

Differential Performance in High vs. Low Stakes Tests: Evidence from the GRE test

Yigal Attali
Educational Testing Service
Rosedale Rd.
MS-16-R
Princeton, NJ 08541
USA
Voice: 609-734-1747
Fax: 609-734-1755
e-mail: yattali@ets.org

Zvika Neeman
The Eitan Berglas School of Economics
Tel Aviv University
P.O.B. 39040
Ramat Aviv, Tel Aviv, 69978
ISRAEL
Office: +972-3-6409488
Fax: +972-3-6409908
e-mail: zvika@post.tau.ac.il

Analia Schlosser
The Eitan Berglas School of Economics
Tel Aviv University
P.O.B. 39040
Ramat Aviv, Tel Aviv, 69978
ISRAEL
Office: +972-3-6409064
Cel: +972-54-4902414
Fax: +972-3-6409908
e-mail: analias@post.tau.ac.il

Differential Performance in High vs. Low Stakes Tests: Evidence from the GRE test¹

Yigal Attali

Educational Testing Service

Zvika Neeman

Tel Aviv University

Analia Schlosser

Tel Aviv University

July, 2018

Abstract

We study how different demographic groups respond to incentives by comparing their performance in “high” and “low” stakes situations. The high stakes situation is the GRE examination and the low stakes situation is a voluntary experimental section of the GRE that examinees were invited to participate in after completing the GRE. We find that Males exhibit a larger drop in performance between the high and low stakes examinations than females, and Whites exhibit a larger drop in performance compared to Asians, Blacks, and Hispanics. Differences in performance between high and low stakes tests are partly explained by the fact that males and whites exert lower effort in low stakes tests compared to females and minorities.

¹ We thank comments received at the SOLE meetings, “Discrimination at Work” and “Frontiers in Economics of Education” workshops, and seminar participants at the The Federal Reserve Bank of Chicago, CESifo, Norwegian Business School, University of Zurich, Bar Ilan University, Ben Gurion University, and University of Haifa. This research was supported by the Israeli Science Foundation (grant No. 1035/12).

1. Introduction

Recently, there has been much interest in the question of whether different demographic groups respond differently to incentives and competitive pressure. Interest in this subject stems from attempts to explain gender, racial, and ethnic differences in human capital accumulation and labor market performance, and is further motivated by the increased use of aptitude tests for college admissions and job screening and the growing use of standardized tests for the assessment of students' learning. While it is clear that motivation affects performance, less attention has been given to demographic group differences in response to performance based incentives.

In this paper, we examine whether individuals respond differently to incentives by analyzing their performance in the Graduate Record Examination General Test (GRE).¹ We examine differences in response to incentives between males and females as well as differences among Whites, Asians, Blacks, and Hispanics. Specifically, we compare performance in the GRE examination in "high" and "low" stakes situations. The high stakes situation is the real GRE examination and the low stakes situation is a voluntary experimental section of the GRE test that examinees were invited to take part in immediately after they finished the real GRE examination.

A unique characteristic of our study is that we observe individuals' performance in a "real" high stakes situation that has important implications for success in life and that is administered to a very large and easily characterizable population, namely the population of applicants to graduate programs in arts and sciences the US. This feature distinguishes our work from most of the literature, which is usually based on controlled experiments that require individuals to perform tasks that might not bear directly on their everyday life, and that manipulate the stakes, degree of competitiveness, or incentive levels in somewhat artificial ways, and where stakes are not as high as in real-life important event. A second distinctive feature of our research is that we are able to observe performance of the same individual in high and low stakes situations that involve the exact same task. A third unique feature of our study is the availability of a rich data on individuals' characteristics that includes information on family background, college major and academic performance, and intended graduate field of studies. These comprehensive data allow us to compare individuals of similar academic and family backgrounds and examine the persistence of our

¹ The GRE test is a commercially-run psychometric examination that is part of the requirements for admission into most graduate programs in arts and sciences in the US and other English speaking countries. Each year, more than 600,000 prospective graduate school applicants from approximately 230 countries take the GRE General Test. The exam measures verbal reasoning, quantitative reasoning, critical thinking, and analytical writing skills that have been acquired over a long period of time and that are not related to any specific field of study. For more information, see the ETS website: <http://www.ets.org/gre/general/about/>.

results across different subgroups. A fourth important advantage of our study is that we are able to observe the selection of individuals into the experiment and examine the extent of differential selection within and across groups. Notably, we do not find any evidence of differential selection into the experiment, neither according to gender, race or ethnicity, nor according to individual's scores in the "real" GRE exam.

Our results show that males exhibit a larger difference in performance between the high and low stakes GRE test than females and that Whites exhibit a larger difference in performance between the high and low stakes GRE test compared to Asians, Blacks, and Hispanics. A direct consequence of our findings is that test score gaps between males and females or between Whites and Blacks or Hispanics are larger in a high stakes test than in a low stakes test, while the test score gap between Asians and Whites is larger in the low stakes test. Specifically, while males outperform females in the high stakes quantitative section of the GRE by .55 standard deviations (SD), the gender gap in performance in the low stakes section is only .30 SD. Similarly, males' advantage in the high stakes verbal section is .26 SD while the gender gap in the low stakes section is only .07 SD. Whites outperform Blacks and Hispanics in the high stakes quantitative section by 1.1 SD and .42 SD, respectively, but the gaps are significantly reduced in the low stakes section to .63 and .14 SD. This pattern is reversed for Asians because they outperform whites by .51 SD in the high stakes quantitative section, so that the gap increases to .55 SD in the low stakes section. These group differences in performance between high and low stakes tests appear across all undergraduate GPA levels, family backgrounds (measured by mother's education), and even among students with similar orientation towards math and sciences (identified by their undergraduate major or intended graduate field of studies).

We explore various alternative explanations for the differential response to incentives across demographic groups and show that the higher differential performance of males and whites between the high and the low stakes test is partially explained by lower levels of effort exerted by these groups in the low stakes situations compared to women and minorities, respectively. We do not find evidence supporting alternative explanations such as test anxiety or stereotype threat.

Our findings imply that inference of ability from cognitive test scores is not straightforward: differences in the perceived importance of the test can significantly affect the ranking of individuals by performance and may have important implications for the analysis of performance gaps by gender, race, and ethnicity. The results from our paper have two main implications:

- (1) Stakes have to be taken into account when analyzing performance gaps between groups

(2) Some groups are mostly driven by incentives while other groups exert high effort even if stakes are low or “nearly zero”.

While these two implications do not, in themselves, amount to direct policy recommendations, they are nevertheless highly relevant for policy. For example, they imply that any analysis of gender or race test score gaps, or studies that examine the effect of a specific educational intervention by gender or race, should take into account the stakes of the test involved in order to interpret the results and effectiveness of the intervention. In addition, our results highlight the fact that university or job admission policies that use standardized aptitude tests should take into account that such tests measure only performance under a high stakes setup and are less informative about individuals’ performance in low stakes or zero stakes situations, which may be as important at the university or job.

Most of the experimental literature about gender differences in performance focuses on a comparison of performance between a competitive setting where the best performer receives a higher payment and a non-competitive environment where subjects are paid according to their own performance (using a piece-rate schedule). A common finding in these studies is that while the performance of men improves under competition, women’s performance is unchanged or even declines slightly (see, e.g. Gneezy et al., 2003, and Gneezy and Rustichini, 2004). A second finding is that women “shy away from competition.” Namely, given the choice, women prefer to be compensated according to a non-competitive piece-rate compensation schedule over participation in competitive tournaments (see, e.g., Datta Gupta et al., 2005; Niederle and Vesterlund, 2007; Dohmen and Falk, 2011).

There are several variations and extensions to these studies that examine whether the results vary by: (a) the gender composition of the group involved in the tournament; (b) the type of task involved (tasks requiring effort vs. skills, or tasks where males or females have a stereotypical or real advantage); (c) the information provided about own and others’ performance during the experiment; (d) the use of priming; (e) letting participants choose the gender of their competitors; (f) manipulating the risk associated with the payments; and (g) the number of iterations involved. For recent reviews of this literature, see Croson and Gneezy (2009), Azmat and Petrongolo (2014), and Niederle (2016).

Our paper differs from these previous studies in several aspects: first, we compare performance between a high stakes setting that has important consequences for life and a task that has almost zero stakes. In a sense, this is more similar to a comparison between performance under a piece-rate and a flat-rate payment scheme. Second, even though GRE scores are also reported in percentiles, the exam is not presented as a direct tournament between subjects (certainly not among those tested in a specific

date and test center).² Accordingly, the focus of our study is not a comparison between a competitive and a non-competitive environment but rather a contrast between a high stakes and a very low stakes setting. As our results show, males invest less effort than females when stakes are low. We therefore add new insights to the experimental literature cited above by suggesting that gender differences found in these lab experiments may significantly understate differences in important real life situations given that stakes levels of lab experiments are relatively low.

Evidence on gender differences in real world situations is limited to a small number of recent studies and remains an important empirical open question. Paserman (2010) studies performance of professional tennis players and finds that performance decreases under high competitive pressure but this result is similar for both men and women. Similarly, Lavy (2008) finds no gender differences in performance of high school teachers who participated in a performance-based tournament. On the other hand, in a field experiment among administrative job seekers, Flory et al. (2010) find that women are indeed less likely to apply for jobs that include performance based payment schemes but this gender gap disappears when the framing of the job is switched from being male- to female-oriented.³

A number of studies within the educational measurement literature demonstrate that high stakes situations induce stronger motivation and higher effort.⁴ However, high stakes also increase test anxiety and so might harm performance (Cassaday and Johnson, 2002). Indeed, Ariely et al. (2009) found that strong incentives can lead to “choking under pressure” both in cognitive and physical tasks, although they did not find gender differences. Performance in tests is also affected by noncognitive skills as shown by Heckman and Rubinstein (2001), Cunha and Heckman (2007), Borghans et al. (2008), and Segal (2010).⁵

Levitt et al. (2016) examine how timing, type of rewards, and framing of rewards affect performance in a series of field experiments involving primary and secondary school students in Chicago. They report that in most cases, boys were more likely to respond to incentives than girls were. Azmat et al. (2016) is the closest paper to ours. They exploited the variation in the stakes of tests administered to students

² While GRE test scores are relative to other students, the competition between students is less salient on the day of the exam as the pool of competitors is very large and not directly visible or known ex ante to GRE test takers.

³ Other studies that compare gender performance by degree of competitiveness include Jurajda and Munich (2011) and Ors et al. (2008).

⁴ For example, Cole et al. (2008) show that students’ effort is positively related to their self reports about the interest, usefulness, and importance of the test; and that effort is, in turn, positively related to performance. For a review of the literature on the effects of incentives and test taking motivation see O’Neil, Surgue, and Baker (1996).

⁵ Several studies (see e.g., Duckworth and Seligman, 2006; and the references therein) suggest that girls outperform boys in school because they are more serious, diligent, studious, and self-disciplined than boys. Other important noncognitive dimensions that affect test performance are discussed by the literature on stereotype threat that suggests that performance of a group is likely to be affected by exposure to stereotypes that characterize the group (see Steele, 1997; Steele and Aronson, 1995; and Spencer et al., 1999).

attending a Spanish private school and show that performance of female students declines as the stakes become higher while males' performance improves. Their finding is consistent with ours, but we examine the performance of a much larger population (GRE test takers) and show gender differences in response to incentives across a wide range of students' background characteristics, fields of study, and ability levels. In addition, we are able to explore the role played by students' effort in explaining our findings, and rule out some alternative explanations (including females' choking under pressure). Our study also expands the literature by examining differential performance by race and ethnicity. To the best of our knowledge, no other study has examined differences in response to incentives among ethnic groups.

Our paper is also related to Babcock et al. (2017) who find that women, more than men, volunteer, are asked to volunteer, and accept requests to volunteer for "low promotability" tasks. Their results suggest that women's higher tendency to volunteer seems to be shaped by women's beliefs rather than preferences. Accordingly, Babcock et al. suggest several alternative assignment schemes to reduce the gender gap in participation in low stakes activities such as turn-taking or random assignment.

In our study, the decision to participate in the low stakes task, which is analogous to "volunteering," does not generate a group benefit as in Babcock et al. However, we examine not just willingness to participate in the low stakes task, but also effort exerted conditional upon participation. That is, our setting contains both the binary decision of whether to volunteer or not, as well as a continuous decision with respect to how much effort to exert after volunteering. Our results show that while men and women are equally likely to volunteer, the performance of men is significantly lower. Our results therefore suggest that even if men and women are randomly assigned to participate in a certain committee, women might invest more time and effort conditional on participation. Consequently, a random assignment mechanism might not overcome the problem of inequality in investment in "*low promotability*" tasks.

The rest of the paper proceeds as follows. In the next section we describe the experimental setup and data. In Section 3, we present the empirical framework. In Section 4 we present the results and in Section 5 we discuss alternative explanations for our findings as well as other related observations. Section 6 concludes.

2. Experimental Set-up and Data

We use data from a previous study conducted by Bridgeman et al. (2004), whose purpose was to examine the effect of time limits on performance in the GRE Computer Adaptive Test (CAT) examination. All examinees who took the GRE CAT General Test during October-November 2001 were invited to participate in an experiment. At the end of the regular test, a screen appeared that invited examinees to

voluntarily participate in a research project that would require them to take an additional test section for experimental purposes.⁶ GRE examinees who agreed to participate in the experiment were promised a monetary reward if they perform well compared to their performance in the real examination.⁷

Participants in the experiment were randomly assigned into one of four groups: one group was administered a quantitative section (Q-section) with standard time limit (45 minutes), a second group was administered a verbal section (V-section) with standard time limit (30 minutes), the third group was administered a quantitative section with extended time limit (68 minutes) and the fourth group was administered a verbal section with extended time limit (45 minutes). The research sections were taken from regular CAT pools (over 300 items each) that did not overlap with the pools used for the real examination. The only difference between the experimental section and the real sections was the appearance of a screen that indicated that performance on the experimental section did not contribute to the examinee's official test score. We therefore consider performance in the real section to be performance in a high stakes situation and performance in the experimental section to be performance in a low stakes (or almost zero stakes) situation. Even though a monetary reward based on performance was offered to those who participated in the experiment, it is clear that success in the experimental section was less significant to examinees and involved less pressure. More importantly, since the monetary reward was conditional on performance relative to one's own achievement in the high stakes section rather than on absolute performance, incentives to perform well in the experimental section were similar for all participants in the experiment.

Appendix Table A1 shows details of the construction process of our analysis sample. From a total of 81,231 GRE examinees in all centers (including overseas), 46,038 were US citizens who took the GRE test in centers located in the US. We focus on US citizens tested in the US to avoid dealing with a more heterogeneous population and to control for a similar testing environment. In addition, we want to abstract from differences in performance that are due to language difficulties. 15,945 out of the 46,038 US examinees agreed to participate in the experiment. About half of them (8,232) were randomized into the regular time limit sections and were administered either an extra Q-section (3,922) or an extra V-

⁶ Students saw their score in the regular test only after the experimental section. They were never told their score in the experimental section.

