

Do high stakes muddle the information from standardized tests? Evidence from Brazil's ENEM exam

Germán Reyes
Evan Riehl
Ruqing Xu

July 2023

ABSTRACT. This paper evaluates two criticisms of high-stakes college admission exams: 1) High-stakes exams help wealthy students “game the system” through test prep; and 2) This gaming is not informative about academic potential. To evaluate these claims, we exploit the introduction of Brazil’s national standardized admission exam, the ENEM. The staggered adoption of the ENEM by selective universities meant that high school seniors took either a low-stakes (from their perspective) school accountability test or a high-stakes exam that determined admission to top universities in their state, depending on their location and cohort. We find that test score gaps between high- and lower-income students expanded when selective universities adopted the ENEM in admissions, but the increase in exam stakes also made ENEM scores *more* informative for students’ college outcomes. Our results show that high-stakes admission exams increase socioeconomic inequality in access to selective colleges, but they also help universities identify students who are likely to succeed academically.

Reyes: briq Institute on Behavior & Inequality, Schaumburg-Lippe-Str. 5-9, 53113 Bonn, Germany (email: german.reyes@briq-institute.org). Riehl (corresponding author): Department of Economics, Cornell University, 266 Ives Hall, Ithaca, NY 14853 (email: eriehl@cornell.edu). Xu: Department of Economics, Cornell University, 449B Uris Hall, Ithaca, NY 14853 (email: rx24@cornell.edu).

1. INTRODUCTION

The sorting of students to colleges has important education and labor market consequences for individuals and for society (Hoxby, 2009; MacLeod et al., 2017; Chetty et al., 2020). Standardized admission exams are the primary—and, in some cases, the only—instrument that colleges use to select students. These exams are perceived as having very high stakes in countries with stratified college systems such as Brazil, China, South Korea, and the United States. Correspondingly, a large industry of courses, tutors, and materials has arisen in these countries to help students prepare for admission exams (Ramey and Ramey, 2010).

Concerns over test prep have led a growing number of U.S. colleges to reduce the importance of high-stakes exams in admissions in favor of lower-stakes admission signals such as high school grades or socioeconomic factors.¹ This movement is motivated by two common criticisms of high-stakes admission exams. First, critics argue that high-stakes exams help wealthy student “game the system” because they have greater access to test prep services (Buchmann et al., 2010). Second, critics contend that this gaming creates “bias” in scores in the sense that it is not informative about individuals’ academic potential (Soares, 2015).

In this paper, we empirically evaluate these two criticisms by asking how an increase in the stakes of a college admission exam impacts the distribution of scores and the exam’s informativeness for students’ potential to succeed in college. Our research questions are motivated by recent theoretical work showing that signaling concerns can make high-stakes exams more informative about an individual’s “gaming ability” and less informative about the “natural” signal of a student’s aptitude that would arise from a low-stakes exam (Frankel and Kartik, 2019). The degree to which gaming increases inequality in college access is unclear since there are many dimensions of inequality that contribute to test score gaps between high- and lower-income students. Further, as Frankel and Kartik note, it is unclear whether colleges would prefer to receive a signal with or without gaming. On the one hand, gaming may make the exam a worse measure of the material that students learned in high school. On the other hand, gaming could reflect effort, the capacity to learn, or other characteristics that help students succeed in college.

To examine how the stakes of an exam impact the distribution and informativeness of scores, we exploit a unique natural experiment in Brazil. From 2009–2017, Brazil’s system of highly-selective federal universities transitioned from their own admission exams to a common standardized test called the ENEM (*Exame Nacional do Ensino Médio*). Federal universities in different states varied in the timing at which they adopted the ENEM in admissions (Machado and Szerman, 2021; Mello, 2022). Yet the ENEM also serves a second

¹ For example, in April 2021, the University of California removed ACT/SAT scores from the formula that determines which students are eligible for admission to the UC system. The current formula is based only on students’ grades in certain high school courses.

purpose as a high school accountability test, so many high school seniors took the exam regardless of its role in college admissions. Since most Brazilian students attend college close to home, this meant that ENEM participants were either taking a low-stakes (from their perspective) school accountability test or a high-stakes exam that governed admission to nearby federal universities, depending on their location and cohort. We define a sample of 2009–2017 high school seniors who were likely to take the ENEM regardless of its stakes, and then use a difference-in-differences design that exploits changes in the exam’s stakes across states and cohorts.

To implement our design, we link administrative records from the ENEM exam to nationwide college and labor market data. The ENEM data include individuals’ responses to each exam question, which allows us to ask how the increase in exam stakes affected students’ overall scores in each subject as well as their performance on different types of exam questions. We use the 2010–2019 waves of Brazil’s higher education census to examine whether the higher-stakes ENEM scores were more or less informative for students’ college persistence and graduation outcomes. Lastly, we examine impacts on the informativeness of scores for labor market earnings using Brazil’s national employer-employee data for the years 2016–2018.

We have two main findings. First, test score gaps between high- and lower-income students widened when federal universities adopted the exam in admissions. The increase in the ENEM’s stakes expanded test score gaps between private and public high school students by roughly 10 percent (relative to the mean gaps in pre-adoption cohorts). This increase was driven by private school students earning higher scores on the high-stakes exam, with roughly similar point estimates in each subject test. Racial and other socioeconomic test score gaps expanded by a similar percentage. We show that these point estimates imply quantitatively significant increases in the selectivity of programs that private school students could gain admission to.

Second, the increase in exam stakes caused the ENEM scores to become *more* informative for students’ academic potential. Specifically, the adoption of the ENEM exam by federal universities increased the correlation coefficients between ENEM scores and students’ college persistence and graduation outcomes by roughly 10–30 percent, depending on the outcome measure. These correlation coefficients increased both overall and measured within college programs, which shows that our findings are due to an increase in the informativeness of scores rather than an impact on the programs students attended. We find some evidence that the high-stakes ENEM scores were also more informative for earnings outcomes, although our labor market data is measured too early in students’ careers to be conclusive. Overall, our results show that the high-stakes ENEM exam improved colleges’ ability to identify students who were likely to succeed.

