

*Brave Boys and Play-it-Safe Girls:*

Gender Differences in Willingness to Guess in a Large Scale Natural Field  
Experiment\*

Nagore Iriberry<sup>+</sup> and Pedro Rey-Biel<sup>\*\*</sup>

October 2020

**Abstract**

Multiple-choice tests are extensively used to measure individuals' knowledge and aptitudes. We study gender differences in willingness to guess using approximately 10,000 multiple-choice math tests, where, for all participants, in half of the questions, omitted answers were rewarded while for the other half they scored the same as wrong answers. Using a within-participant regression analysis, we show that female participants leave significantly more omitted questions than males when there is a reward for omitted questions. This gender difference, which is stronger among high ability and older participants, hurts female performance as measured by the final score and position in the ranking. We conclude that it is important to use gender neutral scoring rules that do not differentiate between wrong answers and omitted questions in order to accurately measure individuals' knowledge and aptitudes.

Keywords: gender differences, willingness to guess, natural field experiment, perceived ability in math, risk preferences, confidence.

JEL classification: C93, D81, I20, J16.

---

\* We thank researchers at various institutions and conferences for their helpful comments and feedback. We are thankful to the organizers of *Concurso de Primavera de Matemáticas*, who collaborated throughout the project. In particular, we would like to honor the memory of Joaquín Hernández, who put great effort and passion into this initiative and who, unfortunately, left us too soon. Iriberry acknowledges financial support from the Departamento Vasco de Educación, Política Lingüística y Cultura (IT1367-19), Ministerio de Economía y Competitividad and Fondo Europeo de Desarrollo Regional (ECO 2015-66027-P and PID2019-106146GB-I00). Rey-Biel acknowledges funding from Programa Ramón y Cajal and from Universidad Ramón Llull (ESADE).

<sup>+</sup> University of the Basque Country, IKERBASQUE, Basque Foundation for Science. E-mail nagore.iriberri@gmail.com.

<sup>\*\*</sup> Universitat Ramon Llull, ESADE. Email: pedro.rey@esade.edu.

## 1. Introduction

Multiple-choice tests are one of the most frequently used means to measure individuals' knowledge and aptitudes. On top of their common use in everyday academic life, performance on some multiple-choice tests plays a crucial role in shaping labor market outcomes. For example, the Scholastic Aptitude Test (SAT) in the USA and the Graduate Record Examination (GRE), used worldwide, are standardized multiple-choice tests that play a key role in shaping students' future outcomes. Similarly, licensing exams for many professions, such as in Medicine and Law, are also based on multiple-choice tests: United States Medical Licensing Examination (USMLE), Medico Interno Residente (MIR) in Spain, and bar exams in Law. One crucial decision a multiple-choice test designer needs to make is whether wrong answers and omitted questions (questions with no answer) are scored the same. GRE and USMLE do not have differential grading, while MIR in Spain does. An interesting example is the SAT, which used to have differential grading but has changed recently to non-differential scoring rule.

The main motivation for scoring wrong answers and omitted questions differently is to prevent test takers from getting the question right by chance, which would add noise to the measure of knowledge and aptitudes. However, one important concern is that multiple-choice tests in which incorrect answers are scored differently than omitted questions may lead individuals with different degrees of confidence and/or risk aversion to follow different strategies when answering, which might also misrepresent those students' knowledge and aptitudes. An extensive body of literature has documented that women are on average more risk averse (Eckel and Grossman, 2008, Croson and Gneezy, 2009, and Filippin and Crosetto, 2016) and less confident (Beyer, 1999, and Barber and Odean, 2001) than men. Hence, an informed decision on the optimal scoring rule regarding omitted questions and wrong answers requires studying its effect on gender differences in willingness to guess and ultimately on performance.

In collaboration with the organizers of *Concurso de Primavera de Matemáticas*, we conduct a large-scale within-subject natural field experiment to test for and understand the mechanisms behind gender differences in willingness to guess. *Concurso de Primavera de Matemáticas* is a regional math contest in which primary education, secondary education and high school students from the region of Madrid participate annually. In the 2016, 2017 and 2018 contests, which had a total of approximately 10,000 test takers, we designed tests with no differential score between omitted questions and

wrong answers for the first 13 of 25 test questions (no reward for omitted), and for the last 12 test questions, wrong answers were scored 0 and omitted questions +1 (reward for omitted). Correct answers always yielded +5. There were no other systematic differences between the two parts of the exam. We compare within-participant willingness to guess and overall performance across both parts of the test.<sup>1</sup>

We find that in the no reward section female participants leave on average more omitted questions than males (0.17 standard deviations of the mean), even though the dominant strategy is to answer all questions when there is no reward. Most importantly, this difference increases (by an additional 0.14 standard deviations of the mean) in the reward section of the test. Moreover, the gender differential in willingness to guess has important consequences for the gender differences in the final scores and the ranking of the participants. Females tend to show lower performance on the math test (0.22 standard deviations of the mean), which leads females to lag approximately 52 positions behind in the ranking of test takers (with an average of approximately 1000 total positions). This female underperformance increases by 0.04 standard deviations of the mean and women lose approximately 11 additional positions in the ranking under the differential scoring rule for omitted questions and wrong answers.

We explore three heterogeneity effects. First, using two different measures of ability (the number of correct answers when there is no reward and math grade at school), we test whether the gender gap in willingness to guess varies with ability. As expected, high-ability participants leave fewer omitted questions than low-ability ones. However, we find that the gender differential for willingness to guess is indeed stronger among high-ability participants (0.26 standard deviations of the mean), while we find no significant gender differential for the low-ability participants. This finding has important implications for competitive settings, as those who will be selected will be at the top of the distribution. Second, using four different age categories, we explore the differential gender effect of the scoring rule across different ages. Participants in their final years of high school (16-17 years old) show a significantly higher gender differential between the reward and no-reward parts of the test than younger participants (10-11 years old). This again has

---

<sup>1</sup> Notice that, prior to our intervention, the contest organizers were already using a grading rule rewarding omitted questions in all 25 questions. Upon our request, they agreed on using a non-differential scoring rule for half of the questions but did not allow us to have treatments differing on whether it was the first or the second part of the test the one with differential grading.

important implications as it suggests that these differences tend to increase when they matter most, for example, in ages in which students take tests that determine important outcomes in education, such as university entry tests (Coffman and Klinowski, 2020). Third, we test if gender differences in willingness to guess vary between public and non-public schools, where the non-public is a proxy of higher socio-economic status. We find no evidence that supports any heterogeneity regarding this dimension.

We show that the main findings remain intact when performing a series of robustness tests regarding different econometrics specification, additional controls for math ability, alternative identification specifications and different samples. We also contemplate alternative explanations, such as gender differences in the understanding of the scoring rule, gender differences in time management and experimenter demand effects, providing evidence against them.

Previous literature has shown that women omit more questions than men when there is a differential scoring rule for wrong answers and omitted questions, mostly based on observational data (Swineford, 1941; Anderson, 1989; Atkins et al., 1991; Ramos and Lambating, 1996; Tannenbaum D., 2012, Pekarinen (2014), Akyol, Key, and Krishna, 2016, Riener and Wagner, 2017). Only recently there have been important advances in pursuing randomized controlled trials in the laboratory (Baldiga, 2014), in the field (Ben-Shakhar and Sinai, 1991, Espinosa and Gardeazábal, 2013, and Funk and Perrone, 2016, Karle et al., 2019 and Atwater and Saygin, 2020) and using before-after quasi-controlled studies (Coffman and Klinowski, 2020) to test for the causal effect of differential scoring rules on male and female test takers' willingness to guess and performance. Although all these studies find that female students leave more omitted questions than males when there is differential grading of wrong answers and omitted questions, there is disagreement over whether this differential grading hurts females or not. On the one hand, Funk and Perrone (2016) do not find any harmful effect for female students, arguing that females in their sample have higher average ability than males. Akyol, Key and Krishna (2016) estimate negative effects for females and for risk averse students but conclude that the effects are small, making a case for differential scoring of omitted questions and wrong answers. On the other hand, Baldiga (2014), and Coffman and Klinowski (2020) find that a differential scoring rule for omitted questions and wrong answers has a significant negative impact on the gender gap in performance. Karle et al. (2019) match data from multiple-choice exams with students' risk preferences elicited in a classroom

experiment, and find similar gender effects explained mainly by differences in loss aversion.

Our study contributes to the existing randomized controlled trials in the following ways. First, we complement existing laboratory experiments providing the cleanest identification strategy using a within-subject design. Our within-subject design allows us to attribute differential gender responses in answering behavior to the scoring rule and not to possible gender differences in motivation to perform well in the exam, as discussed in Gneezy et al., (2019). Second, we provide external validity, carrying out a natural experiment with contest participants who are unaware of being studied. The study of settings that involve a competitive component might be more informative about behavior on high-stakes tests that determine entry to university and the attainment of professional licenses, as they also have a competitive component. Further, most existing studies use samples of Economics undergraduate students, while our sample includes a larger number of students from both public and non-public schools, who vary in age and ability. Most recent studies have shown interesting heterogeneity effects regarding ability differences (Funk and Perrone, 2016; Akyol, Key and Krishna, 2016, and Coffman and Klinowski, 2020). We find that high-ability female participants are indeed more affected, which resonates with the results of Akyol, Key and Krishna (2016) and Coffman and Klinowski (2020). In addition, we explore heterogeneity effects regarding age, and unlike Riener and Wagner (2018) for a large sample of German students, we find that the gender differential increases with age. Finally, we also study heterogeneity with respect to the socio-economic status of the students, measured by whether students attend a public or non-public school. Third, similar in spirit to the laboratory experiment by Baldiga (2014), we also contribute to the understanding of the underlying mechanism and test how much the gender differences in willingness to guess are due to confidence, overconfidence and risk aversion. We indeed also find suggestive evidence that gender differences in risk aversion are the main factor adding external validity to this finding. Finally, the differential scoring rule *rewards* omitted questions rather than *penalizing* wrong answers. Espinosa and Gardeazábal (2010 and 2013) show that these two approaches are only strategically equivalent under risk neutrality and that under risk aversion, penalties will lead to more omitted questions than rewards.<sup>2</sup> Therefore, this study generalizes our

---

<sup>2</sup> Balart et al. (2020) show that framing the scoring rule either in positive or negative marking affects men and women differently in their willingness to guess in multiple-choice tests.

understanding of the gender differences in willingness to guess showing that they are still an important concern in a more favorable scoring rule for females than when using penalties for wrong answers.

In summary, our results generalize existing results from laboratory and field experiments by confirming, using a within-subject design in a natural competitive setting with a large sample of participants and a reward-for-omitted scoring rule, that gender differences do exist in reaction to differential scoring rules and that these do harm female performance. Our paper complements both Baldiga (2014) and Coffman and Klinowski (2020), with a within-subject design from a natural experiment.

The paper has the following structure. Section 2 describes the setting, the data and the main descriptive statistics. Section 3 includes the main results, robustness tests, and heterogeneity results. Section 4 summarizes a discussion on the underlying mechanism. Section 5 concludes.

## **2 The Data**

### **2.1 The Setting: Mathematics Test**

The Mathematics Department of Universidad Complutense de Madrid has been organizing annually since 1996 a regional math contest, *Concurso de Primavera de Matemáticas*, in the Madrid region of Spain.<sup>3</sup> As explained on the website, the contest has two main goals: to “motivate a large number of students by showing them that thinking and studying math can be fun,” and “to promote thinking outside the box and textbooks when solving problems, using logical reasoning, class geometry, parity issues, the properties of numbers, and probability.” It is a two-stage elimination math contest in which the best first-stage participants at each school are selected to participate in a second and final stage, which takes place at a single location. Each year, approximately 40,000 students participate in the first stage math test and over 3,000 students in the second stage math test. Iriberry and Rey-Biel (2019) analyzed within-subject gender differences between the two stages, which differ in terms of competitive pressure, using the 2014 data. In this study, we use data from the second-stage math test from the 2016, 2017 and 2018 contests which, in all cases, lasted for one and half hours.

---

<sup>3</sup> See the organization’s website at <https://www.concursoprivavera.es/#concurso> for more details.

A large number of schools in Madrid participate in this initiative. As shown in Iriberry and Rey-Biel (2019), the sample of participating schools ranges between 30% of the primary education schools and 50% of the secondary education schools in the region (see Table A1 in Iriberry and Rey-Biel, 2019). Regarding the school characteristics, they are more likely to be private schools, they tend to have a relatively large number of students and, as expected, show better results in mathematics, as measured by the standardized test administered and evaluated by the Department of Education in the region of Madrid.<sup>4</sup>

The rules of the math test we study are clearly established. First, there are four different tests, one for each age group. These are referred to as levels 1 to 4 and are grouped such that students from two consecutive school years take the same math test. Thus, level 1 includes children in their fifth and sixth academic years of primary school, and participants are therefore aged 10 and 11. Similarly, level 2 includes 12-13-year-olds, level 3 includes 14-15-year-olds and level 4 includes 16-17-year-olds. Second, the math test takes place on the campus of Universidad Complutense de Madrid on a pre-specified day in April. Third, the top three contestants in each level obtain prizes. Additionally, the top 5% of participants receive a diploma and a small gift in a public ceremony.<sup>5</sup> Fourth, the test for each level consists of 25 multiple-choice questions, all of which are set by the organizers. The questions for each level are designed so that students in the lower school year in each level have already seen the necessary material to answer the questions correctly.

Each question has 5 possible answers, only one of which is correct. Up to 2015, the scoring rule was the same for all 25 questions: 0 for wrong answers, +1 point for omitted questions, and +5 points for correctly answered questions. For the 2016, 2017 and 2018 contests, we collaborated with the organizers to create a math test with two parts that would differ in terms of the scoring rule. For the first 13 questions, the grading system awards 0 points for both omitted questions and wrong answers and +5 points for questions answered correctly. For the remaining 12 questions, questions 14-25, the grading system

---

<sup>4</sup> In particular, we use the standardized test called “Conocimientos y Destrezas Indispensables” (CDI – “Essential Knowledge & Skills”), which includes the subjects of Math, Spanish Language and General Culture. For more information, see <http://www.educa2.madrid.org/web/cdi/pruebas-cdi>

<sup>5</sup> As indicated on the website, what the main prizes will be is not revealed ex-ante. In past years, prizes were scientific calculators or iPads, and the gifts for the top 5% in stage 2 were books. The most important reward is the prestige associated with being among the top 5% of all contestants, which is publicly announced on the website and in a public award ceremony.

awards 0 points for wrong answers, +1 point for omitted questions and +5 points for questions answered correctly. Figure A1 in the Appendix shows how the scoring rule was described to participants.

We explicitly instructed the organizers to keep other things the same, i.e., the content or difficulty of the questions. The mean values of correct answers per question for all the questions on the math test are presented in Figure A2, separately for years in which there was no intervention on the scoring rule between the first and the second parts (2013-2015) and for years in which there was such intervention (2016-2018). As expected, in the years in which there was an intervention, there are more correct answers in the first part than in the second part of the test, as in the first part participants answer more questions than in the second and by luck there will be more correct answers. As additional evidence that the findings can be attributed to the change in the scoring rule and not to the first/sooner and second/later parts of the tests, we perform two robustness checks. First, we use the previous three editions of the exam (2013-2015, with no differential scoring rule across the first and second parts of the tests) as placebo (see Table A2 in the Appendix). Second, we perform a triple difference regression between male/female, first/second parts and 2013-2015/2016-2018 editions (see Table A3 in the Appendix).

Finally, after studying the performance results for the 2016 contest and to better understand the underlying mechanism, we administered a questionnaire immediately after the end of the math test to the participants in 2017 and 2018. Figure A3 in the Appendix includes an English version of the questionnaire. The first five questions listed were used in Iriberry and Rey-Biel (2019), as they focused on the differences between the stage 1 and stage 2 tests. We included questions 6 to 10 to understand whether gender differences in hours of preparation, confidence, overconfidence, risk preferences and perceived math ability can explain any of the gender differences observed in the number of omitted questions. We will use these controls to shed light on the underlying mechanism, which will be summarized in Section 4 and extended in Appendix A.

## **2.2 Descriptive Statistics**

The database consists of the participants who took the 2016, 2017 and 2018 editions of the test. Table 1 shows the descriptive statistics of the main outcome and control variables, overall and by gender. The last column shows the  $p$ -values for the F-

test of equality of variances across gender for the continuous variables and Fisher's exact test for categorical values.

[Table 1 about here]

This database contains a total of 9,907 math tests from 7,833 different participants. It is not a gender balanced sample, as 66% of the test takers are male. Looking at the control variables, we see that some students participate in multiple editions of the contest (*Participation Time*). In particular, 198 participants take the math test in all three years, 991 participate twice, and the remaining 7,331 take the math test just once. Female participants are less likely to participate more than once. The three different contests do not show large differences in overall participation or female participation. Regarding participation in different levels, level 2 is the most popular, and level 4 has the lowest number of participants. Female representation is also lowest in the last level, which is partially explained by female students being less likely to choose the math-science track in high school.

The performance data include the rank, score, and number of correct and omitted questions for each part of the test. When students register to take the math test, schools are asked to provide participants' math grade at school, which is available for approximately 90% of participants. For regression analysis, we will use the standardized math grade at school level in order to control for softer and more stringent schools. As expected, participants have on average high grades in math (*Math at School*), with an average of 8.40 out of 10, and female students indeed show higher performance than males (8.55 for females and 8.32 for males). However, the gender differences reverse when looking at the score on the math test we study, as on both parts of the test, female participants obtain a lower score than male participants. On the first part of the test, when there is no reward for omitted questions, male participants obtain an average score of 29.50 points, and females obtain an average score of 26.47 (out of the maximum score of 65). On the second part, when there is a reward for omitted questions, males obtain 23.30 points on average, while females obtain 20.67 points (out of the maximum score of 60). The slight difference in score between the first and the second part of the test is because the first part has 13 questions while the second part has 12 questions (see Table A5 for robustness). This also carries over into the ranking between male and female participants. Females rank lower than males, on average approximately 54 positions behind (with an

average of approximately 1000 positions), and this difference increases for the math test with the reward for omitted questions, where female participants rank on average 64 positions behind.

[Figure 1 about here]

The number of omitted questions, which is the focus of this paper, shows clear gender differences between the no reward and the reward parts of the test, which is consistent with the mainstream literature. Figure 1 shows the cumulative distribution of the *No. of Omitted* by gender when there is no reward (top) and when there is a reward for omitted questions (bottom), which complements the descriptive statistics in Table 1 well. Note that when there is no reward, the optimum behavior is to answer all questions, but when there is a reward for omitted questions, the optimum behavior depends on one's knowledge, confidence and risk aversion. Although participants should answer all questions when there is no reward for omitted questions, participants indeed omit on average 0.65 questions. In addition, women leave slightly more questions unanswered, at 0.86 questions; thus, while 80% of male participants indeed answer all questions, only 74% of female participants do. More importantly, when there is a reward, participants on average leave 4.82 questions unanswered, with male participants leaving 4.51 questions unanswered and females 5.40. In both panels of Figure 1, the distribution of female participants stochastically dominates that of male participants, and the differences on the reward part of the test are larger. In Table 1, we see that male participants also have a higher number of correct answers and a higher proportion of correct answers than females, but these differences are not large across the two parts.

### 3. Results

#### 3.1. Do Female Participants Leave More Unanswered Questions than Males When There is a Reward for Omitting Questions Compared to When There is No Reward?

Given our within-subject design, our identification strategy relies on comparing gender differences in performance between the no reward and rewarded parts of the multiple choice test. The outcome variables of interest are the number of omitted questions ("*zomitted*"), the proportion of correct answers over the number of answered questions ("*zprop\_correct*") and the final score ("*zscore*") and ranking ("*zrank*"). We use

standardized values by contest year, level and part of the test for all outcome variables, as these refer to different tests and questions. Table 2 shows the estimation results.