⁷ Specifically, the instructions stated "It is important for our research that you try to do your best in this section. The sum of \$250 will be awarded to each of 100 individuals testing from September 1 to October 31. These awards will recognize the efforts of the 100 test takers who score the highest on questions in the research section relative to how well they did on the preceding sections. In this way, test takers at all ability levels will be eligible for the award. Award recipients will be notified by mail." See Bridgeman et al. (2004) for more details about the experiment design and implementation.

section (4,310).⁸ We select only experiment participants who were randomized into the regular time limit experimental groups because we are interested in examining differences in performance in the exact same task that differs only by the stake examinees associated with it.⁹

A unique feature of our research design that distinguishes our study from most of the experimental literature is that we are able to identify and characterize the experiment participants out of the full population of interest (i.e., GRE examinees in our case). Table 1 compares the characteristics of the full sample of US GRE test takers and the sample of experiment participants.¹⁰ The two populations are virtually identical in terms of proportions of females, males, and minorities. For example, women comprise 66 percent of the full population of US domestic examinees while the share of women among those who agreed to participate in the Q or the V section was 65 and 66 respectively. Likewise, whites make up about 78 percent of GRE US domestic examinees and they are equally represented among experiment participants. The shares of Blacks, Hispanics, and Asians range between 6 and 5.5 percent in both the full sample and the sample of experiment participants.¹¹

Participants in the experiment also have similar GRE test scores to those in the full relevant sub-population from which they were drawn. For example, males are located, on average, at the 56 percentile rank of the Q-score distribution, which is equal to the average performance of male participants in the experiment. The median score (57 percentile rank) and standard deviation (27 points) are also identical for the full sample of GRE US male test takers, the sample of experiment participants randomized to the Q-section, and the sample of experiment participants randomized to the V-section. The test score distribution of female GRE test takers is also identical to that of female experiment participants. We observe also the same result when comparing test score distributions within each race/ethnicity. Overall, the results presented in Table 1 show that there is no differential selection into the experiment according

⁸ Since the experimental sections were randomized among the full sample of experiment participants, which included all students (US and international) tested in all centers around the world, the proportion of US participants assigned to each section is not exactly 50 percent.

⁹ One limitation of our study is that we were not able to randomize the order of the tests, so that all examinees received the low stakes test after the high stakes test. As we discuss below, we believe this constraint does not affect our main results or interpretation.

¹⁰ Due to data restrictions we cannot compare experiment participants to non-participants because we received the data on experiment participants and the data on the full population of GRE examinees in two separate datasets that lacked individual identifiers.

¹¹ Reported proportions by race/ethnicity do not add up to one because the following additional groups are not reported in the table: American Indian, Alaskan, and examinees with missing race/ethnicity.

to gender, race/ethnicity or GRE test scores, nor do we find any evidence of differential selection within each gender or race/ethnic group.¹²

GRE test takers are required to fill out a form upon registration to the exam. The form collects information on basic background characteristics, college studies, and intended graduate field of studies.¹³ Appendix Table A2 reports descriptive statistics of these background characteristics for the sample of experiment participants stratified by gender, race, and ethnicity. Note that the comparisons presented here are across the population of GRE test takers, which is a selected sample of college students, and therefore they do not represent group differences across the population of college students but rather differences across college students who intend to pursue graduate studies.

Averages reported in columns 2 and 3 of Table A2 show that males and females come from similar family backgrounds as measured by both mother's and father's educational levels and by the proportion of native English speakers. Females and males have also similar distributions of undergraduate GPA (UGPA). Nevertheless, males are more likely to come from undergraduate majors in math, computer science, physics or engineering and they are also more likely to intend to pursue graduate studies in these fields (26 percent for males versus 5 percent for females).

Columns 3 through 6 in Table A2 report descriptive statistics of the analysis sample stratified by race/ethnicity. Maternal education is similar among Whites and Asians but Asians are more likely to have a father with at least some graduate studies or a professional degree relative to Whites (45 versus 35 percent). Hispanics and Blacks come from less educated families. Asians are less likely to be native English speakers (86 percent) relative to Whites (93 percent), Blacks (95 percent), and Hispanics (90 percent). In terms of undergraduate achievement, we observe that Whites and Asians have similar UGPAs distributions but Hispanics and Blacks have, on average, lower UGPAs. Asians are more likely to do math, science, and engineering either as an undergraduate major or as an intended field of graduate studies (30 percent) relative to Whites (11 percent), Blacks (8 percent), or Hispanics (12 percent).

¹² While we do not find differences in observable characteristics, there could still be differences in unobserved characteristics. Nevertheless, for the purpose of our study, we should worry about differential selection into the experiment by unobservables across demographic groups. The fact that we did not find evidence for differential selection across groups according to observables suggests that the presence of large differences in selection by unobservables across groups is very unlikely.

¹³ We obtained the background information on experiment participants only so we only analyze selection in the experiment according to gender, race, ethnicity, and GRE scores in the high stakes section.

3. Empirical Framework

Our main objective is to examine how performance of different demographic groups changes as a function of the stakes of the test (high stakes: real GRE exam and low stakes: experimental section). We summarize our main finding in Figure 1 using an ordinal metric, which is free of the specific scale of test scores. We ranked individuals according to their performance in each test and plot the rank change distribution (in percentile points) between the high and low stake test by gender and race for each test. Panels (a) and (b) show that men's ranking declines by 4 percentile points in the low stakes test relative to the high stakes test while women's ranking improves by 2 percentile points. Panels (c) and (d) show that ranking of whites declines while the ranking of minorities improves when switching from the high to the low stakes test in both the Q- and the V-sections. Focusing on the Q-section, which is less likely to be affected by language problems of minorities we see that whites' ranking declined by almost one percentile points while that of minorities improved by about 5 percentile points.¹⁴ The rank changes between men and women and between whites and minorities are statistically different (p-values of Mann-Whitney tests <0.0001).

We now turn to measure individuals' change in performance using a simple regression model to control for additional individuals' characteristics and quantify the average change in performance between the high and low stakes test for each group. We estimate the following first difference equation for each of the experimental samples (i.e. individuals randomized to the experimental Q or V section):¹⁵

$$(1) Y_{iHS} - Y_{iLS} = \beta_0 + \beta_1 Female_i + \beta_2 Black_i + \beta_3 Hispanic_i + \beta_4 Asian_i + \beta_5 Other_i + x_i' \gamma + u_i$$

where Y_{iHS} denotes the test score of individual i in the high stakes section; Y_{iLS} is the test score of individual i in the low stakes section; x is vector of individual characteristics that includes the following covariates: mother's and father's education, dummies for UGPA, undergraduate major, intended graduate field of studies, and disability status. *Female*, *Black*, *Hispanic*, *Asian*, and *Other* are dummy variables for the gender and race/ethnicity of the examinee.¹⁶ Whites and males are the omitted categories. The coefficients of interest are $\beta_1, \beta_2, \beta_3, \beta_4$ that denote the difference in performance gap between the high and the low stakes test of the relevant group (Females or Blacks/Hispanics/Asian) relative to the omitted category (Males or Whites). To simplify the exposition, we reverse the sign of the coefficients and report in all tables differences between males and females and differences between Whites and Blacks/Hispanics/Asians.

¹⁴ Minorities include Asians, Hispanics, and Blacks. We excluded students who defined themselves as American Indian or Alaskan Native (43) or other race (271).

¹⁵ Note that at that time, there was only one Q/V section. The high stake GRE score was based on all items in that section.

¹⁶ Race/ethnicity categories in the GRE form are exclusive (i.e., it is not possible to check more than one option).

Note that by using a first difference specification we are differencing out an individual's fixed effect that accounts for all factors that affect examinee's performance in both the low stakes and the high stakes test. By including a vector of covariates we allow for individual's characteristics to affect the change in performance between the high and low stakes situation.¹⁷

GRE scores in the quantitative and verbal sections range between 200 and 800, in 10-point increments. To ease the interpretation of the results, we transformed these raw scores into percentile ranks using the GRE official percentile rank tables.¹⁸ All results presented below are based on GRE percentile ranks. As we show below, we obtain similar results when using raw scores, log of raw scores or z-scores.

4. Results

4.1. Differences in Performance by gender, race, and ethnicity

Panel A of Table 2 exhibits examinees' performance in the high stakes test for males, females, whites, blacks, Hispanics, and Asians and the gaps between groups.¹⁹ Similar to other comparisons of GRE scores by gender, males outperform females in both the quantitative and verbal sections among the participants in our experiment. On average, Males are placed about 15.3 percentile points higher in the test score distribution of the Q-section relative to females. The gender gap in the V-section is smaller but still sizable, with males scoring about 6.5 percentile points higher than females. Asians have the highest achievements among all ethnic/racial groups in the Q-section. Their test scores are about 15 percentile points above Whites. Hispanics lag behind Whites by an average of 10.6 percentile points. Q-scores of Blacks are lower and they are placed, on average, about 25 percentile points below Whites in the test score distribution.

¹⁷ An alternative approach is to estimate a conditional model that regresses the score in the low stakes test on the score in the high stakes test. The score change model described in equation (1) and the conditional regression model both attempt to adjust for baseline outcomes but they answer different questions. The score change model examines how groups, on average, differ in score changes between the high and the low stakes test. The conditional regression model asks whether the score change of an individual who belongs to one group differs from the score change of an individual who belongs to another group under the assumption that the two had come from a population with the same baseline level. The two approaches are expected to provide equivalent answers when the groups have similar baseline outcomes. However, as discussed by Cribbie and Jamieson (2000), when baseline means differ between groups, conditional regression suffers from directional bias. Namely, conditional regression augments differences when groups start at different levels and then remain parallel or diverge (see Lord's Paradox - Lord, 1967) and attenuates differences when groups start at different levels and then converge. Because the demographic groups we examine have different baseline GRE performance, we choose to estimate models of score change.

¹⁸ For more information regarding on the interpretation of GRE scores, exam administration and validity see Educational Testing Service (2007).

¹⁹ The percentile scores of males and females do not add to 100 since they are constructed using the official GRE tables, which include also international examinees and are based on several years of data.

In the verbal section, Whites outperform Asians although the difference between groups is not statistically significant. The gap between Whites and Blacks is a bit smaller (23 percentile points) while the gap between Whites and Hispanics is about 12 percentile points. With the exception of Whites vs. Asians in the verbal section, all gaps between groups in the high stakes section are statistically significant.

Panel B of Table 2 reports students' performance in the experimental section and gaps by gender and race/ethnicity. On average, performance in the low stakes test is lower than in the high stakes test for all groups. Notably, gaps between males and females or whites and blacks or Hispanics are narrower in the experimental section (even though they are still statistically significant). For example, the score gap between males and females shrinks from 15 to 11 percentile points in the Q-section and from 7 to 2 percentile points in the V-section. The score gap between Whites and Blacks shrinks from 25 to 19 percentile points in the Q-section and from 23 to 18 in the V-section and the gap between Whites and Hispanics shrinks from 11 to 5 percentile points in the Q-section and from 12 to 11 percentile points in the V-section. The gap between Asians and whites in the Q-section widens between the high and the low stake test (from 15 to 18 percentile points) because Asians outperform whites in this exam.

Table 3 reports the change in performance between the high and the low stakes section for each demographic group (first row of each panel) and the difference (second and third row) in the drop in performance between males and females or between whites and Blacks/Hispanics/Asians. Males' performance drops by 11.6 percentile points from the high to the low stakes Q-sections while females' performance drops by only 7.1 points. The gap in the drop in performance between males and females is significant and stands at 4.5 percentile points (s.e.=0.784). That is, a switch from the high to the low stakes situation narrows the gender gap in the quantitative test by about 4.5 percentile points (although is still significant), which is equivalent to a 30 percent drop in the gender gap of the high stakes test. The differential change in performance remains almost unchanged after controlling for individual's background characteristics and academic achievement. This finding is important as it suggests that our results are unlikely to be driven by differences in family background and academic achievement.

We also find a similar gender gap in the V-section. Males' scores drop by 10.4 percentile points, on average, while females' scores drop by a smaller magnitude of 6.1 percentile points. That is, males' scores drop by 4.3 percentile points (s.e.=0.783) more relative to females. Note that the proportional drop in males' performance is also larger than females'. Namely, males' scores drop by 21 percent while females' scores drop by 18 percent in the Q-section. Similarly, we find that males' scores in the V-section drop by 17 percent while females' scores drop by 11 percent.

The stratification by race/ethnicity shows that whites exhibit the largest drop in performance between the high and the low stakes Q-section. Whites' performance drops by 9.4 percentile points, while that of Asians drops by 7 percentile points, Blacks' performance drops by 3 percentile points, and Hispanics' performance drops by 3.8 percentile points. Differences in the performance drop between Whites and each of the minority groups are all significant. The controlled difference between Whites and Blacks, after accounting for individual's characteristics, is of 4.3 percentile points (s.e.=1.05). The equivalent difference between Whites and Hispanics is 5.21 (s.e.=1.40) and the difference between Whites and Asians is 3.2 (s.e.=1.70). In the verbal section, the performance drop from the high to the low stakes section is larger for Whites than for Blacks (7.8 percentile points versus 2.3 percentile points). But Hispanics and Asians exhibit a similar drop in performance to that of Whites. We suspect that the different pattern obtained for Asians and Hispanics in the V-section could be related to language dominance.

Overall, the evidence presented in Table 3 shows that males and Whites exhibit the largest drop in performance between the high and the low stakes tests compared to females and minorities. Our results are robust to nonlinear transformations and alternative definitions of the dependent variable as reported in Appendix Table A3. In the first row of panels A and B, we report differences in performance in the quantitative and verbal sections using raw scores (scaled between 200 and 800). In the second row of each panel, we show differences in performance using the natural logarithm of raw scores. In the third row, we report results based on z-scores.²⁰ All alternative metrics yield results that are equivalent to our main findings: males' drop in performance between the high and low stakes section is 5 percent or .17 SD larger than the drop of females; whites' drop in performance in the Q-section is 8 percent or .23 SD larger than the drop of blacks; 7 percent or .23 SD larger than the drop of Hispanics and 7 percent or .19 SD larger than the drop of Asians. These additional results show that our findings are not driven by a specific scale used to measure achievement. Furthermore, as we show in Figure 1, we obtain the same results when we rely only on the ordinal information embedded in scores.

The fourth row of each panel in Table A3 replicates our main results using the samples of examinees randomized into experimental sections with extended time limit (67.5 minutes for the Q-section and 45 minutes for the V-section). Estimates are similar to our main results showing that our findings are replicable in additional settings. In addition, they demonstrate that our results are not sensitive to time constraints or differential responses by gender or ethnicity to the length of the exam.

²⁰ Z-scores are computed using the mean and standard deviation of the high stakes test.

We also examine how the change in performance varies by students' performance in the high stakes exam. To examine this issue we divide the high stakes score distribution for each group into deciles and define for each individual his/her score decile in the high and low stakes section. We plot in Figure 2 the average score decile of the low stakes section as a function of the score decile in the high stakes section by gender and race. Overall, with the exception of those located at the bottom of the test score distribution in the high stakes section, there is a similar drop in performance (in percentage terms) in all parts of the high stakes score distribution with males having a larger drop relative to females and whites having a larger drop in performance relative to minorities.

Another relevant question is whether the results are driven by a small group of males or whites that has a large performance drop or are evident among most individuals who belong to those demographic groups. Figure 3 plots the CDF of the difference in score (measured in percentiles) between the high and low stakes section by gender/race and section. For most individuals the change in performance is of a few percentile points but males have a larger drop in performance than females. In addition, a larger proportion of males has a substantial drop in performance relative to females. The same pattern is observed for whites versus minorities in the Q-section: whites have a larger drop in performance relative to minorities and those who have a very large drop in performance are disproportionately represented by whites.

We further explore this issue by re-estimating our main model after dropping from each demographic group those individuals with the largest drop in performance (i.e., those in the top 10-percentile distribution of the performance change in their demographic group). Results from this subsample (reported in the last row of Appendix table A3) show that differences between demographic groups in performance change are very similar to differences observed for the full sample. Again, males and whites have the larger drop in performance relative to females and minorities. This implies that results are not only driven by a few extreme values of a specific demographic group.

4.2 Within Race/Ethnicity and Gender Differences in Performance

We check for gender and race/ethnicity interactions by examining whether differences between males and females appear across all race/ethnic groups and whether differences between Whites and minorities show up for males and for females.²¹

²¹ The conclusions described in this subsection rely on samples that are stratified by gender and race/ethnicity and that are relatively small for Blacks, Hispanics, and Asians so the results should be taken with caution.

Table 4 reports performance in the high and low stakes section for each gender and ethnicity/race as well as differences in performance between males and females within each race/ethnicity and between Whites and minorities for males and females separately. We focus in the Q-section as performance is less influenced by language constraints among Hispanics and Asians. The results show that White males have the largest differential performance between the high and the low stakes test compared to Black, Asian, and, Hispanic males. We obtain a similar result for females with the exception of Asian females who behave similarly to White females.

