



The role of student effort on performance in PISA: Revisiting the gender gap in achievement

Lina Anaya

University of Arkansas

Gema Zamarro

University of Arkansas

International assessments are important to benchmark the quality of education across countries. However, on low-stakes tests, students' incentives to invest their maximum effort may be minimal. Research stresses that ignoring students' effort when interpreting results from low-stakes assessments can lead to biased interpretations of test performance across groups of examinees. We use data from the Programme for International Student Assessment (PISA), a low-stakes test, to analyze the extent to which student effort helps to explain test scores heterogeneity across countries and by gender groups. Our results highlight the importance of accounting for differences in student effort to understand cross-country heterogeneity in performance and variations in gender achievement gaps across nations. We find that, once we account for differential student effort across gender groups, the estimated gender achievement gap in math and science could be up to 12 and 6 times wider, respectively, and up to 49 percent narrower in reading, in favor of boys. In math and science, the gap widens in most countries, even among some of the top 20 most gender-equal countries. Altogether, our effort measures on average explain between 36 and 40 percent of the cross-country variation in test scores.

VERSION: September 2020

Suggested citation: Anaya, Lina M., and Gema Zamarro. (2020). The role of student effort on performance in PISA: Revisiting the gender gap in achievement . (EdWorkingPaper: 20-295). Retrieved from Annenberg Institute at Brown University: <https://doi.org/10.26300/f056-5t92>

The role of student effort on performance in PISA: Revisiting the gender gap in achievement

Lina Anaya*
University of Arkansas

Gema Zamarro
University of Arkansas

September, 2020

*Corresponding Author: Lina Anaya, University of Arkansas, 211 Graduate Education Building,
Fayetteville, Arkansas, 72701, email: lanaya@aurk.edu

Abstract:

International assessments are important to benchmark the quality of education across countries. However, on low-stakes tests, students' incentives to invest their maximum effort may be minimal. Research stresses that ignoring students' effort when interpreting results from low-stakes assessments can lead to biased interpretations of test performance across groups of examinees. We use data from the Programme for International Student Assessment (PISA), a low-stakes test, to analyze the extent to which student effort helps to explain test scores heterogeneity across countries and by gender groups. Our results highlight the importance of accounting for differences in student effort to understand cross-country heterogeneity in performance and variations in gender achievement gaps across nations. We find that, once we account for differential student effort across gender groups, the estimated gender achievement gap in math and science could be up to 12 and 6 times wider, respectively, and up to 49 percent narrower in reading, in favor of boys. In math and science, the gap widens in most countries, even among some of the top 20 most gender-equal countries. Altogether, our effort measures on average explain between 36 and 40 percent of the cross-country variation in test scores.

JEL Codes: I20, J16, C83

Keywords: Student Effort, Gender Gaps, Rapid guessing, PISA 2015.

1. Introduction

Understanding how well a school or an educational system educates its students is important for stakeholders such as parents, teachers, and governments. Standardized assessments help policymakers to benchmark the quality of schools or a country's educational system relative to other nations. However, when students do not face the consequences for high or low performance, their incentives to invest their maximum effort on the test may be minimal. Thus, differences in test performance may not just reflect variations in actual content knowledge but also differences in other non-content-knowledge factors, such as student effort. One such example is low-stakes international assessments, such as PISA (Programme for International Student Assessment) or TIMSS (Trends in International Mathematics Science Study), in which differences in student effort may be essential for explaining part of the observed differences in student achievement across and within countries by gender.

Several studies find that ignoring student effort may lead to biased conclusions about the test performance of a group of examinees (Demars, 2007; Swerdzewski, Harmes, & Finney, 2011; Wise & Kong, 2005; Wise & DeMars, 2010). This problem can worsen when making international comparisons of achievement. Evidence from international assessments shows that student effort is essential to understand differences in test performance within and across countries (Boe, May, & Boruch, 2002; Debeer, Buchholz, Hartig, & Janssen, 2014; Zamarro, Hitt, & Mendez, 2019).

In this paper, we revisit the prior literature studying the role of effort on explaining differences in test scores to analyze the extent to which student effort contributes to explain variation in test performance in math, reading, and science, across countries, as well as within countries by gender. We use data from the PISA 2015 computer assessment and student computer-

based survey to construct measures of student effort based on the instances of rapid-guessing responses in the test and the effort students put forward in the survey (i.e. item non-response rates), respectively. Prior research from PISA suggests that student item non-response rates contribute to explain a significant part of the variation across countries in test scores (Zamarro et al., 2019).

To compute student rapid-guessing rates, we use the inverse response-time-effort (RTE) score as introduced by Wise & Kong (2005). Following Wise & Kong (2005), we use the information on response times for each question to calculate the proportion of questions of the assessment in which the examinee does not engage in solution behavior (i.e., the examinee does not take the time to analyze the question [Schnipke, 1995; Schnipke and Scrams, 1997]).

Differences in student effort could help explain differences in student performance across countries, as well as test score gender gaps within countries. Obtaining a better understanding of the role of effort on gender achievement gaps is important given the predictive power of math and science performance on explaining women's underrepresentation in science occupations (Anaya, Stafford, & Zamarro, 2017; Ceci, Ginther, Kahn, & Williams, 2014; Nix, Perez-Felkner, & Thomas, 2015; Perez-Felkner, Nix, & Thomas, 2017).

If student effort varies by gender, differences in effort could affect our understanding of gender gaps in test performance. Along these lines, Balart and Oosterveen (2019) use measures of decline in performance throughout the PISA test and find that girls are better at sustaining test performance than boys. According to the authors, this result has consequences for the measurement of the gender achievement gap because in longer assessments, the gap in math and science is smaller compared to shorter assessments. Using data from the U.S., Soland (2019) obtains similar findings. Soland (2019) measures effort based on response times of test questions and finds that

after removing the effect of effort in test scores, the gender gap in math achievement would be wider, and it is more sensitive to effort-adjustment than the reading gap.

Our study advances the current state of knowledge in two ways:

1. We contribute to prior literature about student effort in international assessments (Balart & Oosterveen, 2019; Boe et al., 2002; Debeer et al., 2014; Zamarro et al., 2019) that, to our knowledge, mostly uses data from paper-based assessments, by exploiting response times of test questions to build a measure of rapid-guessing in an international test such as PISA. We also contribute to this literature by replicating effort measures traditionally used in paper based-surveys, such as item non-response rates, on a computer-based assessment.
2. We contribute to the rapid-guessing literature by studying instances of rapid-guessing on a large international representative sample. Most of the research using this technique focuses on U.S. samples, and some of them are small convenient samples. Additionally, few studies analyze whether or not there are gender differences in test effort (DeMars, Bashkov, & Socha, 2013; Soland, 2018; Soland, 2019; Soland, Jensen, Keys, Bi, & Wolk, 2019; Swerdzewski et al., 2011; Wise & Ma, 2012; Wise & DeMars, 2005; Wise, Pastor, & Kong, 2009).

We find evidence of variation of rapid-guessing behavior in PISA. In line with prior research, we find that student effort explains a significant part of the variation in PISA scores across countries. Altogether, our effort measures represent, on average, between 36 and 40 percent of the variation in test performance across countries. Also, the probability of engaging in rapid-guessing behavior is higher for boys than for girls, which has implications for estimated gender

gaps in performance. Accounting for student effort affects the estimated gender gaps in achievement. We find that the gender achievement gap could be up to 6 and 11 times wider in science and math, respectively, and up to 50 percent narrower in reading, in favor of boys.

The remaining parts of this document are organized as follows: Section 2 presents the literature review; section 3 explains in more detail the data we use in this study; section 4 describes the measures of student effort in PISA that we use in the paper; section 5 shows the methodology and results; and, section 6 presents our conclusions.

2. Literature review

Student motivation or effort is an essential element to understand student achievement in low-stakes assessments. Wise & DeMars (2005) define student motivation as the amount of effort or energy that a student invests towards achieving the highest possible score on a test. When students do not face consequences for performance, their incentives to invest their maximum effort on the test may be minimal. As a result, ignoring the role of students' motivation in the interpretation of test scores may lead to biased conclusions given that the resulting scores may not be an accurate indicator of students' ability (Kane, 2006; Swerdzewski et al., 2011; Wise & DeMars, 2005).

A significant first step to take student effort into account when interpreting test scores is to identify who the low-effort examinees are. Researchers who analyze student effort using large representative samples from international assessments have developed several methods to calculate student effort using paper-based assessments (Boe et al., 2002; Borghans & Schils, 2012; Debeer et al., 2014; Zamarro et al., 2019). For example, Debeer et al. (2014) focuses on the reading

achievement data from PISA 2009 and defines effort as the difference in test performance due to the different positions a group of reading questions occupy on the test.

In contrast, Borghans and Schils (2012) employ the rate of decline in performance as the test progresses, while Zamarro et al. (2019) not only employ the rate of decline in performance but also measure the careless answering patterns and item non-response rates on the survey students take after the PISA 2009 test, in order to measure student effort. The authors find that item non-response in the survey has the highest predictive power in explaining differences in test scores across countries. Previous work also highlights the importance of item non-response rates, as a proxy for non-cognitive skills, to understand how differences in student effort can explain cross-country differences in achievement (Boe et al., 2002).