[Table 2 about here]

Columns 1-4 show the estimation results for the OLS specification with the standard errors clustered at the participant level. All regressions control for year, level and school fixed effects. The coefficients of interest are *Female* and, in particular, the interaction between *Female* and *Reward*. Female participants leave on average more omitted questions than males on the no reward part of the test (0.17 standard deviations of the mean), but most importantly, this difference increases 0.14 additional standard deviations of the mean on the reward part of the test. This is not the case for the proportion of correct answers. Though female participants show a lower proportion of correct answers (0.17 standard deviations of the mean), this difference does not increase on the reward part. Female participants leaving more questions unanswered than men on the reward than the no reward part has important consequences for how male and female participants perform under different reward systems. Female participants score on average worse than males (0.22 standard deviations of the mean) and receive lower rankings (52 positions behind) for the no reward part of the test. More importantly, this gap increases when there is a reward for omitted questions. Regarding the score, the gender gap increases by 0.04 standard deviations of the mean. Regarding the ranking, the gender gap increases by approximately 11 more positions.

Columns 5-8 and columns 9-12 in Table 2 show the equivalent estimation results for the random effects and individual fixed effects model specifications. Random effects and individual fixed effects models assume different specifications regarding the error term and therefore allow testing for robustness to the specification of the main effects. The variable of interest, the interaction between *Female* and *Reward*, maintains the same magnitude and significance levels. Hereafter, we will use the OLS estimation, with the standard errors clustered at the participant level.

We comment on the effect of the two main control variables: *Participation Time* for experience with the math test and *Math at School* for ability. The more experienced the participant, as one would expect, the higher the score, the lower the number of omitted questions and the higher the proportion of correct answers. Regarding the control for ability, we find that the higher the math grade at school, as expected, the better the score

and the higher the proportion of correct answers. Unexpectedly, the higher the math grade at school, the higher the number of omitted questions. However, note that in the fixed effects specification (column 9), the math grade is, as one would expect, negatively correlated with the number of omitted questions. As a robustness test, Table A1 in the Appendix shows the exact same table but with an alternative measure for ability; instead of *Math at School*, we control for individual ability by the number of correct answers on the no reward part of the test. The results for the main variable of interest, the interaction between *Female* and *Reward*, are very similar in terms of both the magnitude and the significance levels. Finally, note that, as shown in Figure 1, not controlling for ability in any way, leaves the main result on the coefficient of *Female\*Reward* unaffected.

To summarize, female underperformance increases when moving from the no reward to the reward part of the test, showing that differential grading of omitted questions and wrong answers hurts women more than non-differential grading for omitted questions and wrong answers.

### **3.2. Robustness Tests**

This section shows a series of robustness tests of the main gender differential result shown in Section 3.1.

First, choosing the non-differential scoring rule in the first part of the test and the differential scoring rule in the second and not varying this specific order is not ideal. It is possible that the observed effect is due to male and female participants reacting differently to fatigue over time when taking the test, although we could not find any evidence for this type of behavior. If female participants tend to get tired or lose interest in the math test *before* males do, the observed effect would be confounded with gender differences in performance due to fatigue over time. To rule this possibility out, we perform two additional analyses. First, we use the exact equivalent data from the 2013, 2014, and 2015 contests as placebo. Remember that in these three years, the scoring rule did not change across the test; thus, we can measure whether male and female participants show differential performance and willingness to guess between the first and second parts of the test. Table A2 in the Appendix shows the results: columns 1 to 4 for the 2013 contest, columns 5-8 for the 2014 contest and columns 9-12 for the 2015 contest. We find no evidence of any gender difference in performance between the different parts of 2013 and 2014 tests. In 2015, females showed better performance on the second half of the test.

Second, we perform a triple difference regression between male/female, first/second parts and 2013-2015/2016-2018 editions. Table A3 in the Appendix shows the results. Columns 1-4 do not control for individual math ability, while columns 5-8 control for math grade at school and columns 9-12 control for the number of correct answers in the non-rewarded part of the test. In this new specification the variable of interest is the triple interaction *Female\*Second Half\*Editions 2016-2018*. We can observe that all the results of Table 2 are unaltered in this new specification. These additional identification strategies show that the results are robust.

Second, we already observed that about 25% of participants do not choose an optimal strategy in the no-reward part of the test leaving one or more questions unanswered. Furthermore, we observed that women tend to show higher frequencies for this non-optimal behavior. This poses the question: is the main finding on gender difference in willingness to guess driven mainly by gender differences in confusion and/or non-optimal reaction to incentives? We therefore proceed to replicate the analysis restricting the overall sample to those participants who showed optimal behavior in the no-reward part of the test, that is, those participants who left 0 questions omitted (7,740 out of 9,900 participants). Estimation results in Table A4 show that the effect is not driven by gender differences in confusion and that the magnitude of the effect is even larger (0.33 standard deviations of the mean) if we restrict our main analysis to participants who showed optimal behavior in the first part of the test.

Third, although we standardized the first and the second parts of the test, the first part included 13 questions while the second included 12 questions. To test if this unbalance number of questions across the two parts affected in any way the results, we performed the regression analysis 13 times, taking one question from the first part out each of the times. Table A5 shows the estimated main coefficient (*Female\*Reward*) for the 13 different regressions for each of the relevant outcome variables. As can be seen by these estimated coefficients, the gender differential is very robust both in terms of magnitude and significance.

Next, one issue that we have not considered so far is gender differences in time management during the test. It could be the case that female participants need more time and that the gender difference in the number of omitted questions is mostly driven by participants' behavior in the very last part of the test. Given we have performance in all

questions in the test, we can indeed split the second part of the test, the last 12 questions, into two different parts: questions 14-19 (labeled, Reward\_14-19) and questions 20-25 (labeled, Reward\_20-25). In Table A6 in the Appendix, we can see the estimation results, when we split the reward part into two. The gender difference in the number of omitted questions seems to be equally present in both the first half and the second half of the reward part of the test, such that we can rule out the gender difference in time management as a plausible alternative explanation.

Fifth, male and female participants may differ in their knowledge of math. We do not find any support for this when looking at the math grades from school, as female participants indeed outperform males in this domain (see Table 1). However, if we look at the number of correct responses on the no reward part of the test, we do see that while male participants obtain approximately 5.90 correct answers, females obtain approximately 5.29 correct answers. Note that, in all our analyses so far, we do control for math ability (using either math grade at school or the number of correct answers on the no reward part of the test), so any differences due to ability are being controlled for. As a further robustness test to account for ability differences between male and female participants, we have replicated the main regression analysis, shown in Table 2, using a subsample where we matched male and female participants with the same number of correct answers on the no reward part of the test. While column 1 in Table A7 in the Appendix shows that in the overall sample, female participants show worse performance, both on the reward and no reward parts of the test, column 2, by construction, shows that in the matched sample, females are comparable to males in terms of performance on the no reward part of the test, although on the reward part of the test, the gender gap persists. Columns 3-6 in Table A7 in the Appendix replicate the main results for the rest of the outcome variables with the matched sample. Though these male and female participants do not differ in their ability, the gender differential is still significant for the reward part of the test.

Finally, one may argue that the use of a within-subject design makes participants prone to experimenter demand effects, which may potentially explain differential performance across the two parts of the test by the same subject. This could also explain our main result if women were more prone to experimenter demand effects than men. De Quidt et al. (2018) actually find, in a laboratory setting in a different context in which subjects are aware of their participation in an experiment, that when inducing strong

experimenter demand effects women in fact react more strongly to them than men. We believe that our results are unlikely driven by pure experimenter demand effects. The strongest evidence for this is that, other studies, using between-subject designs and therefore not affected by experimenter demand effects, find results that are consistent with ours (such as in Baldiga, 2014).

### **3.3. Analysis along the Ability Distribution: Are High-Ability Female Participants Particularly Affected?**

An important source of variation when looking at a large sample of math test takers is ability. There are two possible proxies for ability. First, if we take the number of correct answers on the part with no reward as a proxy for ability, we can observe in Figure 2 that there is large variation. The number of correct answers varies between 0 and 13, with a median of 6. Second, the math grade at school also shows some but definitely less variation, since those selected to participate in the second stage of the contest tend to be the best math students at each school and thus obtain similar grades. Although both measures show a positive correlation (0.22), due to its larger variation, we use the number of correct answers on the no reward part of the test as a proxy for ability and use the variation in math grade at school level as a robustness test, which we will discuss as a robustness check.

[Figure 2 over here]

We now study whether the gender differential in the number of omitted questions between the no reward part of the test and the reward part of the test varies substantially by participants' ability.

[Figure 3 over here]

Figure 3 displays the gender differences by ability graphically. Figure 3a shows the number of omitted questions on the non-differential scoring rule and the differential scoring rule parts of the test by low and high ability and by gender. We define low ability as the standardized number of correct answers on the no-reward part being below 0 and high ability as the standardized number of correct answers being above 0. As expected, high-ability participants leave fewer omitted questions on both parts of the test. Additionally, female participants always leave more questions unanswered. However, the

gender difference between the two parts is larger among the high-ability participants. Figures 3b for low ability and 3c for high ability take a closer look at the number of omitted questions on the reward part of the test by gender. Low- and high-ability female participants behave similarly, although as expected, high-ability females leave fewer questions unanswered. However, for male participants, low- and high-ability participants' behavior differs substantially, particularly with significantly more participants omitting no questions at all, which is less evident for female participants.

[Table 3 about here]

Table 3 shows the regression analysis results for the number of omitted questions. We take two complementary approaches.

First, as shown in columns 1 and 2, we consider a binary category for low and high ability using the standardized value of the number of correct questions on the no reward part of the test. For the low-ability participants, the gender differential is not significantly different from zero, while for the high-ability participants, it is highly significant and the magnitude is much higher (0.26 standard deviations of the mean) than the average effect we found in Section 3.1 (0.14 standard deviations of the mean). As shown in column 3, the triple interaction of *Female*, *Reward* and *High Ability* is highly significant, and the magnitude corresponds to the difference between the female and reward coefficients in columns 1 and 2.

Second, we also consider a continuous variable of ability, looking at the actual number of correct questions on the no reward part of the test. Column 4 shows the interaction among *Female*, *Reward* and the *No. of Correct Answers No Reward*, showing, consistent with the results in previous columns, that the gender differential when moving from the no reward part to the reward part is larger among the participants of higher ability (0.07 standard deviations of the mean). As a robustness test, we also perform the same exercise but use standardized math grade at school level as an alternative proxy for ability. Table A8 in the Appendix shows the results. The conclusions are very similar when looking at the interaction between *Female* and *Reward* for the low- and high-ability students, although the magnitudes are slightly lower when using standardized math grade at school level as a proxy for ability.

Third, and finally, we can also split the sample by the main outcome variable, which is the score participants obtain in the test. In columns 5 and 6 in Table 3, we show the same regression analysis for participants who ranked in the bottom 25% and for those who ranked in the top 25%. We see a very similar pattern by ability. For those who show low performance in the test, gender differences are large but they do not show a differential behavior between the scoring rules, while for those who show high performance in the test, gender differences are small but they show a much larger difference between the two scoring rules.

We conclude that gender differences in reaction to a differential scoring rule are larger among the high ability than among the low ability participants.

### **3.4. Analysis Regarding Age: Are Younger/Older Female Participants Equally Affected?**

An interesting feature of our sample is that we can observe male and female participants from a young age (in their fifth and sixth academic years of primary school) up to their final two years of high school, right before going to university. Exploiting this variation, we test whether gender differences in willingness to guess vary with age.

[Table 4 about here]

Table 4 shows the results. Columns 1 to 4 show the regression analysis for each of the levels separately. The coefficient of interest, the interaction between *Female* and *Reward*, shows increasing magnitudes from the lowest academic level (0.03 standard deviations of the mean for the youngest participants in their 5<sup>th</sup> and 6<sup>th</sup> grade of primary school) to the highest academic level (0.23 among the oldest participants in high school). Column 5 shows the results when all levels are included in one regression to test how different the gender differences are across academic levels. The gender differential among older school participants (Level 2, Level 3 and Level 4) is significantly different from the gender differential among the youngest participants.

We therefore conclude that in our sample the gender differential in willingness to guess increases as participants get older. Acknowledging that in our sample, as participants get older they also get selected differently into the Math subject, this suggests that the gender differences would be if anything larger when tests that determine important educational and future labor markets such as university entry tests.

### **3.5. Analysis Regarding Public and Non-Public Schools: Are Female Participants from Different Socio-economic Backgrounds Equally Affected?**

As a final heterogeneity analysis, we explore the gender differential in willingness to guess that might depend on whether participants attend a public school (relatively lower socio-economic background) or a non-public school (relatively higher socio-economic background). The latter includes fully private schools and schools that use both public and private funding (referred to as mixed).

Table 5 shows the results. Columns (1) to (3) do not differentiate between the fully private and the mixed funding schools, while columns (4) to (6) do. Interestingly, participants from non-public schools show a much higher willingness to guess in the reward part of the test but most importantly, the gender interaction is not significant. We therefore find no evidence that female participants from different socio-economic backgrounds react differently to the grading scheme regarding the omitted questions.

## **4. Discussion on the Underlying Mechanism**

Motivated by the gender difference we found in the 2016 data and to understand the underlying mechanism, we designed a questionnaire that would allow us to measure the effects of confidence, risk aversion, and competitiveness, which we administered with the 2017 and 2018 editions of the test. Regarding confidence, we used two measures: confidence in their perceived math ability and the difference between their guessed number of correct answers and the actual number of correct answers, which we refer to as overconfidence. Risk attitudes were measured by asking participants about how certain they had to be about the answer in order not to skip it. We also measured participants' overall preferences for competition. According to our questionnaire, female participants show lower confidence in their perceived math ability, lower overconfidence, higher risk aversion and a lower competitive attitude. We acknowledge that these measures, and in particular our risk preferences measure, are not standard and that they may entail different aspects of answering strategy, so we use them to look for suggestive evidence on the underlying mechanism. See Figure A4 for graphical representation of all these measures by gender, as well as TA9 for the mean values by gender.

In Appendix A, we lay out our analysis of how much all these control variables contribute to the explanation of the gender differential in the willingness to guess. All of them take values in the expected direction: the more confident and more overconfident,

the more willing they are to provide an answer; the riskier, the more willing they are to provide an answer and the more competitive, the more willing they are to provide an answer. However, when all these controls are interacted with the reward part of the test, the female coefficient is left unchanged with the only exception of the measure on risk. We find that among the controls we include, the measure on risk preferences has the largest explanatory power, even though it does not fully explain the female coefficient.

## **5. Conclusion**

Using performance data from a natural field experiment with approximately 10,000 observations, we test for gender differences in willingness to guess when there is differential grading for omitting questions and providing a wrong answer. We find that women always leave more omitted questions but that this behavior becomes even more prominent when there is differential grading for omitted questions and wrong answers. This has negative consequences for female participants, both in terms of their final score and their ranking, demonstrably hurting female performance on the math test.

We also find that this gender differential is stronger among high-ability and older participants. Finally, we find suggestive evidence that gender differences in risk aversion are the most likely underlying mechanism. When explaining the stronger effect for high ability and older participants, although again risk preferences seem to be important, a significant part of the heterogeneity in the gender differential remains unexplained. We leave this interesting puzzle for future research.

Adding to the increasing amount of evidence, and based on a strong within-subject identification strategy, we conclude that a gender neutral grading rule requires non-differential scoring for omitted questions and wrong answers, at least compared to the alternative of having a mild reward for omitting questions. It is still an open question whether gender differences persist, or even if they change sign or size, when the differential score for omitted questions and wrong answers is very large.

## References:

- Akyol, S.P., Key, J. and Krishna, K. (2016). "Hit or Miss? Test Taking Behavior in Multiple Choice Exams." NBER Working Paper Nr. 22401.
- Anderson, J. (1989). "Sex-Related Differences on Objective Tests among Undergraduates." *Educational Studies in Mathematics*, 20(2):165–177.
- Atkins, W.J., Leder, G.C., O'Halloran, P.J., Pollard, G.H. and Taylor, P. (1991). "Measuring Risk Taking." *Educational Studies in Mathematics*, 22(3), 297-308.
- Atwater, A., Saygin, P. O. (2020), "Gender Differences in Willingness to Guess on High-Stakes Standardized Tests". Mimeo.
- Balart, P., Ezquerra, L., Hernandez-Arenaz, I. (2020). "Framing Effects on Test-Taking Behavior: Evidence from a Field Experiment". Universitat de les Illes Balears, *Mimeo*.
- Baldiga, K. (2014). "Gender Differences in Willingness to Guess." *Management Science*, 60(2): 434-448.
- Barber, B.M., and Odean, T. (2001). "Boys Will Be Boys: Gender, Overconfidence and Common Stock Investment." *Quarterly Journal of Economics*, 116(1), 261-292.
- Beyer, S. (1999). "Gender Differences in the Accuracy of Grade Expectancies and Evaluations." *Sex Roles*, 41:314, 279-296.
- Bordalo, P., Coffman, K., Gennaioli, N., and A. Shleifer. (2019). "Beliefs about Gender." *American Economic Review*, 109 (3): 739-73.
- Coffman, Katherine B., and David Klinowski. (2020). "The Impact of Penalties for Wrong Answers on the Gender Gap in Test Scores." *Proceedings of the National Academy of Sciences of the United States of America*, 117 (forthcoming). (Pre-published online April 6, 2020).
- Crosan R. and Gneezy U. (2009). "Gender Differences in Preferences." *Journal of Economic Literature*, 47(2): 448-474.
- De Quidt, J., Haushofer, J. and Roth, C., (2018). "Measuring and Bounding Experimenters Demand." *The American Economic Review* 108 (11): 3266-3202.
- Dohmen, T, A. Falk, B. Golsteyn, D. Huffman and U. Sunde (2017), "Risk attitudes across the life course." *The Economic Journal* 127(605): 95-116.
- Eckel C. and Grossman P. (2008). "Men, Women and Risk Aversion: Experimental Evidence." *Handbook of Experimental Economics Results*.
- Espinosa M.P. and Gardeazabal J. (2013). "Do Students Behave Rationally in Multiple Choice Tests? Evidence from a Field Experiment." *Journal of Economics and Management*, Vol. 9, No. 2, 107-135.

Espinosa Alejos, María Paz & Gardeazábal, Javier, 2010. "Optimal Correction for Guessing in Multiple-Choice Tests." *Journal of Mathematical Psychology* 54(5), 415-425.

Filippin, A., and Crosetto P. (2016) "A Reconsideration of Gender Differences in Risk Attitudes." *Management Science* 62(11): 3138-3160.

Funk, P., and Perrone, H. (2016). "Gender Differences in Academic Performance: The Role of Negative Marking in Multiple-Choice Exams." Working Paper.

Gneezy, U., List, J. A., Livingston, J. A., Qin, X., Sadoff, S. and Xu, Y. (2019). "Measuring Success in Education: The Role of Effort on the Test Itself." *American Economic Review: Insights* 1(3): 291-308.

Iriberry, N. and Rey-Biel, P. (2019). "Competitive Pressure Widens the Gender Gap in Performance: Evidence from a Two-Stage Competition in Mathematics." *The Economic Journal* 129(620), pp. 1863–1893.

Karle, H., Engelmann, D., and M. Peitz, (2019) "Student Performance and Loss Aversion," *Rationality & Competition Discussion Paper No. 181*.

Niederle, M. and Vesterlund, L., "Gender and Competition", *Annual Review in Economics*, 2011, 3, 601–30.

Pekkarinen, T., (2015). "Gender Differences in Behaviour under Competitive Pressure: Evidence on Omission Patterns in University Entrance Examinations". *Journal of Economic Behavior and Organization*, 115: 94-110.

Ramos, I. and Lambating, J. (1996). "Gender Difference in Risk-Taking Behavior and their Relationship to SAT-Mathematics Performance." *School Science and Mathematics*, 96(4): 202-207.