To shed light on mechanisms, we show that private school students’ performance improved across a wide range of exam competencies on the higher-stakes test. For example, the adoption of the ENEM by federal universities increased private students’ probability of a correct answer in each of seven topic areas that are intended to align with high school math curricula (e.g., Algebra, Geometry, and Statistics). Thus the increase in exam stakes did not lead to narrowly-targeted improvements in student performance, as critics often contend. Rather, our finding of broad-based improvements in performance suggests that the high-stakes exam was a better measure of characteristics that benefit students academically, such as the willingness to exert effort, the ability to focus, or the capacity to learn new material.

In sum, our paper shows that there is a tradeoff in the movement to reduce the use of high-stakes exams in college admissions. Our findings show that high-stakes tests do shift the distribution of scores toward wealthier students, and thus a switch to lower-stakes admission signals can help selective colleges admit more socioeconomically diverse classes. But colleges also want to admit students who can succeed in their programs, and our results show that high-stakes exams help colleges identify such students.

Our paper relates most directly to work on the desirability of using standardized exams in college admissions (Amrein and Berliner, 2002; Rothstein, 2004; Scott-Clayton, 2012; Bettinger et al., 2013; Riehl, 2023). The standard approach to examining the desirability of admission exams is to correlate their scores with measures of college success.² This approach has inherent limitations in that it measures exam informativeness in a static admissions environment. We show that changes in admission criteria alter the informativeness of different instruments as prospective students respond to the new system. Our findings also inform the broader literature on how the design and implementation of admission exams affect inequality in college access (Hoxby and Turner, 2013; Bulman, 2015; Pallais, 2015; Goodman, 2016; Bhattacharya et al., 2017; Goodman et al., 2020; Reyes, 2023).³

Our paper also relates to research on the score effects of high- versus low-stakes testing. Prior research shows that high-stakes exams can increase gender test score gaps due to competitive pressures (Ors et al., 2013; Azmat et al., 2016; Cai et al., 2019) or distort incentives for learning (Jacob, 2005; Neal and Schanzenbach, 2010). Other work emphasizes that mental effort is necessary to perform well on standardized tests, and thus motivation to exert such effort can influence performance (Wise and DeMars, 2005; Finn, 2015). Several

² The College Board promotes its SAT exam by showing that scores are highly correlated with college persistence rates and first-year GPA (e.g., Westrick et al., 2019). On the other hand, Rothstein (2004) finds that SAT scores have limited predictive power once one controls for high school GPA and demographic characteristics, which questions their desirability as an admission instrument.

³ A related literature studies the merits of different school admission mechanisms, including centralized systems (Abdulkadiroğlu et al., 2005; Machado and Szerman, 2021), affirmative action (Durlauf, 2008; Bertrand et al., 2010; Bagde et al., 2016), and “percent plans” (Long, 2004; Kain et al., 2005; Niu and Tienda, 2010; Cullen et al., 2013; Daugherty et al., 2014; Kapor, 2015).

papers show that higher incentives can have heterogeneous impacts by inducing a higher effort response from individuals who have low intrinsic motivation to perform well (Duckworth et al., 2011; Segal, 2012; Jalava et al., 2015; Gneezy et al., 2019). Our main contribution to this work is to show how a change in the exam stakes impacts the informativeness of scores. Our paper is also new in showing that increasing the stakes of college admission exams can exacerbate income-based test score gaps.

Lastly, our results provide empirical evidence that informs the theoretical literature on muddled information (Fischer and Verrecchia, 2000; Frankel and Kartik, 2019, 2022). In these models, test scores reflect two distinct skills: the “natural action” (the test score without incentives) and “gaming ability” (the marginal cost of improving test scores as stakes increase). A key result in this literature is that as stakes rise, test scores become more informative about students’ gaming ability (Frankel and Kartik, 2019). However, these models are silent on whether gaming ability is informative for the individual’s likelihood of succeeding in college, which is a key outcome that colleges care about in determining admissions. Our findings show that, at least in the case of college admissions, test scores that include gaming ability are more informative than those that reflect only the natural action. Our question-level results suggest that gaming ability may largely reflect broad learning capacity or non-cognitive skills that are beneficial in college, such as work ethic.

The paper proceeds as follows. Section 2 provides background on Brazilian federal universities and the ENEM exam. Section 3 describes the data and identification strategy. Section 4 shows how the ENEM adoption by federal universities affected test score gaps and performance on different exam competencies. Section 5 shows how the increase in ENEM stakes affected the exam’s informativeness for college and labor market outcomes. Section 6 concludes.

2. INSTITUTIONAL BACKGROUND

2.1. Colleges and high schools in Brazil. The higher education system in Brazil is heavily privatized, but the most prestigious institutions tend to be in its system of *federal universities*. In 2009, there were 59 federal universities, with a presence in all of Brazil’s 27 states. Together, federal universities account for about 11 percent of total college enrollment. Brazil also has a system of 40 *state universities* managed by the governments of each state. Federal and state universities are tuition-free, highly-selective, and consistently top the national college rankings. The Brazilian higher education system additionally includes over 2,000 private universities and technical colleges that make up roughly 80 percent of total enrollment. While a handful of these private institutions are elite and selective, the majority are moderately selective or follow an open enrollment policy.⁴

⁴ Appendix Table A1 provides summary statistics on Brazilian high schools and colleges in 2009.

The situation is reversed at the secondary level, where private high schools represent a small but socioeconomically-advantaged share of enrollment. In 2009, 14 percent of secondary students attended a private high school, while 85 percent attended a public school managed by the state government.⁵ Importantly, private high school students are vastly overrepresented in the higher education system. They accounted for 40 percent of all incoming college students in 2009 and made up 47 percent of federal university enrollees (Appendix Table A1).

2.2. Federal university admissions and the ENEM exam. Admission to federal universities is highly competitive and relies exclusively on test scores from entrance exams. Before 2009, each federal university designed and administered its own admission test (known as *vestibular* exams). Thus it was burdensome for students to apply to more than one university as they had to prepare for multiple tests and travel to each school on a specific date to sit for the exam.

To centralize federal university admissions, the Ministry of Education developed a national standardized college admission exam called the ENEM (*Exame Nacional do Ensino Médio*). The ENEM exam was initially created in 1998 for the purpose of high school accountability. Between 1998 and 2008, the ENEM exam featured 63 multi-disciplinary questions, and the government published annual lists of school-mean scores.⁶ In 2009, the Ministry redesigned and expanded the exam so that it could serve as a tool for college admissions. The post-2009 ENEM exam resembles the ACT exam in the United States; it contains 180 questions across four distinct subject areas (math, language arts, natural sciences, and social sciences) along with a written essay. The exam spans two days of testing and is taken by over five million students each November, making it the second-largest admission test globally. As part of this effort, the Ministry also created a centralized admission platform called SISU (*Sistema de Seleção Unificada*), which allocates students to colleges on the basis of their preferences and ENEM scores.