Comparisons between males and females within each race/ethnicity group reveal that males exhibit a larger drop in performance relative to females among Whites, Blacks, and Hispanics although differences between genders are only statistically significant among Whites. In contrast, we observe no gender differences among Asians. In fact, the drop observed among females is even larger than the drop observed among males, although the difference is not statistically significant.

4.3 Heterogeneous effects

Table 5 reports the gender gap in students' performance in high and low stakes tests for different subsamples stratified by undergraduate GPA (UGPA), student's major, intended field of graduate studies, and mother's education. We focus on gender gap and not on gap by race/ethnicity since subgroups are too small for that stratification. Panel A reports results for the Q-section and panel B reports results for the V-section. Rows 1 through 5 in both panels present estimates for the samples stratified by UGPA. As expected, students with higher UGPA have higher scores in both the high and the low stakes sections of the quantitative and verbal exams. Males' advantage in the high stakes test appears across all cells of the UGPA distribution both in the quantitative and the verbal sections. Again, we observe that the gender gap in performance is narrower in the low stakes section in each of the cells stratified by UGPAs and is even insignificant when comparing performance in the V-section between male and female students with an UGPA of A, A- or B-.

We see in columns 9 and 10 of the table that all students, regardless of their UGPA exhibit a significant drop in performance between the high and the low stakes sections (both the quantitative and the verbal).²² Males' performance drop is larger than females' drop across all levels of UGPA (see columns 11 and 12) and is evident both in absolute and percentage terms.

²² We use UGPA to stratify the sample (instead of using the score in the high stakes section) because it provides a measure of students' performance that is taken independently and before the realization of the dependent variable.

The next two rows of Table 5 (in both panels A and B) report the gender gap in performance for the sample of students who majored in math, computer science, physics or engineering or who intend to pursue graduate studies in one of these fields (to simplify the discussion we will call them math and science students). We focus on these students to target a population of females that is expected to be highly selected.²³ While females represent the majority among the full population of GRE examinees (65 percent) they are a minority among math and science students (26 percent). It is therefore interesting to examine whether we find the same results in a subsample where selection by gender goes in the opposite direction.

As seen in columns 3 and 4 of table 5, achievement in the GRE Q-section is much higher among math and science students relative to the full sample and even relative to those students whose UGPA is an “A”. Math and science students also attain higher scores in the V-section relative to the full sample but they score slightly lower compared to those students with an “A” UGPA. The gender gap in the high stakes Q-section among math and science students is smaller (8.7 percentile points) than the gender gap in the full sample (15.3 percentile points), although we still observe that males have higher achievement than females. The gender gap among those who intend to pursue graduate studies in these fields is even narrower (7.1 percentile points) although still significant. In contrast, there is no gender gap achievement in the V high stakes section in the subsamples of math and science students.

Achievement of math and science students in the Q low stakes section is lower than in the high stakes section but these students still perform better relative to other students in the low stakes section. Consistent with our previous results, the gender gap in Q performance among math and science students is narrower in the low stakes section relative to the high stakes section and is even insignificant. The pattern for the V section is similar with math and science females even outperforming their male counterparts in the low stakes V-section.

Even in this subsample of math and science students, the drop in performance between the high and the low stakes test is larger for males (who reduce their performance by about 12-13 percentile points in both subjects) compared to females (who reduce their performance by 6-7 percentile points in the Q section and by 4-5 percentile points in the V section). The larger drop in males’ performance is evident both in absolute terms and relative to the outcome means in the high stakes test. The gender differences in relative performance in these subsamples is about 5 percentile points in the Q section and 8 percentile

²³ We focus here in a more limited number of fields than the traditional STEM definition (e.g., we exclude biology) to select those fields that are predominately populated by males. Our results do not change when using the broader definition of STEM fields.

points in the V sections. Both gaps are statistically significant and do not change much after controlling for examinees' observed characteristics. This finding is important because it shows that the larger drop in performance among men is found even in subsamples that exhibit no differences in performance in the high stakes test.

We also look at gender gaps within groups stratified by mother's education. We were curious to check whether female examinees whose mothers attended graduate school would behave more like males and exhibit a larger gap in performance between the high and low stakes situation. This turned out not to be the case. The gender gap in relative performance between high and low stakes test appears across all levels of maternal education in both the quantitative and the verbal sections.

5. Discussion

The evidence presented above shows that men and Whites exhibit a larger difference in performance between high and low stakes tests compared to women and minorities. The larger decline in performance found among men and whites can be due to at least two different reasons: (i) men and Whites do not exert as much effort in low stakes situations compared to women and minorities, respectively; (ii) women and minorities find it relatively more difficult to deal with high stakes and stressful situations.²⁴ We examine below the plausibility of these alternative explanations and discuss some other interpretations. We acknowledge that our data do not allow us to rigorously test the relative contribution of each explanation. Nevertheless, we believe the evidence presented below provides interesting directions for further research.

5.1 Do Men and Whites Exert Less Effort in Low Stakes Situations?

To examine the likelihood of the first explanation, we would ideally like to measure effort invested in the test. More effort could be exerted by trying harder to solve each question (i.e., investment of more mental energy) or by investment of more time. Figure 4 plots the distribution of time spent by examinees in the experimental Q and V-sections by gender, race, and ethnicity.²⁵ The figure shows that there is a significant

²⁴ Alternatively, men and whites are arguably better able to boost their performance when stakes are high or the task is challenging. This explanation is harder to assess as it is impossible to establish an ability baseline that is independent of performance in a given test of a given stake. It is challenging to even conceive of a thought experiment that could possibly answer this question because performance always depends on the perceived importance of the test.

²⁵ Unfortunately, there is no information on time spent in the real GRE test. However, students usually exhaust the time limit.

variation in time invested in the experimental section. Some examinees spent very little time and some exhausted the time limit (45 minutes for the Q-section and 30 minutes for the V-section).

Figure 5 exhibits the relationship between achievement in the experimental section and time invested in that section for males, females, Whites, Blacks, Hispanics, and Asians. The figure shows that achievement increases with time invested in the quantitative section for all gender, racial, and ethnic groups. The relationship between time invested and performance in the verbal section is also positive at the lower values of the distribution but switches sign after about 20 minutes. Overall, it is clear from the figures that it is impossible to receive a high score without investing some minimal amount of time. We therefore conclude that subjects who invested very little time were obviously not exerting much effort. We define an indicator of low effort for individuals who invested less than ten minutes in the experimental section. While the ten minutes cutoff is somewhat arbitrary, we choose a time threshold that clearly suggests low effort and cannot be confounded with the ability to solve a test quickly.²⁶

We plot in Figure 6 the cumulative test score distribution in the high stake section stratifying individuals by time spent in the experimental section (below 10 minutes versus at least 10 minutes). Each quadrant in the figure refers to a specific demographic group and section (Quantitative or Verbal). We also report p-values of Kolmogorov-Smirnov tests of equality between the two distributions and p-values of t-tests of equality of means (assuming unequal variances).

For the quantitative section (panels a through d), we see no differences in the high stakes test score distribution between subjects who invested low effort in the experimental section and those who invested some reasonable amount of time. Indeed, we cannot reject the hypothesis of equality of distributions or equality of means for each demographic group. This finding shows that achievement in the high stakes section is unrelated to effort levels invested in the low stakes section and implies that baseline differences in achievement in the high stakes section between demographic groups are unlikely to explain group differences in effort levels. Given that the chances of improving one's score are probably lower for individuals who obtained higher scores in the high stakes section, the result reported in Figure 6 suggests that individuals were not thinking about the chances of winning the prize when deciding about effort levels in the low stakes section.

For the verbal section (panels e through h) we see no differences in test score distributions or means between those who invested low effort and others among males. We see some differences in the

²⁶ All participants who invested less than 10 minutes in the experimental Q-section were located below the 58th percentile of the test score distribution of that section. 94% of all those who spent less than 10 minutes in the V-section were also located below the 58th percentile.

test score distribution for females (p -value of K-S test=0.04). Nevertheless, differences in the distribution derive from differences in the dispersion around the mean, with a larger variance among those investing low effort. Indeed, we cannot reject the hypothesis of equality of means between the two groups (p -value=0.931). For minorities we find lower effort levels among those with lower scores in the high stakes section (although the difference in distributions is not statistically significant). These differences are the opposite of what we would expect if experiment participants were considering the monetary incentive when deciding about effort levels in the low stakes test. Nevertheless, as discussed above, language difficulties might have affected performance of minorities in the verbal section so we prefer not to put too much weight in the comparison of performance between whites and minorities in this section.

Taken together, the evidence presented in Figure 6, suggests that effort exerted by individuals in the experimental section is not related to performance in the “real” GRE test across all demographic groups in the Q-section and among males, females, and whites in the V-section.

Table 6 reports the share of examinees who invested less than 10 minutes in the experimental Q- and V-sections stratified by gender, race/ethnicity, academic achievement, and parental education. We also report p -values that test for equality of proportions between groups. The results show that males appear to exert less effort in the experimental section compared to females. 17 percent of the males who participated in the Q-experiment spent less than ten minutes in the experimental section while the equivalent among females is 13 percent. Gender differences are similar for the V-section. It is important to recall that, as shown in Table 1, the share of males and females among experiment participants was equal to their share in the full population of GRE test takers. This suggests that gender differences in effort among experiment participants cannot be attributed to a differential selection into the experiment. Statistics by race/ethnicity show that Whites are more likely to invest low effort relative to Blacks and Asians. Whites also appear to invest less effort than Hispanics, although differences in this case are smaller and not statistically significant.

The stratification of the sample by background characteristics and achievement shows that students with more educated parents are more likely to invest less in the exam. In contrast, we find no clear relationship between the likelihood of low effort and students’ achievement, neither when defined by students’ scores in the high stakes section nor when defined by students’ UGPAs. This last finding is important as it shows that the decision to exert low effort in the low stakes section is unrelated to students’ academic performance, suggesting that other factors are likely to play a more important role in determining performance in low stakes situations. The lack of a relationship between students’ academic performance and effort invested in the low stakes section suggests also that our previous results on group

differences in performance drop are unlikely to be explained by differences in academic achievement between groups.

We plot in Figure 7 estimates along with confidence intervals for differences in the change in performance from the high to the low stakes section between males and females or whites and minorities when we limit the sample to individuals who spent at least X minutes in the experimental section (for $X=0-45$ in the Q-section and $X=0-30$ in the V-section).²⁷ The figure shows that there is a larger gap by gender or race among those who spent a short time in the experimental section. Nevertheless, we observe that the larger drop in performance among males and whites relative to females and minorities is evident along the whole distribution of time spent in the experimental section. Appendix Table A4 reports estimates for specific points of the figure (individuals who spent at least ten minutes in the experimental section and those who spent at least three minutes). The last row of the table reports estimates from a model that uses the full sample and controls for a fourth order polynomial of time invested in the low stakes section.²⁸ We observe that differences between groups are reduced when accounting for time spent in the experimental section. Nevertheless, we see that the gap in differential performance between males and females and between whites and blacks or Hispanics is still sizable and significant. Note that while we use time invested in the low stakes section as a proxy for effort, we do not observe mental effort, a factor that might explain the remaining differences in performance change between groups.

To summarize, evidence on time invested in the experimental section suggests that the larger gap in performance between the high and the low stakes section found among men and Whites can be partly explained by a lower level of effort exerted by these groups in the low stakes section.

5.2 Are Women and Minorities More Subject to Stress in High Stakes Situations?

As noted above, a second possible explanation for the larger gap in performance between the high and the low stakes section among men and Whites could be a higher level of stress and test anxiety among females and minorities that hinders their performance in high stakes situations. To examine this explanation, we inspect the distribution of changes in performance between the high and the low stakes test. Although most individuals have lower test scores in the low stakes section, we find that some students do improve their performance. This improvement can be due to the volatility of, or measurement error, in test scores, due to learning or increased familiarity with the test, or due to a lower level of stress

²⁷ The figure reports estimates and confidence intervals obtained from a series of regressions based on equation (1) where we limit the sample to individuals spending at least X minutes in the experimental section.

²⁸ Results are very similar if we use a lower or higher degree of polynomial.

and anxiety involved in the low stakes test. We adjust for score volatility and compare the share of examinees who improved their performance across demographic groups.

Columns 1 and 6 of table 7 report the share of examinees who improved their scores in the quantitative and in the verbal experimental sections. To adjust for score improvement due to score volatility and measurement error, we define a score gain for cases where the difference between the low-stakes score and the high-stakes score divided by the conditional standard error of measurement of difference scores is greater than 1.65.²⁹ Roughly 1.5 percent of examinees have a significant score gain in the experimental Q-section and 5.3 percent in the V-section. Columns 2 through 5 and 7 through 10 report differences in the share of examinees who improve scores by gender and by race/ethnicity. The first row reports raw differences between groups, the second row reports differences after controlling for students' background characteristics, and the third row reports odds ratios between females/minorities and males/whites. Overall, we find very small and insignificant differences in the likelihood of improving the score by gender. Odds ratios are close to one for both sections (i.e. small effect size) meaning that the odds of improving the score for males and females are similar. With the exception of Hispanics in the quantitative section and Blacks in the verbal section, all other differences between whites and minorities are small and insignificant with odds ratios that are close to one.

We further explore the differential impact of test anxiety across groups using an alternative approach that takes advantage of additional information reported by examinees in the background questionnaire. The questionnaire asked examinees to report the reason(s) for taking the GRE test, allowing them to mark various alternatives. About 7 percent marked "practice" as one of the reasons for taking the exam.³⁰ If test anxiety hinders performance of females, blacks or Hispanics relative to males or whites in the high stakes section, we expect to find smaller group differences in the performance drop between the high to the low stakes section among those taking the test for practice.³¹ To examine this, we estimated our basic model of drop in performance (as in Table 3) while adding interactions between

²⁹ We use the conditional standard error of measurement of difference scores reported in Table 6b of the official ETS publication and define an indicator for score improvement following the ETS definition of significant GRE score differences (see ETS, 2007).

³⁰ The main reasons were admission to graduate school (96%) and graduate department admissions requirement (29%). Other reasons include fellowship/scholarship application requirement (23%), undergraduate program exit requirement (1%), and other (3%). Applicants were instructed to select all reasons that apply, so that reasons do not add up to 100%. The background questionnaire is filled by examinees before the test so it is not affected by their performance.

³¹ Students who took the exam for practice might be different from those who took the exam for university admission. However, for the purpose of our comparison, we only need to assume that selection works in a similar direction for all demographic groups.

an indicator for taking the test for practice and the demographic groups. Estimates reported in Table 8 show that the gap between demographic groups among those taking the exam for practice is not smaller than the gap estimated among those who are taking the exam for admission to graduate school or fellowship application and are probably facing a more stressing situation.

Taken together, the evidence presented in Tables 7 and 8 suggests that test anxiety in the high stakes section is unlikely to be the reason for the smaller change in performance between the high and the low stakes tests observed among females and minorities.

5.3 Other Explanations

An additional explanation for our results could be that the monetary prize offered to experiment participants had a differential impact on different demographic groups. While this is possible, we note that the prize consisted of \$250 (1.5 times the GRE cost) paid to 100 individuals out of 30,000 experiment participants. Such an amount distributed to such a small number of participants seems too low to have a significant differential effect in performance. Alternatively, it is arguably the case that differences in performance in the experimental section arises from group differences in their opportunity cost of time. However, as shown in Table 1, participation rates in the experiment were similar across demographic groups, suggesting that there were no group differences in the perceived cost or benefit of participating in the experiment.

To further assess the impact of the monetary prize and the opportunity cost of time on performance in the experimental section, we examined the association between the change in performance (from the high to the low stakes section) and earning levels at the state of residence of the examinee. We use two different measures of earnings: median annual earnings of full time workers and median annual earnings of college graduates computed separately by gender and state.³² If the monetary prize or the opportunity cost of time had any impact on performance at the experimental section, we should expect a smaller reduction in performance in states with lower earnings levels. We report in Appendix Table A7, regression estimates for the association between the change in performance and median earnings for males and females. Columns 1 and 3 report estimates from simple bivariate models and columns 2 and 4 report estimates from regressions that control for examinee characteristics. Overall, we do not find any association between median earnings at the state of residence of the examinee and his/her change in

³² Earnings come from data published by the U.S. Census Bureau based on 5-year average earnings by state and gender from American Community Survey for the years 2005-2008.

performance suggesting that our main results are unlikely to be explained by a differential impact of the monetary prize or the opportunity cost of time.