Riener, G., Wagner, V. (2017). "Shying Away from Demanding Tasks? Experimental Evidence on Gender Differences in Answering Multiple-choice Questions." *Economics of Education Review*, 59, 43-62.

Gerhard R., Valentin W. (2018). "Gender Differences in Willingness to Compete and Answering Multiple-choice Questions: The Role of Age". *Economics Letters*, 164: 86-89.

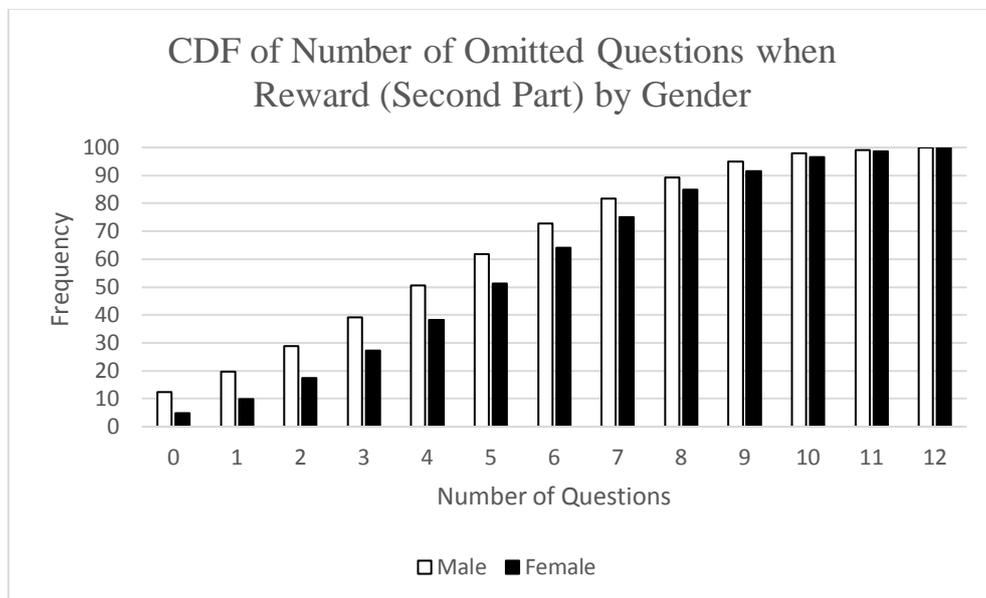
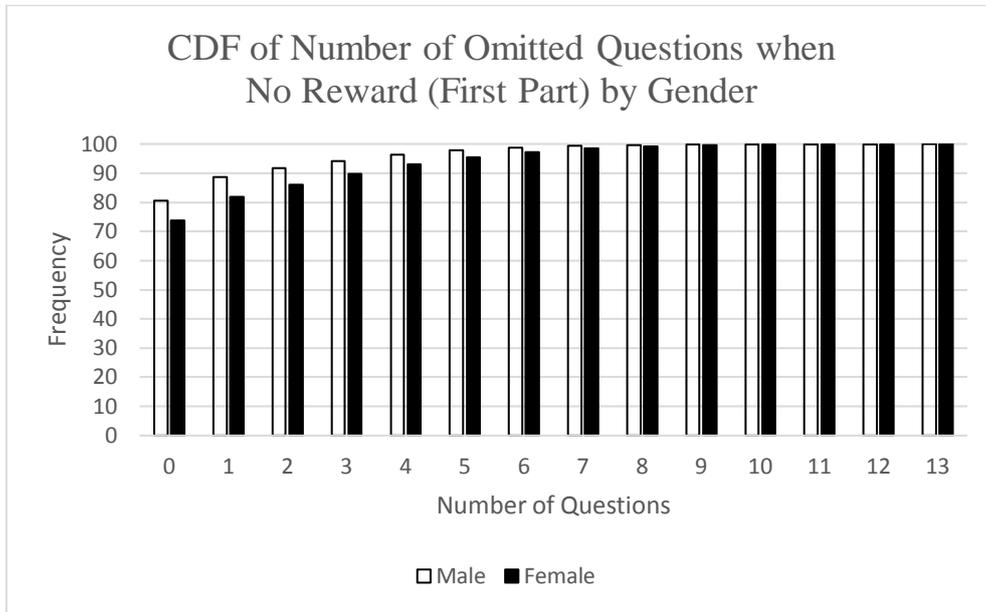
Schildberg-Hörisch, H., (2018). "Are Risk Preferences Stable?". *Journal of Economic Perspectives* 32(2), 135–154.

Swineford, F. (1941). "Analysis of a Personality Trait." *Journal of Educational Psychology*, 32(6):438–444.

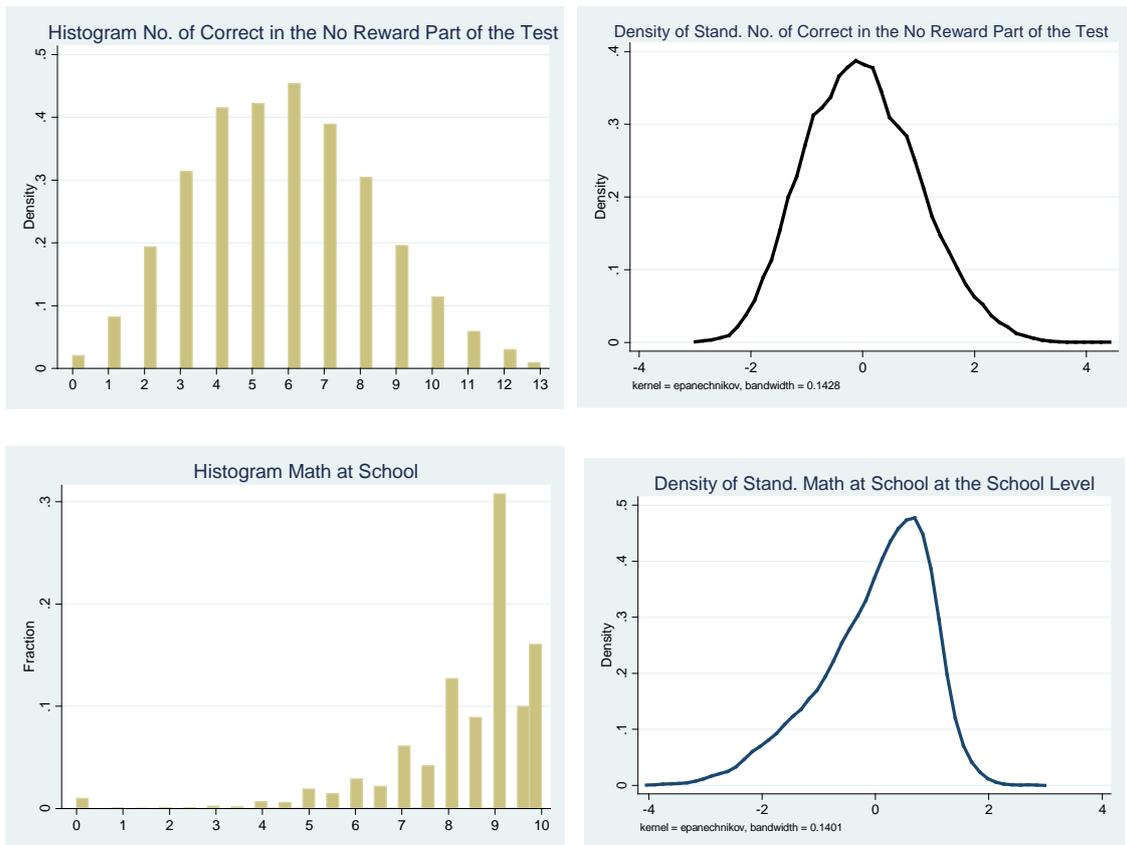
Tannenbaum D. (2012). "Do Gender Differences in Risk Aversion Explain the Gender Gap in SAT Scores? Uncovering Risk Attitudes and the Test Score Gap". Unpublished paper, University of Chicago, Chicago.

**Figures and Tables**

**Figure 1. No. Omitted when No Reward and when Reward by Gender**

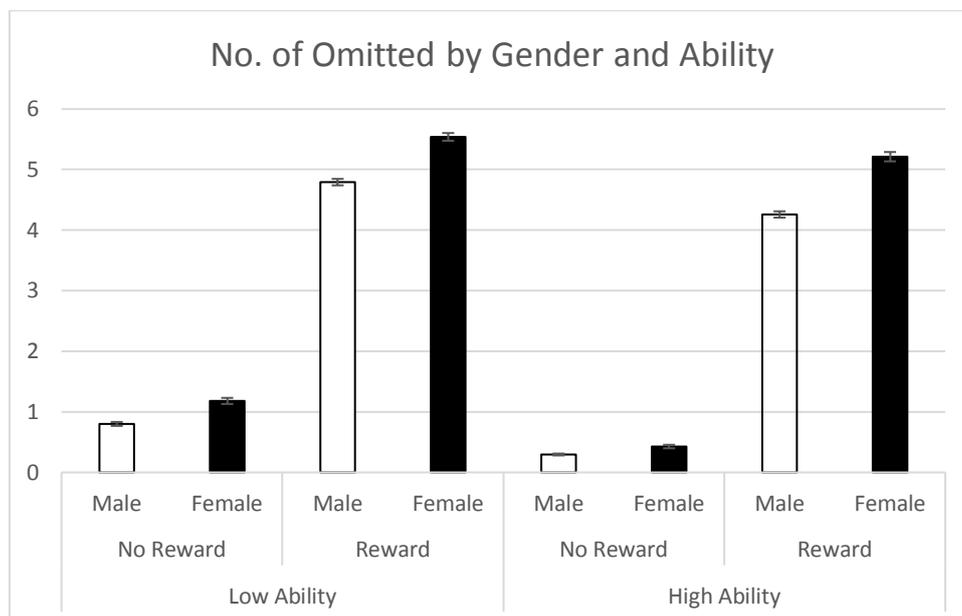


**Figure 2. Variation in No. of Correct No Reward Part of the Test and in Math at School**

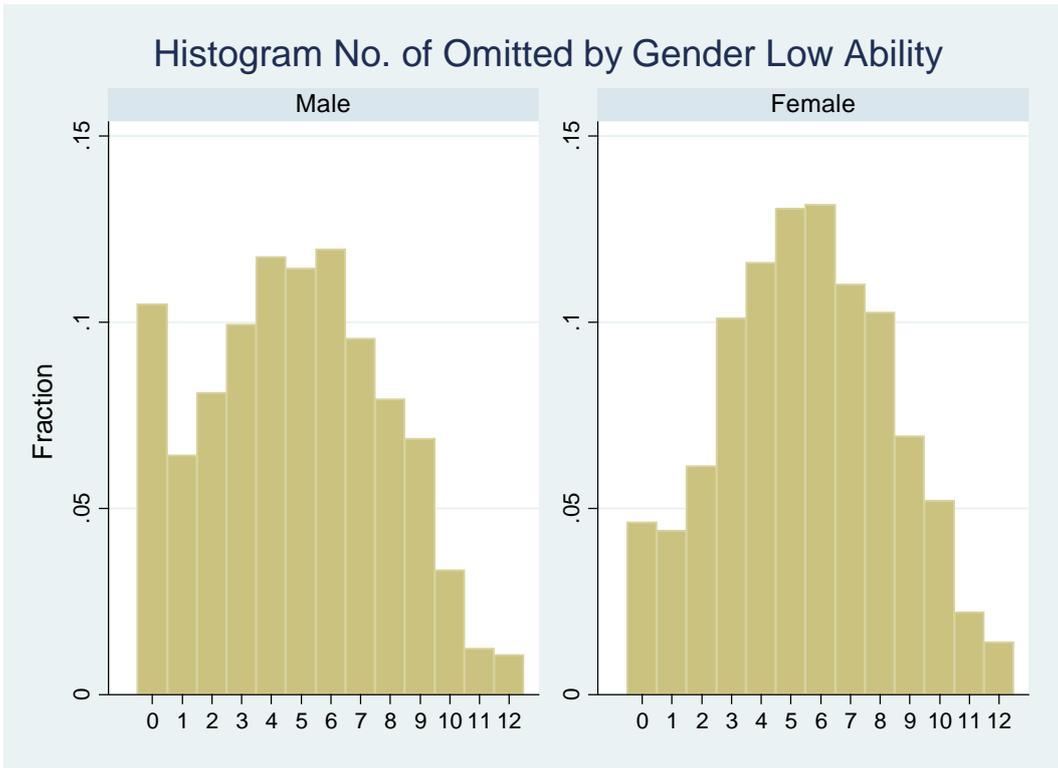


**Figure 3. No. of Omitted by Gender and Ability: Low Ability: Standardized No. of Correct in No Reward < 0 and High Ability: Standardized No. of Correct in No Reward > 0**

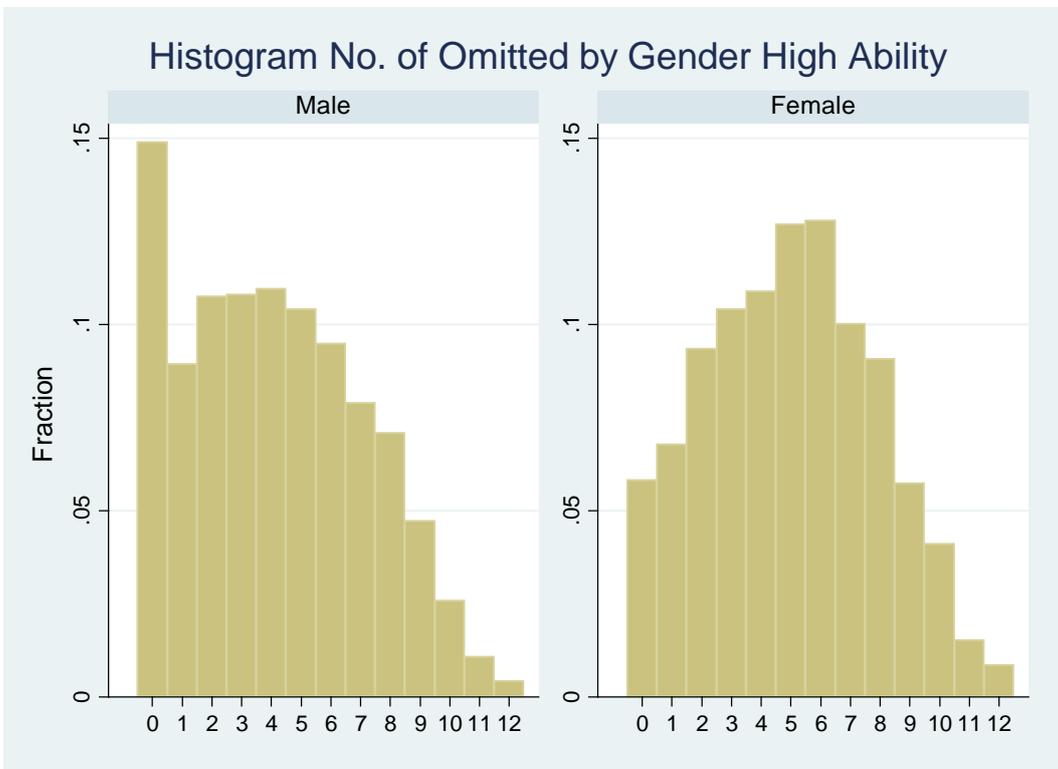
**3a. No. of Omitted by Gender, Scoring Rule and by Ability**



**3b. No. of Omitted when Reward for Omitted by Gender for Low Ability**



**3c. No. of Omitted when Reward for Omitted by Gender for High Ability**





**Table 2. Gender Differences between the No Reward and the Reward Parts of the Test**

	OLS				RE				FE			
	zomitted (1)	zprop_correct (2)	zscore (3)	rank (4)	zomitted (5)	zprop_correct (6)	zscore (7)	rank (8)	zomitted (9)	zprop_correct (10)	zscore (11)	rank (12)
Female	0.170*** (0.0249)	-0.167*** (0.0210)	-0.216*** (0.0204)	-51.91*** (5.443)	0.171*** (0.0252)	-0.170*** (0.0211)	-0.216*** (0.0204)	-52.59*** (5.438)				
Reward	-0.0462*** (0.0162)	-0.00742 (0.0140)	0.0142 (0.0131)	39.65*** (3.469)	-0.0462*** (0.0162)	-0.00742 (0.0140)	0.0142 (0.0131)	39.65*** (3.469)	-0.0462*** (0.0160)	-0.00742 (0.0138)	0.0142 (0.0129)	39.65*** (3.415)
Female*Reward	0.145*** (0.0295)	0.0194 (0.0245)	-0.0449** (0.0218)	-10.74* (5.941)	0.145*** (0.0295)	0.0194 (0.0245)	-0.0449** (0.0218)	-10.74* (5.941)	0.145*** (0.0290)	0.0194 (0.0242)	-0.0449** (0.0215)	-10.74* (5.849)
Math at School	0.0381*** (0.00862)	0.238*** (0.00822)	0.230*** (0.00841)	55.97*** (2.140)	0.0367*** (0.00851)	0.224*** (0.00818)	0.213*** (0.00830)	52.10*** (2.145)	-0.0438** (0.0208)	0.0660*** (0.0248)	0.0753*** (0.0247)	15.96** (6.792)
Participation Time	-0.156*** (0.0195)	0.302*** (0.0202)	0.379*** (0.0216)	97.35*** (5.337)	-0.128*** (0.0183)	0.214*** (0.0198)	0.265*** (0.0204)	68.26*** (5.208)	-0.129* (0.0691)	-0.0132 (0.0764)	0.0524 (0.0755)	1.359 (21.08)
Observations	17,850	17,850	17,850	17,850	17,850	17,850	17,850	17,850	17,850	17,850	17,850	17,850
R-squared	0.096	0.236	0.288	0.334					0.015	0.026	0.040	0.101
Number of participants					7,786	7,786	7,786	7,786	7,786	7,786	7,786	7,786

Notes: Observations are at the Math test's parts level. The first three outcome variables, *No. Omitted*, *Prop. of Correct*, and *Score* are standardized at the edition, level and part of the test levels. *Rank* measures the position in the rank by edition, level and test's parts level, where higher values represent better positions within the rank. *Female* takes the value of 1 if the participant is female. *Reward* takes the value of 1 if the outcome variable refers to the part of the test with reward for omitted questions. *Math at School* measures the standardized Math grade at school level and *Participation Time* takes the values of 1, 2, 3 if it is the first, second or third time that the participant does the Math test. Columns 1-4 show the OLS specification where the standard errors are clustered at the participant level. Columns 5-8 show the RE model specification and columns 9-12 show the FE specification model. All specifications include edition, level and school fixed effects. Standard errors in parenthesis, where \*\*\* p<0.01, \*\* p<0.05, \* p<0.1

**Table 3. Gender Differences between the No Reward and the Reward Parts of the Test:  
Variation along the Ability Distribution**

	Low Ability omitted (1)	High Ability omitted (2)	Interaction omitted (3)	Continuous omitted (4)	Bottom 25% omitted (5)	Top 25% omitted (6)
Female	0.191*** (0.0387)	0.0848*** (0.0226)	0.201*** (0.0380)	0.405*** (0.0737)	0.273*** (0.0608)	0.0928*** (0.0306)
Reward	-0.0929*** (0.0248)	-0.00137 (0.0194)	-0.0929*** (0.0244)	-0.105** (0.0447)	-0.0934** (0.0404)	-0.252*** (0.0233)
Female*Reward	0.0505 (0.0425)	0.264*** (0.0358)	0.0505 (0.0418)	-0.228*** (0.0789)	-0.0451 (0.0651)	0.192*** (0.0471)
High Ability			-0.0172 (0.0297)			
High Ability*Reward			0.0915*** (0.0307)			
Female*High Ability			-0.122*** (0.0439)			
Female*Reward*High Ability			0.213*** (0.0541)			
No. Of Correct No Reward	-0.0887*** (0.0117)	-0.0602*** (0.00540)	-0.0743*** (0.00596)	-0.0712*** (0.00561)	-0.107*** (0.0195)	-0.0139** (0.00657)
Participation Time	-0.0192 (0.0387)	-0.0891*** (0.0195)	-0.0675*** (0.0182)	-0.0683*** (0.0182)	-0.0288 (0.0808)	-0.0994*** (0.0220)
No. Of Correct No Reward*Reward				0.0102 (0.00657)		
Female*No. Of Correct No Reward				-0.0477*** (0.0106)		
Female*Reward*No. Of Correct No Reward				0.0691*** (0.0122)		
Observations	10,048	9,720	19,768	19,768	4,930	4,904
R-squared	0.123	0.153	0.114	0.115	0.203	0.206