Although the college admission version of the ENEM exam began in 2009, federal universities varied in the timing at which they switched from their institution-specific tests to the ENEM exam. The Ministry of Education provided financial incentives to adopt the ENEM, but universities had unilateral control over their admission methods and some were initially uncertain about the content of the new ENEM Machado and Szerman (2021).⁷ Thus some

⁵ Roughly 0.5 percent of Brazilian students attend a high school managed by the federal government (Appendix Table A1). We define “private high schools” to include both private and federal high schools since their students are similar in terms of socioeconomic status and achievement.

⁶ See, for example, *Lista do ENEM 2015: Notas das escolas*, *Globo*, October 5, 2016, available at: <https://especiais.g1.globo.com/educacao/enem/2015/enem-2015-medias-por-escola/>.

⁷ See Otero et al. (2021), Machado and Szerman (2021). and Mello (2022) for details on the implementation of the ENEM/SISU system and its adoption by universities.

federal universities began using the ENEM immediately in 2009, while others adopted five or more years later.⁸ The variation in the timing at which federal universities adopted the ENEM exam is the basis of our empirical strategy, as we describe in the next section.

3. DATA AND IDENTIFICATION

3.1. Data. Our base dataset includes administrative records on all individuals who took the ENEM exam in 2007–2017 (INEP, 2019a). This dataset is compiled by the National Institute of Educational Studies, or INEP (*Instituto Nacional de Estudos e Pesquisas Educacionais*). The data contains scores on each exam subject, demographic characteristics, and information on individuals’ high schools. The data also contains individuals’ responses to each exam question, which allows us to observe which questions individuals got right and wrong. Lastly, we observe information on the content of each question, including the learning objectives, the Item Response Theory (IRT) parameters, and the text of the question.

To measure longer-run outcomes, we link the ENEM data to two other administrative datasets at the individual level.⁹ First, we measure college outcomes by linking to INEP’s higher education census (*Censo da Educação Superior*) for the years 2010–2019 (INEP, 2022). This dataset contains information on the universe of students who were enrolled in the Brazilian higher education system in these years, including each student’s university, major, admission method, enrollment year, and graduation/drop-out outcome.

Second, we measure labor market outcomes by linking to Brazil’s employee-employer dataset, the RAIS (*Relação Anual de Informações Sociais*), for the years 2016–2018 (RAIS, 2022). The RAIS is maintained by the Ministry of Labor and covers the entire population of formal-sector workers in Brazil. Our main labor market outcome is an individual’s hourly wage, which we compute as average monthly earnings divided by average monthly contracted hours. Many of the ENEM participants in our sample were in college during the period of our RAIS data, and even those who had left college were still early in their careers. Thus our earnings outcomes may not capture the long-run returns to their college investments.

3.2. Sample. We begin by defining a sample with a consistent composition of ENEM exam takers over time. The total number of ENEM exam takers increased significantly after the exam was converted into a college admissions test in 2009, as illustrated by the black bars in Panel A of Figure 1. Since our goal is to examine how the increase in the exam’s stakes impacted the distribution of scores, we define a sample in which the number of test takers remained relatively constant over these years. For this, we take advantage of the fact that

⁸ Some state universities also adopted the ENEM as their admission test, but to this date, many still design and administer their own admission exams.

⁹ We linked the three administrative datasets at the individual level in a secure data room at INEP’s facilities in Brasília and extracted results for our analysis. See Appendix C.2 for details on the merge.

many high school students took the ENEM in their senior year regardless of its stakes because of the exam’s legacy as a high school accountability test.

Our analysis sample includes high school seniors at schools that met the criteria to be included in the government’s accountability reports in each year in 2005–2015. To define our sample, we use a dataset that contains school-level mean ENEM scores from 2005–2015, which were computed by INEP and distributed to federal and municipal agencies for publication (INEP, 2019b).¹⁰ Our analysis sample includes only ENEM exam takers who: 1) are in their last year of high school; and 2) attended a school that appears in the INEP school-level dataset in *each* year from 2005 to 2015. The red bars in Panel A of Figure 1 show that our analysis sample contains a small subset of all ENEM participants, but the number of exam takers in our sample remains relatively constant between 2007 and 2017.¹¹ Section 3.5 presents formal tests for balance in our sample.

Table 1 shows that our analysis sample is positively selected on socioeconomic status and academic performance relative to other ENEM test takers. This table reports mean demographic characteristics (Panel A), ENEM scores (Panel B), and college and labor market outcomes (Panel C) for 2009–2017 ENEM participants. Columns (A)–(C) show statistics for all ENEM exam takers, all high school seniors, and high school seniors in our analysis sample, respectively. Our sample contains roughly 2.5 million high school seniors, which is six percent of all ENEM test takers and 22 percent of all high school seniors. On average, students in our sample are four years younger than the typical ENEM participant, and they are roughly 10 percentage points (pp) more likely to be white and to have a college-educated parent. Relative to both the average test taker and the average high school senior, students in our sample score about 0.2–0.3 standard deviations (SD) higher on each of the ENEM subjects.¹²

Despite this positive selection, there is substantial inequality between private and public high school students in our sample. Columns (D)–(F) of Table 1 report statistics for private school students, public school students, and the private/public gap. 32 percent of students in our sample attended a private high school. Relative to public school students, private students were 26pp more likely to be white, 44pp more likely to have a college-educated mother, and 52pp more likely to come from a high-income family. Mean ENEM score gaps

¹⁰ Appendix C.3 provides details on the criteria to be included in the INEP report.

¹¹ Appendix Table A5 shows that our main results are robust to different sample selection criteria.

