup>Department of Biology Teaching and Learning, University of Minnesota, Minneapolis, MN,  
13 USA

14 <sup>3</sup>Department of Ecology & Evolutionary Biology, Cornell University, Ithaca, NY, USA

15 <sup>4</sup>Fuller Evolutionary Biology Program, Cornell Lab of Ornithology, Ithaca, NY, USA

16 <sup>5</sup>Center for Learning and Teaching, The American University in Cairo, Cairo, Egypt

17 <sup>6</sup>Department of Biological Sciences, University of Bergen, Bergen, Norway

18 <sup>7</sup>Department of Computer Science and Engineering, University of Minnesota, Minneapolis,  
19 MN, USA

20 <sup>8</sup>Department of Chemistry, University of Minnesota, Minneapolis, MN, USA

21 <sup>9</sup>Department of Chemical Engineering and Materials Science, University of Minnesota,  
22 Minneapolis, MN, USA

23 <sup>10</sup>Department of Ecology, Evolution and Behavior, University of Minnesota, Minneapolis,  
24 MN, USA

25 <sup>11</sup>School of Mathematics, University of Minnesota, Minneapolis, MN, USA

26 <sup>12</sup>Department of Plant Pathology, Washington State University, Pullman, WA, USA

27 <sup>13</sup>Department of Biological Sciences, Bethel University, St. Paul, MN, USA

28 <sup>14</sup>Department of Biology, University of Puget Sound, Tacoma, WA

29  
30  
31  
32  
33 *This article has been accepted for publication in Bioscience published by Oxford University Press.*

34 Word count: 8,102

35 Keywords: STEM equity; gender; in-class interactions; class size; active learning

36

37 **Abstract**

38 As Science, Technology, Engineering, and Mathematics (STEM) classrooms in higher education  
39 transition from lecturing to active learning, the frequency of student interactions in class  
40 increases. Previous research documents a gender bias in participation, with women  
41 participating less than would be expected based on their numeric proportions. Here we asked  
42 which attributes of the learning environment contribute to decreased female participation:  
43 abundance of in-class interactions, diversity of interactions, proportion of women in class,  
44 instructor gender, class size, and whether the course targeted lower division (first and second  
45 year) or upper division (third or fourth year) students. We calculated likelihood ratios of female  
46 participation from over 5,300 student-instructor interactions observed across multiple  
47 institutions. We falsify several alternative hypotheses and demonstrate that increasing class  
48 size has the largest negative association. We also found that when instructors use a diverse  
49 range of teaching strategies, women are more likely to participate after small-group  
50 discussions.

51

52 **Introduction**

53 Active learning can be distinguished from traditional lecturing through its emphasis on diverse  
54 types of engagement strategies, including structured student-instructor interactions during  
55 activities or guided inquiry (Haak et al., 2011; Smith et al., 2009). Substantial evidence supports  
56 interactive classes as a more effective form of instruction compared to traditional lecture  
57 (Freeman et al., 2014), particularly for at-risk students (Lorenzo *et al.*, 2006; Beichner *et al.*,  
58 2007; Haak *et al.*, 2011; Ballen *et al.*, 2017b). However, the most effective and equitable types  
59 of interactions that support all students in their learning are a subject of current debate. This  
60 question is particularly critical in gateway courses that are required for all students before they  
61 can pursue more specialized coursework. Across the Science, Technology, Engineering, and  
62 Mathematics (STEM) disciplines, students struggle in gateway courses, and failure rates are  
63 high (Freeman et al., 2011; National Academies of Sciences and Medicine, 2016). Thus, it is  
64 critical that gateway courses are systematically assessed to identify which elements within the  
65 classrooms leads to gaps in participation, and provide structure when needed.

66 Previous research demonstrates a pervasive gender gap in participation in  
67 undergraduate STEM courses (Eddy et al., 2014), a trend that persists beyond undergraduate  
68 lecture halls. In fact, it has been shown that women audience members ask fewer questions  
69 than men after academic seminar and conference talks (Carter et al., 2017; Hinsley et al., 2017;  
70 Pritchard et al., 2014). These patterns may contribute to a general tendency to undervalue the  
71 contributions of women, and lead to documented phenomena such as proportionately fewer  
72 women awarded prestigious fellowships (Wold & Wenneras, 2010) and grants (Ledin et al.,  
73 2007), fewer female first (O’Dorchai et al., 2009) and last authors (Holman et al., 2018; Murray  
74 et al., 2018), fewer women invited as speakers at symposia (Isbell et al., 2012), and fewer  
75 women occupying high-status positions in STEM (O’Dorchai *et al.*, 2009; Beede *et al.*, 2011).  
76 Thus, factors that contribute to unequal participation should be identified and proper  
77 interventions should be designed early in STEM education.

78 Variability in female participation across classrooms indicates the presence of  
79 underlying, course-specific factors that create environments more or less encouraging to the  
80 input of women. We selected six course elements from the literature that may impact female

81 participation, and used deductive methods to understand each element’s relative impact on  
82 equitable participation from our sample of observations (Table 1).

83 We examined how the abundance of interactions, diversity of interactions, instructor  
84 gender, proportion of women in the class, class size, and class division affect three specific  
85 types of student participation: (1) **voluntary responses**, when an instructor poses a question  
86 and an individual raises their hand to answer without conferring with their peers; (2) **group**  
87 **responses**, when an instructor poses a question and students have the opportunity to talk to  
88 their peers before answering; (3) **total responses**, or all student-instructor interactions  
89 observed across a class period. A summary of our reasoning for several hypotheses (predictors)  
90 for female participation is provided in Table 1. We addressed the following research question as  
91 it applies across multiple universities: what leads to gendered participation in science lectures  
92 in higher education? We developed a number of alternative hypotheses that might predict why  
93 in some environments we observe individuals of one gender speaking more than another (Table  
94 1).

95  
96 **Table 1.** Alternative hypotheses that may explain, in isolation or in combination, equitable in-  
97 class participation in STEM courses.

Predictor	Reasoning: Students may be more comfortable speaking in class...
<b>Abundance of student-instructor interactions per class period</b>	...if participation is normalized through many different instances of student-instructor interactions throughout class (Kuh and Hu, 2001; Komarraju <i>et al.</i> , 2010).
<b>Diversity of interactions</b>	...if the instructor uses a wide range of teaching strategies, generally involving peer discussions, (e.g., small-group discussions, classroom response systems, think-pair-share) intended to encourage equitable participation (Premo and Cavagnetto, 2018).
<b>Instructor gender</b>	...if the gender of the instructor matches their own (Crombie <i>et al.</i> , 2003; Cotner <i>et al.</i> , 2011).
<b>Proportion of women in the class</b>	...if genders are represented in relatively equitable proportions, so that the under-represented gender does not feel isolated in the larger social setting (Dahlerup, 1988).

---

**Class size**

...if they are in a classroom with fewer students (Kokkelenberg *et al.*, 2008; Schanzenbach, 2014; Ballen *et al.*, 2018a).

**Lower division or upper division**

...if they are in an upper division course, having cleared the hurdle of the introductory, “weed out” courses (Brewer and Smith, 2011). Alternatively, students warmed to instructional methods over time, including in-class activities.

---

98

99 **Data collection**

100 We collected student behavioral data from 44 courses across the United States. As part of the  
101 creation of this larger collaborative research group, we solicited participation through an  
102 existing professional network from instructors from instructors who teach majors, nonmajors,  
103 or both, from a range of institutions. Volunteers represent Bethel University, Cornell University,  
104 University of Minnesota, University of Puget Sound, the American University in Cairo, Egypt,  
105 and University of Bergen, Norway (Table 2). Participating institutions were a convenience  
106 sample chosen from a range of institutional types (public and private, large and small) and  
107 settings (college towns to large metropolitan areas). During the 2-year study period,  
108 approximately 5,200 students enrolled in the sampled courses, and observers categorized over  
109 5,300 interactions between the instructors and students (Research Coordination Network,  
110 National Science Foundation RCN–UBE Incubator: Equity and Diversity in Undergraduate STEM;  
111 #1729935 awarded to S Cotner and CJ Ballen). We included courses from across STEM fields,  
112 including biology, physics, computer science, and chemistry (details in the raw data file).  
113 Demographic information collected by university registrars revealed that on average 53.8% of  
114 the students in these classes identified as female, but this number ranged from 20.4% to 79.6%,  
115 depending on the specific class. All aspects of research were reviewed and approved by each  
116 schools’ respective Institutional Review Boards (Bethel IRB 180518; Cornell IRB 1410005010;  
117 University of Minnesota IRB 00000800; University of Puget Sound IRB 1617-006; American  
118 University in Cairo 2016-2017-0012; University of Bergen NSD 46727).