Computer-based assessments create a new opportunity for researchers to develop new measures of student effort. Wise and Kong (2005) propose using the response-time-effort (RTE) score, which focuses on examinees' response times in computer-based-low-stakes assessments, as a proxy for motivation. This idea comes from Schnipke (1995) and Schnipke and Scrams (1997), who define solution behavior as the situation in which the examinee takes the time to analyze the question in order to find the right answer and, rapid-guessing behavior, when the examinee rapidly chooses a response.

Although in high-stakes evaluations, rapid-guessing may represent the hurry to answer all the questions, when examinees do not have enough time to complete the test using solution behavior (Schnipke, 1995; Schnipke & Scrams, 1997), Wise and Kong (2005) argue that in a low-stakes context, responses given within a short time represent students' low engagement in trying to find the right answer. As a result, the RTE score represents the proportion of test questions for

which the examinee exhibits solution behavior (Wise & Kong, 2005). When the RTE score is close to zero, it represents a low-effort student who rapidly guesses most of the test question answers, while an RTE close to one represents a high-effort examinee who engages in solution behavior in answering most of the questions. Therefore, the rapid-guessing rate is defined as the inverse RTE score.

To develop the RTE scores, Wise and Kong (2005) use data from a low-stakes computer test of a random sample of about 400 college students. To set the time thresholds that separate rapid-guessing from solution behavior, Wise and Kong (2005) conduct a visual inspection of response time distributions and question structure for each question separately. Wise and Kong (2005) show that RTE is then a valid measure of student motivation because of its high reliability, alpha of .97, and its correlation with other measures of motivation such as self-reported test effort. Additionally, their results show that RTE is weakly correlated with SAT scores, which exemplifies that student motivation can be differentiable from ability, a distinction not easily possible using self-reported measures of effort. Finally, the RTE approach evinces that the rate at which rapid guessers choose the right answer is not higher than the probability of getting the question right by chance, which suggests that this method creates a reliable distinction between rapid-guessing and solution behavior.

Although other studies obtain similar findings to Wise and Kong (2005) regarding the RTE score validity (Kong, Wise, & Bhola, 2007; Swerdzewski et al., 2011; Wise, 2006), performing a question-by-question inspection to set time thresholds can be tedious and unfeasible on long assessments such as PISA. Instead, Wise and Ma (2012) propose using the normative threshold (NT) method to set the question-by-question time thresholds. In the NT method, the time threshold

is a percentage of the mean response time of a given question. The threshold should not exceed a maximum value of 10 seconds; thresholds above 10 seconds may not produce a reliable classification of rapid-guessing and solution behavior (Setzer, Wise, van, & Ling, 2013).

Wise and Ma (2012) evaluate the performance of three thresholds, 10, 15, and 20 percent of the mean question-specific response time, on identifying rapid-guessing responses. Using data from a large-scale computer-based assessment that has more than 200 thousand students from the third to the ninth grades in the U.S., the authors find that only the NT at 10 percent of the mean shows accuracy in classifying solution and rapid-guessing behavior. In contrast, the NT at 15 and 20 percent provide evidence of classifying effortful responses as rapid-guessing. The authors recommend using the NT at 10 percent of the mean given its better accuracy in classifying effortful and non-effortful responses.

Concerning how low student effort can potentially distort average test score results, as well as proficiency rates for a group of examinees, Wise and DeMars (2010) exclude from the calculation of group test performance the test score data of low-effort students in order to obtain a cleaner measure of overall achievement. The authors use a sample of about 300 college students who take a low-stakes computer test and then remove from the sample the test scores of low-effort examinees whose RTE score is below 90 percent. Their findings show that the mean test score gains almost doubled, and the percentage of students scoring at or above the proficiency score increased approximately by eight percentage points after adjusting test scores by effort. Our paper contributes to this literature by studying patterns of rapid-guessing in PISA and studying their importance on observed differences in test performance across countries, as well as differences in test score gender gaps within each country.

There is little research available that explicitly studies the effect of student effort on gender differences in test performance (DeMars et al., 2013; Soland, 2018; Soland, 2019; Wise et al., 2009). In this respect, this paper contributes to an emerging literature on this topic. DeMars et al. (2013) study gender differences in test effort using RTE scores of a random sample of about 2,000 college students. The authors find that, on average, male students have a lower RTE score than their female peers. At the lower tail of the RTE score distribution, the gender differences are more significant given that a higher percentage of male students engage in rapid-guessing behavior. However, the limitation of this study is that the sample size hinders generalizing the findings.

Along these lines, Soland (2018) and Soland (2019) extend the analysis from DeMars et al. (2013) and Wise et al. (2009) by not only studying gender differences in the RTE scores but also assessing how accounting for student effort may change the measured achievement gaps in math and reading. Soland (2018) and Soland (2019) use student data from five and seven states in the U.S., respectively, that come from the Measures of Academic Progress (MAP) test. The findings suggest that although the male-female differences in rapid-guessing rates do not change the interpretations of achievement gaps in a significant way, the gender gap in math is more sensitive to effort-adjustment than the reading gap. Soland (2019) calls into question whether or not recent progress in narrowing the gap in math may reflect differences in effort rather than test score gains by female students.

A related work that connects student effort with gender achievement gaps, but using data from international assessments, also highlights the implications of effort in the measurement of gender gaps in test scores. Balart and Oosterveen (2019) employ the rate of decline in performance throughout the PISA 2015 test to study gender differences in sustaining performance and its

implications for the gender achievement gap. The authors find that in longer assessments, the gender gap in math and science decreases, which occurs because, in most countries, girls are better able to sustain performance throughout the test relative to boys, even in math and science subjects.

In this paper, we use data from the computer-based assessment PISA 2015 to examine to what extent student effort helps explain cross-country variation in test performance, as well as gender gaps in achievement, within each country, in the subjects of math, reading, and science. Our study builds upon the previous work we present in this literature review, especially on previous work from Soland (2018), Soland (2019), Balart and Oosterveen (2019), DeMars et al. (2013), Deeber et al. (2014), Zamorro et al. (2019), Wise and Ma (2012), Weinstein and Roediger (2012), Anaya et al. (2019), and Bard and Weinstein (2017). Our study advances the current state of knowledge in two ways:

First, we contribute to the student effort literature in international assessments such as PISA (Balart & Oosterveen, 2019; Debeer et al., 2014; Zamorro et al., 2019) by using the NT method and RTE approach to measure student motivation. To our knowledge, this method has not been applied to the full PISA achievement sample given that assessments before 2015 are paper-based assessments. Therefore, studies that use earlier versions of PISA adopt other approaches to define student effort because it is not possible to obtain response times for a paper-based test.

We find two studies that use the NT, or a similar method, to identify low-effort examinees in PISA 2015; however, they focus on only one subject or a subsample of students and do not analyze the consequences of low-effort on gender achievement gaps (Akyol, Krishna, & Wang, 2018; Michaelides, Ivanova, & Nicolaou, 2020). In contrast, Balart and Oosterveen (2019)'s work focuses on gender achievement gaps, but it uses a different measure of effort.

Second, we contribute to the RTE literature by replicating the RTE approach and the NT method in a large international representative sample. Most of the research using this technique focuses on U.S. samples, and some of them are based on small convenient samples (DeMars et al., 2013; Soland, 2018; Soland, 2019; Soland et al., 2019; Swerdzewski et al., 2011; Wise & Ma, 2012; Wise & DeMars, 2005; Wise et al., 2009). Besides, few studies analyze gender differences in student effort using the RTE approach (DeMars et al., 2013; Soland, 2018; Soland, 2019; Wise et al., 2009) and the implications for gender achievement gaps. Only Soland (2018) and Soland (2019) assess the effects of rapid-guessing behavior on the measurement of gender achievement gaps in math and reading; however, these studies only use a sample of students from the U.S.

3. Data

The Programme for International Student Assessment (PISA) is a triannual survey, managed by the Organization for Economic Co-operation and Development (OECD), which evaluates how well 15-year-old students are capable of using their knowledge and skills to meet real-life challenges in the areas of mathematics, reading, and science. The number of participants in 2015 was about 540,000 students from 72 countries and economies¹. In addition to the three core evaluation subjects, PISA 2015 evaluated students on collaborative problem solving and financial literacy. These last two subjects were optional for the participant countries. Every PISA wave focuses on a subject; in 2015, the primary area of assessment was science, and therefore, the evaluation included more questions about this topic.

¹ To simplify, in the rest of this document we use the term countries to refer to countries and economies. See Table 1 for the list of countries and abbreviations.

For the first time, the main form of assessment in PISA 2015 was computer-based. Paper-based assessments were available to countries that had limited access to computers. These two forms of assessments lasted about two hours. After the completion of the test, students answered a background questionnaire about 30 minutes long that collected information about home environment, school, and learning experiences.

For this study, we restrict our sample to those countries and economies that took the computer-based test. We also exclude the test booklets that have clusters about cooperative problem solving, financial literacy, or that were designed for students with special needs. Our final sample contains 55 countries/economies². We only focus on the computer-based assessment because this form includes response times for each student on each question, which we use later in order to define rapid-guessing behavior.

In the PISA 2015 assessment, the test booklets are randomly assigned to students within each country. The total number of questions in these booklets ranges from 47 to 71 questions with an average of 60 total questions.