¹² ENEM scores, as reported to the public, are scaled to have a mean of 500 and a SD of 100 in the population of 2009 high school seniors who took the exam. Throughout the paper, we report ENEM scores in SD units relative to this population. For ENEM subject scores, our transformation is: Transformed score = (Scale score – 500)/100. Our transformation is different for writing and overall scores since they are on different scales. In all cases, a score of zero in our paper is equivalent to the performance of the average high school senior who took the ENEM in 2009, and a score of one is 1 SD higher within this population. These transformations preserve the comparability of test scores across cohorts. See Appendix C.1 for more details.

are on the order of 1 SD; the test score gap is largest in math, with private students scoring 1.4 SDs higher than public students on average.¹³ There is also substantial inequality in college and labor market outcomes. Private students were 27pp more likely to go to college and 15pp more likely to attend a federal university. Mean hourly wages during our data period are 68 percent higher for private students.

3.3. ENEM exam stakes. Our identification strategy exploits the gradual adoption of the ENEM exam by federal universities. The solid red line in Panel B of Figure 1 plots the proportion of all federal university enrollees in each year who were admitted using the ENEM exam. Although the college admissions version of the ENEM exam was first administered in November 2009, only 28 percent of federal university students nationwide were admitted using the ENEM in the following year.¹⁴ The proportion of federal university seats that were allocated using the ENEM grew over subsequent years as more institutions switched from their own tests to the ENEM, reaching a peak of 72 percent in 2016.

This gradual adoption created geographic variation in the stakes of the exam because Brazilian students typically attend college in their home state. The black dashed line in Panel B of Figure 1 plots the proportion of federal university enrollees who attended college in the state where they were born. On average, 81 percent of federal university students are from in state. Although there is evidence that the ENEM exam increased geographic mobility (Machado and Szerman, 2021), these effects were modest; the proportion of in-state students at federal universities remained above 80 percent throughout 2010–2018. Thus, the stakes of the ENEM exam varied across states and cohorts for students who wished to attend a federal university in their home state.

We use this variation to define two measures of ENEM stakes at the state \times year level. Our benchmark measure, which we denote by $\text{ProportionENEM}_{st}$, is a continuous variable that equals the proportion of new federal university enrollees in state s and year t who were admitted using the ENEM exam. Our continuous treatment variable has two advantages: it is simple to define, and it is more-powered because it includes all variation in ENEM adoption timing. For example, $\text{ProportionENEM}_{st}$ incorporates variation in ENEM adoption across federal universities within the same state, as well as variation in the use of the ENEM across programs within the same university.

Second, we define a binary treatment variable that equals one in years after each state “adopted” the ENEM exam. For this we follow research on tipping points (e.g., Card et al., 2008) in identifying structural breaks in the time series of federal universities’ use of the ENEM. For each state s , we regress an annual time series of the proportion of federal

¹³ For reference, the white/Black gap in the 2017 U.S. SAT math exam was 0.85 SDs (College Board, 2017).

¹⁴ Since the ENEM is administered in November, scores are used for admission to university cohorts that begin in the following calendar year.

university enrollees who were admitted using the ENEM on a linear trend break function for each possible candidate adoption year τ_s .¹⁵ We define the state’s *ENEM adoption year* as the value τ_s^* that yields the highest R^2 across these regressions. Our binary measure, which we denote by HighStakes_{st} , is an indicator for years equal to or after the state’s ENEM adoption year, τ_s^* . Our binary treatment variable allows us to present our results using event study graphs, and it helps to address potential concerns about two-way fixed effects models with treatment effect heterogeneity (discussed below).

Figure 2 shows the relationship between our continuous and binary measures of ENEM stakes. In Panel A, we categorize Brazil’s 27 states into ten groups based on their year of ENEM adoption, τ_s^* . The graph plots the mean of $\text{ProportionENEM}_{st}$ in these groups (y -axis) for each ENEM exam year (x -axis). In each group, the proportion of federal university students who were admitted using the ENEM increases sharply in the state’s ENEM adoption year. Panel B presents an event-study version of Panel A, in which the x -axis denotes years relative to the state’s ENEM adoption year. On average, the share of a state’s federal university admissions that used the ENEM exam increased by 56 percent in the adoption year, and this share remains at a high level in subsequent years. Appendix Table A2 shows the values of $\text{ProportionENEM}_{st}$ and HighStakes_{st} in each state and exam year.

3.4. Regression models. Our benchmark regression model is a two-way fixed effects specification estimated at the high school \times year level:

$$(1) \quad Y_{ht} = \gamma_{s(h)} + \gamma_t + \beta \text{ProportionENEM}_{s(h)t} + \epsilon_{ht}.$$

Y_{ht} is an average outcome for students who attended high school h and took the ENEM exam in year t . We include fixed effects for years, γ_t , and for the states in which each high school is located, $\gamma_{s(h)}$. The variable of interest is our continuous treatment variable, $\text{ProportionENEM}_{s(h)t}$, which measures the stakes of the ENEM exam in state $s(h)$ and cohort t . In alternate specifications, we replace $\text{ProportionENEM}_{s(h)t}$ with our binary treatment variable, $\text{HighStakes}_{s(h)t}$. We weight our regressions by the number of individuals in each ht cell to recover population estimates within our sample. Our benchmark regressions include high school seniors who took the college admissions version of the ENEM exam in 2009–2017,

¹⁵ Specifically, we estimate the following regression for each state s :

$$\text{ProportionENEM}_{st} = \delta_s^0 + \delta_s^1 \mathbb{1}\{t \geq \tau_s\} + \delta_s^2 \mathbb{1}\{t \geq \tau_s\}(t - \tau_s) + \delta_s^3 \mathbb{1}\{t < \tau_s\}(t - \tau_s) + \epsilon_{st},$$

where $\text{ProportionENEM}_{st}$ is our continuous treatment variable. We estimate this regression for all candidate adoption years $\tau_s \in \{2008, \dots, 2016\}$ and pick the value τ_s^* that yields the highest R^2 value. Lastly, we define our binary treatment variable to be $\text{HighStakes}_{st} = \mathbb{1}\{t \geq \tau_s^*\}$. We define one state (Sergipe) as a “never adopter” since the value of $\text{ProportionENEM}_{st}$ never exceeds 0.06.

which holds the structure of the ENEM exam fixed over time.¹⁶ We cluster standard errors at the state level.