119

120 **Table 2.** Six universities participated in the current study, representing diverse geographic  
121 locations across the world.

Institution	Location	Undergraduate enrollment	Institution type	# of courses sampled
American University in Cairo	Cairo, Egypt	5,474	Private	4
Bethel University	St Paul, MN, US	2,800	Faith-based, private	1
Cornell University	Ithaca, NY, US	14,907	Public and private	2
University of Bergen	Bergen, Norway	17,000	Public	2
University of Minnesota	Minneapolis, MN, US	30,511	Public	32
University of Puget Sound	Tacoma, WA, US	2,553	Private	3

122

## 123 **Research methods**

### 124 *Measuring In-Class Participation*

125 We conducted ~1 hour training sessions for observers to characterize classroom participation  
126 as broad types of interactions that occur over a class period, which were further characterized  
127 as either ‘voluntary responses’ or ‘group responses.’ For each type of interaction that takes  
128 place during a class period, an observer recorded the gender of the student participant (1 =  
129 male or 0 = female). The complete (not collapsed) list of categories included: (1) ‘voluntary  
130 response,’ when an instructor poses a question, and an individual raises their hand to answer  
131 without conferring with their group; (2) ‘individual spontaneous question,’ in which a student  
132 asks an instructor an unprompted question or is only very generally prompted (e.g. ‘does  
133 anyone have a question?’); (3) ‘individual spontaneous call,’ when a student makes a comment  
134 not prompted by the instructor; (4) ‘cold call,’ a non-voluntary response after the instructor  
135 calls randomly on an individual (in this scenario, students have not conferred with a group); (5)  
136 ‘spontaneous call post-Think Pair Share (TPS),’ a non-voluntary response after the instructor  
137 calls randomly on a group after they discuss a posed question; (6) ‘voluntary response post-  
138 TPS,’ a voluntary response after the instructor poses a question, students confer, and a student  
139 volunteers to answer the question; (7) ‘voluntary response post-TPS and clicker,’ a voluntary  
140 response after the instructor poses a question, students confer, students answer the question  
141 using a personal response system (e.g., iclicker, TopHat, ChimeIn), and then a student



142 volunteers to answer the question (either after the instructor shows the answer or before; this  
143 category is different from voluntary response post-TPS (#6) in that students have committed to  
144 an answer before responding); and (8) 'circulating instructor question or comment,' when the  
145 instructor is circulating around the classroom, and a student calls them over with a question or  
146 comment (note: we do not distinguish based on content of the interaction because it is often  
147 difficult to identify what is said from the observer's perspective).

148

149 To increase power of analyses, we focus on the most robust categories or combined relevant  
150 values to create broader categories. The final values we included in analyses were (A) **voluntary**  
151 **responses**, the most common type of interaction in which an instructor poses a question, and  
152 an individual raises their hand to answer without conferring with their group (#1 above), and  
153 (B) **group responses**, or any interactions that occur between the student and the instructor  
154 after students have some opportunity to discuss a topic with group members (combination of  
155 #5-7 described above), and (C) **total responses**, or all interactions between the student and  
156 instructor. To clarify, while (C) is not exclusive to (A) and (B), (A) and (B) are exclusive to one  
157 another. Category (C) is the sum of (A) and (B), in addition to a small number of additional  
158 interactions from the original categories described above. Across the two years of observations,  
159 inter-observer reliability at the University of Minnesota was consistently well within acceptable  
160 range among observers' ability to identify voluntary responses and group responses (Cohen's  
161 kappa > 0.90; Hallgren, 2012).

162

163 Because some interactions in our observations were not strictly content related (e.g. instructor  
164 and student discuss current event not related to class) or used only a few times across all  
165 observations, categories 2-4 and 8 were excluded from our analysis (but note they are included  
166 in the total responses variable). For example, students asked individual spontaneous questions  
167 in the beginning of class more often than any other point during lecture, and these rarely  
168 related to the material. Instead we prioritized categories 1, 5-7 because these reliably produced  
169 content-related interactions between instructor and student. We included courses with at least  
170 two full-class observations (minimum 2, maximum 20, average 9.6 observations per course).

171 Only categories that had a total of five or more student–instructor interactions across observed  
172 class sessions for a given course were included in the analyses.

173

#### 174 *Quantifying Predictor Variables*

175 To measure the abundance of instructor-student interactions in class, we calculated the  
176 average number of student-instructor interactions per class period across all observed class  
177 periods. Class period duration varied, so when appropriate, we scaled the average number of  
178 interactions to fit a 50-minute class period. To measure the diversity of these interactions, we  
179 applied Simpson’s diversity index to calculate equitability, or evenness, of teaching strategies  
180 per class (Simpson, 1949). Classically, Simpson’s diversity index is calculated using the number  
181 and abundance of biological species observed, and is used in ecology to quantify the  
182 biodiversity within a habitat. By considering relative abundances, a diversity index depends not  
183 only on species richness but the evenness of individuals distributed among species. Here, we  
184 use the number of interaction types, and how often instructors use each interaction type, to  
185 quantify Simpson’s diversity index of teaching strategies within a classroom (see Supplementary  
186 materials 1 for details and equation). Values range from 0 to 1, with 1 being complete evenness  
187 of teaching strategies. In an education context, low values reflect classrooms with little  
188 variation in instructor-student interaction types; high values reflect classrooms with lots of  
189 different types of instructor-student interactions used frequently.

190 We measured the proportion of women in the class using institutional data when  
191 possible, and information from survey data obtained at the beginning of the semester that  
192 asked “Which pronoun do you prefer to describe yourself?” Students could choose between  
193 she/her, he/him, they/them, or other. Instructor gender was estimated at three levels: man (or  
194 men), woman (or women), or both (both men and women). This is because some classes were  
195 taught **by a man or woman, or** co-taught by men only, women only, or both men and women,  
196 for which we obtained measurements from each instructor. We obtained class size information  
197 from the institution or directly from the instructor.

198 We categorized classes at two levels — those that primarily enrolled first and second  
199 year students (lower division), or classes enrolling third and fourth year students (upper

200 division). We acknowledge that students in upper division courses do not represent a random  
201 sample of students from lower division courses, and multiple selective forces may shape  
202 student samples.

203

#### 204 ***Statistical analyses***

205 We measured outcomes as likelihood ratios,  $LR_w$ , or the likelihood that a participant is a  
206 woman compared to the likelihood that a participant is a man in a given category of interaction,  
207 such that a value of one means that the likelihood of a woman participating is the same as that  
208 of a man. To calculate likelihood ratios, we divided the proportion of instructor-student  
209 interactions with women,  $I_w$ , by the proportion of women in the class,  $C_w$ . We then took this  
210 value and divided it by the proportion of instructor-student interactions with men,  $I_m$ , over  
211 proportion of men in the class,  $C_m$ .

$$212 \quad \quad \quad LR_w = (I_w/C_w)/(I_m/C_m)$$

213

214 For example, consider a semester over which we observed student participation in one  
215 class. We found that of all student-instructor interactions observed, 30% involved female  
216 students and 70% involved male students. In this example, the class composition was 80  
217 women and 120 men (in other words, 40% women and 60% men). With these values, our  
218 outcome would be  $((0.30 / 0.40)/(0.70/0.60))=0.64$  (i.e., in this class women participate 0.64  
219 times as much as men participate). Values less than 1 indicate that women were less likely to  
220 participate relative to men, and values above one indicate that women were more likely to  
221 participate. We used linear mixed-effects models with the LME4 package in R (Bates et al.,  
222 2014; R Core Team, 2014) to test the impact of predictors on the following outcome percentage  
223 differentials across institutions: voluntary responses, group responses, and total responses. We  
224 used the number of classroom observations as a weighted variable because it encodes how  
225 many original observations were conducted in each classroom, and therefore larger weights are  
226 assigned to courses with more ‘reliable’ estimates. A model that treated all of the classroom  
227 data sets equally would give less observed classes more influence and highly observed classes  
228 too little influence. Weighting variables gives each data point the appropriate amount of  
229 influence over the parameter estimates, and is particularly useful in smaller datasets.