Notes: Observations are at the Math test's parts level. *No. Omitted* is standardized at the edition, level and part of the test levels. *Female* takes the value of 1 if the participant is female. *Reward* takes the value of 1 if the outcome variable refers to the part of the test with reward for omitted questions. *No. of Correct No Reward* measures the number of correct questions in the part of the test without any reward for omitted questions, and *Participation Time* takes the values of 1, 2, 3 if it is the first, second or third time that the participant does the Math test. *High Ability* takes value 1 if the participant's standardized number of correct answers in the no reward part is > 0. Bottom25 and Top25 are selected samples in which the total zscore take values that are among the bottom 25 and top 25%. All columns show the OLS specification where the standard errors are clustered at the participant level and include edition, level and school fixed effects. Standard errors in parenthesis, where \*\*\* p<0.01, \*\* p<0.05, \* p<0.1

**Table 4. Gender Differences between the No Reward and the Reward Parts of the Test:  
Variation along Age**

	Level 1 zomitted (1)	Level 2 zomitted (2)	Level 3 zomitted (3)	Level 4 zomitted (4)	Overall zomitted (5)
Female	0.214*** (0.0502)	0.199*** (0.0453)	0.0809* (0.0459)	0.179** (0.0700)	0.213*** (0.0488)
Reward	-0.000431 (0.0320)	-0.0614** (0.0289)	-0.0583* (0.0313)	-0.0613 (0.0424)	-0.000431 (0.0314)
Female*Reward	0.0273 (0.0587)	0.161*** (0.0521)	0.184*** (0.0549)	0.234*** (0.0861)	0.0273 (0.0576)
Math at School	-0.0118 (0.0201)	0.00452 (0.0171)	0.0978*** (0.0147)	0.0235 (0.0207)	0.0383*** (0.00863)
Participation Time	-0.272*** (0.0585)	-0.171*** (0.0379)	-0.166*** (0.0338)	-0.0741* (0.0440)	-0.156*** (0.0195)
Level 2					0.0758** (0.0371)
Level 3					0.114*** (0.0385)
Level 4					0.170*** (0.0442)
Female*Level 2					-0.0130 (0.0651)
Female*Level 3					-0.0949 (0.0656)
Female*Level 4					-0.0777 (0.0851)
Level 2*Reward					-0.0610 (0.0420)
Level 3*Reward					-0.0578 (0.0437)
Level 4*Reward					-0.0609 (0.0513)
Level 2*Female*Reward					0.134* (0.0765)
Level 3*Female*Reward					0.156** (0.0784)
Level 4*Female*Reward					0.207** (0.101)
Observations	4,250	5,764	5,044	2,792	17,850
R-squared	0.167	0.150	0.178	0.212	0.096

*Notes:* Observations are at the Math test's parts level. *No. Omitted* is standardized at the edition, level and part of the test levels. *Female* takes the value of 1 if the participant is female. *Reward* takes the value of 1 if the outcome variable refers to the part of the test with reward for omitted questions. *Math at School* measures the standardized Math grade at school level, and *Participation Time* takes the values of 1, 2, 3 if it is the first, second or third time that the participant does the Math test. All columns show the OLS specification where the standard errors are clustered at the participant level and include edition, level and school fixed effects. Standard errors in parenthesis, where \*\*\* p<0.01, \*\* p<0.05, \* p<0.1

**Table 5. Gender Differences between the No Reward and the Reward Parts of the Test:  
Variation along Public and Non-Public Schools**

	Public omitted (1)	Non-Public omitted (2)	omitted (3)	Public omitted (4)	Mixed omitted (5)	Private omitted (6)
Female	0.201*** (0.0366)	0.147*** (0.0340)	0.200*** (0.0365)	0.201*** (0.0366)	0.0970** (0.0384)	0.274*** (0.0713)
Reward	0.0328 (0.0239)	-0.122*** (0.0220)	0.0328 (0.0238)	0.0328 (0.0239)	-0.0930*** (0.0267)	-0.191*** (0.0385)
Female*Reward	0.138*** (0.0437)	0.152*** (0.0399)	0.138*** (0.0435)	0.138*** (0.0437)	0.140*** (0.0461)	0.174** (0.0802)
Math at School	0.0382*** (0.0125)	0.0366*** (0.0116)	0.0383*** (0.00854)	0.0382*** (0.0125)	0.0656*** (0.0135)	-0.0237 (0.0224)
Participation Time	-0.129*** (0.0292)	-0.188*** (0.0271)	-0.161*** (0.0200)	-0.129*** (0.0292)	-0.171*** (0.0319)	-0.224*** (0.0509)
Non-Public			-0.0932 (0.557)			
Female*Non-Public			-0.0506 (0.0500)			
Non-Public*Reward			-0.155*** (0.0324)			
Female*Reward*Non-Public			0.0143 (0.0591)			
Observations	8,826	8,896	17,722	8,826	6,348	2,548
R-squared	0.102	0.087	0.098	0.102	0.079	0.119

*Notes* : Observations are at the Math test's parts level. *No. Omitted* is standardized at the edition, level and part of the test levels. *Female* takes the value of 1 if the participant is female. *Reward* takes the value of 1 if the outcome variable refers to the part of the test with reward for omitted questions. *Math at School* measures the standardized Math grade at school level, and *Participation Time* takes the values of 1, 2, 3 if it is the first, second or third time that the participant does the Math test. *Non-Public* takes the value of 1 when the school is not publicly financed, which can be mixed or fully private. All columns show the OLS specification where the standard errors are clustered at the participant level and include edition, level and school fixed effects. Standard errors in parenthesis, where \*\*\* p<0.01, \*\* p<0.05, \* p<0.1

## Appendix A: Underlying Mechanism: Confidence, Risk, and Competitiveness

In this section we present the analysis on the underlying mechanism. We first describe the control variables we collected in the questionnaire (see Figure A3) and then present the main results when these control variables are included.

Table A9 shows the descriptive statistics of the control variables we collected in the survey administered with the 2017 and 2018 test editions. Figure A4 show the mean values of all these control variables by gender. The variables of interest are *No. of Preparation Hours*, *Overconfidence*, *Perceived Math Ability*, *Perceived Gender Nature of Math* and *Risk*. All these variables show significant gender differences with one exception (see *p*-value in the last column of Table A9).

Male and female participants show a very similar number of reported hours devoted to preparing for the test (see question 5 in Figure A3 in the Appendix).<sup>1</sup> Figure A4a shows the probability density function of the number of preparation hours by gender, truncated at the value of 30 because most participants' answers lie below that value, which again shows that both male and female participants devote a similar amount of time to preparing for the test.

Overconfidence is measured by the difference between the guessed number of correct answers (see question 7 in Figure A3 in the Appendix) and the actual number of correct answers; thus, the more positive the value, the higher the overconfidence. Figure A4b plots the observations; the x-axis shows the number of correct answers and the y-axis shows the number of guessed correct answers. Both male and female participants are overconfident (also found by Beyer, 1999). However, also consistent with other findings, as in Barber and Odean (2001), female participants in our setting show lower values of overconfidence than male participants. Note that overconfidence is measuring a lower bound of the gender difference, as it is restricted to the questions that were actually answered and women answer fewer questions. Related to confidence, male participants also show higher agreement with the statement "I am good at math" than female participants, as shown in Figure A4c (see question 9 in Figure A3 in the Appendix), so perceived math ability is higher for male than for female participants, which is also consistent with Bordalo et al. (2019). Finally, also related to confidence, we measure

---

<sup>1</sup> Fifteen participants reported very high numbers of preparation hours. We replaced those values with missing to avoid outliers.

participants' perception of the gendered nature of the math task (see question 10 in Figure A3 in the Appendix). As shown in Figure A4d, the large majority of participants (94.44% of them) believe math to be gender neutral, such that men and women are equally good at/knowledgeable about math. Interestingly, this is the case for participants of different ages, who show only minor differences across age: among the youngest participants, 96% of both male and female participants believe that math is gender neutral, while among the oldest participants, 91% of males and 93% of females believe that math is gender neutral. However, both genders show some type of bias: a small fraction of male participants believe that men are better at math than women, and a small fraction of female participants believe that women are better at math than men.

We measure risk by the following question (see question 8 in Figure A3 in the Appendix): “When omitting a question was worth 1 point I answered the question ....” There are 5 possible answers (from 1 for “When I was Absolutely Sure” to 5 for “Always”); the higher the number, the more risk-loving the participant is. Figure A4e shows the histogram of all possible answers by gender. Clearly, more female participants than males answer the test question when absolutely or almost sure, so male and female participants differ in their risk preferences, consistent with Eckel and Grossman (2008) Croson and Gneezy (2009) and Filippin and Crosetto (2016). Given our non-standard method of eliciting risk preferences, two observations are noteworthy.<sup>2</sup> First, we cannot distinguish risk aversion from loss aversion (Karle et al, 2019). Second, as mentioned in the introduction, our risk measure may also be influenced by confidence, as perceived probability of knowing the answer might also be affected by participants' ability and confidence in their own ability.

Finally, given that our setting is competitive, we also measure participants' non-competitiveness, using their answers to question 1 in the questionnaire: “It is more important for me to be selected for Stage 2 than being among the winners in Stage 2.” The more they agree with this statement, the less competitive they are. As expected, and consistent with the literature on gender and competitiveness (Niederle and Vesterlund,

---

<sup>2</sup> Due to institutional restrictions, we were not allowed to elicit risk preferences in a more standard manner, for example, offering participants monetary lotteries.

2011), female participants show lower degrees of agreement with this statement, as shown in Figure A4f.<sup>3</sup>

One might be concerned about the correlation between the different measures of confidence, overconfidence, non-competitiveness and risk, as well as how all these measures correlate with ability (standardized Math at school level). Regarding the correlations among confidence, overconfidence, non-competitiveness and risk, all four of them show low correlations, which suggests that they are indeed measuring different dimensions of personality.<sup>4</sup> Regarding their correlations with ability, as one would expect, confidence correlates with ability positively (0.20); overconfidence correlates negatively with ability (-0.14), as does risk-loving preferences, although the correlation is very low (-0.06). Non-competitiveness correlates also positively with ability (0.008). With the exception of the overconfidence measure, all other measures show low correlations with ability, which indicates that our risk and confidence measures are independent of ability. Nevertheless, we acknowledge that our risk preferences measure is confounded with answering strategy.

We now proceed to test if competitiveness, confidence and/or risk have explanatory power using the answers to the survey about answering strategy that participants in the 2017 and 2018 editions of the test filled up after the exam. In particular, we check whether the gender differential in the number of questions omitted between the no reward and the reward part of the test remains significant when these variables are controlled for.<sup>5</sup> Table A10 shows the estimation results from this exercise. Columns 1 and 2 show the main specification, as in Table 2, but in the sample for which we have control variables collected via the questionnaire. In column 2, we find that, as in the main sample, female participants omit more questions than males when moving from the no

---

<sup>3</sup> Our measure of competitiveness is also different from the standardly used measure of competitive preferences due to limitations imposed by the organization on what we could include in the questionnaire.

<sup>4</sup> In particular, confidence correlates positively with overconfidence (0.09), and negatively with non-competitiveness (-0.007) and risk (-0.04). In addition, overconfidence correlates positively with non-competitiveness (0.06) and negatively with risk (-0.02). Finally, non-competitiveness correlates negatively with risk (-0.11). All these correlations are nevertheless very low. Although correlation values differ slightly, we observe very similar patterns when looking at correlations between our preference measures separated by gender.

<sup>5</sup> Given our setting is competitive by nature, we also checked how much of the gender differential in willingness to guess might be explained by participants' attitudes toward competition. Using question 1 in the questionnaire, where participants assess how much they agree with the statement "It is more important to participate in the competition than to win the competition", we find that, although female participants show a less competitive attitude than males, this does not show any explanatory power in how male and female participants decide on their willingness to guess. These results are available upon request.

reward to the reward part of the test, although the magnitude is slightly lower than in the main specification. In column 3, we add the three main control variables: perceived math ability, overconfidence and risk, and the three of them have the expected sign. The more confident and the more risk-loving the participant is, the fewer the number of omitted questions. The female coefficient decreases, but the interaction of *Female* and *Reward* remains exactly the same as in column 2.

In columns 4 to 7, we interact each of the control variables with the variable *Reward*, as these control variables can have a differential effect under the two scoring rules. As for competitiveness, the more competitive the participant is the fewer the omitted questions. Importantly, when adding these interactions with respect to the competitiveness measure, which is insignificant as one might expect, the main result on *Female* and *Reward* remains unchanged, which suggests that competitiveness plays no role in explaining the gender difference in willingness to guess (column 4). Similarly, when adding these interactions with respect to the two confidence measures, the main coefficient of interest, the interaction between *Female* and *Reward*, changes very little, suggesting that confidence does not explain why female participants leave more questions unanswered (column 5 and 6). However, when interacting the risk measure with reward, we clearly see that the coefficient of *Female* and *Reward* decreases substantially such that it is no longer significant. This shows that differences in risk preference between male and female participants explain 40% of the gender gap and thus, it is indeed an important factor explaining why male and female participants differ in their behavior in omitting questions (column 7).

What about the gender differences found between the low- and high-ability participants and older participants? In Sections 3.3 and 3.4, we found that the gender difference in the number of omitted questions punished high-ability and older females in particular. Could it be that gender differences in risk and overconfidence are different between the low- and high-ability participants or among the younger/older participants?

We must first examine gender differences in confidence, overconfidence and risk by ability. Figure A5 shows the graphs. Gender differences are present among both the high- and low-ability participants, and they always follow the same pattern: female participants show lower perceived math ability, lower levels of overconfidence and higher risk aversion. However, the gender differences in these control variables between the low- and high-ability participants are not striking.

We perform a similar exercise as we do in Table A10 but in two sub-samples, the low- and high-ability participants. Table A11 shows the results. Columns 1 and 2 reproduce the main results found in the first two columns in Table 3, and columns 3 and 4 replicate the same results for the sample of participants for whom we have questionnaire responses. The estimated values of the main variable of interest, the interaction between *Female* and *Reward*, are very similar in the overall sample and the sample for which we have questionnaire responses. Columns 5 and 6 add the main control variables of competitiveness, confidence and risk, and the results do not change significantly. However, when we add the interaction between each of the control variables of confidence and risk, we again see that the interaction between *Female* and *Reward* changes the most when the risk measure is interacted with *Reward*. This again suggests that gender differences in risk preferences underlie the greater gender differences seen among the high-ability participants, explaining 30% of the gender gap. Nevertheless, it is also important to note that, contrary to the main analysis in Table 6, in Table 7, the *Female* and *Reward* interaction remains significant for the high-ability participants when adding risk measures, so some of the differences remain unexplained.

We then look at gender differences in confidence, overconfidence and risk by age, focusing on the most distant age groups: participants in levels 1 and 4. Figure A6 shows the graphs. Interestingly, older participants are less confident, more calibrated (less overconfident) and more risk-loving than the youngest participants.<sup>6</sup> More importantly, gender differences are present both among the younger and older participants. Furthermore, these differences always follow the same pattern: female participants show lower perceived math ability, lower levels of overconfidence and higher risk aversion. However, the gender differences among youngest and oldest participants again are not striking.

We perform a similar exercise as we do in Table A10 but in two sub-samples, the level 1 and level 4 participants. Table A12 shows the results. Columns 1 and 2 reproduce the main results found in the first two columns in Table 4, and columns 3 and 4 replicate the same results for the sample of participants for whom we questionnaire responses. The estimated values of the main variable of interest, the interaction between *Female* and

---

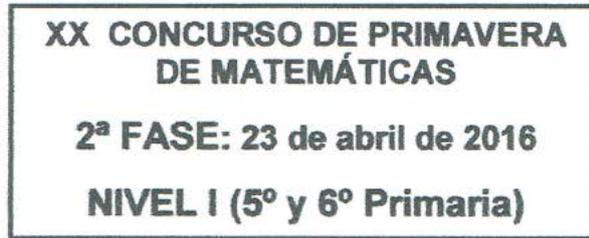
<sup>6</sup> Schildberg-Hörisch (2018) survey the evidence on whether risk preferences remain stable over time concluding they do. Dohmen et al. (2017) on the other hand find that, if anything, individuals become more risk averse as they grow older.

*Reward*, are of very similar magnitude in the overall and the sample for which we have the questionnaire responses, although in the sample with questionnaire responses the value is not significantly different from zero. Columns 5 and 6 add the main control variables of competitiveness, confidence and risk, and the results do not change significantly. However, when we add the interaction between each of the control variables of confidence and risk, we again see that the interaction between *Female* and *Reward* changes the most when the risk measure is interacted with *Reward*. This again suggests that gender differences in risk preferences underlie the greater gender differences among the older participants, explaining up to 50% (but not all) of the gender gap.

We conclude that among all the controls we can include, our risk preferences measure is the one that offers the highest explanatory power in explaining the gender differential. Nevertheless, as we acknowledge that our risk preferences measure might be confounded with measures of answering strategy, we interpret this as suggestive evidence on the underlying mechanism.

## Figures and Tables in the Appendix

Figure A1. Description of Grading System in the Math Test



!!! Lee detenidamente estas instrucciones!!!

Escribe tu nombre y los datos que se te piden en la hoja de respuestas. No pases la página hasta que se te indique.

La prueba tiene una duración de 1 HORA 30 MINUTOS.

No está permitido el uso de calculadoras, reglas graduadas, ni ningún otro instrumento de medida.

Es difícil contestar bien a todas las preguntas en el tiempo indicado. Concéntrate en las que veas más asequibles. Cuando hayas contestado a esas, inténtalo con las restantes.

### PUNTUACIÓN

En los problemas 1 a 13:

Cada respuesta correcta te aportará	5 puntos
Cada pregunta en blanco o errónea	0 puntos

En los problemas 14 a 25:

Cada respuesta correcta te aportará	5 puntos
Cada pregunta que dejes en blanco	1 punto
Cada respuesta errónea	0 puntos

In Questions 1-13:  
Correct answers will score 5 points  
Omitted questions and wrong answers will score 0 point

In Questions 14-25:  
Correct answers will score 5 points  
Omitted questions will score 1 point  
Wrong answers will score 0 points

EN LA HOJA DE RESPUESTAS, MARCA CON UNA ASPA  LA QUE CONSIDERES CORRECTA.