4. Measuring student effort in PISA

4.1. Rapid-guessing in the entire assessment

We define rapid-guessing as the inverse RTE score ($1 - RTE$), and it represents the proportion of responses, out of all test questions, in which an examinee engages in rapid-guessing behavior. To create the rapid-guessing variable, we first calculate average response time for each question across all test booklets within each country. Second, we use the NT method at 10 percent

² We restrict our analytical sample to countries and economies. We exclude the adjudicated regions of USA Massachusetts, USA North Carolina, and the adjudicated regions from Spain.

of the mean to set time thresholds for each question within each country. We focus on 10 percent of the mean because prior evidence suggests that this threshold has better accuracy in classifying rapid-guessing and solution behavior (Wise & Ma, 2012). Finally, we identify the number of questions in which an examinee's response time is below the 10 percent of the mean³ to calculate the inverse RTE score (i.e., the proportion of rapid-guessing responses) on the complete test for each student within each country.

When calculating the rapid-guessing rate on the test, we exclude response times from students whose total time in completing the test exceeds 120 minutes⁴, which represents 2,492 observations. Although the test was expected to last two hours, we are unsure of whether or not some students obtained extra time. Total time above 120 minutes could also occur because test proctors had to log off the computer assessment one by one. According to what we see in the data, it seems that in some cases, the proctor did not end the session, or there was a technical problem in the data collection because we find some records of total time spent on the assessment of up to 14 hours.

Tables 2 and 3 show descriptive statistics of rapid-guessing behavior in the complete assessment, as well as other variables of interests that we describe in the following sections. Students in the estimation sample take, on average, 75 minutes to complete the assessment (see table 2). Approximately four observations have total times of less than one minute, which may occur because of a technical problem in data collection or because the students decided not to

³ We also performed a sensitivity analysis using a more conservative threshold of 5% of the mean response time and our findings do not change significantly. Results are available from the authors upon request.

⁴ We also conducted our estimations without excluding outliers in total time and the results do not change meaningfully. Estimates excluding outliers are the ones presented in the paper since they are more conservative. The results that did not exclude outliers are available upon request.

complete the assessment. The variation in total time is lower between countries than within countries, which suggests that the distribution of total time across countries probably does not vary meaningfully.

Although the proportion of rapid-guessing on the test ranges from 0 to 100 percent, students across countries on average rapidly guess 5 percent of all test questions (see table 2). Since the average number of questions in PISA booklets is 60, a 5 percent rapid-guessing rate on the test is equivalent to rapidly guessing about 3 questions on average. Table 2 also shows that the variation in rapid-guessing behavior is higher across all students, regardless of country, and within countries rather than between countries. The standard deviations for the whole sample show that, overall, the average dispersion in the proportion of rapid-guessing responses is about 8 percentage points. When comparing students within each country, the variation is slightly lower, showing that the dispersion of rapid-guessing proportions is, on average, 7 percentage points above or below the mean. In contrast, the variation between countries is roughly a third lower, with a standard deviation of about 2 percentage points.

When we look at the average rapid-guessing rate for boys and girls (see table 3), their rates differ roughly by one percentage point. Girls have a slightly lower probability of engaging in rapid-guessing behavior than boys. This result is similar to prior research which finds that female students, on average, have lower rapid-guessing rates than boys have (DeMars et al., 2013; Soland, 2018; Soland, 2019). This result is consistent with the difference in total time between girls and boys. Girls, on average, take 5 minutes longer than boys do in completing the assessment.

In summary, we find descriptive evidence of rapid-guessing behavior in PISA 2015. The dispersion of this variable is higher when we compare all students, regardless of country, and when

we compare examinees within each country. The variation is lower across countries, which suggests that across countries, the distributions of rapid-guessing behavior probably are not very different from each other. The latter does not necessarily imply that student effort is not relevant to explain cross-country variations in achievement. Zamarro et al. (2019) find that despite the fact that cross-country variation in student effort is lower than the within-country variation, the differences in student effort across countries are still relevant in explaining cross-country heterogeneity in test scores. Finally, we observe that girls, on average, exhibit more effort and take more time to complete the test than boys do.

4.2. Item non-response on the student background survey

We replicate the Zamarro et al. (2019) approach by calculating the item non-response rate in the student survey, but this time by using a computer-based survey from PISA. This rate corresponds to the proportion of questions that a student skips or does not complete on the survey.⁵ We focus on the item non-response rate since previous research finds that this indicator has the highest predictive power in explaining cross-country variation in performance on paper-based assessments (Boe et al., 2002; Zamarro et al., 2019). According to table 2, students do not respond to between 0 and 98 percent of survey items, and on average, they leave blank 7 percent of the questions. The variation between and within countries on the item non-response rate is almost twice the variation on the rapid-guessing rate on the test. Girls on average have a roughly 2-percentage-points lower item non-response rate than boys have (see table 3). Overall, girls consistently show higher effort than boys do both in the test and in the survey.

⁵ Although we have response times for this questionnaire, we do not construct rapid-guessing rate for the background survey because PISA does not report response times for each question but for a group of items.

5. Estimating the role of student effort in explaining cross-country differences in achievement and within-country differences in gender achievement gaps

We follow a similar methodological approach to that of Zamarro et al. (2019) and conduct a country-random-effects estimation for each tested subject in PISA to assess the role that student effort may have in explaining cross-country differences in performance and within-country gender achievement gaps. Our dependent variable in model 1 corresponds to the plausible value j (i.e., test score) that student i from country c obtained on the subject s . The variables *INRsurvey* and *RGtest* represent the item non-response rate on the student background survey and the proportion of rapid-guessing responses on the entire assessment, respectively. The terms α and ε represent the country random-effect for the subject s and the error term, respectively.

$$TestScore_{ic}^{sj} = \beta_0^{sj} + \beta_1^{sj} INRsurvey_{ic} + \beta_2^{sj} RGtest_{ic} + \alpha_c^{sj} + \varepsilon_{ic}^{sj} \quad (1)$$

PISA reports test scores as plausible values. These scores are calculated using a multiple imputation method that aims to increase accuracy in measuring students' skills⁶. Each student has 30 possible values in total; ten plausible values for each subject. We estimate model 1 using as dependent variable each of the 10 plausible values on each subject, and we report the average estimated coefficients for each subject in table 4. We first examine effort measures separately and estimate equation 1 for each effort measure. We replicate Zamarro et al. (2019) results and find that item non-response is also a statistically significant predictor of test performance in this computer-based assessment.

⁶ For further information about plausible values and multiple imputation method, see chapter 9 of the PISA 2015 technical report.

From equation 1, we follow Zamarro et al.'s (2019) approach and obtain effort-adjusted test scores ($\overline{Adjusted Score}_{ic}^s$) for each student and subject by obtaining the average of the sum of the estimated coefficients of the intercept, the country random-effect, and the residuals ($\hat{\beta}_0^{sj} + \hat{\alpha}_c^{sj} + \hat{\varepsilon}_{ic}^{sj}$) using the following formula:

$$\overline{Adjusted Score}_{ic}^s = \sum_{j=1}^{10} \frac{\hat{\beta}_0^{sj} + \hat{\alpha}_c^{sj} + \hat{\varepsilon}_{ic}^{sj}}{10} \quad (2)$$

We then compute the average adjusted score for each subject across the 10 plausible values. Then we calculate the average effort-adjusted gender gap \overline{GAP}_c^s for each country and subject by subtracting the average effort-adjusted test score of girls minus that of boys using the formula:

$$\overline{GAP}_c^s = \sum_{g=1}^{G_c} \frac{\overline{Adjusted Score}_{G_c}^s}{G_c} - \sum_{b=1}^{B_c} \frac{\overline{Adjusted Score}_{B_c}^s}{B_c} \quad (3)$$

Where G_c and B_c represent the sample sizes of girls (G) and boys (B) from country c , respectively.

To calculate the effort-unadjusted test scores, we use formula 2 and replace the numerator with the actual plausible values that each student on the estimation sample obtained on each subject. Then we calculate the average effort-unadjusted achievement gap \overline{GAP}_c^s for each subject and country using formula 3 but replacing the numerator with the effort-unadjusted score that boys and girls in the estimation sample obtained on each subject. On average, students score before effort-adjustment 471, 474, and 476 points on the subjects of math, reading, and science, respectively (see table 2). Before effort-adjustment, girls score on average, 25 points higher on

reading than boys do, whereas in math and science, girls score 9 and 4 points lower than boys do, respectively (see table 3).

After calculating the average achievement gap for each subject and country using test scores, we compare the effort-adjusted and unadjusted gap using the following percentage change formula:

$$\Delta\%GAP_c^S = \frac{\overline{GAP_c^S} - \overline{GAP_c^S}}{\overline{GAP_c^S}} * 100 \quad (4)$$

Formula 4 represents the percentage change of the achievement gap relative to the average effort-unadjusted gap. In other words, formula 4 shows, compared to the unadjusted gap, what would be the expected percentage change in the average gender achievement gap for each country, and subject, in the absence of student effort heterogeneity. We adjust the signs of the calculated percentage changes such that negative signs represent a widening of the gender achievement gap, and positive signs represent a reduction of the gap.

5.1. Results of the role of student effort in explaining cross-country differences in student achievement

When we analyze to what extent our effort measures explain the variation in performance, we find that both item non-response rates and rapid-guessing are relevant predictors of test scores (see table 4). A one standard deviation increase in the proportion of rapid-guessing responses in the test is associated with a decrease of 0.13, 0.15, and 0.16 standard deviations on the math, science, and reading test scores, respectively (see columns 3, 6, and 9). Regarding the item non-response variable, a one standard deviation increase on this variable is associated with a decrease of 0.17, 0.18, and 0.21 standard deviations on the math, science, and reading test scores,

respectively (see columns 3, 6, and 9). These findings suggest that low-effort students often experience lower test performance.