230 For the multi-university analyses, we included schools ('uni') as a random variable in the  
 231 mixed effects model. Starting with a null model, we used Akaike's information criterion to  
 232 assess model fit (Table 3). We chose the most parsimonious model that best fit the data by  
 233 calculating AIC differences ( $\Delta_i$ ), and Akaike weights ( $w_i$ ) which both represent different ways to  
 234 assess of strength of each model as the best model. We only included data that included all  
 235 predictor variables (Supplementary material 2: Model selection summary tables).

236 Because the majority of classes observed were from the University of Minnesota (UMN),  
 237 we were also interested in whether apparent trends persisted across the non-UMN institutions  
 238 (N = 12). We ran post hoc analyses on non-UMN institutions to address this question.

239

240 **Table 3.** Best fit models for analyses of total responses, voluntary response, and group response  
 241 across all institutions.

Outcome variable	Best fit model
Total responses	~class size + (1   uni)
Voluntary Response	~class size + (1   uni)
Group Response	~class size + Simpson's diversity index + (1   uni)

242

## 243 **Results**

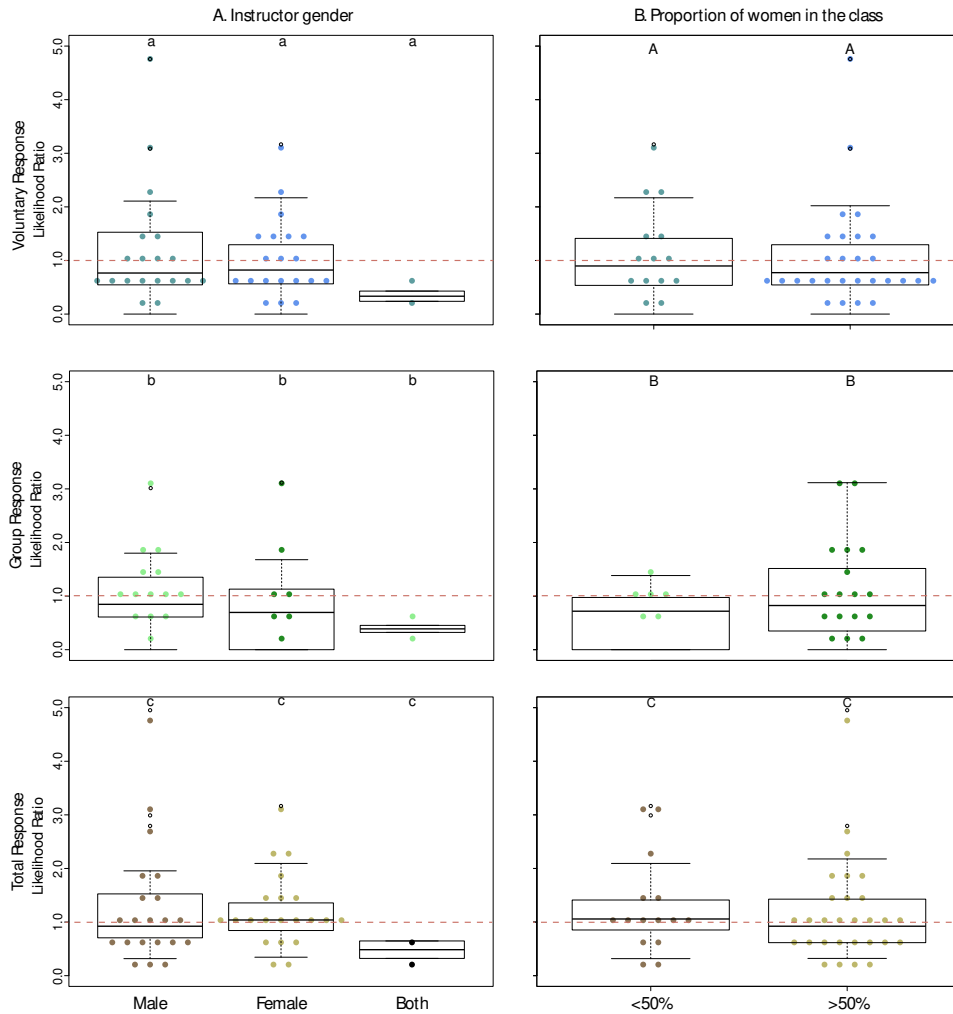
### 244 ***Analyses of courses across six universities with mixed-effects models***

245 Overall, across all classes, the average likelihood ratio for voluntary, group, and total  
 246 interactions were 1.03 (0.92 SD), 0.86 (0.81 SD), and 1.2 (0.91 SD), respectively. To examine  
 247 factors that explain observed variation in the data, we used linear mixed-effects models across  
 248 the 44 classes. Our multilevel model accounted for fixed and random effects to explain  
 249 variation in the data (e.g., instructor gender as a fixed effect, and school as a random effect).  
 250 This approach controls for nonindependence in sampling due to the nested nature of our data  
 251 (Theobald, 2018). We present data to falsify a number of alternative hypotheses: in our sample  
 252 of observed classes, gender bias in participation was not predicted by:

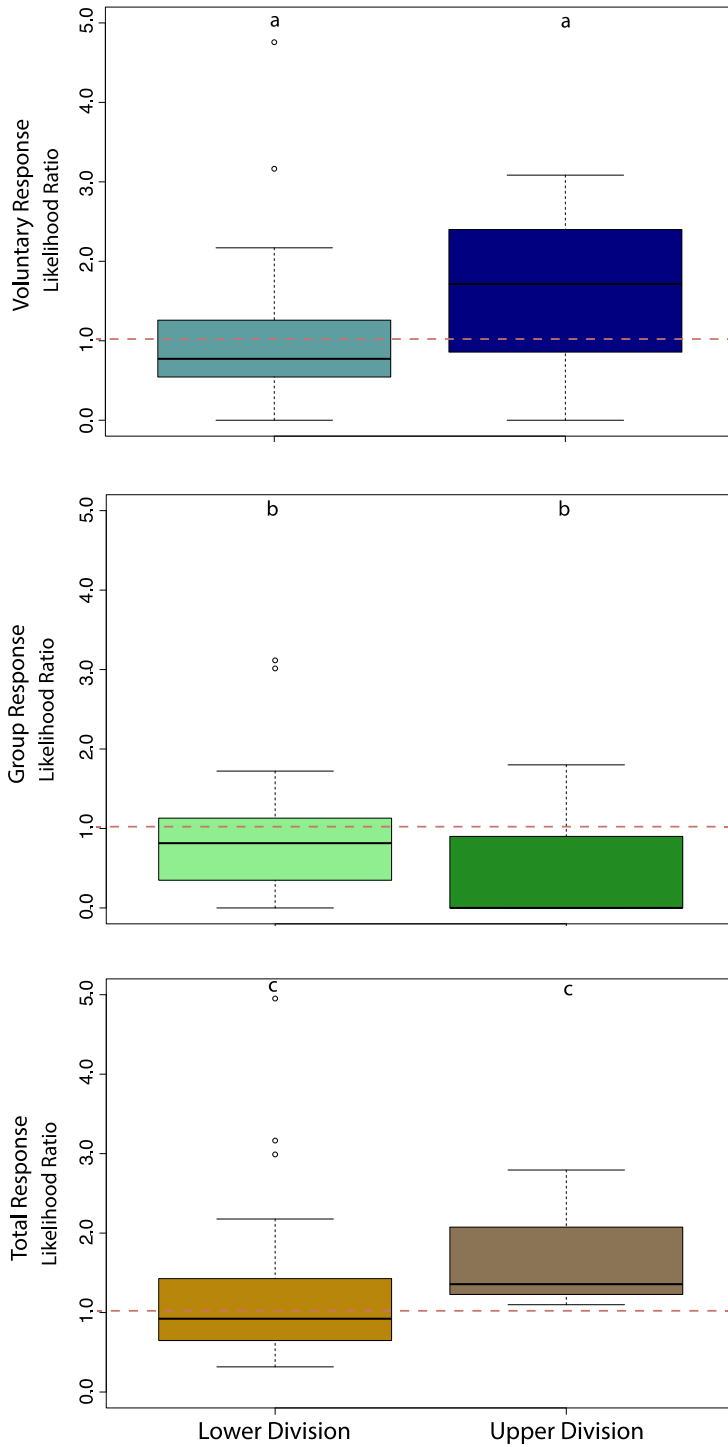
- 253 • the abundance of interactions in the class (Supplementary material 1)
- 254 • the gender(s) of the instructors (Figure 1A)