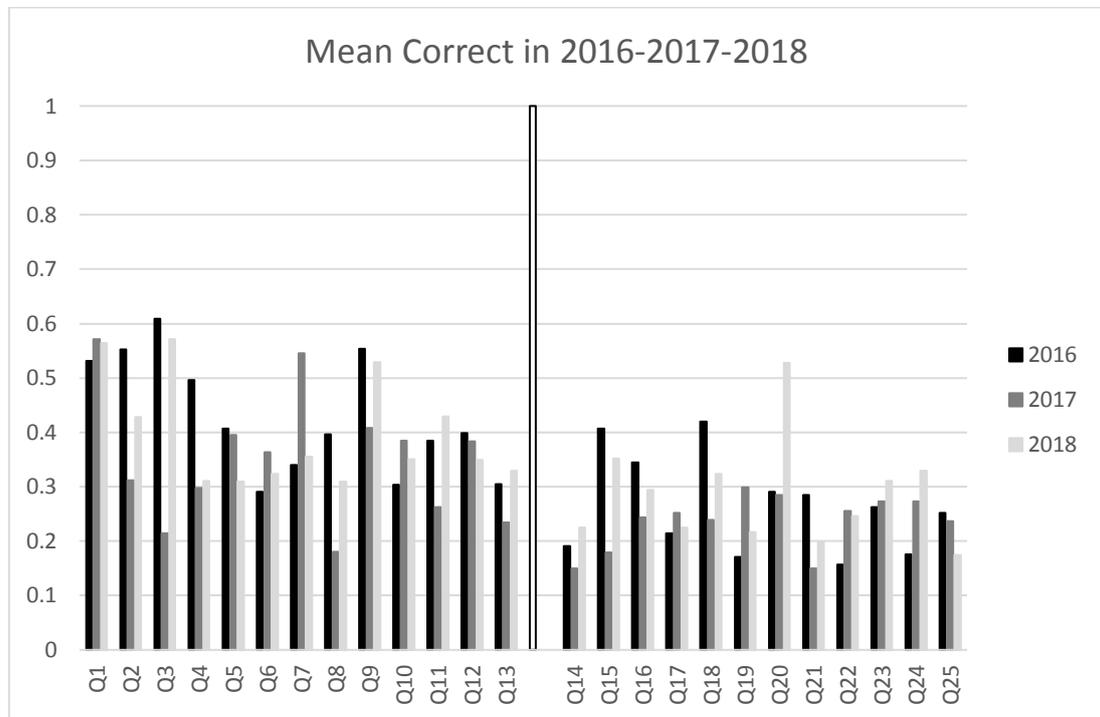
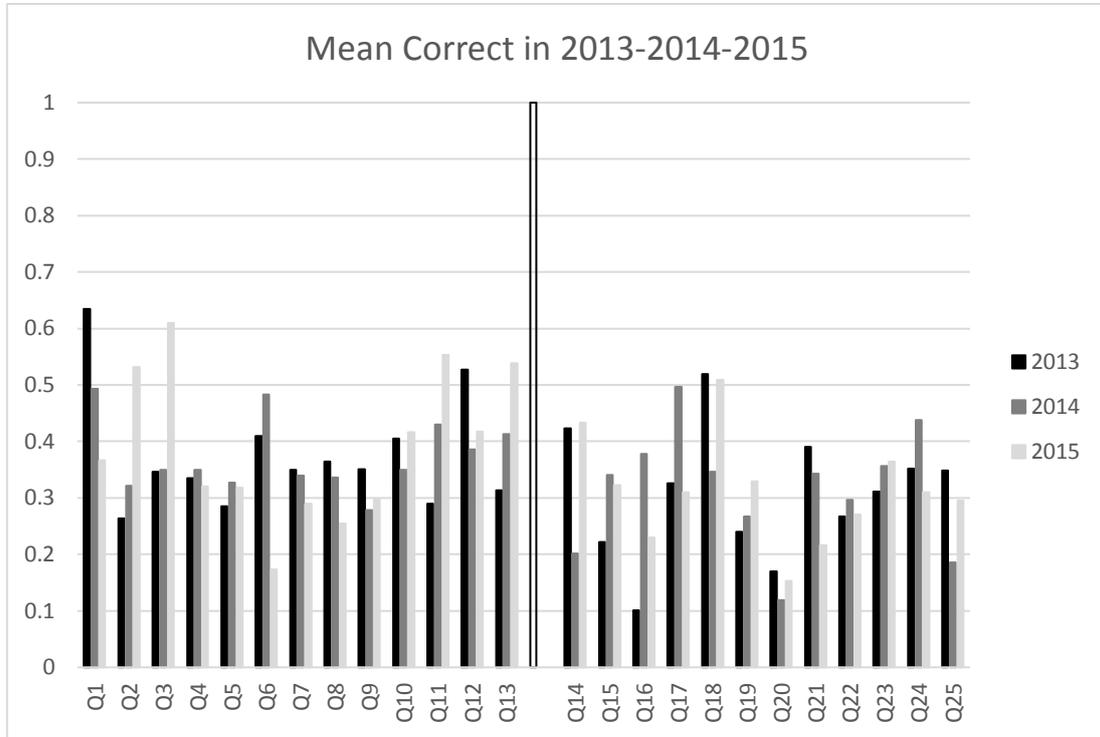
SI TE EQUIVOCAS, ESCRIBE "NO" EN LA EQUIVOCADA Y MARCA LA QUE CREAS CORRECTA.

CONVOCA  
Facultad de Matemáticas de la UCM

ORGANIZA  
Asociación Matemática  
Concurso de Primavera

COLABORAN  
Universidad Complutense de Madrid  
Consejería de Educación de la Comunidad de Madrid  
El Corte Inglés  
Grupo ANAYA  
Grupo SM  
Smartick

**Figure A2. Mean Values of Correct Per Question: First Part (Questions 1-13) and Second Part (Questions 14-25) for years 2013-2015 and 2016-2018**

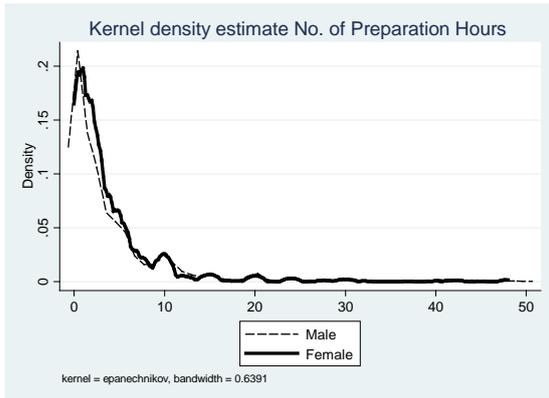


### Figure A3. Questionnaire at the end of the Test

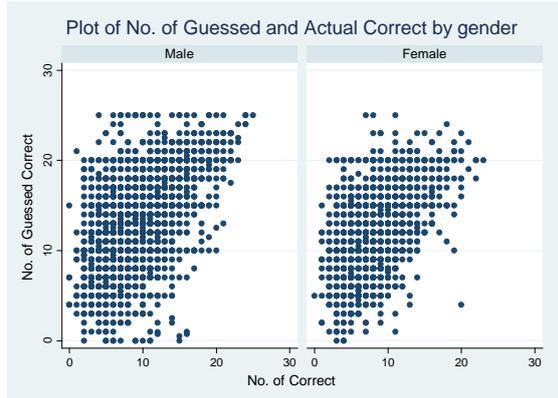
For the following statements please, say your agreement level (1 referring to Strongly Disagree and 5 to Strongly Agree):

1. "It is more important for me to be selected for Stage 2 than being among the winners in Stage 2."
2. "It is more important to my parents being selected for Stage 2 than being among the winners in Stage 2."
3. "It is more important to do well in Stage 2 than in Stage 1."
4. "I have devoted more hours to prepare Stage 2 test than Stage 1 test."
6. "While doing the test I felt more pressure during Stage 2 than in Stage 1"
9. "I am good at Mathematics"
  
5. How many hours did you devote to prepare Stage 2 test?
7. How many questions do you expect to get right?
8. When omitting a question was worth 1 point I answered the question \_\_\_\_\_
  - a. when I was absolutely sure.
  - b. when I was almost sure.
  - c. when I was uncertain between 2 answers.
  - d. when I was uncertain between 3 answers.
  - e. always.
  
10. I believe \_\_\_\_\_ at Math
  - a. men are better than women
  - b. men and women are equally good
  - c. women are better than men

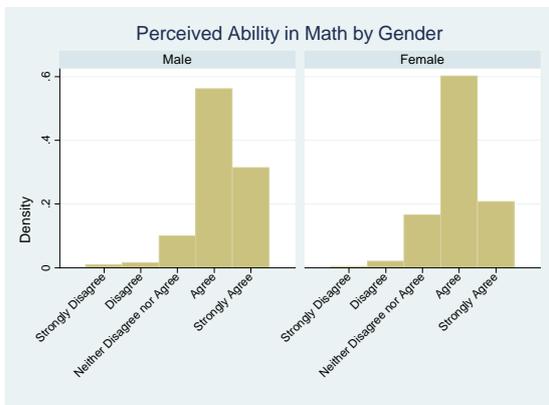
**Figure A4. Descriptive Statistics on the Control Variables from the Questionnaire**



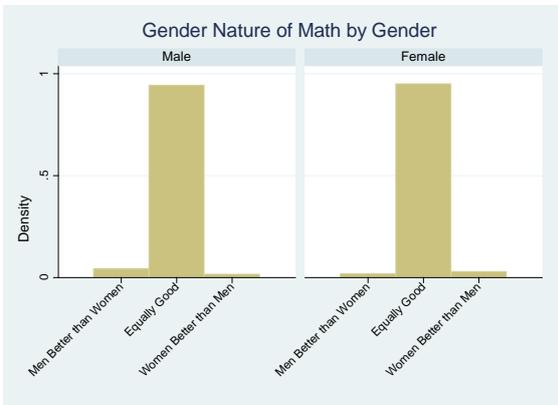
**A4a**



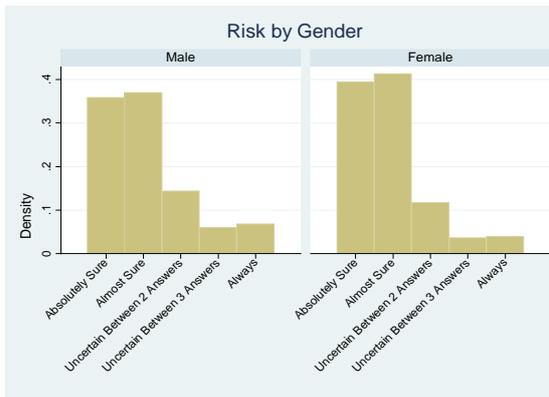
**A4b**



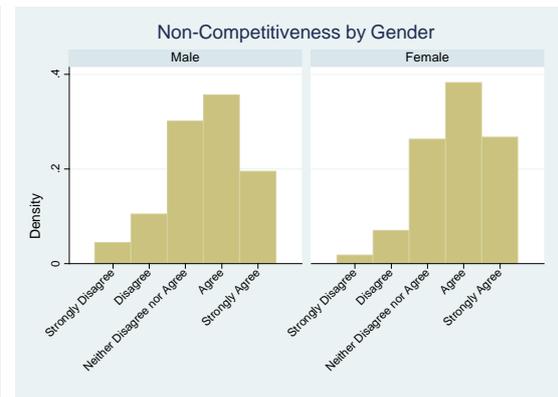
**A4c**



**A4d**

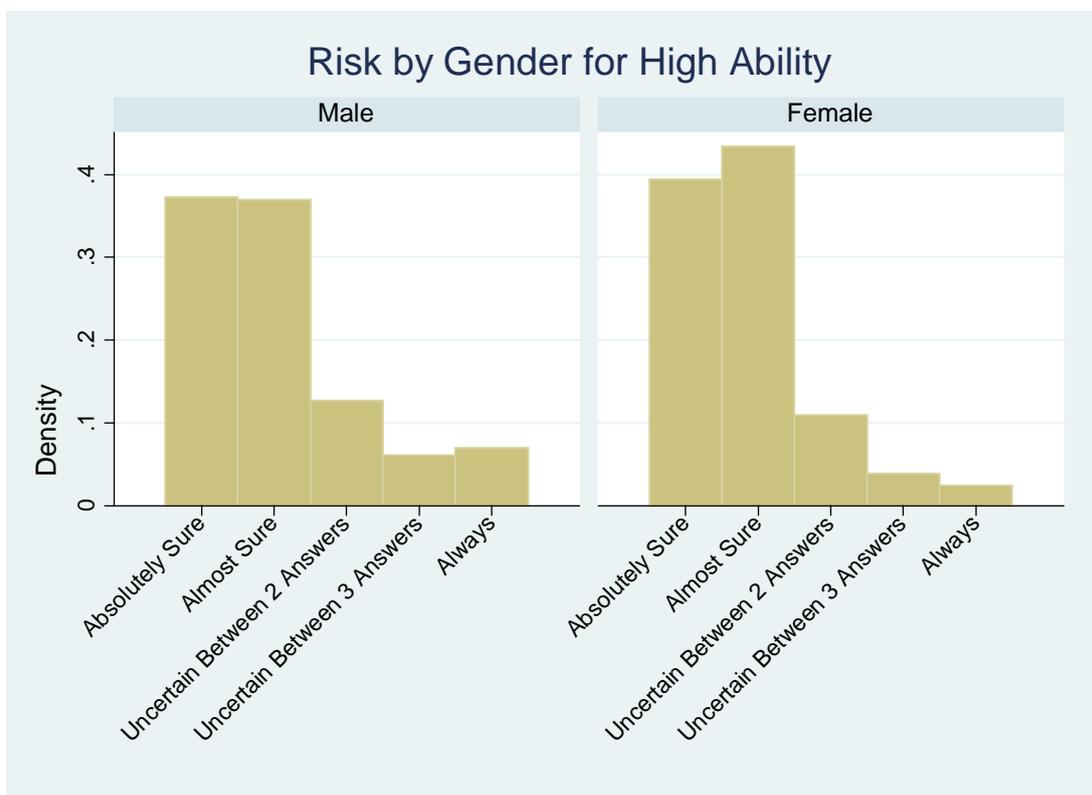
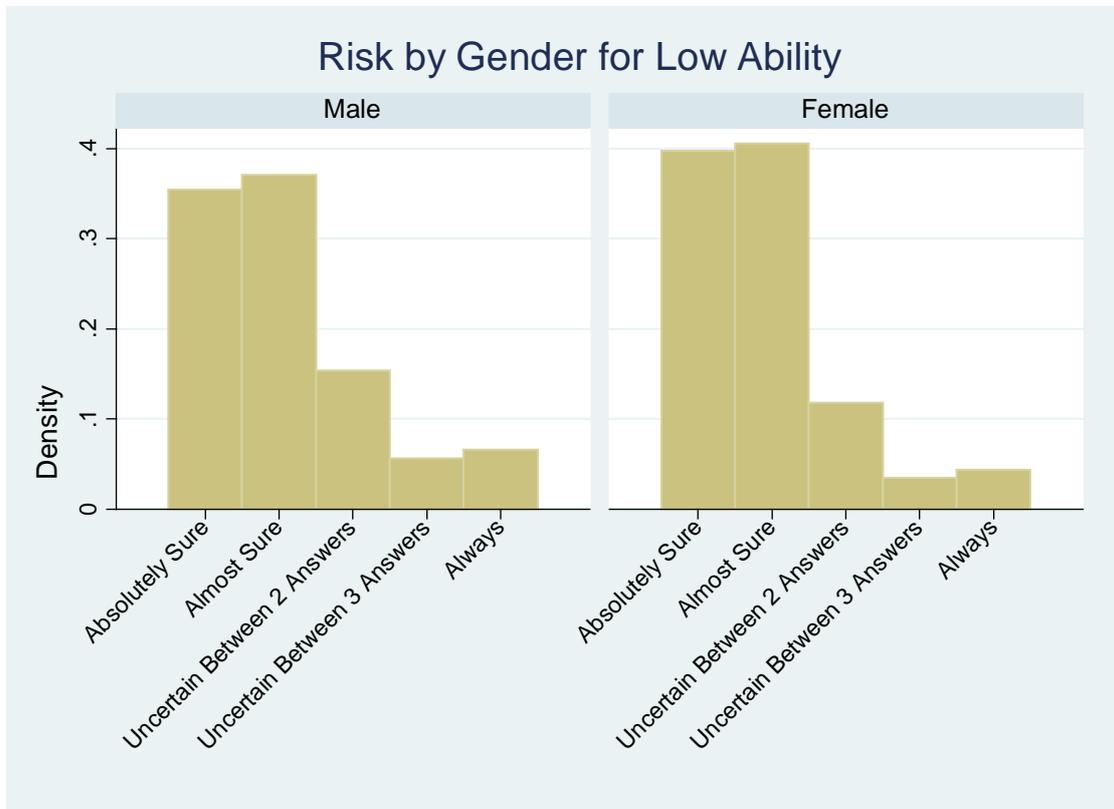


**A4e**

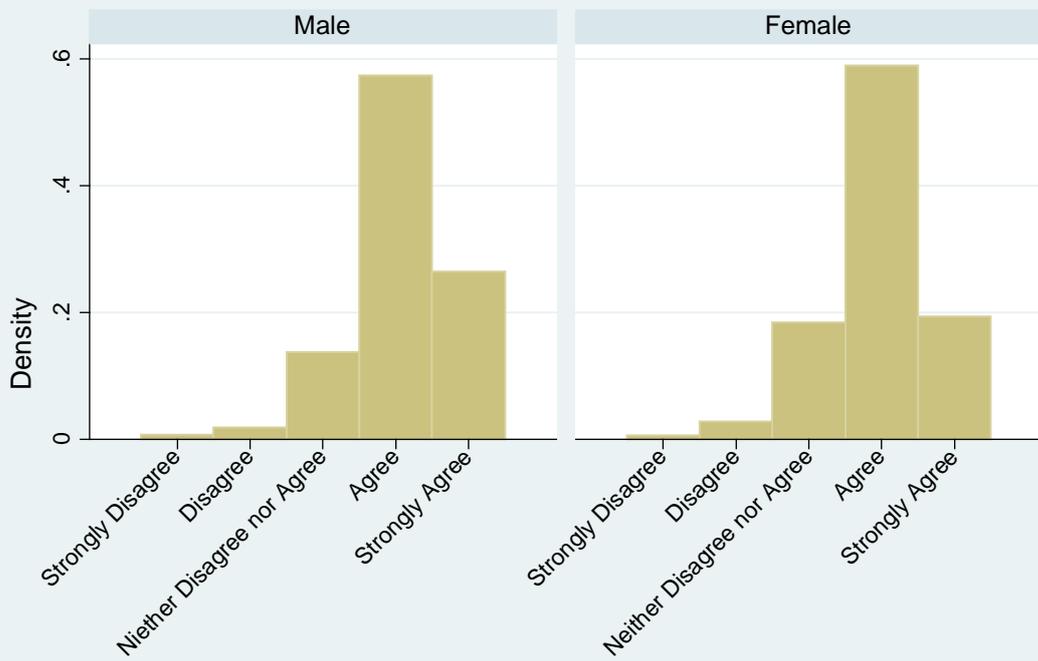


**A4f**

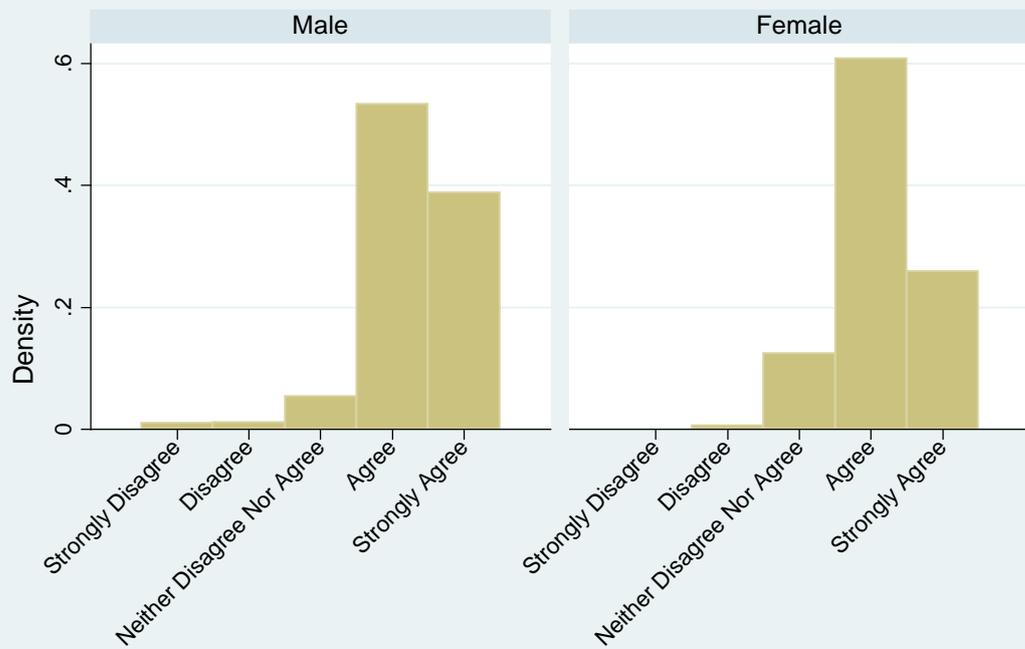
**Figure A5. Risk, Confidence and Overconfidence by Gender: Low Ability:  
Standardized No. of Correct in No Reward<0 and High Ability:  
Standardized No. of Correct in No Reward>0**