Additionally, we find that our effort measures have more explanatory power across countries than within countries. Altogether, our effort measures explain between 36 and 40 percent of the variation in test performance across countries, which is similar to Zamarro et al.'s (2019) findings, versus about 3 to 4 percent of the within-country variation in test scores (see table 4). This finding is not very surprising. Previous work by Wise, Soland & Bo (2020) examine the distortive effect of effort heterogeneity in test scores at the school level using data from a pilot computer-based assessment from PISA in the U.S. Although the authors find variation in effort across schools, the mean test scores for each school after effort-adjustment do not significantly change compared to the effort-unadjusted scores. These effort measures may perform better at capturing differences in effort across different contexts or cultures than within similar environments, such as schools or countries.

5.2. Results of the role of student effort on gender achievement gaps

In this section, figures 1, 2, and 3 present the percentage change of the gender achievement gap in the absence of student effort heterogeneity, relative to the effort-unadjusted gap. Countries in the green color correspond to a reduction of the gap, represented by a positive percentage change. In contrast, the remaining colors correspond to a widening of the gap represented by a negative sign; the darker the color of a country is, the wider the gap becomes. Tables 5, 6, and 7 show the effort-adjusted and unadjusted gaps, as well as the percentage change for each country and subject.

The widening of the gap in math achievement occurs in 48 out of 55 countries and ranges from 2.9 to up to 1,158 percent (see figure 1 and table 5). The smallest increase occurs in Costa Rica, whereas the highest increase occurs in Norway. The latter means that, relative to the unadjusted gap, in Norway, the gap in math achievement could be up to 12 times wider in favor of boys in the absence of variation in student effort. The size of the effort-unadjusted gap in Norway is about 0.3 points in favor of girls, while after adjustment, girls fall behind boys by about 3.4 points, which represents a difference of about 3.7 points between the two gaps (see table 5). Another meaningful change occurs in Qatar. Before the adjustment, the gap is about 12 points in favor of girls, but after effort-adjustment, it becomes 8.6 points in favor of boys, which represents a widening of the gap by roughly 20.5 points, or 172 percent, favoring boys (see table 5).

In contrast, only in 7 out of 55 countries, the gap in math achievement becomes lower in the absence of student effort heterogeneity, according to figure 1. The decrease in the gap ranges from 12 to up to 76 percent (see table 5). The smallest decline occurs in the province of Macao in China, whereas the highest decline occurs in Thailand. In the latter case, the size of the effort-unadjusted gap is about 1.9 points in favor of girls, and after adjustment, its size is about 0.5 points, which represents a reduction of 1.4 points (or 75 percent) in the math achievement gap.

We obtain similar results when we look at the change in the science achievement gap in figure 2. In 43 out of 55 countries, the widening of the gap ranges from 3.5 percent up to 645 percent. The smallest increase in the science gap occurs in Costa Rica, whereas the highest increase occurs in Iceland (see table 6). The latter means that in Iceland, the gap becomes about 6 times wider after effort-adjustment, relative to the unadjusted gap. The effort-unadjusted gap in Iceland

is roughly 0.8 points in favor of girls, whereas after adjustment, girls fall behind boys by roughly 4.5 points, which represents a widening of the gap of about 5 points (see table 6).

When we analyze the percentage change in the reading achievement gap (see figure 3), the results are very different from those in math and science since most countries now appear in the green color. In 54 out of 55 countries, the reading achievement gap in the absence of variation in student effort narrows from 2 to up to 49 percent (see table 7). The smallest reduction of the gap occurs in the Dominican Republic, whereas the highest reduction occurs in Qatar. In the latter country, the effort-unadjusted reading gap is about 54 points in favor of girls; after adjustment, it is about 27 points. Although the effort-adjusted gap in Qatar still favors girls, the gap experiences a reduction of roughly 26 points, or 49 percent, favoring boys relative to the unadjusted. Only in Peru, the reading gap widens by 38 percent in the absence of student effort variation.

Overall, in most PISA countries that took the computer assessment, the gender achievement gap in math and science could be up to 6 and 11 times wider in favor of boys, respectively, in the absence of variation in student effort. Surprisingly, this widening of the gap in these two subjects is the highest among some of the top 20 gender-equal countries, according to the 2015 Global Gender Gap Index (GGGI), such as Iceland, Norway, Sweden, Netherlands, Latvia, and France. In contrast, the gender gap in reading could narrow up to 49 percent in favor of boys in the absence of variation in student effort. Our findings are consistent with Soland (2018) and Soland (2019), who find that the male-female gap in math is more sensitive to test effort compared to the reading gap.

Finally, we analyze the relationship between the 2015 Global Gender Gap Index (GGGI) and the percentage change of the gap after adjusting for student effort and we find that, relative to

the original gap, the effort-adjusted gap is even wider in countries with a higher gender equality as measured in the GGGI. More gender-equal countries tend to have a higher widening of the gap in math and science (correlations of -0.4 and -0,35, respectively), relative to the original gap.

6. Conclusions

In this paper, we use data from PISA 2015, a triannual survey that evaluates 15-year-old students from 74 countries in math, reading, and science to study the effect of student effort on cross-country differences in performance as well as within-country gender gaps in achievement. We restrict our sample to the countries which take the computer-based test and use innovative measures of effort based on rapid-guessing on the test and item non-response on the survey.

Altogether, our effort measures, on average, explain between 36 and 40 percent of the variation in test scores across countries. Our results also suggest that the estimated gender achievement gap in math and science could be up to 12 and 6 times wider, respectively, in favor of boys in the absence of variation in student effort. The gap in these two subjects widens in most of the countries in our sample, even among some of the top 20 gender-equal countries according to the GGGI in 2015, such as Iceland, Norway, Sweden, Netherlands, Latvia, and France. In contrast, the estimated gender gap in reading could narrow up to 49 percent in favor of boys. Our results highlight the importance of accounting for student effort to understand not only cross-country differences in performance but also variations in the measurement of the achievement gaps across nations.

References

- Akyol, Ş P., Krishna, K., & Wang, J. (2018). Taking PISA seriously: How accurate are low stakes exams? *Taking PISA Seriously: How Accurate are Low Stakes Exams?*, Retrieved from <https://www.nber.org/papers/w24930>
- Anaya, L., Stafford, F., & Zamarro, G. (2017). Gender gaps in math performance, perceived mathematical ability and college STEM education: The role of parental occupation. *EDRE Working Paper, 2017-21* doi:10.2139/ssrn.3068971
- Balart, P., & Oosterveen, M. (2019). Females show more sustained performance during test-taking than males. *Nature Communications, 10*(1), 3798. doi:10.1038/s41467-019-11691-y
- Boe, E. E., May, H., & Boruch, R. F. (2002). Student task persistence in the third international mathematics and science study: A major source of achievement differences at the national, classroom, and student levels. Retrieved from <https://eric.ed.gov/?id=ED478493>
- Borghans, L., & Schils, T. (2012). The leaning tower of pisa: Decomposing achievement test scores into cognitive and noncognitive components. *Unpublished Manuscript*, Retrieved from https://pdfs.semanticscholar.org/add9/e3d2a408bf1758e5cb3774c91e7f26b8d0b9.pdf?_ga=2.134769521.1849659902.1599322793-289877629.1599322793
- Ceci, S. J., Ginther, D. K., Kahn, S., & Williams, W. M. (2014). Women in academic science: A changing landscape. *Psychol Sci Public Interest, 15*(3), 75-141. doi:10.1177/1529100614541236

- Debeer, D., Buchholz, J., Hartig, J., & Janssen, R. (2014). Student, school, and country differences in sustained test-taking effort in the 2009 PISA reading assessment. *Journal of Educational and Behavioral Statistics*, 39(6), 502-523. doi:10.3102/1076998614558485
- Demars, C. E. (2007). Changes in rapid-guessing behavior over a series of assessments. *Educational Assessment*, 12(1), 23-45. doi:10.1080/10627190709336946
- DeMars, C. E., Bashkov, B. M., & Socha, A. B. (2013). The role of gender in test-taking motivation under low-stakes conditions. *Research & Practice in Assessment*, 8, 69-82. Retrieved from <http://www.rpajournal.com/dev/wp-content/uploads/2013/11/A4.pdf>
- Kane, M. (2006). Content-related validity evidence in test development. (pp. 131-153). Mahwah, NJ, US: Lawrence Erlbaum Associates Publishers. Retrieved from <https://psycnet.apa.org/record/2006-01815-007>
- Kong, X. J., Wise, S. L., & Bhola, D. S. (2007). Setting the response time threshold parameter to differentiate solution behavior from rapid-guessing behavior. *Educational and Psychological Measurement*, 67(4), 606-619. doi:10.1177/0013164406294779
- Michaelides, M. P., Ivanova, M., & Nicolaou, C. (2020). The relationship between response-time effort and accuracy in PISA science multiple choice items. *International Journal of Testing*, 1-19. doi:10.1080/15305058.2019.1706529
- Nix, S., Perez-Felkner, L., & Thomas, K. (2015). Perceived mathematical ability under challenge: A longitudinal perspective on sex segregation among STEM degree fields. *Frontiers in Psychology*, 6, 530. doi:10.3389/fpsyg.2015.00530