- 255 • the proportion of women sitting in the classroom (i.e., 'critical mass effect'; Figure 1B)
- 256 • whether courses were lower (first and second year) or upper division (third or fourth
- 257 year) (Figure 2A).

258 During the model selection process, all of these variables were eliminated because they did not  
259 significantly improve the fit of the model to the data (Supplementary material 2: Model  
260 selection summary tables; Results tables). The classroom trait that had the largest impact on  
261 equitable participation was class size, with women demonstrating higher levels of voluntary  
262 responses and total responses in smaller classes across six institutions (voluntary responses  $B =$   
263  $-0.005$ ,  $t(24.810) = -3.483$ ,  $P = 0.002$ ,  $SE = 0.001$ ; total responses  $B = -0.004$ ,  $t(25.274) = -2.890$   $P$   
264  $= 0.008$ ,  $SE = 0.001$ ; Figure 3). Based on these estimates, as class size increased, fewer women  
265 were likely to voluntarily respond to questions posed by the instructor. Based on the estimated  
266 effect size, an increase in class size from 50 to 150 students decreased the likelihood of a  
267 woman participating relative to a man by 50%. Class size did not have a significant impact on  
268 gender-specific group responses across six institutions ( $B = -0.004$ ,  $t(17.805) = -1.643$ ,  $P = 0.118$ ,  
269  $SE = 0.002$ ). The Simpson's diversity index, which considers the variety of interactions, and how  
270 often instructors used each type of interaction, significantly predicted group response  
271 likelihood ratios ( $B = 2.114$ ,  $t(26.897) = 2.473$ ,  $P = 0.020$ ,  $SE = 0.855$ ; Figure 4A), with increasing  
272 likelihood of female participation as teaching methods varied. Future research will profit from  
273 an explicit focus on this course component to clarify the full impact of group discussions on  
274 equitable participation.

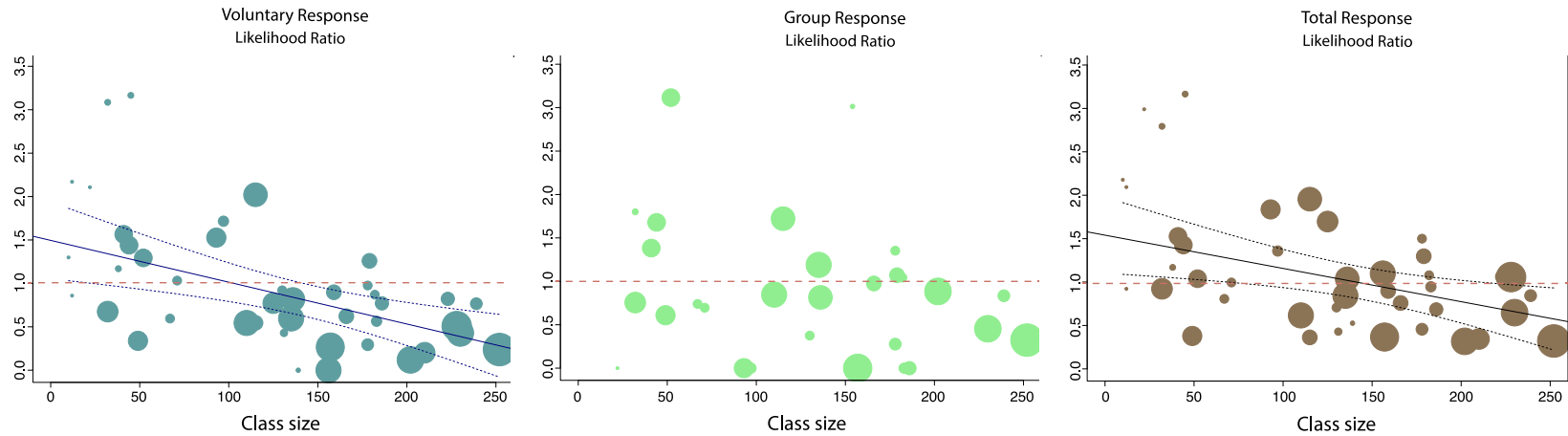


275  
 276 **Figure 1.** A. Instructor gender: Likelihood of female voluntary responses (blue), group responses  
 277 (green), and total responses (brown) based on the instructor gender. B. Proportion of women in  
 278 the classroom: Likelihood of female voluntary responses (blue), group responses (green), and  
 279 total responses (brown) based on the proportion of women in the classroom (either under 50%  
 280 or over 50%). Letters at the top of each panel indicate non-significant differences ( $P > 0.05$ ).  
 281 Values less than 1 indicate fewer women participated relative to men, and values above one  
 282 indicate more women participated. The dashed line indicates parity in participation.  
 283  
 284