Another alternative explanation for differential changes in performance could be that performance of females and minorities is lower than expected in the high stakes section due to stereotype threat (e.g. Steele, 1997 and Steel and Aronson, 1995). However, it is unclear why gender and race/ethnicity stereotypes would be more pronounced in the high stakes section. In addition, the fact that we find similar gender differences in both the quantitative and the verbal sections suggest that stereotype threat is unlikely to explain our main results as the theory would predict that females would respond negatively only to the quantitative section. Moreover, stereotype threat theory implies that Asians should respond differently than Blacks and Hispanics in the quantitative section but our findings are similar for the three groups.

We further assess the likelihood of stereotype threat explanation by examining the relationship between gender stereotypes in math and verbal achievement at the state of residence of the examinees and the differential change in performance. To proxy for gender stereotypes at the state of residence of the examinee we use the stereotype adherence index developed by Pope and Sydnor (2010) which reflects gender disparities in test scores favoring boys in math and science and favoring girls in reading and was shown by the authors to be positively associated with other measures of gender stereotype attitudes at the state level.³³ Higher values in this index mean a stronger gender stereotype. To facilitate interpretation of the results, we transform this index into a z-score. We hypothesize that stereotype threat plays a more important role in states with higher values in the stereotype index. Therefore, for our results to be consistent with stereotype threat, we should observe a larger gender differential in the Q-section and a smaller gender differential in the V-section in states with a higher stereotype index. In Appendix Table A6, we examine this hypothesis by regressing the score difference between the high and the low stakes section on a female indicator, the gender stereotype index and an interaction between these two variables. Estimates for the interaction term between female and the stereotype index are all small and insignificant, meaning that there is no apparent relationship between state gender stereotypes and the gender gap in differential performance between the high and the low stakes section. Moreover, their sign goes in the opposite direction than would be expected by the stereotype threat theory.

³³ Pope and Sydnor (2010) use test score data from the National Assessment of Educational Progress (NAEP) and show that states that have larger gender disparities in stereotypically male-dominated tests of math and science also tend to have larger gender disparities (of the opposite sign) in stereotypically female-dominated tests of reading. The authors develop a state stereotype adherence index that is defined as the average of the male-female ratio in math and science and female-male ratio in reading for the top 5 percent of the students.

An additional alternative interpretation of our findings could be that group differences in underlying ability might generate differential drop in performance. However, as we note above, we observe the same pattern of gender and race/ethnic differences across different subsamples and even in subsamples that exhibit similar performance in the high or the low stakes section.

It could also be the case that women and minorities become less fatigued by the GRE examination than men and Whites, respectively and therefore exhibit a smaller drop in performance in the experimental section. This argument seems unlikely as it goes against recent psychological and medical literature that claims that, if anything, females appear to exhibit a higher level of fatigue after performance of cognitive tasks (see, e.g., Yoon et al., 2009). In addition, we are not aware of any studies that show that Whites exhibit a higher level of fatigue in response to cognitive tasks compared to Blacks, Hispanics, or Asians. Furthermore, in the context of aptitude tests, Ackerman and Kanfer (2009) and Liu et al. (2004) show no evidence for a decline in test performance in the longer test conditions. Moreover, the fact that we find similar participation rates in the experiment among males and females and whites, blacks, Hispanics, and Asians, provides further evidence that a differential effect of fatigue is unlikely to explain our findings. Lastly, as shown in Appendix table A4, the fact that we can replicate our results in the samples of students randomized into the extended time limit sections, provides strong evidence that mitigates this concern.³⁴

One could argue that group differences in performance change between the low and the high stakes section can be explained by differences in learning or test familiarization. To assess this conjecture, we took advantage of one additional piece of information at our disposal. The background questionnaire collected information on examinees' preparation methods for the GRE exam (e.g., use of software or books published by the ETS or other providers, coaching courses offered by commercial companies, coaching courses offered by educational institutions, no preparation, etc.). We coded this information in a vector of dummy variables and re-estimated our main models while controlling for these additional covariates. Results of these expanded models are reported in Appendix Table A7 together with results from our main specification. All estimates are highly similar to our main results suggesting that learning or test familiarization cannot explain our findings.

³⁴ Cotton et al. (2013) find that male's advantage over females erodes over a sequence of quizzes of equal importance. However, they note that this cannot be attributed to differential fatigue because the effect was also found when there was a two week gap between the different quizzes. In fact, they attribute the difference to the fact that boys seem to become less excited about the quizzes over time, which lowers the stakes of these quizzes and is thus consistent with our findings.

Finally, our results might also be explained by the fact that some groups might get “bored” faster than others, and this might harm their performance. This would imply that there are group differences in the likelihood of getting bored when incentives are low. We believe that the implications of our results are still relevant under this alternative interpretation.

5.4 Relation to other assessment tests

Our findings demonstrate that test score gaps between males and females or between Whites and minorities might vary according to the stakes of the test, as each group appears to respond differently to level of stakes. Therefore, it is important to consider the stakes of a test and the differential performance of each group according to the stakes level when analyzing test score gaps. We show this in Appendix Table A8 where we compare the (low stakes) NAEP test scores in mathematics that was administered to 12th grade students in 2015 to the (high stakes) SAT scores at 2015. Consistent with our results, we see that test score gaps between males and females, and Whites and Blacks, Whites and Hispanics and Whites and Asians are larger at the SAT compared to the NAEP. While many explanations are offered for the differences in performance in SAT and NAEP exams, results from our study imply that males might just exert lower effort in NAEP exams. This is supported by Freund and Rock’s (1992) finding of a higher prevalence of pattern-marking behavior among males at the NAEP exam, which they interpret as an attempt to complete the exam as quickly and with as little effort as possible. Similar gender differences in response patterns were also found by Chiacchio et al. (2016) who analyzed students’ responses at the Italian Pisa Science exam of 2006. Our results are also consistent with the claim that standardized tests usually underpredict college and graduate school performance for women and overpredict performance for men (see, e.g., Willingham and Cole, 1997 and Rothstein, 2004).³⁵

5.5 Nature or Nurture

It is interesting to try to determine to what extent differences in performance between high and low stakes situations are socially constructed or innate.³⁶ While this question is beyond the scope of the current study, we speculate that the similarity between Asian males and females suggests that part of the source for the gender differences observed among other ethnic and racial groups might be explained by acquired rather than innate skills. Women are generally perceived as more conscientious, and/or altruistic

³⁵ Our findings also suggest the same pattern for Whites compared to minorities, but this is not the case in practice (see, e.g., Mattern et al., 2008), presumably because the lower performance of minority students in college can be explained by other factors such as their relatively disadvantaged background (Rothstein, 2004).

³⁶ See, e.g., Gneezy et al. (2009), Booth and Nolen (2011), and Booth and Nolen (2012).

than men.³⁷ Other possible explanations may be associated with experimenter demand effects, and or higher salience of the reward for women and non-whites.

A curious finding that relates to this question is presented in Figure 8, where we plot differences in achievement between the high and the low stakes Q-section by students' undergraduate major. Interestingly, those who exhibit the largest gap in achievement between the high and the low stakes section are economics majors. This finding could be either a result of self-selection into economic majors or skills acquired during undergraduate studies. Be that as it may, it is consistent with Rubinstein (2006) who finds that economics majors have a much stronger tendency to maximize profits relative to other undergraduate majors.

6. Conclusions

In this study, we examine the change in performance of females, males, Whites, and minorities between high and low stakes situations by comparing the performance of GRE examinees in the real and in an experimental section of the test. Our results show that males and Whites have the highest change in performance relative to females, Asians, Hispanics, and Blacks. Males' drop in performance between the high and low stakes section is .16 SD larger than the drop of females; whites' drop in performance in the Q-section is .23 SD larger than the drop of blacks and the drop of Hispanics, and .19 SD larger than the drop of Asians. We show that the larger change in performance observed among males and Whites is partially due to the fact that these two groups invest less effort in the low stakes test. We rule out alternative explanations for these findings such as stereotype threat, differences in stress levels, learning or alternative cost of time.

Our findings suggest that men and Whites who perform well in high stake tests might not perform as well in ordinary assignments, and that women and minorities who do not perform so well in high stake tests may do relatively better in low stakes tasks. Accordingly, our results may also have implications for admission policies that are intended to achieve demographic diversity in educational institutions and the workplace. If different groups perform differently in low and high stakes situations, then policymakers may be able to diversify the population admitted to colleges, universities, specific study fields, and workplaces by putting more weight on "low stake" measures of achievement and performance such as high school performance and grades, extra-curricular activities, and recommendation letters (as opposed

³⁷ Interestingly, these stereotypes do not have firm empirical support. In a recent review of various meta-analysis Hyde (2014) finds only small or no gender differences in the Big Five personality traits. The only exception is a higher score of women on "tender-mindedness," which is part of the agreeableness factor (Cohen's d between .3 and 1.07 for US and Japanese adults, but not for Black South Africans).

to focusing entirely performance in standardized tests).

Finally, our results may also have implications for personnel and incentive policies as they suggest that differences in productivity between workers could vary according to the incentive scheme attached to the job. Often times, job candidates are evaluated using high stake tests. However, most jobs require excellence in tasks that are not directly attached to high-powered incentive schemes. This suggests that consideration of performance in high stake tests should not come at the expense of consideration of other indicators of ability that reflect good performance in low stake situations.

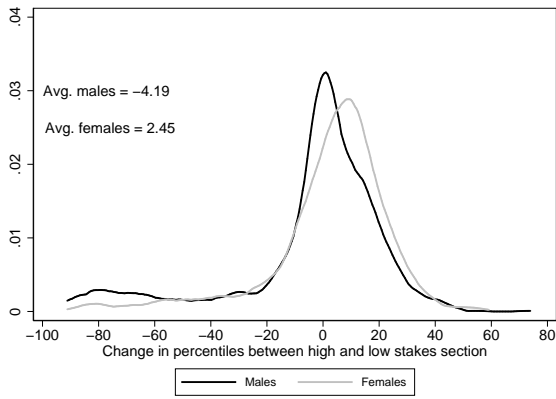
References

- Ackerman P. L. and R. Knafer, (2009) "Test Length and Cognitive Fatigue: an Empirical Examination of Effects of Performance and Test-Taker Reactions", *Journal of Experimental Psychology: Applied*, Vol. 15 No.2, pp. 163-181.
- Azmat, G. and B. Petrongolo (2014) "Gender and the labor market: What have we learned from field and lab experiments?" *Labour Economics*, Vol. 30, pp. 32-40.
- Ariely, D., Gneezy U., and G. Loewenstein (2009) "Large Stakes and Big Mistakes", *Review of Economic Studies*, Vol. 76, pp. 451-469.
- Babcock, Linda, Maria P. Recalde, Lise Vesterlund, and Laurie Weingart. 2017. "Gender Differences in Accepting and Receiving Requests for Tasks with Low Promotability." *American Economic Review*, Vol. 107, No. 3, pp. 714-47.
- Booth, A. and P. Nolen (2011) "Choosing To Compete: How Different Are Girls and Boys?" Forthcoming, *Journal of Economic Behavior and Organization*.
- Booth, A. and P. Nolen (2012) "Gender Differences in Risk Behavior: Does Nurture Matter" Forthcoming, *Economic Journal*.
- Borghans, L. H. Meijers, and B. T. Weel (2008), "The Role of Noncognitive Skills in Explaining Cognitive Test Scores" *Economic Inquiry*, Vol. 46 No. 1, pp. 2-12.
- Bridgeman B., F. Cline, and J. Hessinger (2004) "Effect of Extra Time on Verbal and Quantitative GRE Scores" *Journal Applied Measurement in Education* No. 17, pp. 25-37.
- Cassaday, J. C., and Johnson, R. E. (2002) "Cognitive Test Anxiety and Academic Performance" *Contemporary Educational Psychology*, 27, 270–295.
- Chiacchio, C., S. De Stasio, and C. Fiorilli (2016) "Examining how motivation towards science contributes to omitting behaviors in the Italian PISA 2006 sample, *Learning and Individual Differences* Vol. 50, pp. 56-63.
- Cole, J. S., D. A. Bergin, and T. A. Whittaker (2008) "Predicting student achievement for low stakes test with effort and task value" *Contemporary Educational Psychology* No. 33, pp. 609-624.
- Cotton, Christopher, Frank McIntyre, and Joseph Price (2013) "Gender Differences in Repeated Competition: Evidence from School Math Contests", *Journal of Economic Behavior and Organization*, Vol. 86 .
- Cribbie, R. A. and J. Jamieson, (2000) "Structural Equation and the Regression Bias for Measuring Correlates of Change", *Educational and Psychological Measurement*, Vol. 60 No. 6, pp. 893-907.
- Crosan, Rachel and Uri Gneezy. 2009. "Gender Differences in Preferences." *Journal of Economic Literature*, Vol. 47, No. 2, pp. 448–474.
- Cunha F. and J. J. Heckman (2007), "The Technology of Skill Formation", *American Economic Review*, Vol. 97 No. 2, pp. 31-47.
- Datta Gupta, N., A. Poulsen, and M. C. Villeval (2005) "Male and Female Competitive Behavior: Experimental Evidence" IZA Discussion Paper No. 1833.
- Dohmen, T. J. and A. Falk, (2011) "Performance Pay and Multi-Dimensional Sorting: Productivity, Preferences, and Gender", *American Economic Review*, Vol. 101, No. 2, pp. 493-525.
- Duckworth, A. L., and M. E. P. Seligman (2006) "Self-discipline gives girls the edge: Gender in self-discipline, grades, and achievement test scores" *Journal of Educational Psychology*, Vol. 98 No. 1, pp. 198-208.
- Educational Testing Service, ETS (2007), "Graduate Record Examinations, Guide to the Use of Scores 2007-2008".