- OECD. (2017). PISA 2015 technical report. Retrieved from https://www.oecd.org/pisa/data/2015-technical-report/PISA2015_TechRep_Final.pdf
- Perez-Felkner, L., Nix, S., & Thomas, K. (2017). Gendered pathways: How mathematics ability beliefs shape secondary and postsecondary course and degree field choices. *Frontiers in Psychology*, 8, 386. doi:10.3389/fpsyg.2017.00386
- Schnipke, D. L. (1995). Assessing speededness in computer-based tests using item response times. Retrieved from <https://files.eric.ed.gov/fulltext/ED383742.pdf>
- Schnipke, D. L., & Scrams, D. J. (1997). Modeling item response times with a two-state mixture model: A new method of measuring speededness. *Journal of Educational Measurement*, 34(3), 213-232. Retrieved from <http://www.jstor.org/stable/1435443>
- Setzer, J. C., Wise, S. L., van, d. H., & Ling, G. (2013). An investigation of examinee test-taking effort on a large-scale assessment. *Applied Measurement in Education*, 26(1), 34-49. doi:10.1080/08957347.2013.739453
- Soland, J. (2018). The achievement gap or the engagement gap?investigating the sensitivity of gaps estimates to test motivation. *Applied Measurement in Education*, 31(4), 312-323. doi:10.1080/08957347.2018.1495213
- Soland, J. (2019). Are achievement gap estimates biased by differential student test effort? putting an important policy metric to the test. *Teachers College Record*, 120(12) Retrieved from <https://www.nwea.org/resource-library/research/are-achievement-gap-estimates-biased-by-differential-student-test-effort-3>

- Soland, J., Jensen, N., Keys, T. D., Bi, S. Z., & Wolk, E. (2019). Are test and academic disengagement related? implications for measurement and practice. *Educational Assessment, 24*(2), 119-134. doi:10.1080/10627197.2019.1575723
- Swerdzewski, P. J., Harmes, J. C., & Finney, S. J. (2011). Two approaches for identifying low-motivated students in a low-stakes assessment context. *Applied Measurement in Education, 24*(2), 162-188. doi:10.1080/08957347.2011.555217
- Wise, S. L. (2006). An investigation of the differential effort received by items on a low-stakes computer-based test. *Applied Measurement in Education, 19*(2), 95-114.
doi:10.1207/s15324818ame1902_2
- Wise, S. L., & DeMars, C. E. (2005). Low examinee effort in low-stakes assessment: Problems and potential solutions. *Educational Assessment, 10*(1), 1-17.
doi:10.1207/s15326977ea1001_1
- Wise, S. L., & DeMars, C. E. (2010). Examinee noneffort and the validity of program assessment results. *Educational Assessment, 15*(1), 27-41. doi:10.1080/10627191003673216
- Wise, S. L., & Kong, X. (2005). Response time effort: A new measure of examinee motivation in computer-based tests. *Applied Measurement in Education, 18*(2), 163-183.
doi:10.1207/s15324818ame1802_2
- Wise, S. L., Pastor, D. A., & Kong, X. J. (2009). Correlates of rapid-guessing behavior in low-stakes testing: Implications for test development and measurement practice. *Applied Measurement in Education, 22*(2), 185-205. doi:10.1080/08957340902754650

Wise, S. L., Soland, J., & Bo, Y. (2020). The (non)impact of differential test taker engagement on aggregated scores. *International Journal of Testing*, 20(1), 57-77.

doi:10.1080/15305058.2019.1605999

Wise, S. L., & Ma, L. (2012). Setting response time thresholds for a CAT item pool: The normative threshold method. Retrieved from

<https://www.nwea.org/content/uploads/2012/04/Setting-Response-Time-Thresholds-for-a-CAT-Item-Pool.pdf>

Zamarro, G., Hitt, C., & Mendez, I. (2019). When students don't care: Reexamining international differences in achievement and student effort. *Journal of Human Capital*,