285  
 286  
 287  
 288

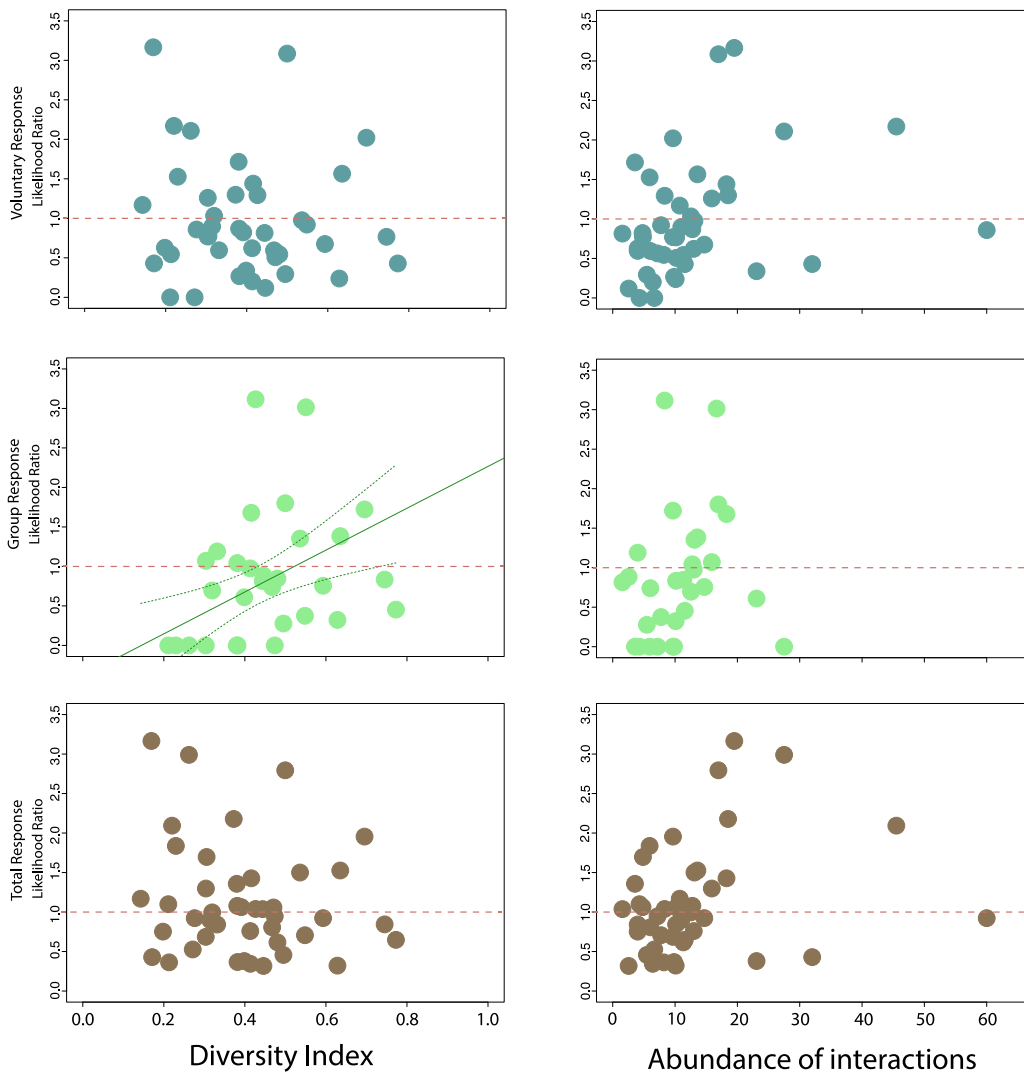
**Figure 2.** Likelihood of female voluntary responses (blue), group responses (green), and total responses (brown) in lower division versus upper division courses across all institutions. Letters above the box plots show statistical non-significance across categories ( $P > 0.05$ ).



290  
291  
292  
293  
294

**Figure 3.** The impact of class size on the likelihood of female voluntary responses (blue), group responses (green), and total responses (brown) across all institutions sampled. Regression lines with confidence intervals denote significant relationships between the likelihood ratio and class size ( $P < 0.05$ ), with values below one indicating women were less likely to participate than men. The size of the symbol is proportional to the number of classes observed.





296

297

298

299

300

301

302

303

304

**Figure 4.** The likelihood of female voluntary responses (blue), group responses (green), and total responses (brown) across all institutions as a function of a calculated in-class ‘Simpson's diversity index’ that measures the amount of varied teaching strategies an instructor uses and the abundance of interactions per 50-minute class period. Regardless of class size, more women participated after group discussions when the instructor used more diverse types of interactions during the class period. Regression lines with confidence intervals denote significant relationships between variables ( $P < 0.05$ ), with values below one indicating women were less likely to participate than men.

305 In order to test whether the relationship between class size and likelihood that women  
306 participate was driven by the data obtained from the University of Minnesota (UMN), we  
307 combined and analyzed all institutions *other than* UMN. Due to the low sample size (N = 12), we  
308 caution readers as they interpret our results. Using Spearman's correlations, we found  
309 significant negative relationships between class size and the likelihood of female participation  
310 with voluntary responses ( $r_s = -0.774$ ,  $P = 0.003$ ) and total responses ( $r_s = -0.770$ ,  $P = 0.003$ ),  
311 but not group responses (N = 9;  $r_s = -0.200$ ,  $P = 0.606$ ) across the 12 non-UMN classes  
312 (Supplementary material 3). For the Simpson's diversity index, we did not observe the same  
313 results when we removed University of Minnesota. We found significant negative relationships  
314 between Simpson's diversity index and the likelihood of female participation across voluntary  
315 responses ( $r_s = -0.755$ ,  $P = 0.005$ ) and total responses ( $r_s = -0.664$ ,  $P = 0.018$ ), but not group  
316 responses (N = 9;  $r_s = -0.050$ ,  $P = 0.898$ ; Supplementary material 3).

317  
318

### 319 **Discussion**

320 We analyzed predictors of female participation as voluntary responses, group responses, and  
321 total responses in lecture, across 44 unique STEM courses (Summary of results, Table 4). We  
322 falsified several alternative hypotheses and demonstrated that gender biased participation  
323 sharply increases in large classes. These results suggest that the reluctance of women to  
324 participate in class is related to traits inherent to large lecture courses. We also used a  
325 modified form of Simpson's diversity index and equitability as a proxy for diverse teaching  
326 strategies in student-instructor interactions (described in Supplementary materials 1). The  
327 Simpson's diversity index measure showed women were more likely to participate after group  
328 work when the instructor employed diverse teaching strategies in the course.

329

330 **Table 4.** Summary of results found in the observational study of student participation across six  
 331 institutions.

Course element tested	Difference?	Notes
Abundance of student-instructor interactions	No	No effect.
Diversity of student-instructor interactions	Yes	More diverse interactions = more female participation after group work.
Proportion of women in the classroom	No	No effect.
Instructor gender	No	No effect.
Class size	Yes	Smaller class size = more female participation in voluntary responses and across all observations.
Lower division or upper division course	No	No effect.

332  
 333 **The impacts of class size**  
 334 Research on the reduction of class size has produced mixed results, largely focused on K-12  
 335 student populations, and on much smaller scales than the data presented here. Despite  
 336 ongoing debates on the effectiveness of reducing class size in K-12 learning spaces, several  
 337 state legislatures have appropriated significant amounts of money to reduce classes to between  
 338 15 to 20 students (summarized in Zinth 2005). For example, in 1990, the Tennessee legislature  
 339 funded a longitudinal study on the impact of reducing the size of K-3 classes on student  
 340 achievement. By following 7,000 students across 79 elementary schools, researchers concluded  
 341 that small class sizes (13-17 students) increased student achievement scores as compared to  
 342 students in regular class sizes (22-25 students). Further, those students who were exposed to

343 small classes early in their education excelled later, after they were re-introduced into regular-  
344 sized classes.

345         Inspired by the results observed in Tennessee, California passed an ambitious education  
346 reform initiative in 1996, committing more than \$1 billion a year to a class-size reduction  
347 program that provided irresistible financial incentives to school districts that reduced the  
348 number of students in K-3 classes. However, California schools confronted unique problems  
349 that did not apply in the Tennessee case study , including a shortage of qualified teachers and  
350 adequate teaching facilities to reduce class size. Additionally, California was more culturally  
351 diverse, with one-third of California’s students living in households in which languages other  
352 than English were primarily spoken. Research into California’s efforts found that class-size  
353 reduction did not benefit school districts serving the state’s most historically underserved  
354 students. This was partly because (1) the effort was more expensive to implement than  
355 expected; (2) in efforts to recruit new staff, they observed a decline in average teacher  
356 qualifications; and (3) in order to create additional classroom spaces, lower-income schools  
357 used facilities and resources at the expense of other programs (Jepsen and Rivkin, 2009). Thus,  
358 impacts of class size reduction efforts can be context-dependent, and care must be taken in  
359 assessing their impacts.

360         Results from studies that focus on the effects of class size in higher education approach  
361 the research on a different scale, and generally with more diverse student populations. Cuseo  
362 (2007) reviewed studies that examined the effects of class size on teaching, learning and  
363 retention. His findings indicate that increasing class size had deleterious impacts on educational  
364 outcomes for students overall, and students enrolled in first year courses in particular. Studies  
365 using big data have echoed these findings, that student achievement declines as class size  
366 increases (Dillon *et al.*, 2002; Kokkelenberg *et al.*, 2008). Maringe and Sing (2014) warn that  
367 increasing class sizes are particularly dangerous when coupled with current national trends  
368 towards increased student mobility, access to higher education, and internationalization of  
369 student composition. They point to impact of the trade-off between individualized instruction  
370 and class size on student participation and engagement, curricular access and interpretation,  
371 opportunities for deep learning for all, and evaluation of student learning and satisfaction.