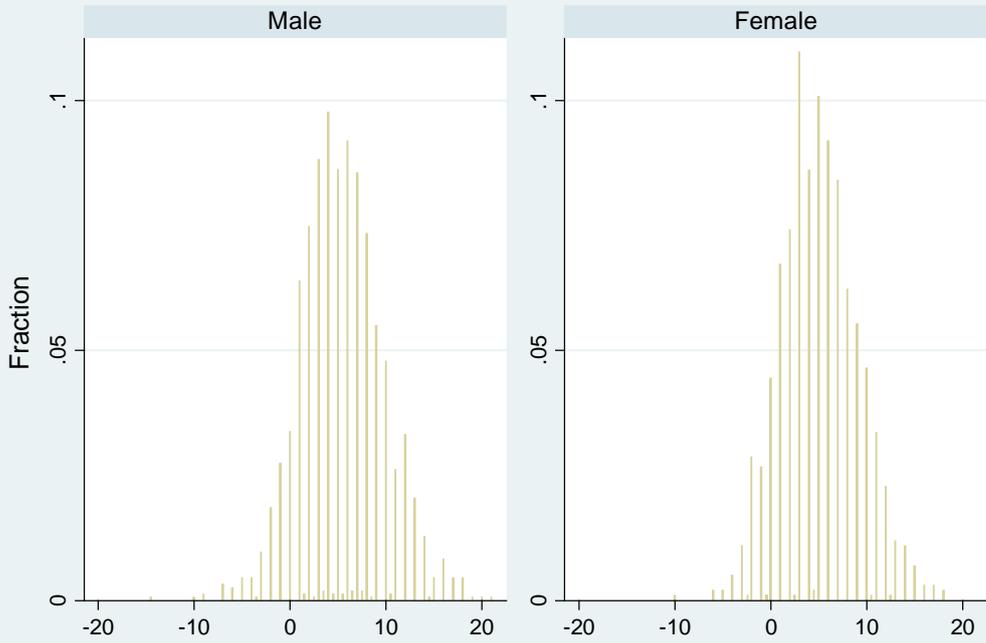
### Perceived Ability Math by Gender for Low Ability



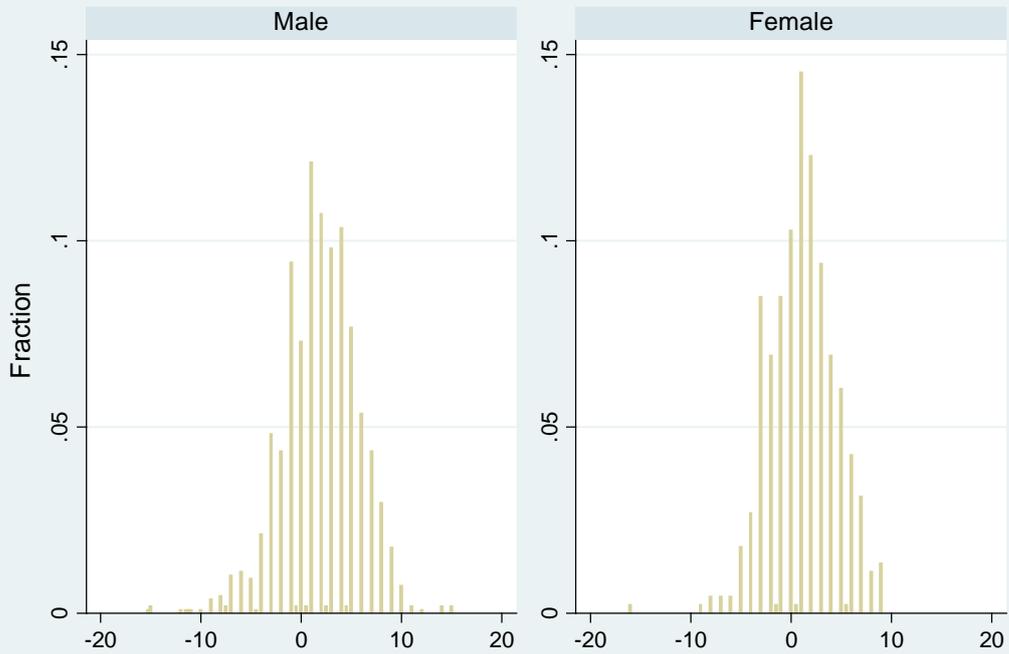
### Perceived Ability Math by Gender for High Ability



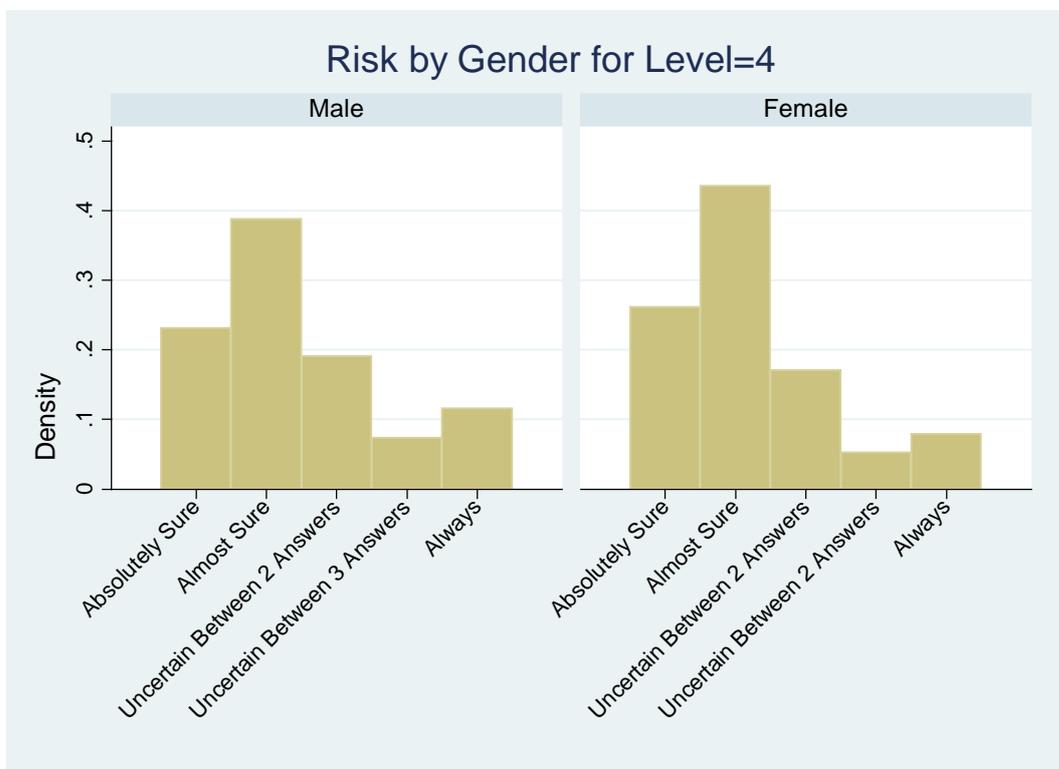
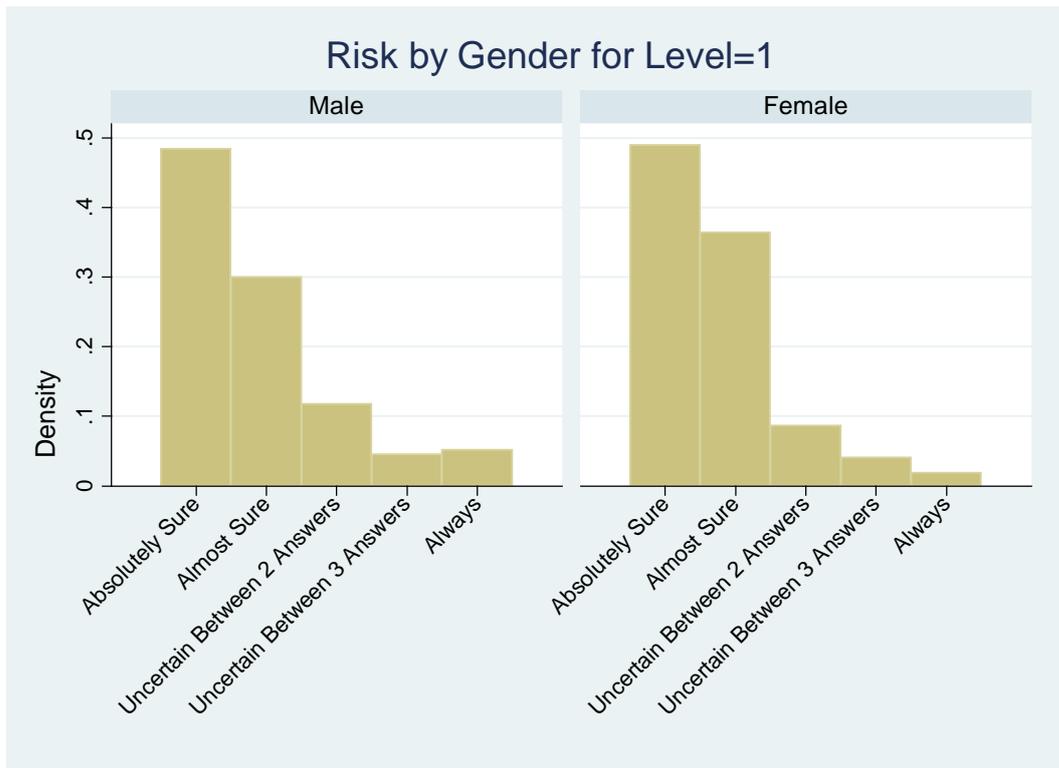
### Histogram Overconfidence by Gender Low Ability



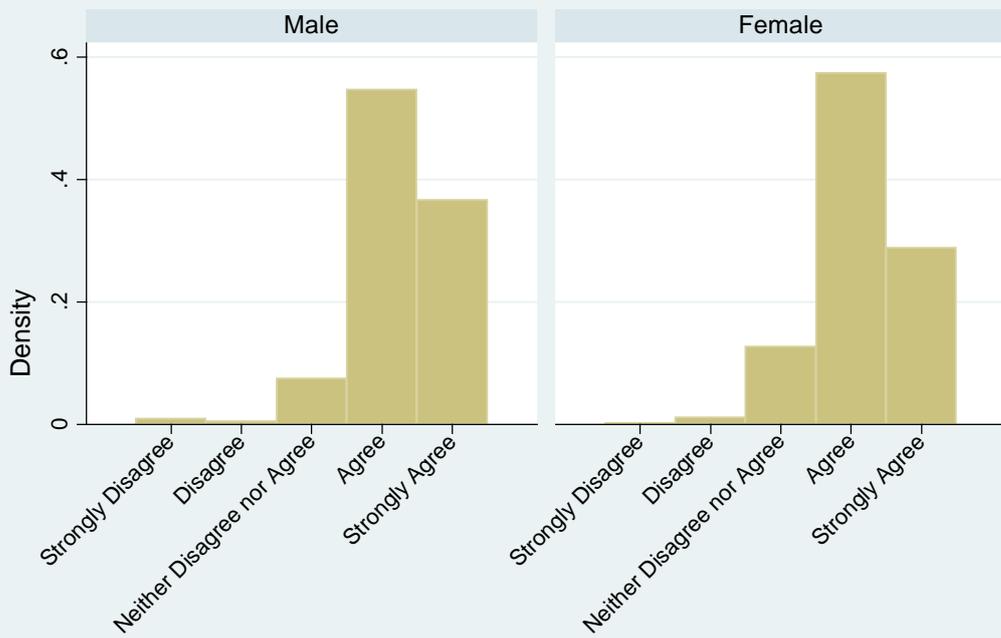
### Histogram Overconfidence by Gender High Ability



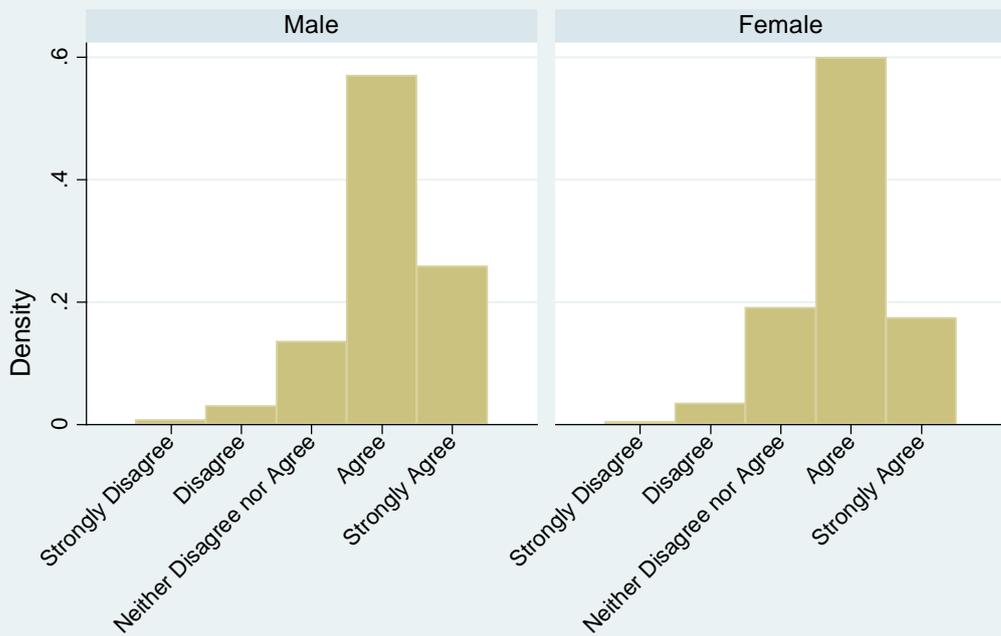
**Figure A6. Risk, Confidence and Overconfidence by Gender: Level 1 and L4**



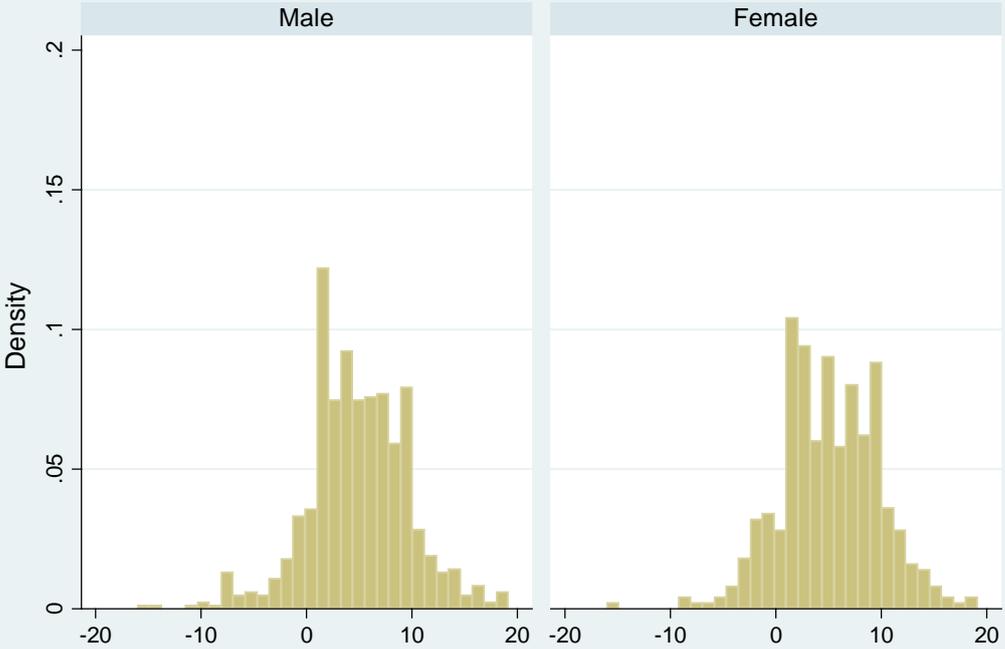
### Perceived Ability Math by Gender for Level=1



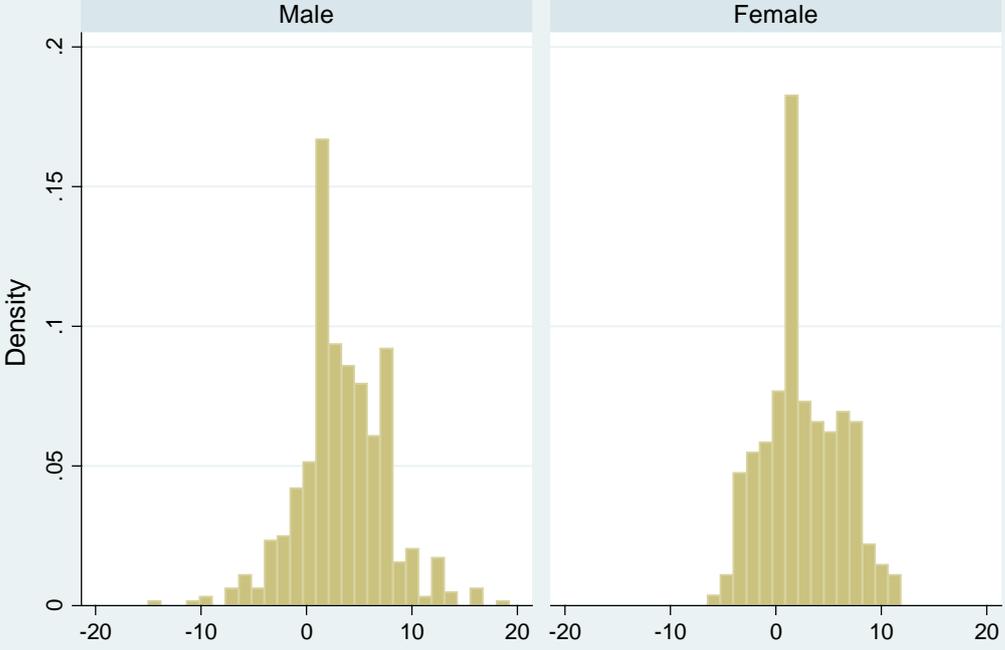
### Perceived Ability Math by Gender for Level==4



Histogram Overconfidence by Gender for Level=1



Histogram Overconfidence by Gender for Level=4



**Table A1. Gender Differences between No Reward and the Reward Parts of the Test with Alternative Control for Ability: No. Of Correct No Reward**