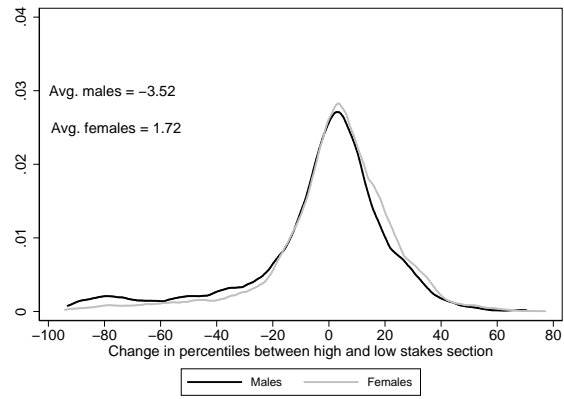
- Flory, J. A., A. Leibbrandt and J. A. List (2010) "Do Competitive Work Places Deter Female Workers? A Large-Scale Natural Experiment on Gender Differences in Job Entry Decisions" NBER Working Paper No. 16546.
- Freund D. S. and D. A. Rock (1992) "A preliminary investigation of pattern-making in 1990 NAEP data," paper presented at the Annual Meeting of the American Educational Research Association (San Francisco, CA, April 20-24).
- Gneezy, U., K. L. Leonard, and J. A. List (2009) "Gender Differences in Competition: Evidence From a Matrilineal and a Patriarchal Society" *Econometrica* Vol. 77 No. 5, pp. 1637-1664.
- Gneezy, U., M. Niederle, and A. Rustichini (2003) "Performance in Competitive Environments: Gender Differences" *Quarterly Journal of Economics*, Vol. 108 No. 3, pp. 1049-1074.
- Gneezy, U., and A. Rustichini (2004) "Gender and Competition at a Young Age" *American Economic Review Papers and Proceedings*, Vol. 94 No. 2, pp. 377-381.
- Gunther, C., N. A. Ekinici, C. Schwieren, and M. Strobel, (2010) "Women Can't Jump? An Experiment on Competitive Attitudes and Stereotype Threat" *Journal of Economic Behavior and Organization*, Vol. 75 No. 3, pp. 395-401.
- Heckman, J. J. and Y. Rubinstein (2001), "The Importance of Noncognitive Skills: Lessons from the GED Testing Program", *American Economic Review Papers and Proceedings*, Vol. 91 No. 2, pp. 145-149.
- Hyde, Janet S., (2014), "Gender Similarities and Differences", *Annual Review of Psychology*, Vol. 65, pp. 373-98.
- Jurajda, Š. and D. Münich (2011), "Gender Gap in Admission Performance under Competitive Pressure", *American Economic Review Papers and Proceedings*, May.
- Lavy, V. (2008) "Gender Differences in Competitiveness in a Real Workplace: Evidence from Performance-Based Pay Tournaments among Teachers." NBER Working Paper 14338.
- Levitt, Steven D., J. List, S. Sadoff, (2016), "The Effect of Performance-Based Incentives on Educational Achievement: Evidence from a Randomized Experiment", NBER Working Paper 22107.
- Liu J., J. R. Allspach, M. Feigenbaum, H. J. Oh, and N. Burton (2004) "A Study of Fatigue Effects from the New SAT," *College Board Research Report No. 2004-2005*.
- Lord F. M. (1967), "A Paradox in the Interpretation of Group Comparisons", *Psychological Bulletin*, Vol. 68 No. 5, pp. 304-305.
- Mattern, K. D., B. F. Patterson, E. J. Shaw, J. L. Kobrin, and S. M. Barbuti (2008) "Differential Validity and Prediction of the SAT" *College Board Research Report No. 2008-4*, The College Board, New York, 2008.
- Niederle, Muriel, (2016) "Gender," *Handbook of Experimental Economics*, second edition, Eds. John Kagel and Alvin E. Roth, Princeton University Press, pp. 481-553.
- Niederle, M., and L. Vesterlund (2007) "Do Women Shy Away From Competition? Do Men Compete Too Much?" *Quarterly Journal of Economics*, Vol. 122 No. 3, pp. 1067-1101.
- Niederle, M., and L. Vesterlund (2010). "Explaining the Gender Gap in Math Test Scores," *Journal of Economic Perspectives*, Vol. 24 No. 2, pp. 124-144.
- O'Neil, H. F., B. Sugrue, and E. L. Baker (1996) "Effects of Motivational Interventions on the National Assessment of Educational Progress Mathematics Performance", *Educational Assessment*, Vol. 2 No. 2, pp. 135-157.
- Örs, E., F. Palomino, and E. Peyrache (2008). "Performance Gender-Gap: Does Competition Matter?" CEPR Discussion Paper No. 6891.
- Paserman, D. (2010) "Gender Differences in Performance in Competitive Environments: Evidence from Professional Tennis Players." Discussion Paper, Boston University.
- Pope, D. G. and J. R. Sydnor (2010), "Geographic Variation in the Gender Differences in Test Scores," *Journal of Economic Perspectives*, Vol. 24 No. 2, pp. 95-108.

- Rothstein, J. (2004), "College Performance Predictions and the SAT," *Journal of Econometrics*, Vol. 121 No.1-2, pp. 123-144.
- Rubinstein, A. (2006), "A Skeptic's Comment on the Study of Economics", *Economic Journal*, Vol. 116, C1-C9.
- Segal, C., (2010), "Motivation, Test Scores, and Economic Success", Working Paper, Universitat Pompeu Fabra.
- Spencer, S., Steele, C. M., and D. M. Quinn (1999) "Stereotype Threat and Women's Math Performance," *Journal of Experimental Social Psychology*, Vol. 35 No. 1, pp. 4-28.
- Steele, C. M., (1997) "A threat in the Air: How Stereotypes Shape the Intellectual Identities and Performance of Women and African-Americans", *American Psychologist*, Vol. 52 No. 6, pp. 613-629.
- Steele, C. M. and J. Aronson (1995) "Stereotype threat and the intellectual test performance of African Americans", *Journal of Personality and Social Psychology*, Vol. 69, pp. 797-811.
- Willingham, W.W., and N. S. Cole, N.S. (1997). *Gender and fair assessment*. Mahwah, NJ: Erlbaum.
- Yoon T., M. Keller, B. Schinder De-Lap, A. Harkins, R. Lepers, and S. Hunter, (2009) "Sex Differences in Response to Cognitive Stress During a Fatiguing Contraction", *Journal of Applied Physiology*, Vol. 107, pp. 1486-1496.

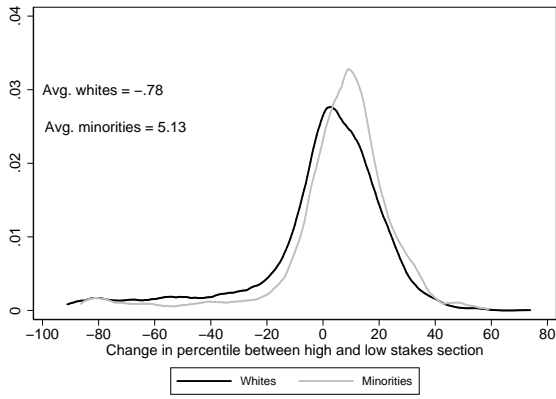
Figure 1: Difference in Ranking Between High and Low Stakes Test



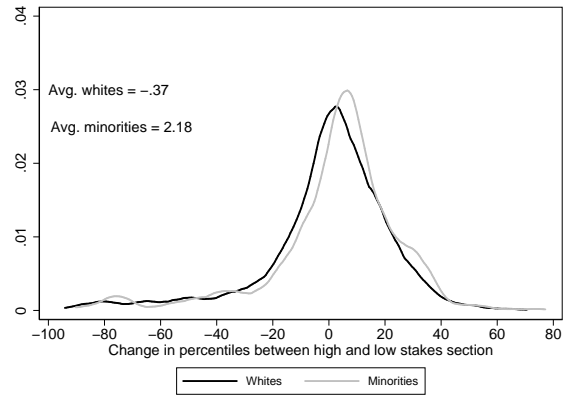
(a) Males vs. Females: Quantitative Section



(b) Males vs. Females: Verbal Section



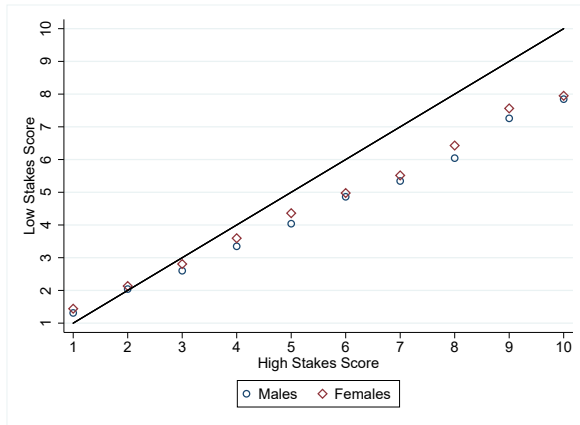
(c) Whites vs. Minorities: Quantitative Section



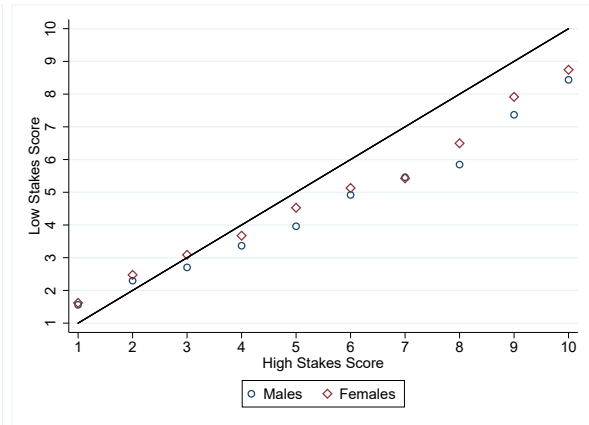
(d) Whites vs. Minorities: Verbal Section

Notes: The figure shows the difference in ranking between the low and the high stakes exam by gender and race ethnicity.

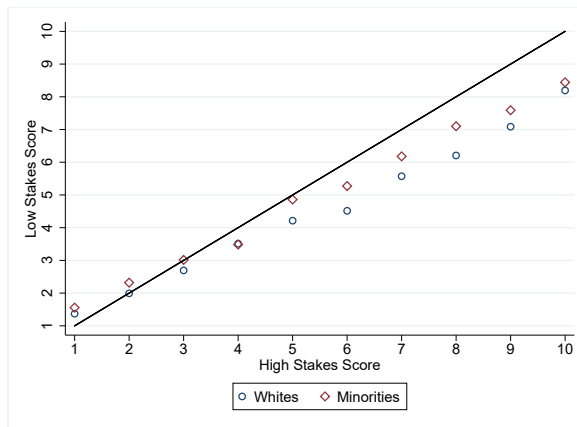
Figure 2: Score Distribution in High and Low Stakes Test



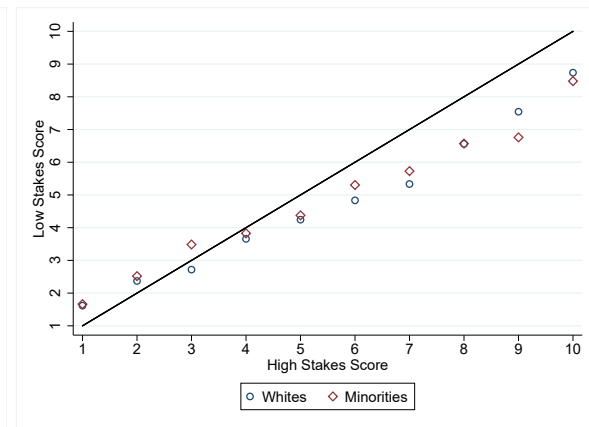
(a) Males vs. Females: Quantitative Section



(b) Males vs. Females: Verbal Section



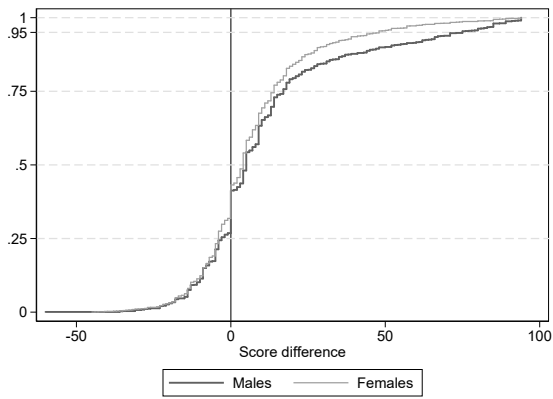
(c) Whites vs. Minorities: Quantitative Section



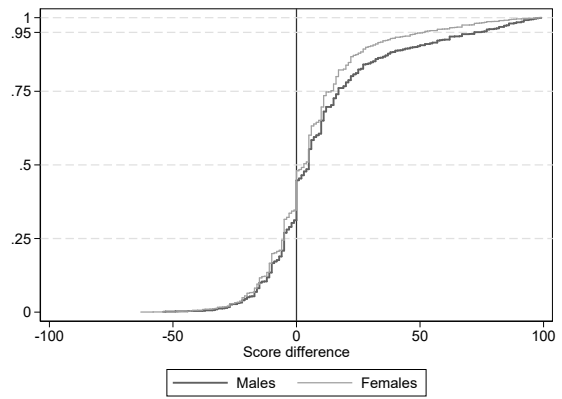
(d) Whites vs. Minorities: Verbal Section

Notes: The figure shows the low stake score as a function of the high stakes score. Scores for each gender, race, and ethnicity, are mapped into deciles using the distribution of the high stakes score of each group.

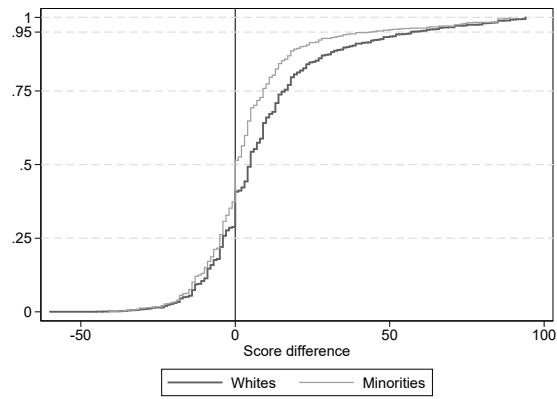
Figure 3: Distribution of Score Difference



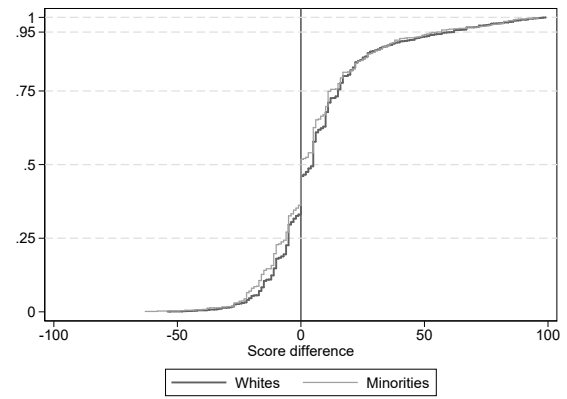
(a) Males vs. Females: Quantitative Section



(b) Males vs. Females: Verbal Section



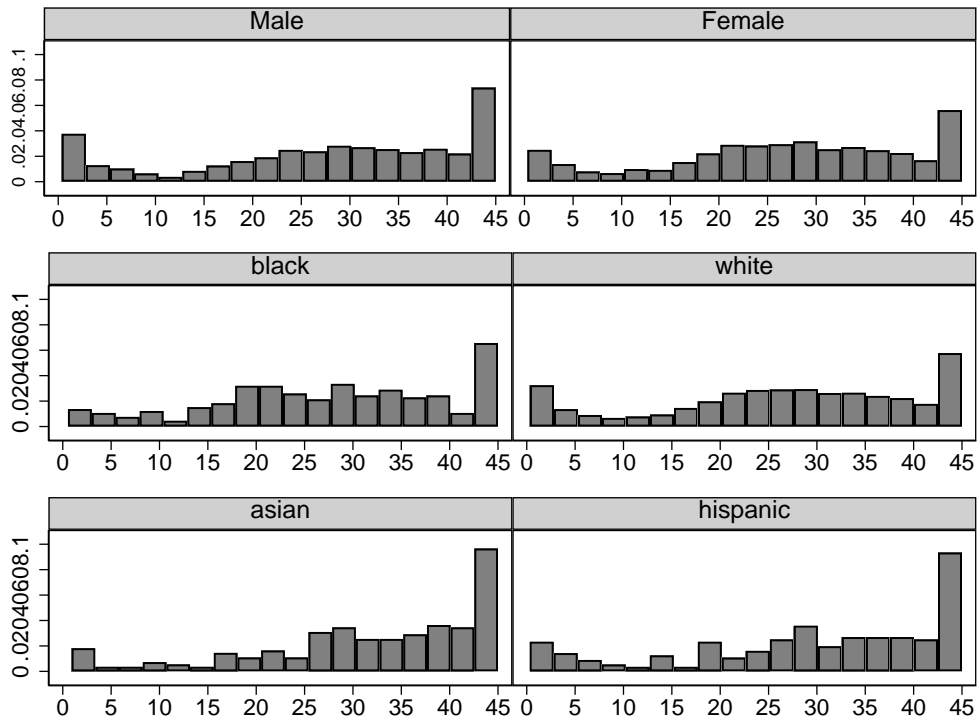
(c) Whites vs. Minorities: Quantitative Section



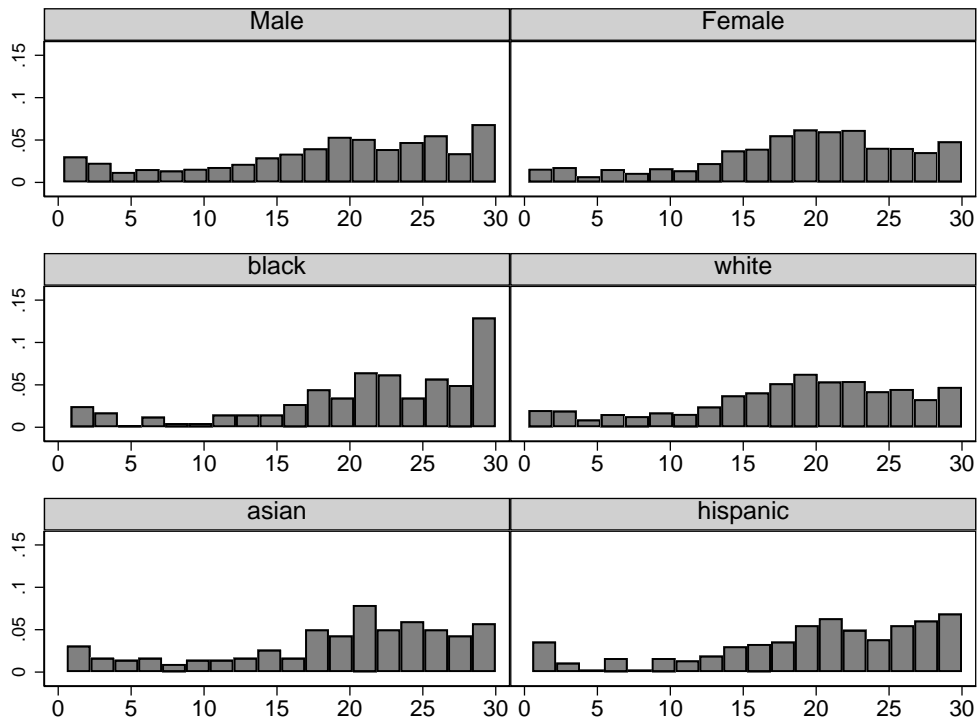
(d) Whites vs. Minorities: Verbal Section

Notes: The figure shows the CDF of the difference in score score (measured in percentiles) between the high and low stakes section by gender/race and section.