372 Renewed focus on this topic is warranted after the recent development of online or  
373 hybrid classes and very large enrollments. For example, students in the University of Central  
374 Florida's College of Business obtained more than 1,800 signatures on a petition criticizing the  
375 college's recent shift to a blended classroom model. Classes that tend to have between 800 and  
376 2,000 students learn through a reduced class time format, which eliminates instructor-led  
377 lectures with the expectation that students spend more time learning with their peers outside  
378 of class to gain more thorough knowledge of the material  
379 ([https://www.insidehighered.com/digital-learning/article/2018/09/21/blended-learning-](https://www.insidehighered.com/digital-learning/article/2018/09/21/blended-learning-model-university-central-florida-draws-business)  
380 [model-university-central-florida-draws-business](https://www.insidehighered.com/digital-learning/article/2018/09/21/blended-learning-model-university-central-florida-draws-business)). From an institutional perspective, while the  
381 additional costs of smaller classes are viewed as prohibitively expensive as enrollment rises,  
382 results such as those presented here should not be ignored. Increased understanding of  
383 qualities that support learning and participation of students in small, medium, and large classes  
384 will improve effectiveness within institutional limitations.

385  
386 **Why do we observe gender differences in participation?**  
387

388 Our data show that the largest gender disparities in participation occur when instructors  
389 elicit voluntary responses from students immediately after asking a question in a large lecture  
390 hall. Previous work suggests that instructors may not provide enough time for most students to  
391 think through a response. Rowe (1974a) reported that when precollege instructors asked  
392 voluntary response questions, the 'wait time' before the instructor rephrased or called on a  
393 student was approximately one second. With approximately one second, students must  
394 formulate a response and decide whether to participate, and many factors unrelated to content  
395 knowledge impact the decision to do so. Some of these factors may differentially affect men  
396 and women. For example, Cooper et al. (2018) showed that men generally have a higher  
397 perception of their own ability in a disciplinary domain. In the context of an interactive  
398 introductory STEM course, this may lead to increased comfort among men in readily  
399 participating in front of a large lecture.

400 Other work shows different factors prevent men and women from participating, with  
401 women citing a central reason as 'not working up the nerve' to ask a question or respond to an

402 answer (Ballen et al., 2018; Carter et al., 2017). Elements of social identity threat may also be at  
403 work, in which a person's social identity (in this case gender), can be, or perceived to be,  
404 negatively stereotyped (Steele et al., 2002). Extensive evidence from the precollege literature  
405 shows that regardless of how girls perform in a subject, they are more concerned about how  
406 instructors will evaluate them (Pomerantz et al., 2002), and are less confident than boys in their  
407 science content knowledge, even after controlling for variation in their performance (Micari et  
408 al., 2007). This difference is apparent in several STEM disciplines at the college level, and likely  
409 plays a role in the observed skewed in-class participation towards males.

410

#### 411 **Limitations**

412 The methods of this study have a number of limitations. We decided to quantify real-time  
413 interactions in classrooms to expand our opportunities to collaborate across universities.  
414 However, this meant that in some classes, observers could not double check whether they  
415 categorized interactions correctly if they were unsure. An advantage of having observers in the  
416 classroom observing in real-time is a reduced uncertainty about student gender of participants,  
417 and observers could move if necessary to better identify students (which is not possible with a  
418 camera). While the person who trained all observers was the same (Ballen), we were only able  
419 to obtain reliability scores across observers at the University of Minnesota. Within the  
420 categories we used (voluntary response or response after group work) we consistently had very  
421 high inter-observer reliability at the University of Minnesota (>0.90), but this was not measured  
422 across all observers. Therefore we cannot rule out the possibility that reliability across other  
423 institutions was lower than at the University of Minnesota. However, for this reason we urge  
424 readers to find analyses of total responses the *most* reliable, which encompasses *all* types of  
425 interactions. Additionally, for responses where the instructor posed a question and selected a  
426 person to answer, there is the possibility that the instructor, being aware of the ongoing study,  
427 would preferentially select women more often than their ratio among those who volunteered.  
428 Instructors report that they did not knowingly do this, and results are similar between  
429 "individual spontaneous question" (i.e., in which a student asks an instructor an unprompted

430 question or is only very generally prompted) where this was not an issue and the other  
431 categories.

432 Another limitation is the binary assignment of gender. Such assignment may not align  
433 with self-identified gender. Gender does not exist as a binary variable but rather along a  
434 continuum (Ainsworth, 2015). In this study we only report male and female genders due to the  
435 limitations of our non-invasive observation methods, and we recognize we are unable to report  
436 more accurate gender identities. While we focused on either lower division (first and second  
437 year) or upper division (third or fourth year) classes, this does not rule out the possibility that  
438 the course level precisely reflects the composition of student experience in those courses.  
439 Specifically, some introductory classes that are required for certain majors can be taken at any  
440 time before graduation, and might include larger proportions of older students than other  
441 introductory classes. We did not examine the composition of students in those classes in this  
442 context specifically. Finally, we removed one class from the analysis because it yielded an  
443 unusually high likelihood ratio. Whereas all other values ranged from zero to four (i.e.,  
444 likelihood of female participation was four times the probability of male participation), in this  
445 class the likelihood of women participating was 18 times higher in two types of participation.  
446 We believe this may have been the impact of one or two very vocal students. While the outlier  
447 did not impact the overall results, it *created* a significant association between outcomes and  
448 whether students were in lower or upper division courses. Because we cannot completely rule  
449 out the possibility that the results which include this data point are a better explanation of  
450 student participation in science, we also provide the model selection and results as they appear  
451 *with* the inclusion of this outlier (Supplementary material 2: Model selection summary,  
452 including outlier; Results tables, including outlier). While the current dataset has limitations,  
453 this kind of collaborative effort among universities still allows us to amass enough data to  
454 assess predictors of behavior and answer larger questions across a broad sample of university  
455 types.

456

457 **What can instructors do to broaden participation?**

458           Instructors who teach large lectures can use many simple, evidence-based strategies to  
459 increase participation. For instance, by simply lengthening ‘wait time’ from one second to  
460 between three to five seconds, Rowe (1974b) found that more students volunteer answers, and  
461 that students' answers were longer and more complex. Additionally, asking students to discuss  
462 questions in pairs or in groups lets students work through problems in a non-threatening  
463 environment, and practice expressing their opinions prior to being called upon (Smith et al.  
464 2009). Our results show that group work mitigated the negative impact of large class size on  
465 female participation. Interdependency theory (Rusbult and Van Lange, 2008) predicts  
466 individuals who are put in positions to invest in and rely on peers for their success will also help  
467 themselves. Previous work demonstrates how increasing interdependency among classroom  
468 peers promotes participation, discussion, and ideas (Brewer and Klein, 2006). In large  
469 classrooms, structured ways to promote interdependency among students is one pathway to  
470 improve equitable participation. Another simple option is to have students respond in *writing*  
471 first rather than out loud, using a student response system that has space for open responses to  
472 questions. After the instructor reports a few anonymous notable answers, they can ask  
473 students to follow-up out loud. To increase the breadth of responses in class, instructors can  
474 ask for multiple volunteers and only call on one or more individuals after a certain number of  
475 students have raised their hands (Tanner, 2013). Instructors can assign student groups a  
476 number, and use a random number generator to spontaneously call on groups. Within student  
477 groups, randomly appointed ‘reporters’ can be responsible for voicing an answer on behalf of  
478 their group, which also takes responsibility off of the individual if the answer is incorrect (Cohen  
479 and Lotan, 2014). Instructors assign reporters based on arbitrary qualities, such as the person  
480 who woke up earliest that morning, or the person sitting closest to the classroom entryway  
481 (Tanner et al. 2013). Critically, our findings suggest that employing a *diversity* of strategies to  
482 promote engagement, rather than simply settling on one or two, is likely to lead to more  
483 equitable participation. We do not explicitly address engagement in this research, but future  
484 research will profit from the study of *engagement equity* as a function of class size. If women  
485 are experiencing large classes differently from men, which contributes to gender gaps in  
486 participation, we may also expect differences in engagement, as well.