	OLS				RE				FE			
	zomitted (1)	zprop_correct (2)	zscore (3)	rank (4)	zomitted (5)	zprop_correct (6)	zscore (7)	rank (8)	zomitted (9)	zprop_correct (10)	zscore (11)	rank (12)
Female	0.153*** (0.0234)	-0.0111 (0.0120)	-0.0539*** (0.00808)	-11.09*** (2.610)	0.153*** (0.0237)	-0.0111 (0.0120)	-0.0539*** (0.00808)	-11.09*** (2.610)				
Reward	-0.0447*** (0.0154)	-0.00382 (0.0134)	0.0148 (0.0125)	40.40*** (3.302)	-0.0447*** (0.0154)	-0.00382 (0.0134)	0.0148 (0.0125)	40.40*** (3.302)	-0.0447*** (0.0152)	-0.00382 (0.0132)	0.0148 (0.0123)	40.40*** (3.254)
Female*Reward	0.132*** (0.0282)	0.0140 (0.0233)	-0.0408** (0.0207)	-10.30* (5.660)	0.132*** (0.0282)	0.0140 (0.0233)	-0.0408** (0.0207)	-10.30* (5.660)	0.132*** (0.0278)	0.0140 (0.0230)	-0.0408** (0.0204)	-10.30* (5.579)
No. Of Correct No Reward	-0.0705*** (0.00401)	0.258*** (0.00245)	0.284*** (0.00231)	70.80*** (0.616)	-0.0680*** (0.00401)	0.258*** (0.00245)	0.284*** (0.00231)	70.80*** (0.616)	-0.0316*** (0.0116)	0.199*** (0.00727)	0.212*** (0.00664)	51.42*** (1.839)
Participation Time	-0.0677*** (0.0182)	0.0794*** (0.0133)	0.128*** (0.0127)	33.22*** (3.465)	-0.0693*** (0.0178)	0.0793*** (0.0133)	0.128*** (0.0127)	33.22*** (3.465)	-0.224** (0.0882)	-0.0615 (0.0547)	0.00445 (0.0490)	-7.554 (14.04)
Observations	19,768	19,768	19,768	19,768	19,768	19,768	19,768	19,768	19,768	19,768	19,768	19,768
R-squared	0.112	0.473	0.591	0.583					0.016	0.070	0.100	0.150
Number of participants					8,537	8,537	8,537	8,537	8,537	8,537	8,537	8,537

Notes: Observations are at the Math test's parts level. The first three outcome variables, *No. Omitted*, *Prop. of Correct* and *Score* are standardized at the edition, level and part of the test levels. *Rank* measures the position in the rank by edition, level and test's parts level, where higher values represent better positions within the rank. *Female* takes the value of 1 if the participant is female. *Reward* takes the value of 1 if the outcome variable refers to the part of the test with reward for omitted questions. *No. of Correct No Reward* measures the number of correct in the part of the test with the reward and *Participation Time* takes the values of 1, 2, 3 if it is the first, second or third time that the participant does the Math test. Columns 1-4 show the OLS specification where the standard errors are clustered at the participant level. Columns 5-8 show the RE model specification and columns 9-12 show the FE specification model. All specifications include edition, level and school fixed effects. Standard errors in parenthesis, where \*\*\* p<0.01, \*\* p<0.05, \* p<0.1

**Table A2. Gender Differences between the First and the Second Parts of the Test, Placebo test with 2013, 2014 and 2015 Editions**

	Placebo 2013 Edition				Placebo 2014 Edition				Placebo 2015 Edition			
	zomitted (1)	zprop_correct (2)	zscore (3)	rank (4)	zomitted (5)	zprop_correct (6)	zscore (7)	rank (8)	zomitted (9)	zprop_correct (10)	zscore (11)	rank (12)
Female	0.150*** (0.0329)	-0.0259 (0.0330)	-0.160*** (0.0313)	-45.79*** (10.48)	0.265*** (0.0500)	-0.151*** (0.0479)	-0.264*** (0.0439)	-28.76*** (6.730)	0.281*** (0.0379)	-0.137*** (0.0370)	-0.255*** (0.0348)	-63.24*** (8.923)
Reward	-0.000267 (0.0124)	-0.0101 (0.0158)	-0.0146 (0.0175)	-8.366 (5.606)	-0.0151 (0.0226)	-0.00447 (0.0279)	-0.00825 (0.0269)	-1.786 (3.448)	0.0289 (0.0182)	-0.0174 (0.0223)	-0.0233 (0.0220)	-8.116 (5.523)
Female*Reward	0.00302 (0.0218)	0.0302 (0.0286)	0.0418 (0.0297)	13.27 (9.683)	0.0460 (0.0403)	0.0137 (0.0502)	0.0252 (0.0439)	1.384 (5.789)	-0.0797*** (0.0308)	0.0481 (0.0395)	0.0642* (0.0359)	17.83* (9.339)
Participation Time	0.565* (0.308)	-0.379** (0.180)	-0.767*** (0.254)	-224.3** (89.23)	-0.0693 (0.0594)	0.431*** (0.0545)	0.469*** (0.0571)	40.12*** (7.654)	-0.135*** (0.0380)	0.257*** (0.0350)	0.348*** (0.0382)	78.96*** (8.491)
Observations	7,794	7,794	7,794	7,794	4,170	4,170	4,170	5,584	6,258	6,258	6,258	6,258
R-squared	0.219	0.224	0.259	0.347	0.257	0.277	0.335	0.498	0.220	0.241	0.297	0.367

Notes: Observations are at the Math test's parts level. The first three outcome variables, *No. Omitted*, *Prop. of Correct* and *Score* are standardized at the level and part of the test levels. Rank measures the position in the rank by level and test's parts level, where higher values represent better positions within the rank. *Female* takes the value of 1 if the participant is female. *Second Half* takes the value of 1 if the outcome variable refers to the questions 14-25. Columns 1-4 show the results for the edition 2013, columns 5-8 for the edition 2014 and columns 9-12 for the edition 2015. In editions 2013, 2014 and 2015 there was differential score for omitted questions and wrong answers for all questions in the test. All regressions include level and school fixed effects. Standard errors, clustered at the participant level, are shown in parenthesis, where \*\*\* p<0.01, \*\* p<0.05, \* p<0.1

**Table A3. Gender Differences between No Reward and the Reward Parts of the Test with Triple Interaction: Editions 2013 to 2018**

	zomitted (1)	zprop_correct (2)	zscore (3)	rank (4)	zomitted (5)	zprop_correct (6)	zscore (7)	rank (8)	zomitted (9)	zprop_correct (10)	zscore (11)	rank (12)
Female	0.238*** (0.0215)	-0.0942*** (0.0211)	-0.216*** (0.0196)	-33.78*** (5.356)	0.256*** (0.0249)	-0.158*** (0.0236)	-0.266*** (0.0222)	-37.23*** (5.532)	0.157*** (0.0191)	0.0526*** (0.0150)	-0.0509*** (0.0101)	2.686 (3.216)
Treatment (Years 2016-2017-2018)	0.0498** (0.0228)	-0.0747*** (0.0220)	-0.107*** (0.0218)	-138.3*** (6.326)	0.226*** (0.0309)	-0.141*** (0.0281)	-0.164*** (0.0295)	-166.6*** (8.756)	0.184*** (0.0206)	-0.318*** (0.0158)	-0.381*** (0.0132)	-216.0*** (4.133)
Second Half	-0.0205** (0.00916)	-0.00737 (0.00888)	9.36e-05 (0.00864)	17.20*** (2.206)	-0.0240** (0.0106)	-0.00754 (0.0102)	0.00449 (0.00977)	20.68*** (2.402)	-0.0205** (0.00916)	-0.00737 (0.00888)	9.36e-05 (0.00864)	17.82*** (2.286)
Female*Treatment	-0.0517* (0.0308)	-0.0246 (0.0280)	0.0448* (0.0266)	-18.67*** (7.171)	-0.0901*** (0.0335)	-0.00754 (0.0297)	0.0546* (0.0284)	-23.12*** (7.316)	-0.0221 (0.0288)	-0.0784*** (0.0190)	-0.0159 (0.0135)	-26.57*** (4.146)
Female*Second Half	0.0102 (0.0162)	0.0289 (0.0194)	0.0298* (0.0180)	-12.17*** (4.639)	0.0111 (0.0200)	0.0269 (0.0241)	0.0203 (0.0218)	-16.98*** (5.143)	0.0102 (0.0162)	0.0289 (0.0194)	0.0298* (0.0180)	-12.37** (4.973)
Female*Second Half*Treatment	0.0974*** (0.0270)	-0.0114 (0.0255)	-0.0558** (0.0227)	25.07*** (6.099)	0.113*** (0.0297)	-0.00768 (0.0294)	-0.0558** (0.0259)	25.15*** (6.571)	0.0974*** (0.0270)	-0.0114 (0.0255)	-0.0558** (0.0227)	24.65*** (6.328)
Math at School					0.0298*** (0.00689)	0.239*** (0.00627)	0.233*** (0.00633)	51.86*** (1.567)				
No. Of Correct No Reward									-0.126*** (0.00269)	0.230*** (0.00186)	0.259*** (0.00161)	68.02*** (0.439)
Participation Time	-0.0969*** (0.0120)	0.281*** (0.0123)	0.343*** (0.0144)	74.96*** (3.039)	-0.106*** (0.0126)	0.221*** (0.0119)	0.283*** (0.0135)	62.80*** (3.063)	0.0116 (0.0111)	0.0837*** (0.00824)	0.121*** (0.00798)	18.80*** (1.982)
Observations	37,988	37,988	37,988	39,402	30,580	30,580	30,580	31,964	37,988	37,988	37,988	37,988
R-squared	0.098	0.163	0.222	0.297	0.102	0.228	0.284	0.345	0.183	0.443	0.578	0.600

Notes: Observations are at the Math test's parts level. The first three outcome variables, *No. Omitted*, *Prop. of Correct* and *Score* are standardized at the edition, level and part of the test levels. *Rank* measures the position in the rank by edition, level and test's parts level, where higher values represent better positions within the rank. *Female* takes the value of 1 if the participant is female. *Treatment* takes the value of 1 if the outcome variable refers to the editions of 2016, 2017 and 2018. *Second-Half* takes the value of 1 if the outcome variables refer to the second half of the test. *Math at School* measures the standardized Math grade at school level and *Participation Time* takes the values of 1, 2, 3 if it is the first, second or third time that the participant does the Math test. *No. of Correct No Reward* measures the number of correct in the part of the test with the reward and *Participation Time* takes the values of 1, 2, 3 if it is the first, second or third time that the participant does the Math test. All columns show the OLS specification where the standard errors are clustered at the participant level. Standard errors in parenthesis, where \*\*\* p<0.01, \*\* p<0.05, \* p<0.1.

**Table A4. Gender Differences between the No Reward and the Reward Parts of the Test for those Participants who Leave 0 Omitted Questions in the No Reward Part of the Test**

	OLS				RE				FE			
	zomitted (1)	zprop_correct (2)	zscore (3)	rank (4)	zomitted (5)	zprop_correct (6)	zscore (7)	rank (8)	zomitted (9)	zprop_correct (10)	zscore (11)	rank (12)
Female	-0.0131* (0.00763)	-0.196*** (0.0226)	-0.195*** (0.0230)	-47.78*** (6.148)	-0.0133* (0.00762)	-0.195*** (0.0226)	-0.190*** (0.0230)	-47.29*** (6.192)				
Reward	0.185*** (0.0158)	0.0284* (0.0152)	-0.0203 (0.0149)	29.32*** (3.897)	0.185*** (0.0158)	0.0284* (0.0152)	-0.0203 (0.0149)	29.32*** (3.897)	0.185*** (0.0155)	0.0284* (0.0149)	-0.0203 (0.0146)	29.32*** (3.817)
Female*Reward	0.332*** (0.0265)	0.0459* (0.0270)	-0.0798*** (0.0252)	-16.22** (6.851)	0.332*** (0.0265)	0.0459* (0.0270)	-0.0798*** (0.0252)	-16.22** (6.851)	0.332*** (0.0260)	0.0459* (0.0264)	-0.0798*** (0.0247)	-16.22** (6.709)
Math at School	0.0224*** (0.00634)	0.239*** (0.00890)	0.246*** (0.00949)	59.47*** (2.425)	0.0225*** (0.00633)	0.227*** (0.00887)	0.229*** (0.00938)	55.71*** (2.418)	-0.0380** (0.0188)	0.0670** (0.0277)	0.0739*** (0.0281)	14.09* (7.423)
Participation Time	-0.0848*** (0.0147)	0.306*** (0.0214)	0.370*** (0.0236)	93.14*** (5.794)	-0.0838*** (0.0146)	0.228*** (0.0209)	0.268*** (0.0224)	67.99*** (5.683)	-0.0668 (0.0524)	0.0294 (0.0884)	0.0570 (0.0906)	10.24 (24.83)
Observations	13,948	13,948	13,948	13,948	13,948	13,948	13,948	13,948	13,948	13,948	13,948	13,948
R-squared	0.139	0.272	0.310	0.354					0.104	0.028	0.040	0.095
Number of participants					6,144	6,144	6,144	6,144	6,144	6,144	6,144	6,144

Notes: Observations are at the Math test's parts level. The first three outcome variables, *No. Omitted*, *Prop. of Correct*, and *Score* are standardized at the edition, level and part of the test levels. *Rank* measures the position in the rank by edition, level and test's parts level, where higher values represent better positions within the rank. *Female* takes the value of 1 if the participant is female. *Reward* takes the value of 1 if the outcome variable refers to the part of the test with reward for omitted questions. *Math at School* measures the standardized Math grade at school level and *Participation Time* takes the values of 1, 2, 3 if it is the first, second or third time that the participant does the Math test. Columns 1-4 show the OLS specification where the standard errors are clustered at the participant level. Columns 5-8 show the RE model specification and columns 9-12 show the FE specification model. All specifications include edition, level and school fixed effects. Standard errors in parenthesis, where \*\*\* p<0.01, \*\* p<0.05, \* p<0.1

**Table A5. Gender Differences between the No Reward and the Reward Parts of the Test with 12 Questions in Each Part of the Test**

		zomitted (1)	zprop_correct (2)	zscore (3)	rank (4)
Main specification (Table 2) Q1-Q13 in the first part	Female*Reward	0.145*** (0.0295)	0.0194 (0.0245)	-0.0449** (0.0218)	-10.74* (5.941)
Q1 out	Female*Reward	0.144*** (0.0295)	0.0152 (0.0247)	-0.0512** (0.0219)	-12.17** (5.980)
Q2 out	Female*Reward	0.150*** (0.0295)	0.0124 (0.0248)	-0.0522** (0.0221)	-12.55** (5.998)
Q3 out	Female*Reward	0.142*** (0.0295)	0.00303 (0.0247)	-0.0556** (0.0222)	-13.10** (6.004)
Q4 out	Female*Reward	0.151*** (0.0295)	0.0165 (0.0247)	-0.0515** (0.0221)	-11.95** (6.001)
Q5 out	Female*Reward	0.144*** (0.0295)	0.0182 (0.0247)	-0.0450** (0.0221)	-10.97* (6.008)
Q6 out	Female*Reward	0.151*** (0.0294)	0.0235 (0.0246)	-0.0438** (0.0218)	-10.26* (5.929)
Q7 out	Female*Reward	0.144*** (0.0296)	0.0206 (0.0247)	-0.0458** (0.0220)	-11.87** (5.967)
Q8 out	Female*Reward	0.146*** (0.0295)	0.0132 (0.0246)	-0.0487** (0.0220)	-12.01** (5.976)
Q9 out	Female*Reward	0.147*** (0.0295)	0.0150 (0.0248)	-0.0521** (0.0221)	-12.31** (6.014)
Q10 out	Female*Reward	0.147*** (0.0296)	0.0188 (0.0248)	-0.0460** (0.0220)	-11.10* (5.987)
Q11 out	Female*Reward	0.145*** (0.0296)	0.0161 (0.0247)	-0.0499** (0.0222)	-13.11** (6.024)
Q12 out	Female*Reward	0.149*** (0.0295)	0.00860 (0.0247)	-0.0539** (0.0220)	-13.03** (5.975)
Q13 out	Female*Reward	0.148*** (0.0295)	0.0229 (0.0246)	-0.0457** (0.0220)	-10.82* (5.977)

*Notes:* Observations are at the Math test's parts level. The first three outcome variables, *No. Omitted*, *Prop. of Correct*, and *Score* are standardized at the edition, level and part of the test levels. *Rank* measures the position in the rank by edition, level and test's parts level, where higher values represent better positions within the rank. All control variables as the same as the ones in Table 2. The first row copies the *Female\*Reward* coefficient from the main Table 2. The additional rows show the same coefficient taking each time one of the questions out from the first part, such that there are always 12 questions from the first part and 12 questions from the second part of the test. All estimations show the OLS specification where the standard errors are clustered at the participant level. Standard errors in parenthesis, where \*\*\* p<0.01, \*\* p<0.05, \* p<0.1

**Table A6. Gender Differences between the No Reward and the Reward Parts of the Test:  
Splitting the Reward part into two (Q14-Q19 and Q20-Q25)**

	zomitted (1)	zprop_correct (2)	zscore (3)
Female	0.170*** (0.0251)	-0.171*** (0.0211)	-0.220*** (0.0204)
Reward_Q14-Q19	-0.0418*** (0.0162)	-0.0122 (0.0147)	0.00112 (0.0138)
Reward_Q20-Q25	0.00849 (0.0178)	-0.0160 (0.0156)	-0.0111 (0.0147)
Female*Reward_Q14-Q19	0.132*** (0.0298)	0.0313 (0.0254)	-0.00999 (0.0231)
Female*Reward_Q20-Q25	0.0740** (0.0329)	0.0415 (0.0270)	0.0277 (0.0244)
Math at School	0.0378*** (0.00831)	0.198*** (0.00718)	0.201*** (0.00754)
Participation Time	-0.135*** (0.0194)	0.263*** (0.0178)	0.336*** (0.0196)
Observations	26,775	26,775	26,775
R-squared	0.099	0.172	0.224

*Notes:* Observations are at the Math test's parts level. The first three outcome variables, *No. Omitted*, *Prop. of Correct*, and *Score* are standardized at the edition, level and part of the test levels. *Female* takes the value of 1 if the participant is female. *Reward\_Q14-Q19* takes the value of 1 if the outcome variable refers to the part of the test with reward for omitted questions Q14 to Q19. *Reward\_Q20-Q25* takes the value of 1 if the outcome variable refers to the part of the test with reward for omitted questions Q20 to Q25. *Math at School* measures the standardized Math grade at school level and *Participation Time* takes the values of 1, 2, 3 if it is the first, second or third time that the participant does the Math test. Columns 1-3 show the OLS specification where the standard errors are clustered at the participant level. All specifications include edition, level and school fixed effects. Standard errors in parenthesis, where \*\*\* p<0.01, \*\* p<0.05, \* p<0.1

**Table A7. Gender Differences between the No Reward and the Reward Parts of the Test in the Matched Sample based on the No. Of Correct in the No Reward Part**

	zcorrect (1)	zcorrect (2)	zomitted (3)	zprop_correct (4)	zscore (5)	rank (6)
Female	-0.215*** (0.0199)	-0.0193 (0.0230)	0.107*** (0.0289)	0.00406 (0.0237)	-0.0216 (0.0229)	-6.003 (6.158)
Reward	0.0319** (0.0134)	0.160*** (0.0185)	-0.0680*** (0.0230)	0.120*** (0.0196)	0.136*** (0.0182)	72.52*** (4.976)
Female*Reward	-0.0980*** (0.0223)	-0.226*** (0.0258)	0.167*** (0.0334)	-0.108*** (0.0282)	-0.166*** (0.0253)	-43.61*** (6.943)
Math at School	0.208*** (0.00813)	0.171*** (0.00977)	0.0635*** (0.0110)	0.217*** (0.0101)	0.195*** (0.00986)	49.73*** (2.622)
Participation Time	0.379*** (0.0212)	0.355*** (0.0243)	-0.142*** (0.0247)	0.311*** (0.0244)	0.361*** (0.0243)	98.46*** (6.317)
Observations	17,850	12,174	12,174	12,174	12,174	12,174
R-squared	0.279	0.266	0.108	0.233	0.277	0.332