Figure 4: Distribution of Time Invested in the Experimental Section



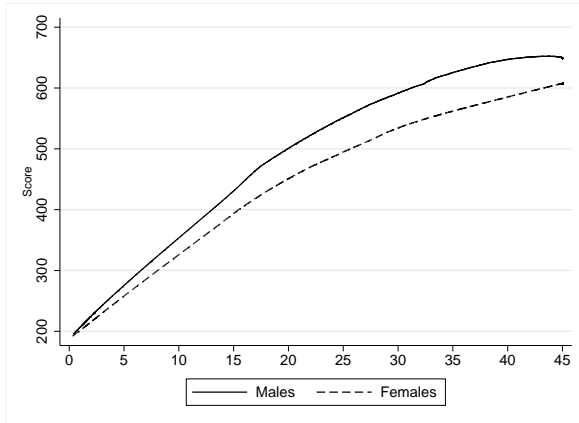
(a) Quantitative Section



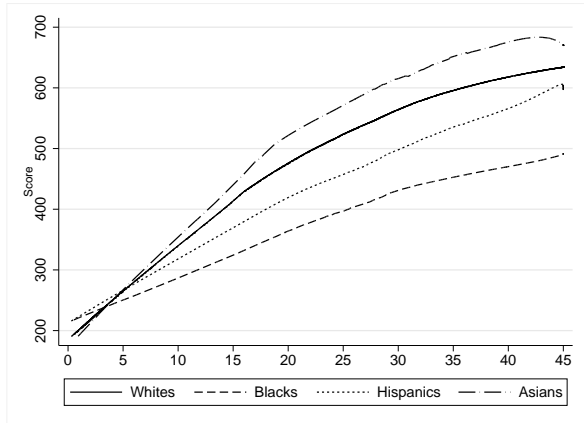
(b) Verbal Section

Notes: The histograms plot the distribution of time spent by examinees in the experimental Q and V-sections by gender, race, and ethnicity

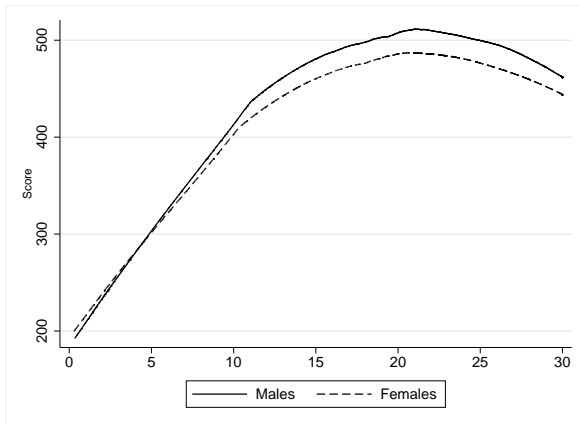
Figure 5: Relationship Between Time Invested in the Experimental Section and Test Score Achieved in that Section



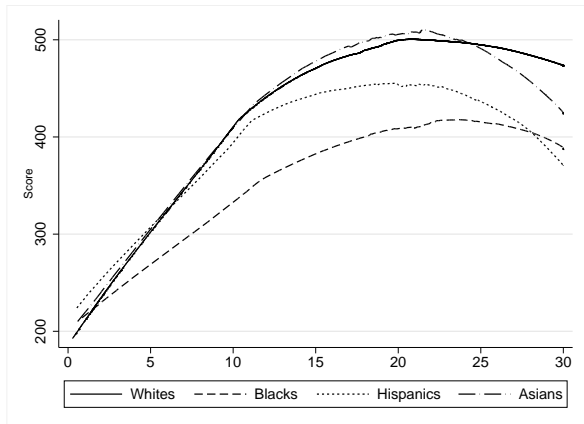
(a) Males vs. Females: Quantitative Section



(b) Whites vs. Minorities: Quantitative Section



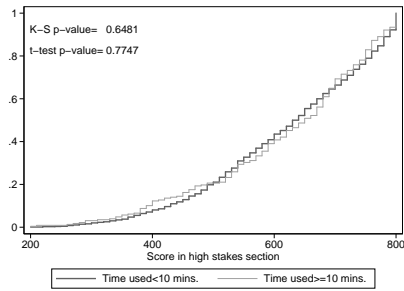
(c) Males vs. Females: Verbal Section



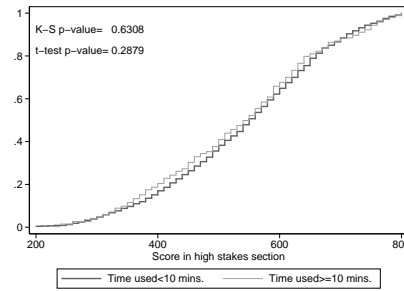
(d) Whites vs. Minorities: Verbal Section

Notes: The figure exhibits the relationship between achievement in the experimental section and time invested in that section using a local weighted regression.

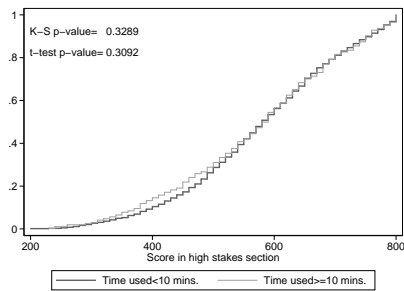
Figure 6: CDFs of Test Score in High Stake Section by Effort Invested in Experimental Section



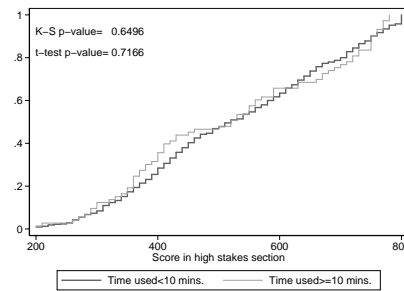
(a) Males: Quantitative Section



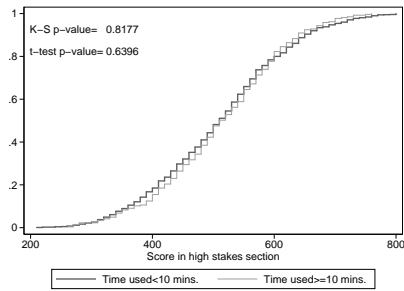
(b) Females: Quantitative Section



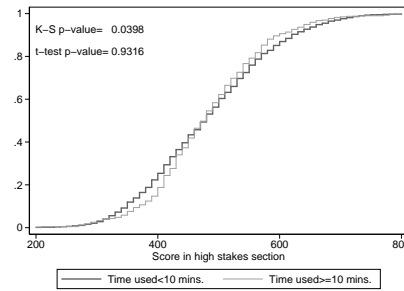
(c) Whites: Quantitative Section



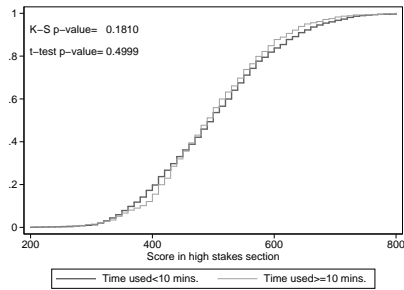
(d) Minorities: Quantitative Section



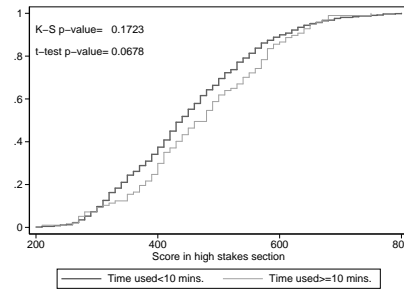
(e) Males: Verbal Section



(f) Females: Verbal Section



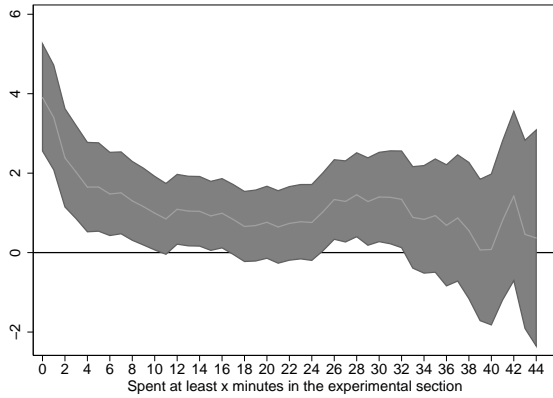
(g) Whites: Verbal Section



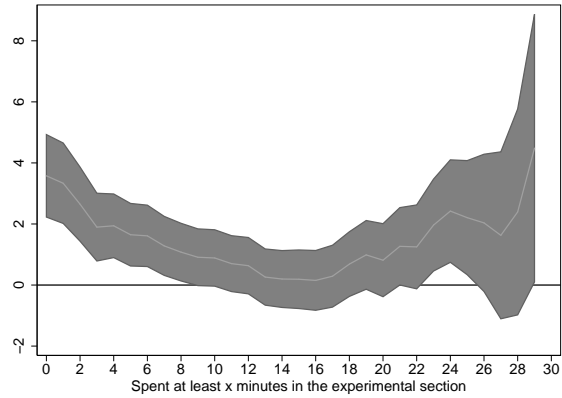
(h) Minorities: Verbal Section

Notes: The figure plots the cumulative test score distribution in the high stake section for individuals who spent below 10 minutes versus at least 10 minutes in the experimental section.

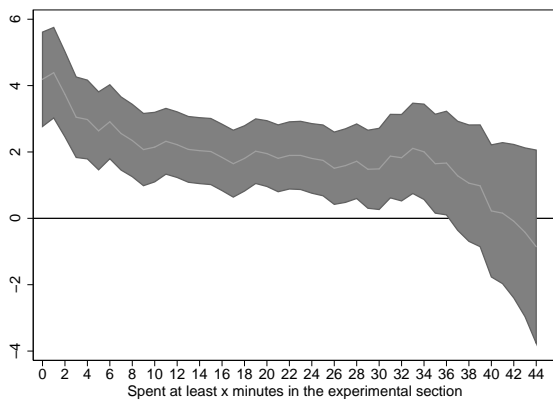
Figure 7: Gap in differential performance for those spending at least x minutes in the experimental section



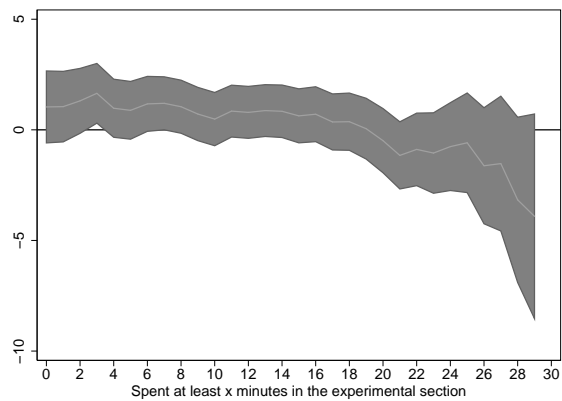
(a) Males vs. Females: Quantitative Section



(b) Males vs. Females: Verbal Section



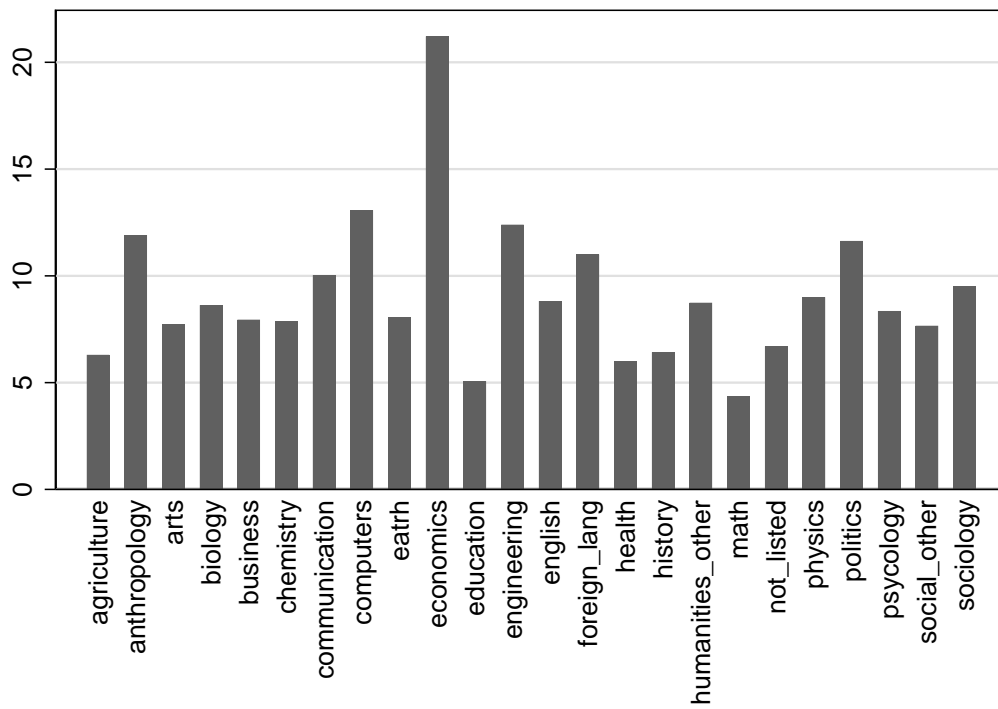
(c) Whites vs. Minorities: Quantitative Section



(d) Whites vs. Minorities: Verbal Section

Notes: The figure plots estimates along with confidence intervals for differential performance by gender and race from a series of regressions that limit the sample to individuals who spent at least X minutes in the experimental section (for X=0-45 in the Q-section and X=0-30 in the V-section).

Figure 8: Performance Gap Between High and Low Stakes Test by Undergraduate Major: Quantitative Section



Notes: The figure plots the performance gap between the high and the low stakes Q-section by subjects' undergraduate major. We include only majors that have at least 30 observations. Test scores are measured in percentile score ranks.