doi:10.1086/705799

Figure 1: Percentage change in the gender gap in math achievement

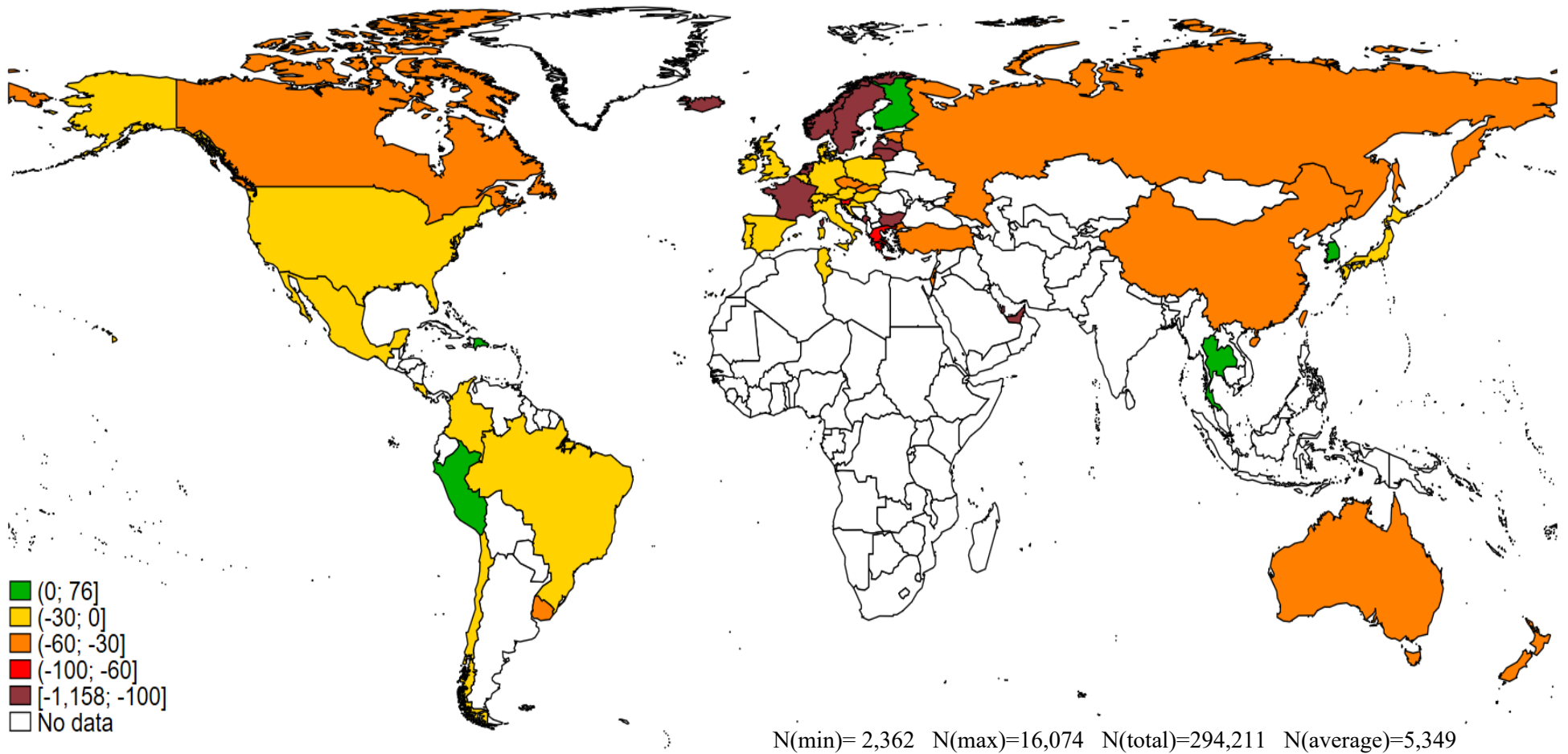


Figure 2: Percentage change in the gender gap in science achievement

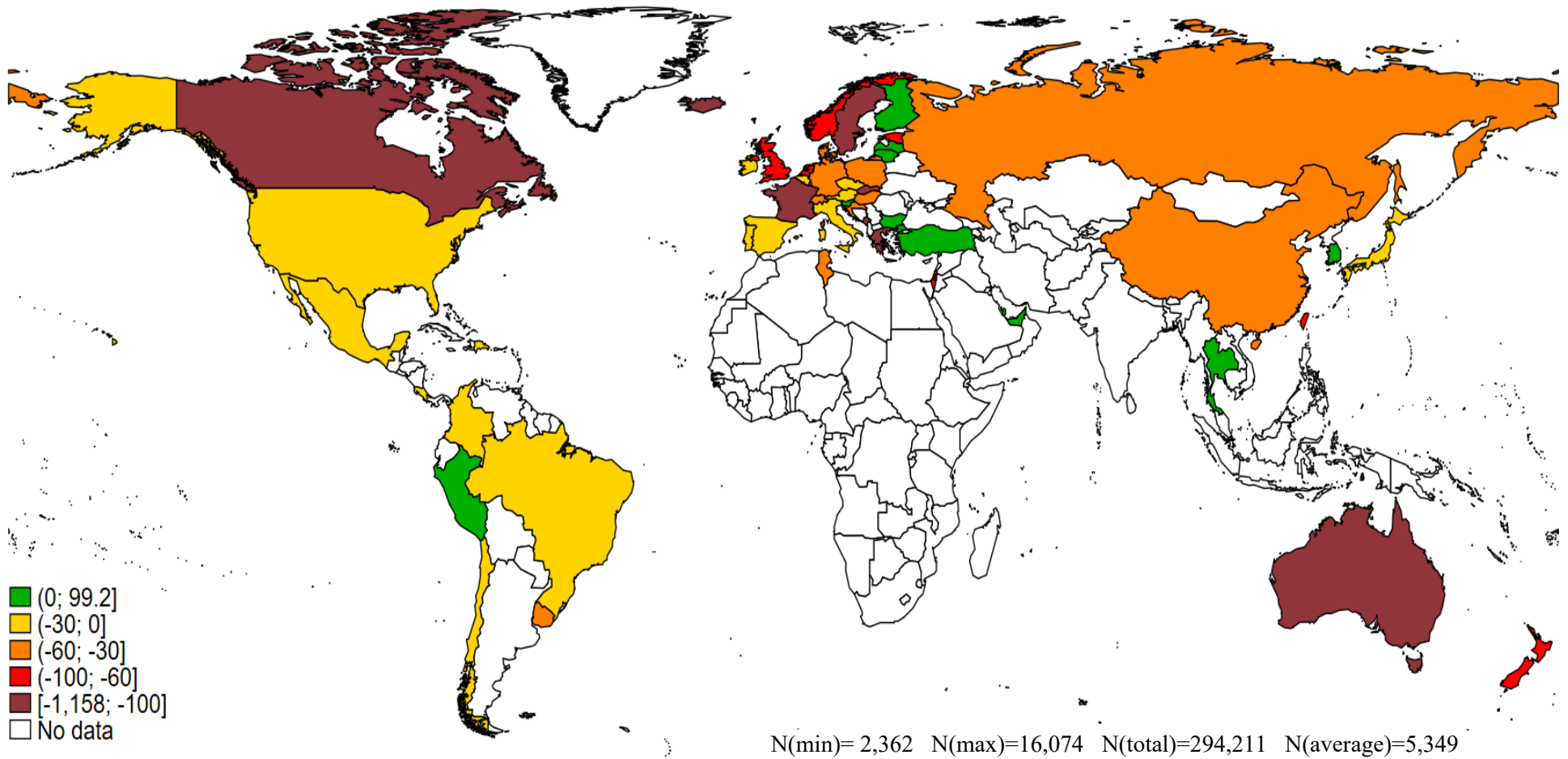


Figure 3: Percentage change in the gender gap in reading achievement

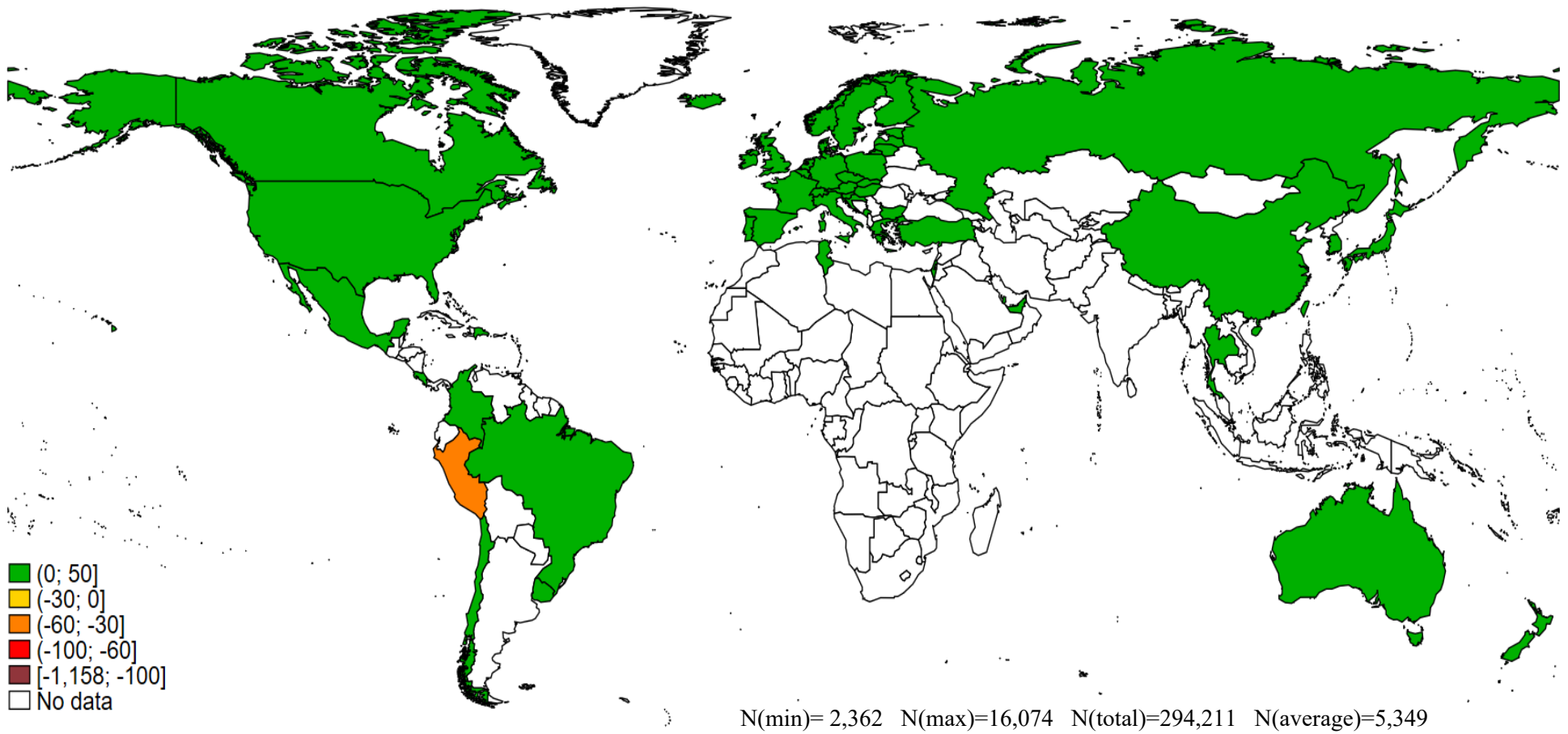


Table 1: Country names and abbreviations in PISA 2015

Abbreviation	Country Name	Abbreviation	Country Name
SGP	Singapore	ESP	Spain
JPN	Japan	LVA	Latvia
EST	Estonia	RUS	Russia
TAP	Chinese Taipei	LUX	Luxembourg
FIN	Finland	ITA	Italy
MAC	Macao	HUN	Hungary
CAN	Canada	LTU	Lithuania
HKG	Hong Kong	HRV	Croatia
QCH	B-S-J-G (China)	ISL	Iceland
KOR	Korea	ISR	Israel
NZL	New Zealand	SVK	Slovak Republic
SVN	Slovenia	GRC	Greece
AUS	Australia	CHL	Chile
GBR	United Kingdom	BGR	Bulgaria
DEU	Germany	ARE	Arab Emirates
NLD	Netherlands	URY	Uruguay
CHE	Switzerland	TUR	Turkey
IRL	Ireland	THA	Thailand
BEL	Belgium	CRI	Costa Rica
DNK	Denmark	QAT	Qatar
POL	Poland	COL	Colombia
PRT	Portugal	MEX	Mexico
NOR	Norway	MNE	Montenegro
USA	United States	BRA	Brazil
AUT	Austria	PER	Peru
FRA	France	TUN	Tunisia
SWE	Sweden	DOM	Dominican Republic
CZE	Czech Republic		

Table 2: Summary statistics of the variables of interest

Variable		Mean	Std. Dev.	Min	Max
Rapid-guessing % - test	Overall	4.6	7.6	0.0	100.0
	Between		2.3	1.0	16.0
	Within		7.1	-11.4	100.9
Item non-response % - survey	Overall	7.3	17.1	0.0	97.9
	Between		4.9	0.6	26.0
	Within		16.1	-18.7	104.5
Total time - test (min)	Overall	74.7	18.6	0.05	120.0
	Between		5.8	55.6	89.9
	Within		17.7	-8.8	133.9
Math score	Overall	470.9	97.9	113.4	826.3
	Between		50.7	331.3	557.6
	Within		82.5	87.5	807.4
Reading score	Overall	473.9	99.2	54.3	812.0
	Between		41.9	361.0	530.5
	Within		89.3	15.0	822.6
Science score	Overall	476.0	99.6	133.4	831.3
	Between		45.1	334.8	547.9
	Within		88.2	132.1	816.9
Observations	Overall student sample			N = 294,211	
	Between countries			n = 55	
	Within-country average sample			Tbar = 5,349.29	

Note: excludes observations with total time above 120 min

Table 3: Descriptive gender differences on the variables of interest

Average	Boys	Girls	Difference
Rapid-guessing % - test	5.1	4.1	1.0***
Item non-response % - survey	8.2	6.4	1.7***
Total time - test (min)	72.4	77.0	-4.6***
Math score	475.5	466.3	9.2***
Reading score	461.3	486.4	-25.1***
Science score	478.1	473.9	4.1***
Total observations	146,741	147,470	

Note: excludes observations with total time above 120 min; ***
p<0.01.