487 For students, the opportunity to reflect on, interact with, and come to a deep  
488 understanding of scientific ideas is central to learning. Providing explicit guidance for  
489 instructors requires a careful investigation of underlying factors that contribute to observed  
490 classroom disparities.

491

## 492 **Conclusion**

493 Our results align with previous work that calls for a halt on the continued expansion of  
494 large introductory ‘gateway’ courses in science (Achilles, 2012; Baker et al., 2016; Cuseo, 2007),  
495 and underscores the importance of continued empirical measurement of factors that either  
496 promote or counter equity in undergraduate STEM (Brewer & Smith, 2011; National Academies  
497 of Sciences and Medicine, 2016). In practice, the gender gap in participation means women in  
498 large STEM courses systematically miss out on opportunities to rehearse articulating their  
499 answers aloud to a science community, in an environment where wrong answers rarely have  
500 negative impacts on consequential outcomes such as grades. These formative experiences are  
501 bound to influence future interactions (e.g. in seminars and conferences; Carter et al., 2017;  
502 Hinsley et al., 2017; Pritchard et al., 2014; Schmidt et al., 2017; Schmidt & Davenport, 2017),  
503 possibly contributing to a general tendency to undervalue the input of women in STEM (e.g., as  
504 grant recipients or speakers; Grunspan et al., 2016; Isbell et al., 2012).

505 Fortunately, while large lectures do pose a clear challenge to student success overall,  
506 and to equitable performance (Ballen et al., 2018) and participation specifically, instructors can  
507 employ simple strategies to minimize some of these challenges. In fact, many evidence-based  
508 active-learning techniques appear to work by making large classes function like smaller classes.  
509 Our results show females were more likely to participate after small group discussions and this  
510 effect was more pronounced when diverse teaching approaches were employed. Further, these  
511 findings support the “course deficit model,” whereby overt instructional choices can minimize  
512 gaps—in this case, in participation—that may contribute to inequalities in STEM (Cotner and  
513 Ballen, 2017). By placing some of the burden of responsibility on instructors, we are in a better  
514 position to be proactive in our classrooms with respect to these inequities.

515           We realize that ultimately, administrators and legislators must grapple with the  
516 problems associated with large classes, and we hope this work can be part of that conversation.  
517 Based on our results, large classes begin to negatively impact students when they are  
518 comprised of more than approximately 120 students. This may be because class size is strongly  
519 associated with the kinds of assignments given and the level of student involvement in class.  
520 Instructors can play an active role in minimizing the problems associated with large classes by  
521 drawing on the active learning literature and exploring which strategies, from an array of  
522 possibilities, are most effective in their own courses. Our results suggest that the best way to  
523 ameliorate the negative impact of large class sizes on female participation is to use diverse  
524 teaching strategies and small group interactions.

525

526 **Acknowledgements.** This work would not have been possible without undergraduate and  
527 graduate research assistants including Neelam Chandiramani, Melissa Khaw, Shivonne  
528 McCarthy, Sergio Molina, Jake Peterson, Meredith Song, Morgan Burkhart, Grace Leland,  
529 Brandon Vanderbush, Mai Vang, Christine Lian, Connor Neil, Tiare Gill, Kylie Young, Irene Deng,  
530 Ellica Spiut, and Gregor-Fausto Siegmund, for collecting and entering participation data; Lucy  
531 Fortson and Mos Kaveh for connecting us with instructors; Jennifer Kroschel, Geri Grosinger,  
532 Oddfrid T. Kårstad Førland, and Rachael Shelden for invaluable logistical support; Aaron  
533 Wynveen, Annika Moe, Mark Decker, and Cheryl Scott for allowing us to observe their classes,  
534 and J.D. Walker for statistical advice and valuable insights. This work was funded in part by a  
535 Research Coordination Network grant from the National Science Foundation, RCN–UBE  
536 Incubator: Equity and Diversity in Undergraduate STEM; #1729935 awarded to S Cotner and CJ  
537 Ballen.

538 **References**

- 539 Achilles CM. 2012. Class-Size Policy: The STAR Experiment and Related Class-Size Studies.  
540 NCPEA Policy Brief. Volume 1, Number 2. *NCPEA Publications*
- 541 Ainsworth C. 2015. Sex redefined. *Nature* **518** (7539): 288
- 542 Baker BD, Farrie D, Sciarra DG. 2016. Mind the Gap: 20 Years of Progress and Retrenchment in  
543 School Funding and Achievement Gaps. *ETS Research Report Series* **2016** (1): 1–37
- 544 Ballen CJ, Aguilon SM, Brunelli R, Drake AG, Wassenberg D, Weiss SL, Zamudio KR, Cotner S.  
545 2018a. Do Small Classes in Higher Education Reduce Performance Gaps in STEM?  
546 *BioScience*
- 547 Ballen CJ, Danielsen M, Jørgensen C, Grytnes J-A, Cotner S. 2017a. Norway’s gender gap:  
548 classroom participation in undergraduate introductory science. *Nordic Journal of STEM*  
549 *Education* **1** (1): 179–186
- 550 Ballen CJ, Lee D, Rakner L, Cotner S. 2018b. Politics a “Chilly” Environment for Undergraduate  
551 Women in Norway. *PS: Political Science & Politics*: 1–6
- 552 Ballen CJ, Wieman C, Salehi S, Searle JB, Zamudio KR. 2017b. Enhancing diversity in  
553 undergraduate science: Self-efficacy drives performance gains with active learning. *CBE-*  
554 *Life Sciences Education* **16** (4): ar56
- 555 Bates D, Mächler M, Bolker B, Walker S. 2014. Fitting linear mixed-effects models using lme4.  
556 *arXiv preprint arXiv:1406.5823*
- 557 Beede DN, Julian TA, Langdon D, McKittrick G, Khan B, Doms ME. 2011. Women in STEM: A  
558 gender gap to innovation. *Economics and Statistics Administration Issue Brief* (04-11)
- 559 Beichner RJ, Saul JM, Abbott DS, Morse JJ, Deardorff D, Allain RJ, Bonham SW, Dancy MH, Risley  
560 JS. 2007. The student-centered activities for large enrollment undergraduate programs  
561 (SCALE-UP) project. *Research-based reform of university physics* **1** (1): 2–39
- 562 Bettinger EP, Long BT. 2005. Do faculty serve as role models? The impact of instructor gender  
563 on female students. *American Economic Review* **95** (2): 152–157
- 564 Bornmann L, Mutz R, Daniel H-D. 2007. Gender differences in grant peer review: A meta-  
565 analysis. *Journal of Informetrics* **1** (3): 226–238
- 566 Brewer CA, Smith D. 2011. Vision and change in undergraduate biology education: a call to

567 action. *American Association for the Advancement of Science, Washington, DC*

568 Brewer S, Klein JD. 2006. Type of positive interdependence and affiliation motive in an  
569 asynchronous, collaborative learning environment. *Educational Technology Research and*  
570 *Development* **54** (4): 331–354

571 Carter A, Croft A, Lukas D, Sandstrom G. 2017. Women’s visibility in academic seminars: women  
572 ask fewer questions than men. *arXiv preprint arXiv:1711.10985*

573 Cohen EG, Lotan RA. 2014. *Designing Groupwork: Strategies for the Heterogeneous Classroom*  
574 *Third Edition*. Teachers College Press.