Table 1. Comparison Between Full Population of GRE Test Takers and Experiment Participants

		A. By gender							B. By Race/Ethnicity											
		Males			Females				Whites			Blacks			Hispanics			Asians		
		Experiment Participants			Experiment Participants				Full Sample			Experiment Participants			Experiment Participants			Experiment Participants		
		Full Sample	Q section	V section	Full Sample	Q section	V section	Full Sample	Q section	V section	Full Sample	Q section	V section	Full Sample	Q section	V section	Full Sample	Q section	V section	
N		15,749	1,369	1,465	30,160	2,553	2,845													
Share		0.34	0.35	0.34	0.66	0.65	0.66													
Quantitative score																				
Mean		55.8	55.6	56.8	40.7	40.3	41.2													
S.D		26.7	27.4	27.0	23.9	24.4	23.9													
Median		57	57	57	39	39	39													
Verbal score																				
Mean		64.1	62.4	62.9	57.0	56.2	56.5													
S.D		24.5	25.0	25.0	24.8	25.0	24.5													
Median		67	67	67	57	57	57													

Notes: The table reports students' performance (in percentile score ranks) of the full sample of GRE test takers and performance of experiment participants stratified by gender and race/ethnicity. The samples are restricted to US citizens tested in the US.

Table 2. Performance in GRE Test by Gender, Race and Ethnicity

	Males (M)	Females (F)	M-F (5)	Whites (W)	Blacks (B)	Hispanics (H)	Asians (A)	W-B (9)	W-H (10)	W-A (11)
	(3)	(4)	(5)	(5)	(6)	(7)	(8)	(9)	(10)	(11)
A. High Stakes Score										
Quantitative Section	55.58 (27.43)	40.28 (24.38)	15.30 (0.85)	46.99 (25.46)	21.85 (21.80)	36.39 (25.33)	62.30 (26.76)	25.13 (1.62)	10.59 (1.75)	-15.32 (1.75)
Number of observations	1,368	2,553		3,026	265	224	224			
Verbal Section	62.90 (24.96)	56.45 (24.54)	6.45 (0.79)	60.55 (23.69)	37.37 (24.23)	48.73 (26.20)	60.84 (26.85)	23.18 (1.58)	11.82 (1.67)	-0.30 (1.56)
Number of observations	1,465	2,845		3,380	248	221	255			
B. Low Stakes Score										
Quantitative Section	43.93 (31.34)	33.16 (25.48)	10.77 (0.93)	37.55 (27.78)	18.90 (19.72)	32.58 (26.39)	55.20 (30.38)	18.65 (1.75)	4.97 (1.90)	-17.64 (1.90)
Number of observations	1,368	2,553		3,026	265	224	224			
Verbal Section	52.48 (30.53)	50.34 (27.65)	2.14 (0.92)	52.79 (28.17)	35.08 (24.08)	42.22 (27.87)	51.78 (31.42)	17.71 (1.85)	10.57 (1.95)	1.01 (1.83)
Number of observations	1,465	2,845		3,380	248	221	255			

Notes: The table reports students test scores in the high stakes and low stakes sections of the GRE and the gaps between males and females and whites and minorities. Test scores are reported in percentile ranks. Standard deviations are reported in parenthesis.

Table 3. Difference in Performance Between High and Low Stakes Test by Gender, Race and Ethnicity

	Males (1)	Females (2)	Whites (3)	Blacks (4)	Hispanics (5)	Asians (6)
A. Quantitative Section						
High Stakes - Low Stakes	11.644 (0.683)	7.115 (0.385)	9.431 (0.399)	2.951 (0.863)	3.808 (1.346)	7.107 (1.561)
Raw Difference between Males and Females or Whites and minority group		4.529 (0.784)		6.480 (0.949)	5.623 (1.402)	2.323 (1.609)
Controlled Difference		3.905 (0.820)		4.276 (1.050)	5.205 (1.402)	3.145 (1.701)
B. Verbal Section						
High Stakes - Low Stakes	10.421 (0.673)	6.108 (0.400)	7.755 (0.390)	2.282 (1.316)	6.511 (1.457)	9.067 (1.625)
Raw Difference between Males and Females or Whites and minority group		4.313 (0.783)		5.473 (1.371)	1.244 (1.506)	-1.312 (1.669)
Controlled Difference		3.577 (0.821)		3.150 (1.472)	0.629 (1.533)	-0.555 (1.706)

Notes: The first row of each panel reports differences in individual's performance between the high and the low stakes section of the GRE by gender, race, and ethnicity. The second row of each panel reports the differences in the drop in performance between males and females or Whites and Blacks/Hispanics/Asians. The third row of each panel reports differences between groups controlling for the following individual covariates: mother's and father's education, indicators for gender or race/ethnicity, UGPA, undergraduate major, intended graduate field of studies, and disability status. Test scores are reported in percentile ranks. Robust standard deviations and standard errors of the differences are reported in parenthesis.

Table 4. Performance in High versus Low Stakes Tests by Gender and Race/Ethnicity - Quantitative Section

	High Stakes		Low Stakes		High-Low Stakes		Controlled Difference (Males-Females) (7)
	Males (1)	Females (2)	Males (3)	Females (4)	Males (5)	Females (6)	
Whites	56.701 (26.403)	41.800 (23.342)	43.914 (25.179)	34.161 (31.132)	12.787 (0.793)	7.639 (0.437)	4.904 (0.945)
Blacks	28.769 (27.739)	19.605 (19.039)	24.215 (16.851)	17.175 (26.150)	4.554 (2.146)	2.430 (0.906)	0.186 (2.531)
Controlled Difference (Whites-Blacks)					5.485 (2.405)	3.568 (1.153)	
Hispanics	44.022 (27.048)	31.363 (22.875)	38.405 (23.230)	28.748 (29.775)	5.618 (2.422)	2.615 (1.561)	0.181 (3.502)
Controlled Difference (Whites-Hispanics)					7.464 (2.608)	4.071 (1.663)	
Asians	72.167 (23.589)	56.386 (26.875)	66.071 (29.090)	48.671 (29.509)	6.095 (2.603)	7.714 (1.955)	-1.307 (4.678)
Controlled Difference (Whites-Asians)					9.266 (2.955)	-0.399 (2.055)	

Notes: The table reports test scores in the Q-section of the GRE exam. Columns 1-2 report mean performance in the high stakes test for each gender-race/ethnicity group. Columns 3-4 report mean performance in the low stakes test for each gender-race/ethnicity group. Performance change between the high and the low stakes tests are reported in columns 5 and 6. Controlled differences in performance change between males and females stratified by race/ethnicity are reported in bold in column 7. Controlled differences in performance change between whites and minorities stratified by gender are reported in bold in columns 5 and 6. Test scores are reported in percentile ranks. Standard deviations and robust standard errors are reported in parenthesis.

Table 5. Performance in High and Low Stakes Tests by Gender and Examinee Characteristics

	Number of Obs.		High Stakes Score			Low Stakes Score			High Stakes - Low Stakes			
	Males (1)	Females (2)	Males (3)	Females (4)	Diff. (5)	Males (6)	Females (7)	Diff. (8)	Males (9)	Females (10)	Raw Diff. (11)	Controlled Diff. (12)
A. Quantitative Section												
<i>Undergraduate GPA</i>												
C or C-	102	134	39.784 (24.462)	21.157 (18.445)	18.628 (2.793)	30.461 (17.397)	18.590 (25.557)	11.871 (2.800)	9.324 (1.947)	2.567 (0.851)	6.756 (2.124)	7.103 (2.320)
B-	144	266	43.028 (25.528)	28.267 (19.377)	14.761 (2.248)	34.458 (19.386)	24.034 (26.841)	10.425 (2.306)	8.569 (1.939)	4.233 (0.837)	4.336 (2.111)	2.295 (2.294)
B	426	855	48.962 (25.942)	36.063 (22.755)	12.899 (1.415)	38.418 (23.056)	29.958 (28.660)	8.460 (1.486)	10.545 (1.152)	6.105 (0.613)	4.439 (1.305)	3.492 (1.375)
A-	393	717	63.237 (24.906)	46.815 (23.935)	16.422 (1.524)	51.438 (27.150)	37.756 (31.765)	13.682 (1.812)	11.799 (1.273)	9.059 (0.823)	2.740 (1.516)	3.109 (1.641)
A	251	490	69.821 (25.227)	50.700 (23.462)	19.121 (1.869)	53.801 (27.321)	42.382 (34.295)	11.419 (2.318)	16.020 (1.908)	8.318 (0.959)	7.702 (2.135)	7.980 (2.529)
Undergrad major in Physics, Math, Comp. or Eng.	362	132	78.644 (17.321)	69.955 (23.107)	8.689 (1.935)	65.870 (27.074)	63.295 (31.352)	2.575 (3.078)	12.773 (1.549)	6.659 (2.121)	6.114 (2.624)	4.244 (2.829)
Grad intended studies in Physics, Math, Comp. or Eng.	340	122	77.674 (18.191)	70.574 (21.707)	7.100 (2.024)	65.515 (25.909)	64.369 (31.265)	1.146 (3.161)	12.159 (1.596)	6.205 (2.167)	5.954 (2.689)	4.457 (2.875)
<i>Maternal Education</i>												
High School or less	320	582	43.903 (26.374)	32.973 (22.986)	10.931 (1.687)	35.581 (23.117)	27.038 (27.255)	8.543 (1.716)	8.322 (1.235)	5.935 (0.672)	2.387 (1.405)	2.091 (1.497)
College or some college	621	1228	58.097 (26.830)	39.965 (23.495)	18.132 (1.214)	46.018 (24.850)	33.800 (32.199)	12.218 (1.356)	12.079 (1.013)	6.165 (0.529)	5.914 (1.142)	5.732 (1.218)
At least some graduate studies or professional degree	357	619	63.588 (25.921)	48.724 (25.125)	14.864 (1.689)	49.952 (27.697)	39.069 (32.106)	10.883 (1.953)	13.636 (1.455)	9.654 (0.929)	3.982 (1.725)	2.829 (1.879)

Table 5 (cont.). Performance in High and Low Stakes Tests by Gender and Examinee Characteristics

	Number of Obs.		High Stakes Score			Low Stakes Score			High Stakes - Low Stakes			
	Males	Females	Males	Females	Diff.	Males	Females	Diff.	Males	Females	Raw Diff.	Controlled Diff.
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)
B. Verbal Section												
<i>Undergraduate GPA</i>												
C or C-	106	161	48.689 (23.915)	38.441 (22.205)	10.248 (2.864)	43.208 (24.116)	35.435 (26.541)	7.773 (3.140)	5.481 (2.036)	3.006 (1.513)	2.475 (2.536)	1.121 (3.641)
B-	167	275	53.695 (26.025)	47.949 (23.273)	5.746 (2.389)	46.144 (25.274)	44.447 (27.002)	1.696 (2.545)	7.551 (1.719)	3.502 (1.129)	4.049 (2.056)	0.677 (2.583)
B	436	945	58.690 (23.905)	51.935 (23.512)	6.755 (1.368)	50.197 (25.740)	46.309 (29.117)	3.888 (1.555)	8.493 (1.165)	5.626 (0.664)	2.867 (1.340)	2.514 (1.392)
A-	405	799	68.225 (22.888)	62.016 (23.097)	6.208 (1.405)	54.138 (27.634)	55.253 (32.032)	-1.115 (1.780)	14.086 (1.391)	6.763 (0.793)	7.323 (1.600)	7.098 (1.738)
A	292	560	74.137 (20.914)	66.366 (22.573)	7.771 (1.589)	61.709 (28.622)	58.664 (31.125)	3.045 (2.130)	12.428 (1.598)	7.702 (0.933)	4.726 (1.850)	3.388 (2.064)
Undergrad major in Physics, Math, Comp. or Eng.	388	161	66.781 (24.124)	65.839 (25.365)	0.942 (2.296)	54.036 (25.708)	62.012 (31.769)	-7.976 (2.824)	12.745 (1.424)	3.826 (1.301)	8.919 (1.929)	7.547 (2.063)
Grad intended studies in Physics, Math, Comp. or Eng.	378	142	66.341 (23.796)	66.056 (24.881)	0.285 (2.372)	53.643 (27.411)	60.535 (31.356)	-6.892 (2.986)	12.698 (1.445)	5.521 (1.340)	7.177 (1.970)	7.506 (2.135)
<i>Maternal Education</i>												
High School or less	344	628	54.302 (26.892)	49.244 (23.959)	5.059 (1.679)	45.959 (25.717)	45.051 (29.148)	0.908 (1.810)	8.343 (1.305)	4.193 (0.745)	4.150 (1.502)	4.197 (1.611)
College or some college	658	1354	64.114 (23.671)	56.078 (23.942)	8.036 (1.134)	53.157 (27.139)	49.908 (30.420)	3.249 (1.343)	10.957 (1.033)	6.171 (0.591)	4.787 (1.190)	4.750 (1.281)
At least some graduate studies or professional degree	376	731	68.830 (22.931)	63.848 (24.094)	4.982 (1.504)	58.495 (28.787)	56.791 (30.521)	1.704 (1.865)	10.335 (1.318)	7.057 (0.827)	3.278 (1.556)	3.614 (1.702)

Notes: The table reports gender differences in performance in the low and the high stakes sections of the GRE test for different subsamples. Panel A reports results for experiment participants in the Q-Section Panel B reports results for experiment participants in the V-Section. Controlled differences in column 12 include the covariates detailed in Table 2. Test scores are reported in percentile ranks. Robust standard deviations and standard errors of the differences are reported in parenthesis. Sample sizes are reported in columns 1 and 2.

Table 6. Share of Experiment Participants who Spent Less than Ten Minutes in the Experimental Section

Share who spent less than ten minutes among	Q-section (1)	V-section (2)
<i>Gender</i>		
Males	0.167	0.181
Females	0.132	0.138
<i>p-value of difference: Males-Females</i>	<i>0.0032</i>	<i>0.0002</i>
<i>Race/ethnicity</i>		
Whites	0.152	0.154
Blacks	0.106	0.101
<i>p-value of difference: Whites-Blacks</i>	<i>0.0405</i>	<i>0.0227</i>
Hispanics	0.129	0.140
<i>p-value of difference: Whites-Hispanics</i>	<i>0.3557</i>	<i>0.5714</i>
Asians	0.071	0.161
<i>p-value of difference: Whites-Asians</i>	<i>0.0010</i>	<i>0.7871</i>
<i>Maternal Education</i>		
High School or less	0.134	0.133
College or some college	0.134	0.155
At least some graduate studies or professional degree	0.163	0.157
<i>p-value of difference</i>	<i>0.0860</i>	<i>0.2100</i>
<i>Paternal Education</i>		
High School or less	0.145	0.136
College or some college	0.130	0.151
At least some graduate studies or professional degree	0.161	0.166
<i>p-value of difference</i>	<i>0.0580</i>	<i>0.1160</i>
<i>Undergraduate GPA</i>		
C or C-	0.148	0.161
B-	0.120	0.122
B	0.128	0.136
A-	0.159	0.176
A	0.151	0.155
<i>p-value of difference</i>	<i>0.1300</i>	<i>0.0220</i>
<i>Achievement decile in high stakes test</i>		
1	0.166	0.160
2	0.147	0.092
3	0.128	0.103
4	0.128	0.152
5	0.153	0.174
6	0.150	0.177
7	0.132	0.170
8	0.137	0.147
9	0.166	0.169
10	0.137	0.133
<i>p-value of difference</i>	<i>0.7220</i>	<i>0.0080</i>
Number of Observations	565	659

Notes: Columns 1 and 2 report the share of examinees that spent less than 10 minutes in the experimental Q or V sections respectively out of their relevant group. The p-values reported in italics test for equality of the coefficients of the different subgroups. P-values for comparisons by gender and race are based on tests for equality of proportions. P-values for other categories are based on chi-squared tests.