Table 4: Average estimated coefficients of the role of student effort on PISA test scores

	(1) Math	(2) Math	(3) Math	(4) Science	(5) Science	(6) Science	(7) Reading	(8) Reading	(9) Reading
Item non-response survey	-0.19*** (0.002)		-0.17*** (0.002)	-0.21*** (0.002)		-0.18*** (0.002)	-0.24*** (0.002)		-0.21*** (0.002)
Rapid-guessing test		-0.16*** (0.002)	-0.13*** (0.002)		-0.18*** (0.002)	-0.15*** (0.002)		-0.20*** (0.002)	-0.16*** (0.002)
Constant	0.05 (0.057)	0.06 (0.063)	0.05 (0.054)	0.03 (0.051)	0.04 (0.055)	0.03 (0.047)	0.03 (0.048)	0.04 (0.050)	0.03 (0.044)
Observations	296,832	294,211	294,211	296,832	294,211	294,211	296,832	294,211	294,211
Number of countries	55	55	55	55	55	55	55	55	55
R-squared within	0.04	0.03	0.06	0.05	0.04	0.07	0.06	0.04	0.09
R-squared overall	0.07	0.03	0.09	0.07	0.04	0.10	0.08	0.05	0.12
R-squared between	0.28	0.11	0.37	0.25	0.17	0.40	0.23	0.15	0.36
Min student sample size	2,368	2,362	2,362	2,368	2,362	2,362	2,368	2,362	2,362
Max student sample size	16,224	16,074	16,074	16,224	16,074	16,074	16,224	16,074	16,074
Average student sample size	5,397	5,349	5,349	5,397	5,349	5,349	5,397	5,349	5,349

Standard errors in parentheses

*** p<0.01, ** p<0.05, * p<0.1

All coefficients are standardized

Table 5: Effort-adjusted and unadjusted math scores and percentage change in the gap

Country	Effort-unadjusted score		Effort-adjusted score		Unadjusted gap	Adjusted gap	Absolute difference	% change in gap
	Girls	Boys	Girls	Boys				
Norway	501.8	501.5	514.6	518.0	0.3	-3.4	3.7	-1157.8
Latvia	485.8	486.2	492.8	496.7	-0.3	-3.9	3.5	-1112.0
Iceland	489.8	490.4	502.2	507.4	-0.6	-5.2	4.7	-815.9
Bulgaria	444.0	442.8	462.8	466.7	1.2	-3.9	5.1	-432.3
Montenegro	416.7	418.4	431.7	438.9	-1.7	-7.2	5.5	-321.7
Hong Kong	551.9	553.1	562.3	566.6	-1.2	-4.3	3.1	-261.3
Sweden	496.6	494.4	510.0	513.0	2.1	-3.0	5.1	-239.3
Lithuania	475.5	473.1	484.5	487.0	2.3	-2.5	4.8	-207.2
Qatar	407.5	395.6	427.6	436.2	11.9	-8.6	20.5	-172.0
United Arab Emirates	432.4	425.0	446.5	449.4	7.4	-3.0	10.4	-140.3
France	496.1	500.1	509.1	518.5	-4.1	-9.4	5.3	-131.6
Netherlands	518.9	521.6	527.4	532.8	-2.7	-5.4	2.7	-102.2
Greece	459.6	466.1	472.8	485.0	-6.5	-12.2	5.7	-86.7
Slovenia	496.7	501.0	505.2	512.6	-4.3	-7.4	3.1	-73.3
Taiwan	538.8	542.8	548.1	554.3	-4.1	-6.2	2.1	-52.3
B-S-J-G (China)	540.4	545.1	552.9	560.1	-4.7	-7.2	2.5	-52.0
Israel	469.6	477.9	487.5	499.9	-8.2	-12.4	4.2	-51.0
Slovak Republic	476.9	484.9	485.5	497.4	-8.0	-11.9	3.9	-48.5
Australia	481.6	486.2	493.3	499.9	-4.6	-6.6	2.1	-45.3
Estonia	518.4	525.2	525.7	535.5	-6.8	-9.8	3.0	-43.7
Turkey	413.7	421.7	426.5	437.9	-8.0	-11.3	3.3	-41.8
New Zealand	491.5	501.2	502.5	516.3	-9.7	-13.7	4.0	-41.0
Canada	500.1	509.0	510.5	523.0	-8.9	-12.5	3.6	-39.9
Czech Republic	501.8	509.1	510.6	520.1	-7.3	-9.6	2.3	-31.2
Russian Federation	491.4	499.1	502.6	512.6	-7.7	-10.0	2.3	-30.3
Uruguay	414.2	429.2	433.5	453.0	-15.0	-19.5	4.5	-29.7
Germany	502.3	522.3	531.6	557.2	-20.0	-25.7	5.7	-28.5
United Kingdom	487.8	498.2	499.3	512.6	-10.4	-13.3	2.9	-28.3
Luxembourg	480.4	494.3	494.4	511.9	-13.9	-17.5	3.6	-25.5

(continued on next page)

	Effort-unadjusted score		Effort-adjusted score					
Country	Girls	Boys	Girls	Boys	Unadjusted gap	Adjusted gap	Absolute difference	% change in gap
Denmark	496.3	507.9	508.6	522.9	-11.6	-14.3	2.7	-22.9
Switzerland	513.1	525.1	528.5	543.0	-12.0	-14.6	2.6	-21.9
Poland	499.0	511.4	508.5	523.5	-12.4	-15.0	2.6	-20.7
Tunisia	363.0	370.9	384.7	394.1	-7.8	-9.4	1.6	-20.3
Colombia	392.2	405.9	406.1	422.1	-13.7	-16.0	2.3	-17.1
Croatia	460.3	472.6	468.6	483.0	-12.3	-14.4	2.1	-16.8
Mexico	411.0	419.1	422.8	432.2	-8.1	-9.4	1.3	-16.6
United States	465.4	475.7	477.0	489.1	-10.4	-12.0	1.7	-16.2
Hungary	481.8	490.1	492.1	501.7	-8.3	-9.7	1.3	-16.1
Austria	487.3	512.5	496.8	526.0	-25.2	-29.2	4.0	-16.1
Chile	434.2	452.5	447.4	468.1	-18.2	-20.7	2.4	-13.2
Japan	526.0	541.2	534.6	551.8	-15.2	-17.1	1.9	-12.8
Ireland	494.9	511.7	505.2	524.1	-16.9	-18.9	2.0	-12.0
Portugal	475.6	487.0	482.2	494.9	-11.5	-12.7	1.2	-10.6
Spain	481.7	500.1	492.0	512.0	-18.4	-20.0	1.6	-8.4
Belgium	506.3	525.4	517.5	538.0	-19.1	-20.5	1.4	-7.1
Brazil	367.8	383.2	394.3	410.5	-15.4	-16.2	0.8	-5.4
Italy	489.7	510.9	499.4	521.7	-21.2	-22.3	1.1	-5.1
Costa Rica	394.8	412.0	408.7	426.4	-17.2	-17.7	0.5	-2.9
Macau	548.1	542.2	559.0	553.9	5.9	5.1	0.7	12.3
Dominican Republic	332.9	329.6	375.3	372.5	3.3	2.8	0.5	15.8
Peru	382.3	392.9	401.9	410.5	-10.6	-8.6	2.0	19.3
Finland	515.7	508.5	524.0	519.5	7.2	4.4	2.8	38.4
Singapore	559.6	555.8	568.5	566.6	3.7	1.9	1.8	48.8
Korea	528.2	521.7	537.1	534.0	6.5	3.1	3.4	52.8
Thailand	429.7	427.8	436.0	435.5	1.9	0.5	1.4	75.6

Note: excludes outliers in total time above 120 minutes

Table 6: Effort-adjusted and unadjusted science scores and percentage change in the gap

Country	Effort-unadjusted score		Effort-adjusted score		Unadjusted gap	Adjusted gap	Absolute difference	% change in gap
	Girls	Boys	Girls	Boys				
Iceland	475.5	474.7	489.5	494.0	0.8	-4.5	5.3	-645.2
France	500.9	499.8	515.8	520.7	1.1	-4.9	6.1	-540.2
Slovak Republic	466.8	468.2	476.7	482.4	-1.3	-5.8	4.4	-332.4
Greece	466.8	464.6	482.0	486.2	2.2	-4.2	6.4	-293.1
Australia	499.5	500.4	512.6	515.8	-0.9	-3.2	2.3	-270.2
Canada	515.7	517.5	527.6	533.4	-1.8	-5.8	4.0	-222.3
Hong Kong	528.4	526.4	540.3	542.0	1.9	-1.7	3.6	-187.1
Sweden	497.3	493.6	512.6	514.8	3.7	-2.2	5.8	-159.1
Montenegro	412.7	408.1	429.7	431.3	4.6	-1.6	6.2	-134.5
Norway	495.7	500.0	510.2	518.8	-4.4	-8.6	4.2	-96.6
United Kingdom	502.7	506.5	515.8	522.9	-3.8	-7.1	3.3	-87.5
Netherlands	514.6	518.4	524.4	531.2	-3.8	-6.9	3.1	-82.7
Singapore	546.6	549.2	556.9	561.6	-2.6	-4.7	2.1	-81.7
Israel	468.0	474.1	488.6	499.4	-6.1	-10.8	4.8	-78.6
Taiwan	529.8	533.3	540.6	546.6	-3.5	-6.0	2.5	-69.9
New Zealand	511.3	518.0	523.8	535.0	-6.6	-11.1	4.5	-67.7
Estonia	533.8	539.0	542.2	550.9	-5.3	-8.7	3.4	-64.1
Uruguay	434.0	443.2	455.9	470.1	-9.2	-14.2	5.1	-55.0
Switzerland	498.9	504.9	516.3	525.3	-6.0	-9.0	3.0	-50.2
Russian Federation	484.0	490.0	496.7	505.2	-6.0	-8.6	2.6	-43.7
Croatia	475.4	480.8	484.9	492.6	-5.4	-7.7	2.3	-43.7
Germany	508.4	523.3	541.1	562.5	-14.9	-21.3	6.4	-42.8
Poland	498.7	505.6	509.6	519.6	-7.0	-10.0	3.0	-42.4
Tunisia	384.1	388.4	408.7	414.7	-4.3	-6.1	1.8	-40.8
Luxembourg	479.9	490.4	495.9	510.5	-10.5	-14.6	4.1	-38.7
Denmark	487.2	495.3	501.3	512.4	-8.1	-11.1	3.0	-37.2
B-S-J-G (China)	525.5	533.4	540.0	550.7	-7.9	-10.7	2.8	-36.1
Hungary	483.4	488.1	495.1	501.3	-4.7	-6.2	1.5	-32.5
Austria	490.6	506.6	501.4	522.0	-16.0	-20.6	4.6	-28.4