The coefficient of interest, β , measures how outcomes changed in a school when the stakes of the ENEM exam increased. We estimate equation (1) separately for public and private high school students to examine how the increase in exam stakes affected scores in these two populations. In addition, we estimate regressions that fully interact the covariates in equation (1) with an indicator for private high schools, Private_h :

$$(2) \quad Y_{ht} = \gamma_{s(h)} + \gamma_t + \beta \text{ProportionENEM}_{s(h)t} + [\tilde{\gamma}_{s(h)} + \tilde{\gamma}_t + \beta^{\text{gap}} \text{ProportionENEM}_{s(h)t}] \text{Private}_h + \nu_{ht}.$$

The β^{gap} coefficient in equation (2) shows how the increase in exam stakes impacted the private/public gap in ENEM scores.

To examine the robustness of our results to concerns about two-way fixed effects models with treatment effect heterogeneity (De Chaisemartin and d’Haultfoeuille, 2020), we use a specification that restricts identification to clean comparisons based on states’ ENEM adoption years. Our approach follows Callaway and Sant’Anna (2021) in estimating separate treatment effects for each pair of ENEM adoption years, τ_s^* and $\tau_{s'}^*$, and then averaging the pairwise treatment effects to recover a single point estimate. For example, one of our pairs contains states that adopted the ENEM in 2010 and 2011, and we restrict the sample to students who took the exam in 2009–2010. In this pair, the 2010 adopters are our treated group since ENEM adoption “switches on” in 2010. The 2011 adopters are our control group since these states had not yet adopted the ENEM in these years. We define groups for all pairwise combinations of ENEM adoption years, and in each pair we restrict the sample to students who took the ENEM prior to the control group’s adoption year. We create a stacked dataset of these pairwise samples and estimate a version of equation (2) that uses our binary treatment variable, HighStakes_{st} , and includes interactions with dummies for the pairwise groups.¹⁷ The resulting β^{gap} coefficients are regression-weighted averages of the pairwise treatment effects. Appendix Table A4 shows the pairwise groups and the structure of our stacked dataset.

3.5. Identification assumptions and balance tests. Our identification relies on an assumption of parallel trends across Brazilian states. This assumption requires that the timing

¹⁶ In robustness analyses, we add in students from the same high schools who took the old 63-question version of the ENEM in 2007–2008.

¹⁷ Our stacked regression specification is:

$$(3) \quad Y_{htg} = \gamma_{s(h)g} + \gamma_{tg} + \beta \text{HighStakes}_{s(h)t} + [\tilde{\gamma}_{s(h)g} + \tilde{\gamma}_{tg} + \beta^{\text{gap}} \text{HighStakes}_{s(h)t}] \text{Private}_h + \epsilon_{htg}.$$

This specification differs from equation (2) in three ways: 1) the dataset is at the high school (h) \times year (t) \times pairwise group (g) level; 2) we include state \times group dummies, $\gamma_{s(h)g}$ and $\tilde{\gamma}_{s(h)g}$, and year \times group dummies, γ_{tg} and $\tilde{\gamma}_{tg}$; and 3) we replace $\text{ProportionENEM}_{s(h)t}$ with $\text{HighStakes}_{s(h)t}$.

of federal universities’ switch to the ENEM exam is unrelated to state-level trends in potential test score outcomes. There are two main ways in which this assumption could be violated. First, the increase in ENEM stakes may have affected the characteristics of students who took the exam within our high school senior sample. Second, the timing at which federal universities adopted the ENEM could be related to trends in student achievement.

Table 2 presents balance tests for the composition of our analysis sample. The dependent variables are the number of exam takers per high school (Panel A), the demographic characteristics of exam takers (Panel B), and an individual’s *predicted* overall score based on demographic characteristics (Panel C). Column (A) shows the mean of each dependent variable in cohorts prior to the state’s ENEM adoption year. Columns (B)–(D) present the β coefficients from equation (1), which we estimate separately for all schools, private schools, and public schools. Column (E) reports β^{gap} coefficients from equation (2), which are equivalent to the difference between the β coefficients in columns (C) and (D).

Our balance tests suggest that the stakes of the ENEM exam are not related to the composition of our sample. In Panel A of Table 2, we do not find significant effects on the number of ENEM takers in our full sample or in the private school subsample. We find that a 100 percentage point increase in the proportion of federal university students who were admitted using the ENEM exam is associated with a roughly 10 percent increase in the number of public school exam takers in our sample (column D); this effect is marginally significant in levels but not in logs. In Panel B, we find no systematic relationship between ENEM stakes and the age, race, parental education, or family income of exam takers in our sample. The increase in exam stakes is associated with a statistically significant but small decrease in the proportion of female exam takers. We cannot reject the hypothesis that the coefficients on all demographic characteristics are jointly equal to zero in any subsample (last row of Panel B). Similarly, we find small and insignificant effects on the demographic-based index of predicted test scores (Panel C). Overall, these tests suggest that the composition of our high school senior sample did not change significantly when the ENEM stakes increased.

Appendix Table A3 shows that there are no systematic differences between states with federal universities that were early- and late-adopters of the ENEM exam, which lends further support to our parallel trends assumption. For example, universities in the most populous state, São Paulo, adopted the ENEM immediately in 2009, while universities in the next two largest states, Minas Gerais and Rio de Janeiro, did not adopt until 2013. There are no clear trends in the size, selectivity, or student body characteristics of federal universities in early- vs. late-adopting states. This argues against the hypothesis that ENEM adoption was correlated with concurrent trends in student achievement.¹⁸

¹⁸ It is likely that the first universities to adopt the ENEM were those that saw the greatest benefits to doing so, e.g., those that found it especially costly to administer their own tests. But we do not have a strong

4. EXAM STAKES AND THE DISTRIBUTION OF SCORES

4.1. **Effects on test scores.** Table 3 presents our main results on how the stakes of the ENEM impacted mean test scores. Column (A) displays the mean private/public school gap in test scores in cohorts prior to the state’s ENEM adoption year. Columns (B)–(D) present β coefficients from equation (1) estimated separately by high school type. Column (E) displays β^{gap} coefficients from equation (2). Our dependent variables are individuals’ test scores in SD units. We examine scores on each of the four multiple choice tests (math, language arts, natural science, social science), average scores across these four core subjects, and scores on the writing component.

We find that the increase in the stakes of the ENEM exam led to a widening of private/public test score gaps. Private school students’ scores increased on the higher-stakes exam in each of the four core subjects (column C), with the largest effect in math (0.143 SDs). Public school students’ scores did not change significantly on the core subjects (column D). Thus test score gaps between private and public school students increased with the stakes of the exam (column E). Our point estimate implies that a 100 percentage point increase in the use of the ENEM by federal universities is associated with a 0.11 SD increase in the private/public test score gap on core subjects. This effect is nine percent of the mean test score gap in lower-stakes cohorts (column A). We also find that the increase in ENEM stakes widened the private/public gap in writing scores by 0.10 SDs.