Table 7. Share of Experiment Participants who Improved their Score in the Low Stakes Section Relative to the High Stakes Section

	Q-section					V-section				
	Mean (1)	Males - Females (2)	Whites- Blacks (3)	Whites- Hispanics (4)	Whites- Asians (5)	Mean (6)	Males- Females (7)	Whites- Blacks (8)	Whites- Hispanics (9)	Whites- Asians (10)
Raw difference	0.015	-0.004 (0.004)	-0.006 (0.009)	-0.018 (0.012)	-0.005 (0.009)	0.053	0.000 (0.007)	-0.038 (0.018)	-0.016 (0.017)	-0.008 (0.015)
Controlled difference		-0.005 (0.005)	-0.009 (0.009)	-0.020 (0.012)	-0.002 (0.010)		-0.004 (0.008)	-0.032 (0.019)	-0.016 (0.017)	-0.008 (0.015)
Odds ratio		1.254	1.473	2.471	1.393		1.003	1.874	1.370	1.177

Notes: Columns 1 and 6 report the share of examinees who improved their score in the experimental Q or V sections respectively relative to the real GRE section. A score gain is defined for cases where the score difference between the low and the high stakes section divided by the standard error of measurement of difference in scores is greater than 1.65. Columns 2-5 and 7-10 report differences between males and females and between whites and minorities in the share of examinees who improve their scores. The first row reports raw differences between groups. The second row reports differences between groups after controlling for examinee's covariates detailed in Table 2. Robust standard errors are reported in parenthesis. The third row reports odds ratios relative to males/whites.

Table 8. Differences in Gap in Performance Drop Between Students Taking Test for Practice and Other Students

	Gaps by Gender		Gaps by Race/Ethnicity					
	Female x		Black x		Hispanic x		Asian x	
	Female	Practice	Black	Practice	Hispanic	Practice	Asian	Practice
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Quantitative section	-3.595 (0.848)	-4.948 (3.436)	-4.314 (1.014)	0.536 (6.174)	-5.002 (1.470)	-2.262 (4.469)	-2.480 (1.781)	-10.524 (4.053)
Verbal section	-3.696 (0.852)	1.580 (2.922)	-2.671 (1.580)	-3.883 (3.738)	-0.761 (1.607)	2.017 (4.769)	0.352 (1.755)	3.178 (7.309)

Notes: The table reports estimates from a regression of test score change on indicators for the different demographic groups and the interaction between demographic groups and practice exam. The model controls also for an indicator of practice exam and student's background characteristics detailed in Table 2. Robust standard errors are reported in parenthesis.

Online appendix

Table A1. Sample Construction Process

	Total	Gender			Race/ethnicity				
		Males	Females	Missing	Whites	Blacks	Hispanics	Asians	Other/ Missing
Population (all GRE tested 9/1/2001-10/31/2001)	81,231	34,723	41,617	4,891					
US citizens tested in the US	46,038	15,749	30,160	129	36,042	2,877	2,400	2,584	2,135
Experiment participants (total)	29,962	13,359	14,803	1,800					
US citizens tested in the US	15,945	5,486	10,458	1	12,374	1,024	850	982	715
Participants in regular time limit experiment	8,232	2,834	5,398	0	6,407	513	445	479	388
Participants in Q section	3,922	1,369	2,553		3,027	265	224	224	182
Participants in V section	4,310	1,465	2,845		3,380	248	221	255	206

Notes: The table reports the process we followed to select our analysis samples.

Online appendix

Table A2. Descriptive Statistics of Experiment Participants

	Males (1)	Females (2)	Whites (3)	Blacks (4)	Hispanics (5)	Asians (6)
Females			0.66	0.74	0.65	0.63
<i>Race/Ethnicity</i>						
Whites	0.78	0.78				
Blacks	0.05	0.07				
Hispanics	0.06	0.05				
Asians	0.06	0.06				
American Indian or Alaskan Native	0.00	0.01				
Other	0.05	0.04				
<i>Mother's Education</i>						
High School or less	0.23	0.22	0.21	0.33	0.40	0.24
College or some college	0.45	0.48	0.48	0.41	0.37	0.46
At least some graduate studies or professional degree	0.26	0.25	0.26	0.19	0.19	0.25
Missing	0.06	0.05	0.04	0.07	0.04	0.05
<i>Father's Education</i>						
High School or less	0.21	0.23	0.20	0.43	0.40	0.15
College or some college	0.40	0.44	0.44	0.38	0.33	0.39
At least some graduate studies or professional degree	0.37	0.32	0.35	0.16	0.25	0.45
Missing	0.01	0.01	0.01	0.04	0.02	0.01
Native English speaker	0.93	0.92	0.93	0.95	0.90	0.86
<i>Undergraduate GPA</i>						
C or C-	0.07	0.05	0.05	0.20	0.08	0.05
B-	0.11	0.10	0.10	0.18	0.13	0.07
B	0.30	0.33	0.32	0.36	0.37	0.36
A-	0.28	0.28	0.30	0.13	0.23	0.30
A	0.19	0.19	0.21	0.07	0.13	0.18
Missing	0.04	0.04	0.03	0.06	0.07	0.05
Undergraduate major in Physics, Math, Comp. Science or Engineering	0.26	0.05	0.12	0.10	0.12	0.31
Grad. intended studies in Physics, Math, Comp. Science or Engineering	0.25	0.05	0.11	0.07	0.13	0.30

Notes: The table reports descriptive statistics of participants in the regular time limit experiment. The samples are restricted to US citizens tested in the US.

Table A3. Robustness Check: Differential Performance in High versus Low Stakes Tests

	Difference in performance between high and low stake test						Controlled difference between groups			
	Males (1)	Females (2)	Whites (3)	Blacks (4)	Hispanics (5)	Asians (6)	Males- Females (7)	Whites- Blacks (8)	Whites- Hispanics (9)	Whites- Asians (10)
A. Quantitative Section										
Raw scores	76.952 (4.402)	51.324 (2.596)	65.601 (2.650)	21.019 (6.085)	32.321 (8.818)	43.616 (9.464)	23.018 (5.414)	32.553 (7.288)	31.455 (9.187)	26.256 (10.395)
Ln(raw scores)	0.182 (0.011)	0.130 (0.007)	0.161 (0.007)	0.054 (0.016)	0.091 (0.022)	0.098 (0.022)	0.048 (0.013)	0.083 (0.019)	0.067 (0.023)	0.072 (0.024)
Standardized scores (z-scores)	0.551 (0.031)	0.367 (0.019)	0.469 (0.019)	0.150 (0.044)	0.231 (0.063)	0.312 (0.068)	0.165 (0.039)	0.233 (0.052)	0.225 (0.066)	0.188 (0.074)
Percentile score ranks Extended time limit sample	11.560 (0.733)	6.330 (0.421)	8.708 (0.421)	1.585 (1.071)	3.995 (1.611)	10.573 (1.789)	4.251 (0.931)	5.081 (1.203)	3.440 (1.741)	-0.901 (1.847)
Percentile score excluding largest drop in performance	4.636 (0.407)	2.124 (0.247)	3.590 (0.245)	-0.370 (0.516)	-1.615 (0.772)	0.622 (0.868)	3.133 (0.533)	4.315 (0.643)	5.479 (0.839)	2.365 (0.944)
B. Verbal Section										
Raw scores	45.993 (3.022)	26.882 (1.748)	34.275 (1.734)	11.935 (5.779)	27.059 (6.118)	39.255 (7.073)	15.658 (3.626)	11.646 (6.442)	4.293 (6.480)	-1.577 (7.391)
Ln(raw scores)	0.121 (0.008)	0.072 (0.005)	0.091 (0.005)	0.037 (0.016)	0.072 (0.016)	0.103 (0.018)	0.041 (0.009)	0.029 (0.017)	0.013 (0.017)	-0.004 (0.019)
Standardized scores (z-scores)	0.430 (0.028)	0.251 (0.016)	0.320 (0.016)	0.112 (0.054)	0.253 (0.057)	0.367 (0.066)	0.146 (0.034)	0.109 (0.060)	0.040 (0.061)	-0.015 (0.069)
Percentile score ranks Extended time limit sample	11.380 (0.748)	4.575 (0.413)	7.428 (0.423)	1.240 (1.101)	3.894 (1.539)	8.008 (1.814)	5.781 (0.895)	3.364 (1.261)	2.347 (1.612)	0.181 (1.844)
Percentile score excluding largest drop in performance	3.473 (0.420)	0.699 (0.260)	1.682 (0.242)	-3.151 (0.857)	1.015 (1.053)	2.052 (1.021)	2.873 (0.539)	4.330 (0.950)	0.656 (1.086)	-0.386 (1.080)

Notes: The table reports differences in performance between the high and the low stakes tests by gender and race/ethnicity and the controlled difference between males and females or whites and minorities. The first row of each panel uses difference in raw test scores as a dependent variable. The second row uses difference in the natural logarithm of raw scores as a dependent variable. The third row uses differences in standardized scores (z-scores). The third row uses difference in percentile score ranks (as all our main results table) based on the sample of students who got the experimental GRE section with an extended time limit. The last row excludes from our main sample students in the top 10 percentile drop in performance of each demographic group. Standard deviations and robust standard errors are reported in parenthesis.

Online appendix

Table A4. Performance Gap Between High and Low Stakes Section by Time Spent in Low Stakes Section

Sample	Controlled difference between groups			
	Males-Females (1)	Whites-Blacks (2)	Whites-Hispanics (3)	Whites-Asians (4)
A. Quantitative Section				
Full	3.905 (0.820)	4.276 (1.050)	5.205 (1.402)	3.145 (1.701)
Time spent in experimental section \geq 10 mins.	0.991 (0.563)	2.100 (0.782)	4.385 (0.990)	0.246 (1.210)
Time spent in experimental section \geq 3 mins.	2.065 (0.713)	3.007 (0.871)	4.440 (1.230)	1.831 (1.412)
Full sample - controlling for 4th order polynomial of time spent in experiment	2.749 (0.582)	1.842 (0.941)	2.549 (1.048)	-1.093 (1.201)
B. Verbal Section				
Full	3.577 (0.821)	3.150 (1.472)	0.629 (1.533)	-0.555 (1.706)
Time spent in experimental section \geq 10 mins.	0.885 (0.561)	2.131 (1.108)	-0.703 (1.246)	-0.071 (1.205)
Time spent in experimental section \geq 3 mins.	1.897 (0.675)	4.144 (1.131)	0.269 (1.353)	0.506 (1.402)
Full sample - controlling for 4th order polynomial of time spent in experiment	1.834 (0.556)	1.076 (1.079)	0.152 (1.270)	-0.686 (1.128)

Notes: The table reports differences in performance between the high and the low stakes tests by gender and race/ethnicity. Panel A reports differences in the Q-section and panel B reports differences in the V-section. The first row of each panel reproduces results reported in table 3. The second row of each panel reports results for the subsample of examinees who spent less than 10 minutes in the experimental section. The third row of each panel reports results for the subsample who spent at least 3 minutes in the experimental section. The fourth row of each panel reports results for the full sample from a model that controls for a 4th order polynomial of time spent in the experimental section. Test scores are reported in percentile ranks. Standard errors are reported in parenthesis.

Table A5. Associations Between Median Earnings at the Examinee State of Residence and Differential Performance

	Males		Females	
	(1)	(2)	(3)	(4)
A. Quantitative Section				
Median earnings (in thousand dollars)	-0.169 (0.161)	-0.221 (0.182)	0.124 (0.130)	0.059 (0.142)
Median earnings of college graduates (in thousand dollars)	0.049 (0.151)	0.083 (0.151)	0.012 (0.143)	-0.036 (0.142)
B. Verbal Section				
Median earnings (in thousand dollars)	0.211 (0.162)	0.089 (0.178)	0.158 (0.076)	0.089 (0.075)
Median earnings of college graduates (in thousand dollars)	0.191 (0.111)	0.131 (0.132)	0.109 (0.112)	0.045 (0.106)
Controls for examinee's covariates	--	✓	--	✓

Notes: The table reports regression estimates for the coefficient of annual median earnings (in thousand US\$) of full time workers or college graduates working full time at the state of residence of the examinee. The dependent variable is the score difference (in percentile points) between the high and the low stakes section. Estimates reported in columns (2) and (4) come from regressions that control for race/ethnicity, mother's and father's education, dummies for UGPA, undergraduate major, intended graduate field of studies, and disability status. Standard errors clustered at the state levels are reported in parenthesis.

Online appendix

Table A6. Differential Performance and Stereotype Threat

	High Stakes - Low Stakes			
	Quantitative Section		Verbal Section	
	(1)	(2)	(3)	(4)
Female	-4.434 (0.692)	-3.745 (0.879)	-4.851 (0.749)	-4.136 (0.863)
State stereotype index	-1.108 (0.651)	-0.765 (0.645)	-0.149 (0.702)	-0.054 (0.657)
Female x State stereotype index	0.606 (0.681)	0.498 (0.747)	-0.796 (0.706)	-0.657 (0.712)
Controls for examinee's covariates	--	✓	--	✓

Notes: The table reports estimates from models that regress the score difference (in percentile points) between the high and the low stakes section on a female dummy, the gender stereotype index of the state of residence of the examinee and the interaction between these two variables. Estimates reported in columns (2) and (4) come from regressions that control also for race/ethnicity, mother's and father's education, dummies for UGPA, undergraduate major, intended graduate field of studies, and disability status. The gender stereotype index was developed by Pope and Sydnor (2010) and reflects gender stereotypes in performance in math and science versus reading at the state level. Higher values denote stronger gender stereotypes. The index is standardized to have mean of 0 and a standard deviation of 1. Standard errors clustered at the state levels are reported in parenthesis.

Online appendix

Table A7. Differential Gap Between High and Low Stakes Section After Controlling for Test Preparation Methods

	Quantitative Section				Verbal Section			
	Males- Females	Whites- Blacks	Whites- Hispanics	Whites- Asians	Males- Females	Whites- Blacks	Whites- Hispanics	Whites- Asians
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Basic specification	3.905 (0.820)	4.276 (1.050)	5.205 (1.402)	3.145 (1.701)	3.577 (0.821)	3.150 (1.472)	0.629 (1.533)	-0.555 (1.706)
Controlling for test preparation methods	3.959 (0.828)	4.019 (1.065)	5.216 (1.413)	3.496 (1.683)	3.746 (0.829)	3.402 (1.468)	0.793 (1.540)	-0.217 (1.707)

Notes: The table reports differences in performance between the high and the low stakes tests by gender and race/ethnicity. Columns 1-4 report differences in the Q-section and columns 5-8 report differences in the V-section. The first row of the table reproduces estimates from the full specification reported in table 3. The second row reports results from regressions that control also for indicators for test preparation methods reported by the examinees: no preparation, software from the ETS, books published by the ETS, software from other providers, books from other providers, attended a coaching course offered by a commercial company, attended a coaching course offered by an educational institution, used *ScoreItNow!* online writing practice, used GRE enhanced diagnostic service, other type of preparation. Robust standard errors are reported in parenthesis.

Online appendix

Table A8. Gender Gaps in Low and High Stakes Tests: NAEP vs. SAT Scores 2015

	12th grade NAEP			SAT		
	Males	Females	Standardized mean difference	Males	Females	Standardized mean difference
A. Math						
Whites	161 (33)	159 (30)	0.063	551 (107)	518 (99)	0.320
Blacks	129 (31)	131 (31)	-0.065	435 (104)	422 (96)	0.130
Hispanics	141 (33)	136 (31)	0.156	473 (70)	443 (64)	0.456
Asian	171 (36)	169 (35)	0.056	611 (125)	585 (126)	0.207
B. Reading						
Whites	290 (40)	301 (37)	-0.285	532 (105)	526 (101)	0.058
Blacks	259 (38)	272 (37)	-0.347	428 (102)	434 (99)	-0.060
Hispanics	272 (38)	279 (37)	-0.187	453 (69)	446 (67)	0.108
Asian	290 (41)	304 (38)	-0.354	525 (128)	526 (124)	-0.008

Notes: The table reports students' achievement in NAEP and SAT exams in 2015. The standardized mean difference (also known as effect size - d) is the average of men minus the average of women divided by the within group standard deviations pooled across groups. Data for NAEP scores was downloaded from the National Center for Education Statistics Website. Data from SAT scores come from SAT Total Group Profile Report 2015 (College Board, 2015).