(continued on next page)

Country	Effort-unadjusted score		Effort-adjusted score		Unadjusted gap	Adjusted gap	Absolute difference	% change in gap
	Girls	Boys	Girls	Boys				
Czech Republic	501.7	510.9	511.7	523.5	-9.3	-11.8	2.6	-27.8
Dominican Republic	333.8	335.9	382.3	385.0	-2.2	-2.7	0.6	-26.8
Ireland	496.8	508.2	508.7	522.3	-11.3	-13.6	2.3	-20.2
Brazil	396.9	401.5	426.4	432.0	-4.7	-5.6	0.9	-19.8
Colombia	419.4	432.8	435.3	451.4	-13.4	-16.1	2.7	-19.8
United States	493.3	503.6	506.6	518.8	-10.3	-12.2	1.9	-18.3
Chile	458.7	474.2	473.7	491.9	-15.5	-18.2	2.7	-17.6
Spain	492.2	502.6	504.0	516.3	-10.5	-12.3	1.8	-17.0
Mexico	417.4	427.3	430.8	442.3	-10.0	-11.5	1.5	-15.1
Japan	532.4	547.6	542.3	559.7	-15.2	-17.4	2.2	-14.7
Portugal	484.0	495.8	491.6	504.8	-11.8	-13.2	1.4	-11.7
Belgium	503.1	517.8	515.8	532.1	-14.8	-16.3	1.5	-10.2
Italy	484.3	501.8	495.4	514.1	-17.5	-18.7	1.2	-6.7
Costa Rica	413.0	430.2	428.8	446.6	-17.2	-17.8	0.6	-3.5
Macau	533.5	527.0	546.2	540.5	6.5	5.7	0.8	12.8
Finland	542.1	522.6	551.6	535.1	19.6	16.5	3.1	16.0
Peru	391.9	403.5	414.1	423.4	-11.6	-9.3	2.3	19.6
Thailand	437.7	429.6	444.9	438.4	8.1	6.5	1.6	19.9
Latvia	498.0	488.2	506.1	500.2	9.9	5.9	4.0	40.6
Bulgaria	456.1	441.8	477.3	468.8	14.3	8.5	5.8	40.6
Korea	522.0	512.8	532.2	527.0	9.2	5.2	4.0	43.3
United Arab Emirates	451.1	425.1	467.2	453.0	26.0	14.2	11.8	45.3
Slovenia	501.9	496.5	511.7	509.8	5.4	1.9	3.6	65.4
Lithuania	472.8	465.1	483.1	480.8	7.8	2.3	5.5	70.5
Turkey	426.4	421.8	441.2	440.3	4.7	0.9	3.8	81.2
Qatar	428.9	405.4	452.0	451.8	23.5	0.2	23.3	99.1

Note: excludes outliers in total time above 120 minutes

Table 7: Effort-adjusted and unadjusted reading scores and percentage change in the gap

Country	Effort-unadjusted score		Effort-adjusted score		Unadjusted gap	Adjusted gap	Absolute difference	% change in gap
	Girls	Boys	Girls	Boys				
Peru	402.6	395.7	427.8	418.2	6.9	9.6	2.7	-38.2
Dominican Republic	376.5	347.2	430.8	402.1	29.4	28.7	0.7	2.3
Macau	525.3	494.4	539.3	509.3	31.0	30.0	0.9	3.0
Costa Rica	436.0	419.3	453.9	437.8	16.7	16.1	0.6	3.8
Brazil	415.8	393.5	450.0	428.8	22.3	21.2	1.1	4.8
Thailand	433.4	404.4	441.5	414.3	29.0	27.2	1.8	6.3
Hungary	490.1	466.4	503.3	481.3	23.7	22.0	1.7	7.2
Finland	551.8	504.7	562.5	518.9	47.1	43.6	3.5	7.5
Australia	509.0	476.7	524.0	494.4	32.3	29.6	2.7	8.2
Tunisia	371.6	348.2	399.5	378.1	23.4	21.4	2.0	8.8
Italy	501.1	485.6	513.6	499.5	15.5	14.1	1.4	9.0
Slovenia	513.7	471.9	524.5	486.7	41.8	37.8	4.0	9.6
Croatia	502.3	475.0	512.9	488.2	27.3	24.6	2.7	9.7
Singapore	540.0	517.1	551.3	530.8	22.9	20.6	2.3	10.1
Taiwan	509.7	484.1	521.6	498.6	25.6	22.9	2.7	10.5
Portugal	494.0	479.3	502.6	489.4	14.7	13.1	1.6	10.6
Korea	539.7	499.0	551.1	514.8	40.7	36.3	4.4	10.8
Czech Republic	514.4	488.5	525.6	502.7	25.9	22.9	2.9	11.3
Poland	521.3	492.2	533.4	507.6	29.1	25.8	3.3	11.3
Latvia	510.9	470.7	519.8	484.2	40.2	35.6	4.5	11.3
United States	508.3	489.5	523.3	506.6	18.8	16.7	2.2	11.4
Russian Federation	508.8	482.9	523.2	500.3	25.9	22.9	3.0	11.6
Mexico	438.1	423.4	453.2	440.3	14.6	12.9	1.7	11.7
Spain	508.1	491.6	521.3	506.8	16.4	14.5	2.0	12.0
Norway	532.6	494.3	548.9	515.4	38.3	33.5	4.8	12.5
Switzerland	501.8	477.0	521.5	500.0	24.8	21.4	3.3	13.5
Hong Kong	545.2	516.0	558.5	533.3	29.2	25.2	4.0	13.7
Estonia	534.9	507.7	544.4	520.9	27.2	23.4	3.8	14.0
Bulgaria	459.5	413.1	483.6	443.7	46.4	39.8	6.6	14.1

(continued on next page)

	Effort-unadjusted score		Effort-adjusted score					
Country	Girls	Boys	Girls	Boys	Unadjusted gap	Adjusted gap	Absolute difference	% change in gap
Belgium	513.7	501.4	528.1	517.6	12.3	10.5	1.7	14.2
Slovak Republic	477.6	443.7	488.8	459.9	33.9	28.9	5.0	14.7
Iceland	502.8	462.6	518.5	484.3	40.2	34.2	6.0	14.9
Netherlands	523.5	500.6	534.3	514.9	22.9	19.4	3.5	15.2
New Zealand	526.7	495.1	540.9	514.4	31.6	26.4	5.1	16.3
Lithuania	485.3	447.4	496.8	465.1	37.9	31.7	6.2	16.3
Turkey	440.4	414.3	456.7	434.9	26.1	21.8	4.3	16.4
B-S-J-G (China)	516.9	498.1	532.8	517.1	18.8	15.7	3.1	16.6
Sweden	522.2	483.3	539.4	507.1	38.9	32.3	6.6	16.9
Denmark	502.4	483.7	518.2	503.0	18.7	15.2	3.4	18.4
Canada	526.7	501.9	540.1	519.9	24.8	20.2	4.6	18.5
France	521.3	488.7	538.0	512.3	32.5	25.7	6.8	21.0
United Kingdom	505.9	488.0	520.6	506.5	17.9	14.1	3.8	21.1
Montenegro	441.9	409.7	461.2	436.0	32.2	25.2	7.1	21.9
Japan	523.0	511.7	534.0	525.2	11.3	8.8	2.5	22.0
Israel	494.5	470.8	517.4	499.0	23.7	18.3	5.4	22.7
Ireland	526.6	515.1	539.7	530.9	11.5	8.9	2.6	22.7
Austria	498.9	476.2	511.1	493.5	22.7	17.5	5.2	22.9
Greece	494.3	463.7	511.1	487.8	30.6	23.3	7.2	23.7
Luxembourg	492.8	474.2	510.8	496.7	18.6	14.1	4.5	24.4
Colombia	443.8	431.7	461.6	452.5	12.1	9.1	3.0	24.9
Uruguay	450.6	428.0	475.5	458.6	22.7	16.9	5.7	25.3
United Arab Emirates	459.8	409.0	477.7	440.2	50.9	37.6	13.3	26.1
Chile	482.5	471.0	499.4	491.0	11.5	8.4	3.1	26.8
Germany	525.3	507.7	563.0	552.8	17.6	10.3	7.3	41.6
Qatar	428.5	374.8	454.1	426.8	53.7	27.4	26.3	49.0

Note: excludes outliers in total time above 120 minutes