Figure 3 shows that test score gaps typically widened in the first ENEM exam cohort after its adoption by federal universities. This figure presents estimates from an event study version of equation (2) using our binary treatment variable, HighStakes_{st} , and our stacked dataset of pairwise ENEM adoption years. This yields coefficients β_l^{gap} that show how the private/public score gap changed in each year l relative to the state’s ENEM adoption year, τ_s^* .¹⁹ In most subjects, we do not see significant pre-trends in the private/public score gap prior to the ENEM adoption year. In all subjects, we find increases in the private/public score gap in the first cohort after ENEM adoption on the order of 0.05 to 0.10 SDs. These wider gaps decline only slightly in subsequent cohorts; for example, the average gap between

prior on how variation in the ENEM’s benefits across universities would relate to trends in the achievement of local students.

¹⁹ Figure 3 plots β_l^{gap} coefficients from the high school (h) \times year (t) \times pairwise group (g) level regression

$$(4) \quad Y_{htg} = \gamma_{s(h)g} + \tilde{\gamma}_{s(h)g} \text{Private}_h + \gamma_{tg} + \tilde{\gamma}_{tg} \text{Private}_h + \sum_{l=-7}^7 [\beta_l + \beta_l^{\text{gap}} \text{Private}_h] \mathbb{1}\{t - \tau_{s(h)}^* = l\} + \epsilon_{htg},$$

where l denotes years relative to the state’s ENEM adoption year, $\tau_{s(h)}^*$. We include state \times group dummies, $\gamma_{s(h)g}$, year \times group dummies, γ_{tg} , and dummies for years l , $\mathbb{1}\{t - \tau_{s(h)}^* = l\}$, omitting $l = -1$. We interact all covariates with a dummy for private schools, Private_h , and plot the β_l^{gap} coefficients from $l = -4$ to 4.

private and public students on core subjects increased by 0.08 SDs in the year of ENEM adoption, and it was still 0.05 SDs higher measured four years later (Panel E).

The magnitudes of these estimates represent meaningful increases in private students’ chances of gaining admission to selective federal university programs. In 2016 data from the centralized admission system, the within-state standard deviation of cutoff scores for federal university programs is 0.52 SDs (in the test score units of our paper). Thus our estimate for the private/public gap in average ENEM scores ($\beta^{\text{gap}} = 0.11$ SD) is 21 percent of a standard deviation in the distribution of federal university program cutoffs. To put this in perspective, consider a private school student whose low-stakes ENEM score would have made them barely eligible for admission to a program at the 50th percentile of their state’s distribution of federal university programs. Our estimate of β^{gap} implies that this student’s high-stakes ENEM score would instead make them eligible for a program at the 58th percentile.

Appendix Table A6 shows that test score gaps by race, mother’s education, and family income also expanded on the higher stakes ENEM exam. The gap in average core subject scores between white/non-white (excluding Asian) students expanded by 0.06 SDs on the higher-stakes test. The gap between students having mothers with/without college education expanded by 0.08 SDs, while the gap by family income expanded by 0.09 SDs. These point estimates are similar to those for private/public school students as a percentage of average low-stakes score gaps. We find no significant effect on the score gaps by father’s education and gender.

4.2. Robustness tests. Table 4 examines the robustness of our results on private/public test score gaps. Column (A) reproduces our benchmark estimates of β^{gap} from column (E) of Table 3. Columns (B)–(F) present estimates of β^{gap} from alternative specifications.

Our results are robust to including demographic controls and using our alternative measure of exam stakes. In column (B) of Table 4, we estimate equation (2) including high school \times year averages of age, gender, and dummies for race, parental education, and family income bins. These demographic controls do not significantly alter our point estimates, which is consistent with the findings of our balance tests in Table 2. In column (C), we replace our continuous treatment variable, $\text{ProportionENEM}_{st}$, with our binary measure of ENEM stakes, HighStakes_{st} . This specification reduces the magnitudes of β^{gap} by about 50 percent in each subject, which is expected since $\text{ProportionENEM}_{st}$ increases by roughly 50 percent following a state’s adoption of the ENEM (Figure 2, Panel B). Yet we continue to find that the increase in exam stakes widened private/public test score gaps in each subject, and the coefficient for the average score remains statistically significant at $p < 0.05$.

Next, we examine the robustness of our benchmark estimates to potential concerns about two-way fixed effects models with treatment effect heterogeneity (De Chaisemartin and d’Haultfoeuille, 2020).

Columns (D)–(F) of Table 4 examine the robustness of our results to treatment effect heterogeneity using three different samples from our stacked dataset (see Section 3.3). Column (D) includes all pairwise combinations of ENEM adoption years that we can estimate using 2009–2017 exam takers. Column (E) focuses on a single pairwise comparison between the two most common ENEM adoption years, 2009 and 2013, which together account for 13 states (see Appendix Table A3). We require a pre-period for 2009 adopters to estimate a treatment effect in this pair, so this sample includes 2007–2012 test takers. In column (F), we include all 2007–2017 test takers and all pairwise combinations in our stacked dataset. Note that in columns (E)–(F), the sample includes two cohorts that took the old 63-question version of the ENEM exam (2007–2008), so these estimates may reflect effects of the ENEM redesign in addition to the impacts of the exam’s adoption by federal universities.²⁰ Appendix Table A4 shows the samples for each regression in columns (D)–(F) of Table 4.

The results in all of these specifications are similar to our benchmark estimates. The point estimates in column (D) are similar to those in column (C), which shows that our results are not impacted by restricting identification to clean pairwise comparisons. We continue to find positive and significant estimates of β^{gap} when we restrict to the simple “2×2” difference-in-differences model that compares 2009 vs. 2013 adopters (column E). Lastly, our results are similar in the full stacked dataset with 2007–2017 test takers (column F). The consistency of estimates across specifications shows that our results are not the result of averaging oppositely-signed treatment effects with negative weights.