*Notes:* Observations are at the Math test's parts level. The outcome variables, *No. of Correct*, *No. Omitted*, *Prop. of Correct*, and *Score* are standardized at the edition, level and part of the test levels. *Rank* measures the position in the rank by edition, level and test's parts level, where higher values represent better positions within the rank. *Female* takes the value of 1 if the participant is female. *Reward* takes the value of 1 if the outcome variable refers to the part of the test with reward for omitted questions. *Math at School* measures the standardized Math grade at school level and *Participation Time* takes the values of 1, 2, 3 if it is the first, second or third time that the participant does the Math test. Column 1 shows the estimation for the whole sample and columns 2 to 6 show the estimation results for the matched sample using the *No. of Correct* in the no reward part of the test with 3386 male and 3386 female participants. All regressions show the OLS specification where the standard errors are clustered at the participant level. All specifications include edition, level and school fixed effects. Standard errors in parenthesis, where \*\*\* p<0.01, \*\* p<0.05, \* p<0.1

**Table A8. Gender Differences between the No Reward and the Reward Parts of the Test along the Ability Distribution with Alternative Measure of Ability: Standardized Math at School Level**

	Low Ability omitted (1)	High Ability omitted (2)	Interaction omitted (3)	Continuous omitted (4)
Female	0.160*** (0.0386)	0.143*** (0.0328)	0.185*** (0.0377)	0.173*** (0.0249)
Reward	-0.0433* (0.0246)	-0.0485** (0.0219)	-0.0433* (0.0239)	-0.0460*** (0.0162)
Female*Reward	0.0835* (0.0476)	0.183*** (0.0384)	0.0835* (0.0463)	0.139*** (0.0297)
High Ability			-0.0109 (0.0352)	
High Ability*Reward			-0.00525 (0.0319)	
Female*High Ability			-0.0239 (0.0488)	
Female*Reward*High Ability			0.0999* (0.0595)	
Math at School	-0.0749 (0.0495)	-0.0238 (0.0278)	0.0401*** (0.0154)	0.0296** (0.0122)
Participation Time	-0.0956*** (0.0368)	-0.182*** (0.0239)	-0.156*** (0.0195)	-0.156*** (0.0195)
Math at School*Reward				0.0122 (0.0168)
Female*Math at School				-0.0133 (0.0255)
Female*Reward*Math at School				0.0417 (0.0320)
Observations	7,404	10,446	17,850	17,850
R-squared	0.156	0.134	0.096	0.096

*Notes:* Observations are at the math test's parts level. *No. Omitted* is standardized at the edition, level and part of the test levels. *Female* takes the value of 1 if the participant is female. *Reward* takes the value of 1 if the outcome variable refers to the part of the test with reward for omitted questions. *Math at School* measures the standardized Math grade at school level and *Participation Time* takes the values of 1, 2, 3 if it is the first, second or third time that the participant does the Math test. *High Ability* takes value 1 if the participant's standardized Math grade is >0. All columns show the OLS specification where the standard errors are clustered at the participant level and include edition, level and school fixed effects. Standard errors in parenthesis, where \*\*\* p<0.01, \*\* p<0.05, \* p<0.1

**Table A9. Descriptive Statistics**

	Variables from the Questionnaire (2017-2018)															
	Overall					Men					Female					<i>p</i> -value
	Obs.	Mean	St. Dev	Min	Max	Obs.	Mean	St. Dev	Min	Max	Obs.	Mean	St. Dev	Min	Max	
No. of Preparation Hours	5925	4.36	8.65	0	100	3897	4.40	8.77	0	100	2028	4.28	8.42	0	100	0.63
Non-Competitiveness	6390	3.64	1.03	1	5	4133	3.55	1.05	1	5	2257	3.81	0.97	1	5	0.00
Overconfidence	4800	3.75	4.39	-16	21	3112	3.83	4.48	-15.25	21	1688	3.61	4.20	-16	18	0.11
Risk	5301	2.04	1.11	1	5	3400	2.11	1.16	1	5	1901	1.91	1.00	1	5	0.00
Perceived Math Ability	6105	4.10	0.72	1	5	3941	4.16	0.73	1	5	2164	3.99	0.70	1	5	0.00
Perceived Gender Nature of Math	6118	1.99	0.24	1	3	3945	1.98	0.24	1	3	2173	2.01	0.22	1	3	0.00

*Notes:* *No. of Preparation Hours* measures the total number of hours devoted to prepare the Math test. *Non-Competitiveness* contains the responses to question 1 in the questionnaire. *Perceived Math Ability* contains the responses to question 9 in the questionnaire. *Overconfidence* measures the difference between the guessed number of correct answers and the actual number of correct answers. And *Risk* contains the responses to question 8 in the questionnaire.

**Table A10. Gender Differences between No Reward and the Reward Parts of the Test:  
Confidence, Overconfidence or Risk?**

	Original Sample	Sample with Questionnaire		Sample with Questionnaire			
	zomitted (1)	zomitted (2)	zomitted (3)	zomitted (4)	zomitted (5)	zomitted (6)	zomitted (7)
Female	0.170*** (0.0249)	0.172*** (0.0348)	0.114*** (0.0348)	0.116*** (0.0350)	0.115*** (0.0348)	0.116*** (0.0348)	0.135*** (0.0346)
Reward	-0.0462*** (0.0162)	-0.0265 (0.0240)	-0.0265 (0.0240)	-0.102 (0.0757)	0.00407 (0.119)	0.0283 (0.0304)	0.427*** (0.0455)
Female*Reward	0.145*** (0.0295)	0.103** (0.0430)	0.103** (0.0430)	0.0975** (0.0434)	0.101** (0.0432)	0.0992** (0.0429)	0.0601 (0.0426)
Math at School	0.0381*** (0.00862)	0.0394*** (0.0128)	0.0311** (0.0126)	0.0311** (0.0126)	0.0311** (0.0126)	0.0311** (0.0126)	0.0311** (0.0126)
Participation Time	-0.156*** (0.0195)	-0.163*** (0.0261)	-0.138*** (0.0240)	-0.138*** (0.0240)	-0.138*** (0.0240)	-0.138*** (0.0240)	-0.138*** (0.0240)
Non-Competitiveness			0.0409*** (0.0117)	0.0304* (0.0159)	0.0409*** (0.0117)	0.0409*** (0.0117)	0.0409*** (0.0117)
Perceived Math Ability			-0.0655*** (0.0167)	-0.0655*** (0.0167)	-0.0618*** (0.0215)	-0.0655*** (0.0167)	-0.0655*** (0.0167)
Overconfidence			-0.0137*** (0.00265)	-0.0137*** (0.00265)	-0.0137*** (0.00265)	-0.00661** (0.00319)	-0.0137*** (0.00265)
Risk			-0.202*** (0.00994)	-0.202*** (0.00994)	-0.202*** (0.00994)	-0.202*** (0.00994)	-0.0946*** (0.0117)
Non-Competitiveness*Reward				0.0210 (0.0204)			
Perceived Math Ability*Reward					-0.00732 (0.0277)		
Overconfidence*Reward						-0.0141*** (0.00433)	
Risk*reward							-0.215*** (0.0171)
Observations	17,850	8,312	8,312	8,312	8,312	8,312	8,312
R-squared	0.096	0.136	0.186	0.186	0.186	0.187	0.200

*Notes:* Observations are at the Math test's parts level. Column 1 shows the main estimation result from column 1 in Table 2 for the original sample. For the rest of the columns, observations are at the Math test's parts level in the edition of 2017 and 2018 for the participants whose questionnaire answers are available. *No. Omitted* is standardized at the edition, level and part of the test levels. *Female* takes the value of 1 if the participant is female. *Reward* takes the value of 1 if the outcome variable refers to the part of the test with reward for omitted questions. *Math at School* measures the standardized Math grade at school level and *Participation Time* takes the values of 1, 2, 3 if it is the first, second or third time that the participant does the Math test. *Non-Competitiveness* contains the responses to question 1 in the questionnaire. *Perceived Math Ability* contains the responses to question 9 in the questionnaire. *Overconfidence* measures the difference between the guessed number of correct answers and the actual number of correct answers. And *Risk* contains the responses to question 8 in the questionnaire. All columns show the OLS specification where the standard errors are clustered at the participant level and include edition, level and school fixed effects. Standard errors in parenthesis, where \*\*\* p<0.01, \*\* p<0.05, \* p<0.1

**Table A11. Gender Differences between No Reward and the Reward Part of the Test along the Distribution of Ability: Confidence, Overconfidence or Risk?**

	Sample				Sample with Questionnaire											
	Original Sample		with Questionnaire		Low Ability		High Ability		Low Ability		High Ability		Low Ability		High Ability	
	Low Ability zomitted (1)	High Ability zomitted (2)	Low Ability zomitted (3)	High Ability zomitted (4)	Low Ability zomitted (5)	High Ability zomitted (6)	Low Ability zomitted (7)	High Ability zomitted (8)	Low Ability zomitted (9)	High Ability zomitted (10)	Low Ability zomitted (11)	High Ability zomitted (12)	Low Ability zomitted (13)	High Ability zomitted (14)		
Female	0.191*** (0.0387)	0.0847*** (0.0226)	0.196*** (0.0539)	0.132*** (0.0351)	0.150*** (0.0526)	0.0641* (0.0348)	0.149*** (0.0529)	0.0735** (0.0348)	0.148*** (0.0526)	0.0684** (0.0348)	0.150*** (0.0525)	0.0709** (0.0347)	0.166*** (0.0525)	0.0929*** (0.0343)		
Reward	-0.0929*** (0.0248)	-0.00137 (0.0194)	-0.0493 (0.0363)	0.0118 (0.0289)	-0.0493 (0.0363)	0.0118 (0.0289)	-0.0204 (0.119)	-0.242*** (0.0854)	-0.199 (0.172)	0.234 (0.156)	-0.0558 (0.0553)	0.0575* (0.0336)	0.339*** (0.0683)	0.543*** (0.0539)		
Female*Reward	0.0505 (0.0425)	0.264*** (0.0358)	-0.0290 (0.0614)	0.222*** (0.0527)	-0.0290 (0.0614)	0.222*** (0.0527)	-0.0275 (0.0620)	0.203*** (0.0531)	-0.0241 (0.0614)	0.214*** (0.0531)	-0.0285 (0.0613)	0.209*** (0.0527)	-0.0611 (0.0611)	0.165*** (0.0514)		
No. Of Correct No Reward	-0.0811*** (0.0109)	-0.0576*** (0.00510)	-0.0782*** (0.0171)	-0.0439*** (0.00823)	-0.122*** (0.0172)	-0.0569*** (0.00817)	-0.122*** (0.0172)	-0.0569*** (0.00817)	-0.122*** (0.0172)	-0.0569*** (0.00817)	-0.122*** (0.0172)	-0.0569*** (0.00817)	-0.122*** (0.0172)	-0.0569*** (0.00817)		
Participation Time	-0.0381 (0.0378)	-0.0972*** (0.0185)	-0.00511 (0.0521)	-0.121*** (0.0270)	0.0404 (0.0477)	-0.0960*** (0.0242)	0.0404 (0.0477)	-0.0960*** (0.0242)	0.0404 (0.0477)	-0.0960*** (0.0242)	0.0404 (0.0477)	-0.0960*** (0.0242)	0.0404 (0.0477)	-0.0960*** (0.0242)		
Non-Competitiveness					0.00617 (0.0171)	0.0275** (0.0125)	0.0101 (0.0252)	-0.00906 (0.0137)	0.00617 (0.0171)	0.0275** (0.0125)	0.00617 (0.0171)	0.0275** (0.0125)	0.00617 (0.0171)	0.0275** (0.0125)		
Perceived Math Ability					-0.0266 (0.0255)	-0.0247 (0.0176)	-0.0266 (0.0255)	-0.0247 (0.0176)	-0.0450 (0.0353)	0.00138 (0.0194)	-0.0266 (0.0255)	-0.0247 (0.0176)	-0.0266 (0.0255)	-0.0247 (0.0176)		
Overconfidence					-0.0472*** (0.00462)	-0.0210*** (0.00353)	-0.0472*** (0.00462)	-0.0210*** (0.00353)	-0.0472*** (0.00462)	-0.0210*** (0.00353)	-0.0478*** (0.00595)	-0.0116*** (0.00370)	-0.0472*** (0.00462)	-0.0210*** (0.00353)		
Risk					-0.229*** (0.0156)	-0.184*** (0.0114)	-0.229*** (0.0156)	-0.184*** (0.0114)	-0.229*** (0.0156)	-0.184*** (0.0114)	-0.229*** (0.0156)	-0.184*** (0.0114)	-0.136*** (0.0197)	-0.0590*** (0.0106)		
Non-Competitiveness*Reward							-0.00782 (0.0311)	0.0732*** (0.0230)								
Perceived Math Ability*Reward									0.0367 (0.0414)	-0.0522 (0.0353)						
Overconfidence*Reward											0.00121 (0.00685)	-0.0188*** (0.00596)				
Risk*reward													-0.186*** (0.0248)	-0.250*** (0.0213)		
Observations	10,048	9,720	4,920	4,366	4,920	4,366	4,920	4,366	4,920	4,366	4,920	4,366	4,920	4,366		
R-squared	0.122	0.153	0.185	0.208	0.244	0.266	0.244	0.268	0.244	0.266	0.244	0.268	0.252	0.297		

Notes: Observations are at the Math test's parts level. Columns 1-2 show the main estimation results from columns 1-2 in Table 3 for the original sample. For the rest of the columns, observations are at the Math test's parts level in the edition of 2017 and 2018 for the participants whose questionnaire answers are available. *No. Omitted* is standardized at the edition, level and part of the test levels. *Female* takes the value of 1 if the participant is female. *Reward* takes the value of 1 if the outcome variable refers to the part of the test with reward for omitted questions. *No. of Correct No Reward* measures the number of correct questions in the part of the test without any reward for omitted questions, and *Participation Time* takes the values of 1, 2, 3 if it is the first, second or third time that the participant does the Math test. *Non-Competitiveness* contains the responses to question 1 in the questionnaire. *Perceived Math Ability* contains the responses to question 9 in the questionnaire. *Overconfidence* measures the difference between the guessed number of correct answers and the actual number of correct answers. And *Risk* contains the responses to question 8 in the questionnaire. All columns show the OLS specification where the standard errors are clustered at the participant level and include edition, level and school fixed effects. Standard errors in parenthesis, where \*\*\* p<0.01, \*\* p<0.05, \* p<0.1

**Table A12. Gender Differences between No Reward and the Reward Part of the Test along Age: Confidence, Overconfidence or Risk?**

	Original Sample		Sample with Questionnaire		Level 1		Level 4		Level 1		Level 4		Level 1		Level 4	
	Level 1	Level 4	Level 1	Level 4	Level 1	Level 4	Level 1	Level 4	Level 1	Level 4	Level 1	Level 4	Level 1	Level 4	Level 1	Level 4
	zomitted	zomitted	zomitted	zomitted	zomitted	zomitted	zomitted	zomitted	zomitted	zomitted	zomitted	zomitted	zomitted	zomitted	zomitted	zomitted
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)	(13)	(14)		
Female	0.214*** (0.0502)	0.179** (0.0701)	0.149** (0.0688)	0.338*** (0.116)	0.130* (0.0687)	0.258** (0.112)	0.122* (0.0688)	0.260** (0.112)	0.134* (0.0687)	0.253** (0.111)	0.130* (0.0687)	0.269** (0.111)	0.140** (0.0687)	0.303*** (0.111)		
Reward	-0.000431 (0.0320)	-0.0613 (0.0424)	0.0289 (0.0497)	-0.0971 (0.0663)	0.0289 (0.0497)	-0.0971 (0.0664)	0.355** (0.178)	-0.136 (0.194)	0.472* (0.262)	-0.405 (0.349)	0.0331 (0.0669)	-0.00927 (0.0862)	0.228*** (0.0849)	0.708*** (0.139)		
Female*Reward	0.0273 (0.0587)	0.234*** (0.0861)	-0.0266 (0.0870)	0.167 (0.139)	-0.0266 (0.0871)	0.167 (0.140)	-0.0104 (0.0872)	0.163 (0.141)	-0.0346 (0.0871)	0.178 (0.139)	-0.0263 (0.0870)	0.145 (0.140)	-0.0455 (0.0868)	0.0764 (0.136)		
Math at School	-0.0125 (0.0200)	0.0234 (0.0207)	-0.0294 (0.0290)	0.0261 (0.0382)	-0.0177 (0.0294)	-0.00783 (0.0345)	-0.0177 (0.0294)	-0.00783 (0.0345)	-0.0177 (0.0294)	-0.00783 (0.0345)	-0.0177 (0.0294)	-0.00783 (0.0345)	-0.0177 (0.0294)	-0.00783 (0.0345)		
Participation Time	-0.277*** (0.0570)	-0.0693 (0.0422)	-0.171** (0.0768)	-0.0432 (0.0703)	-0.157** (0.0744)	-0.00290 (0.0618)	-0.157** (0.0744)	-0.00290 (0.0618)	-0.157** (0.0744)	-0.00290 (0.0618)	-0.157** (0.0744)	-0.00290 (0.0618)	-0.157** (0.0744)	-0.00290 (0.0618)		
Non-Competitiveness					-0.00787 (0.0259)	0.0577 (0.0369)	0.0350 (0.0345)	0.0521 (0.0497)	-0.00787 (0.0260)	0.0577 (0.0370)	-0.00787 (0.0260)	0.0577 (0.0370)	-0.00787 (0.0260)	0.0577 (0.0370)		
Perceived Math Ability					-0.139*** (0.0375)	-0.0129 (0.0485)	-0.139*** (0.0375)	-0.0129 (0.0485)	-0.0866* (0.0474)	-0.0508 (0.0669)	-0.139*** (0.0375)	-0.0129 (0.0485)	-0.139*** (0.0375)	-0.0129 (0.0485)		
Overconfidence					-0.00493 (0.00529)	-0.0262*** (0.00811)	-0.00493 (0.00529)	-0.0262*** (0.00812)	-0.00493 (0.00529)	-0.0262*** (0.00812)	-0.00450 (0.00692)	-0.0131 (0.00957)	-0.00493 (0.00529)	-0.0262*** (0.00812)		
Risk					-0.119*** (0.0210)	-0.268*** (0.0285)	-0.119*** (0.0210)	-0.268*** (0.0285)	-0.119*** (0.0210)	-0.268*** (0.0285)	-0.119*** (0.0210)	-0.268*** (0.0285)	-0.119*** (0.0262)	-0.268*** (0.0353)		
Non-Competitiveness*Reward							-0.0857* (0.0452)	0.0113 (0.0557)								
Perceived Math Ability*Reward									-0.104* (0.0593)	0.0758 (0.0847)						
Overconfidence*Reward											-0.000861 (0.00834)	-0.0262** (0.0132)				
Risk*reward														-0.105*** (0.0343)	-0.321*** (0.0435)	
Observations	4,250	2,792	2,080	1,260	2,080	1,260	2,080	1,260	2,080	1,260	2,080	1,260	2,080	1,260		
R-squared	0.167	0.212	0.212	0.305	0.232	0.364	0.232	0.364	0.232	0.364	0.232	0.364	0.232	0.364		

Notes: Observations are at the Math test's parts level. Columns 1-2 show the main estimation results from columns 1-2 in Table 4 for the original sample. For the rest of the columns, observations are at the Math test's parts level in the edition of 2017 and 2018 for the participants whose questionnaire answers are available. *No. Omitted* is standardized at the edition, level and part of the test levels. *Female* takes the value of 1 if the participant is female. *Reward* takes the value of 1 if the outcome variable refers to the part of the test with reward for omitted questions. *Math at School* measures the Math grade at school, and *Participation Time* takes the values of 1, 2, 3 if it is the first, second or third time that the participant does the Math test. *Non-Competitiveness* contains the responses to question 1 in the questionnaire. *Perceived Math Ability* contains the responses to question 9 in the questionnaire. *Overconfidence* measures the difference between the guessed number of correct answers and the actual number of correct answers. And *Risk* contains the responses to question 8 in the questionnaire. All columns show the OLS specification where the standard errors are clustered at the participant level and include edition, level and school fixed effects. Standard errors in parenthesis, where \*\*\* p<0.01, \*\* p<0.05, \* p<0.1