575 Cooper KM, Krieg A, Brownell SE. 2018. Who perceives they are smarter? Exploring the  
576 influence of student characteristics on student academic self-concept in physiology.  
577 *Advances in physiology education* **42** (2): 200–208

578 Cotner S, Ballen CJ. 2017. Can mixed assessment methods make biology classes more  
579 equitable? *PLoS ONE* **12** (12): e0189610

580 Cotner S, Ballen C, Brooks DC, Moore R. 2011. Instructor gender and student confidence in the  
581 sciences: a need for more role models. *Journal of College Science Teaching* **40** (5): 96–101

582 Crombie G, Pyke SW, Silverthorn N, Jones A, Piccinin S. 2003. Students’ perceptions of their  
583 classroom participation and instructor as a function of gender and context. *The journal of*  
584 *higher education* **74** (1): 51–76

585 Cuseo J. 2007. The empirical case against large class size: adverse effects on the teaching,  
586 learning, and retention of first-year students. *The Journal of Faculty Development* **21** (1):  
587 5–21

588 Dahlerup D. 1988. From a small to a large minority: women in Scandinavian politics.  
589 *Scandinavian Political Studies* **11** (4): 275–298

590 Dillon M, Kokkelenberg EC, Christy SM. 2002. The Effects of Class Size on Student Achievement  
591 in Higher Education: Applying an Earnings Function.

592 Eddy SL, Brownell SE, Wenderoth MP. 2014. Gender gaps in achievement and participation in  
593 multiple introductory biology classrooms. *CBE-Life Sciences Education* **13** (3): 478–492

594 Freeman S, Eddy SL, McDonough M, Smith MK, Okoroafor N, Jordt H, Wenderoth MP. 2014.  
595 Active learning increases student performance in science, engineering, and mathematics.

596 *Proceedings of the National Academy of Sciences* **111** (23): 8410–8415

597 Freeman S, Haak D, Wenderoth MP. 2011. Increased course structure improves performance in  
598 introductory biology. *CBE Life Sci Educ* **10** (2): 175–186 DOI: 10.1187/cbe.10-08-0105

599 Grunspan DZ, Eddy SL, Brownell SE, Wiggins BL, Crowe AJ, Goodreau SM. 2016. Males Under-  
600 Estimate Academic Performance of Their Female Peers in Undergraduate Biology  
601 Classrooms. *PLoS ONE* **11** (2): 1–16 DOI: 10.1371/journal.pone.0148405

602 Haak DC, HilleRisLambers J, Pitre E, Freeman S. 2011. Increased structure and active learning  
603 reduce the achievement gap in introductory biology. *Science* **332** (6034): 1213–1216 DOI:  
604 10.1126/science.1204820

605 Hallgren KA. 2012. Computing inter-rater reliability for observational data: an overview and  
606 tutorial. *Tutorials in quantitative methods for psychology* **8** (1): 23

607 Hinsley A, Sutherland WJ, Johnston A. 2017. Men ask more questions than women at a scientific  
608 conference. *PloS one* **12** (10): e0185534

609 Ho DE, Kelman MG. 2014. Does class size affect the gender gap? a natural experiment in law.  
610 *The Journal of Legal Studies* **43** (2): 291–321

611 Hoffmann F, Oreopoulos P. 2009. A professor like me the influence of instructor gender on  
612 college achievement. *Journal of Human Resources* **44** (2): 479–494

613 Isbell LA, Young TP, Harcourt AH. 2012. Stag parties linger: continued gender bias in a female-  
614 rich scientific discipline. *PLoS One* **7** (11): e49682

615 Jepsen C, Rivkin S. 2009. Class size reduction and student achievement the potential tradeoff  
616 between teacher quality and class size. *Journal of human resources* **44** (1): 223–250

617 Kokkelenberg EC, Dillon M, Christy SM. 2008. The effects of class size on student grades at a  
618 public university. *Economics of Education Review* **27** (2): 221–233

619 Komarraju M, Musulkin S, Bhattacharya G. 2010. Role of student–faculty interactions in  
620 developing college students’ academic self-concept, motivation, and achievement. *Journal*  
621 *of College Student Development* **51** (3): 332–342

622 Kuh GD, Hu S. 2001. The effects of student-faculty interaction in the 1990s. *The review of higher*  
623 *education* **24** (3): 309–332

624 Ledin A, Bornmann L, Gannon F, Wallon G. 2007. A persistent problem: Traditional gender roles

625 hold back female scientists. *EMBO reports* **8** (11): 982–987

626 Lorenzo M, Crouch CH, Mazur E. 2006. Reducing the gender gap in the physics classroom.  
627 *American Journal of Physics* **74** (2): 118–122

628 Micari M, Pazos P, Hartmann MJZ. 2007. A matter of confidence: gender differences in attitudes  
629 toward engaging in lab and course work in undergraduate engineering. *Journal of Women*  
630 *and Minorities in Science and Engineering* **13** (3)

631 National Academies of Sciences and Medicine E. 2016. *Barriers and opportunities for 2-year*  
632 *and 4-year stem degrees: systemic change to support students' diverse pathways*. National  
633 Academies Press.

634 O'Dorchai S, Meulders D, Crippa F, Margherita A. 2009. *She figures 2009—Statistics and*  
635 *indicators on gender equality in science*. Publications Office of the European Union.

636 Pomerantz EM, Altermatt ER, Saxon JL. 2002. Making the grade but feeling distressed: Gender  
637 differences in academic performance and internal distress. *Journal of Educational*  
638 *Psychology* **94** (2): 396

639 Premo J, Cavagnetto A. 2018. Priming students for whole-class interaction: Using  
640 interdependence to support behavioral engagement. *Social Psychology of Education*: 1–21

641 Pritchard J, Masters K, Allen J, Contenta F, Huckvale L, Wilkins S, Zocchi A. 2014. Asking gender  
642 questions. *Astronomy & Geophysics* **55** (6): 6–8

643 Rowe MB. 1974a. Wait-time and rewards as instructional variables, their influence on language,  
644 logic, and fate control: Part one-wait-time. *Journal of research in science teaching* **11** (2):  
645 81–94

646 Rowe MB. 1974b. Relation of wait-time and rewards to the development of language, logic, and  
647 fate control: Part II-Rewards. *Journal of research in science teaching* **11** (4): 291–308

648 Rusbult CE, Van Lange PAM. 2008. Why we need interdependence theory. *Social and*  
649 *Personality Psychology Compass* **2** (5): 2049–2070

650 Schanzenbach DW. 2014. Does Class Size Matter?

651 Schmidt SJ, Davenport JRA. 2017. Who asks questions at astronomy meetings? *NATURE* **1**  
652 (0153): 1

653 Schmidt SJ, Douglas S, Gosnell NM, Muirhead PS, Booth RS, Davenport JRA, Mace GN. 2017. The

654 Role Of Gender In Asking Questions At Cool Stars 18 And 19. *arXiv preprint*  
655 *arXiv:1704.05260*

656 Simpson EH. 1949. Measurement of diversity. *nature*

657 Smith MK, Wood WB, Adams WK, Wieman C, Knight JK, Guild N, Su TT. 2009. Why Peer  
658 Discussion Improves Student Performance on In-Class Concept Questions. *Science* **323**  
659 (5910): 122 LP-124 Available at:  
660 <http://science.sciencemag.org/content/323/5910/122.abstract>

661 Steele CM, Spencer SJ, Aronson J. 2002. Contending with group image: The psychology of  
662 stereotype and social identity threat. In *Advances in Experimental Social*  
663 *Psychology* Academic Press; 379–440.

664 Tanner KD. 2013. Structure matters: twenty-one teaching strategies to promote student  
665 engagement and cultivate classroom equity. *CBE-Life Sciences Education* **12** (3): 322–331

666 Team RC. 2013. R: A language and environment for statistical computing

667 Theobald E. 2018. Students are rarely independent: When, why, and how to use random effects  
668 in discipline-based education research. *CBE—Life Sciences Education* **17** (3): rm2

669 Wold A, WENNERÁS C. 2010. Nepotism and sexism in peer-review. In *Women, Science, and*  
670 *Technology* Routledge; 64–70.

671 Zinth K. 2005. State class-size reduction measures. *Denver: Education Commission of the States*  
672