Our results are also robust to controls for the nationwide rollout of affirmative action during our sample period. Many federal and state universities implemented reserved quotas for disadvantaged students during the late 2000s and early 2010s (Mello, 2022), which could have impacted the achievement of high school seniors through a motivational channel (Akhtari et al., 2020). To examine this possibility, we use the higher education census to compute the fraction of new university students in each state × year who enrolled through reserved quotas, and then add this variable as a control in our regressions. Appendix Table A7 shows that private/public test score gaps are not significantly related to the rollout of

²⁰ The 2007–2008 ENEM reported only a single core-component score plus a writing score. To define scores for each subject, we categorized the multiple choice questions into math, language arts, natural science, and social science, and then computed a separate score for each subject using the IRT parameters. Since the reference populations differ for the 2007–2008 and 2009–2017 exams, our regressions in columns (E)–(F) of Table 4 standardize scores to have mean 0 and SD 1 within each year of our sample. See Appendix C.1 for details.

affirmative action, and that our estimates for the impacts of ENEM stakes remain positive and significant with these controls.

4.3. Potential mechanisms. Our finding that private students earned higher scores on the high-stakes ENEM exam may be driven by several mechanisms. On the one hand, students may have exerted more effort while taking the exam. The typical private student had a better chance of gaining admission to federal universities than the typical public student, and thus private students had a stronger incentive to increase effort when the exam stakes increased. There is significant overlap between the distribution of private school ENEM scores and the distribution of admission cutoff scores for federal university programs, while the public school score distribution is shifted well to the left (Appendix Figure A1). Thus moderate increases in ENEM scores were unlikely to significantly affect the admission chances of public students at most federal universities, which can partly explain why we find no significant changes in their scores.²¹

On the other hand, the increase in private students' scores could reflect test prep, broadly defined. Brazil, like other countries with high-stakes admission tests, has a large industry of exam prep tutors, materials, and courses. In particular, there is an industry of for-profit exam prep courses called *cursinhos* that often run for six months or more (Fernandes, 2015). Private students have greater ability to pay for these services than public students, and it is also common for private high schools to incorporate college admission exam prep into their curricula. Yet test prep could also include other types of learning beyond the services offered in this industry. For example, students may have spent more time reviewing material they had learned in high school in advance of the high-stakes ENEM exam.

It is hard to distinguish between these mechanisms since we do not observe how individuals prepared for the exam, but we can examine whether the improvement in private students' ENEM performance was broad or narrowly-targeted. Critics of high stakes exams often argue that they encourage students to learn skills that raise their scores but are not useful outside of the exam. In this case, one might expect that the increase in private students' scores was driven by certain types of questions that are more amenable to test prep. Conversely, an improvement in performance across many types of exam questions would suggest that our results are driven by increases in effort or other types of broad-based learning.

4.4. Heterogeneity by question content. To explore these potential mechanisms, we estimate heterogeneity in the impacts of the high stakes ENEM exam across different types of questions. Our data includes students' responses to each exam question as well as information

²¹ Some individuals in our sample would have been eligible for affirmative action quotas that were reserved for public school students, but these quotas were not fully implemented at many federal universities until 2016, and they often included race and/or SES criteria in addition to a public school criterion.

on the skills that the questions measure. We use this data to estimate our regression model (2) separately for different groups of exam questions. In these regressions, the dataset is at the high school (h) \times year (t) \times question (q) level, and the dependent variable is the proportion of correct answers in each htq cell. We focus on math performance in the main text because it is the subject with the largest increase in ENEM score gaps (Table 3) and because math exams are often thought to be more “preppable” (Riehl and Welch, 2023). Appendix Table A8 presents results for language arts, natural science, and social science.

Panel A of Table 5 shows that the increase in ENEM stakes expanded the private/public gap in the proportion of correct answers, consistent with our results on scale scores (Table 3). This panel shows results pooling across all 405 math questions (9 years \times 45 questions/year). Column (C) reports the mean proportion of correct answers for public school students in cohorts prior to the state’s ENEM adoption year, and column (E) reports the mean private/public gap in these cohorts. The ENEM math exam is challenging for most students; the average public school student answered only 29.1 percent of the questions correctly. Private school students got 46.7 percent of the questions correct, and thus the private/public gap was 17.6pp. Columns (D) and (F) report the β and β^{gap} coefficients from equation (2). We normalize the β and β^{gap} estimates so that they represent a percentage change from the low-stakes means in columns (C) and (E).²² We do not find a significant change in the proportion of correct answers for public students, but the β^{gap} coefficient suggests that a 100 percentage point increase in ENEM adoption by federal universities expanded the private/public gap by 13.6 percent (a 2.4pp increase).

Panels B–C of Table 5 display estimates of β and β^{gap} for groups of questions that cover different math skills. The ENEM is designed to be closely related to high school curriculum, and the questions are based on a “reference matrix” of skills that educators think are important for students to know by the end of high school.²³ On the math exam, these skills are grouped into seven *topic areas* and 30 *competencies*. Topic areas include different branches of math such as algebra, geometry, and statistics. Competencies are specific abilities within each topic area such as identifying concepts, solving problems, and constructing arguments; these competencies are intended to test reasoning skills more than memorization. Panels B–C of Table 5 show results from estimating equation (2) separately for each topic area and

²² Even if the high-stakes exam induced the same amount of skill accumulation for each type of question, one would expect variation in the un-normalized β and β^{gap} coefficients. This is because skill accumulation has different impacts on the probability of a correct answer across questions that vary in difficulty. For example, an increase in the skill of an average exam taker would meaningfully increase the probability that they correctly answered a question that is of average difficulty, but it would have little impact on their probability of correctly answering very hard or very easy questions. Normalizing the β and β^{gap} by the mean probability of a correct answer helps to make the coefficients more comparable in terms of the underlying skill accumulation.

²³ See: https://download.inep.gov.br/download/enem/matriz_referencia.pdf (accessed in June 2023).

competency. For brevity, Panel C reports only competencies with the five largest and five smallest values of β^{gap} as a percentage of the mean gap.

We find increases in the private/public gap in nearly all math skills. The β^{gap} coefficients are positive and statistically significant at $p < 0.10$ in all seven topic areas (Panel B). The estimates at the competency level are less-powered since these regressions typically include only 10–15 questions across all years, but the β^{gap} coefficients are positive and economically meaningful in all but one of the 30 competencies.

Despite the uniformly positive impacts, there is significant variation in the β^{gap} coefficients. At the topic area level, these estimates range from a 7.6 percent increase in the private/public gap in algebra to a 17.7 percent increase in questions on interpreting data. There is even more variation at the competency level, with coefficients as high as 27.5 percent of the mean private/public gap. We reject equality of the β^{gap} coefficients at both the topic area ($p = 0.060$) and competency ($p < 0.001$) levels, suggesting that the variation in these estimates is not solely driven by statistical noise. We find a similar pattern of results for the language arts, natural science, and social science exams, i.e., increases in the private/public gap in most topics areas but significant variation in the β^{gap} coefficients (Appendix Table A8).

The results in Table 5 show that the increase in ENEM stakes induced a broad-based improvement in private students’ performance. In particular, private students performed better across a wide range of math skills that ENEM designers think are important for high school graduates to know. This suggests that our results are at least partly driven by mechanisms that would cause students to perform better on *all* exam questions, such as increased effort or broad learning.

Yet the variation in the β^{gap} coefficients leaves opens the possibility of some skill-specific test prep. To explore this possibility further, Appendix Table 5 presents estimates of β^{gap} for math questions that are related to topics in an ENEM study guide created by the test prep company *Me Salva!*.²⁴ This study guide contains formulas, definitions, and other problem solving tips for questions that are likely to appear on the ENEM exam. We match topics in the *Me Salva!* study guide to the ENEM math questions using text analysis and then estimate equation (2) separately for different study guide topics.²⁵ Some questions that are related to the *Me Salva!* topics have particularly large values of β^{gap} , including questions about proportions that use the “Rule of Three” and questions that require knowledge of how to compute the median of a set with an even number of elements. In one of our two

²⁴ The study guide is called *The Approved Book: One topic per day to pass the ENEM*, and is available at: <https://cdn.mesalva.com/uploads/medium/attachment/MS2018-livro-do-aprovado.pdf> (accessed in June 2023).

²⁵ See Appendix C.4 for details on our match of *Me Salva!* topics and ENEM questions.

matching strategies, we find that the average β^{gap} coefficient is larger for questions that are covered by the study guide than for questions on uncovered topics. But in both strategies, we find significant improvements in private students' performance across both covered and uncovered questions. This suggests that test prep for specific exam skills plays at most a small role in our findings.²⁶

Regardless of the nature of private students' improved test performance, it is *ex ante* unclear whether the effort or skill accumulation that drove these improvements is useful beyond raising individuals' admission scores. To shed light on this, we now turn to our analysis of the predictive power of the ENEM for college and labor market outcomes.

5. EXAM STAKES AND THE INFORMATIVENESS OF SCORES

5.1. Potential information channels. There are several potential channels through which a change in exam stakes could impact the informativeness of the scores, measured by the correlation between one's test scores and college and labor market outcomes.

The high-stakes test could be more informative through a correlation channel. The correlation could be with student's SES, i.e., the high-stakes test scores may be more predictive of student outcomes simply because they are more correlated with SES, which also tends to benefit students in college and labor market. For example, wealthy students may have greater access to test prep services, and family wealth may also help students succeed in college. On the other hand, the correlation could also exist with student's ability. Students who have higher test ability on the low-stakes exam may engage in more test prep when the stakes increase. Additionally, students with higher ability to succeed in college may have stronger incentives to prep for the high-stakes exam. In other words, higher-ability students could self-select into prepping for the high-stakes exam and thus increase the informativeness of the scores. It is hard to distinguish between the ability correlation and the SES correlation since we do not observe a prior measure of ability that is independent of SES. In the following sections, we will consider them together as the *correlation channel*, since the informativeness of the test is improved merely through increased correlations with other student characteristics.

In contrast, a second channel, which we call *skill accumulation*, captures the direct effect of studying for the high-stakes exam on college success. This channel is operative if the test skills that individuals accumulate from prepping directly improve their college outcomes. The sign of this channel is *a priori* less clear. High-stakes exams are often criticized for creating incentives to engage in test-oriented learning that is not useful outside the exam.

²⁶ Consistent with some role for test prep, Appendix Figure A2 shows that state-level Google searches for "ENEM" and the online prep company "*Descomplica*" increased when federal universities adopted the ENEM.

Thus, it is possible that this channel has a null or negative effect if the skills learned from prepping are not useful in college or crowd out other useful learning.

Finally, high-stakes exam scores may have more predictive validity if they increase the “match quality” between individuals and college programs. This channel arises because an individual’s test score could have a causal effect on their outcome through its influence on which college and/or major they attend. Thus, high-stakes exam scores may be more informative for college success if the distribution of these scores leads to better student/college matches, e.g., on the basis of academic preparation. To distinguish between this channel and those mentioned above, we follow the standard practice that testing agencies use to measure predictive validity. Specifically, in addition to estimating raw correlations between test scores and outcomes, our empirical analysis also estimates correlations after de-meaning each variable *within* college programs.

5.2. Results. To assess the informativeness of the ENEM exam, we use a dataset of state-by-year correlation coefficients between the ENEM exam scores and the exam takers’ college and labor market outcomes. Then, we estimate our benchmark DD regression (1) at the state-year level with these correlation coefficients as the dependent variable:

$$(5) \quad Y_{st} = \gamma_s + \gamma_t + \beta \text{ProportionENEM}_{st} + \epsilon_{st}.$$

Y_{st} is the correlation coefficient between the students’ scores in the ENEM and the students’ outcomes in state s and year t . We include fixed effects for states and years. The variable of interest is our continuous treatment variable, $\text{ProportionENEM}_{st}$, which measures the stakes of the ENEM exam in state s and cohort t . In alternate specifications, we replace $\text{ProportionENEM}_{st}$ with our binary treatment variable, HighStakes_{st} . We weight our regressions by the number of individuals in each st cell to recover population estimates within our sample.

Table 6 presents the results on how the increase in exam stakes impacted the informativeness of test scores. We use the ENEM overall score (core subjects) for all estimates in this table. Column (A) displays the mean correlation coefficient in cohorts prior to the state’s ENEM adoption year. Columns (B) and (C) present β coefficients from equation (5) using raw correlations and correlations after de-meaning each variable within college programs, respectively. Columns (D) and (E) present β coefficients for raw and within-program correlations using our binary treatment variable, HighStakes_{st} .

We find that the increase in the stakes of the ENEM exam improved the informativeness of the exam scores. The raw correlations between ENEM scores and student outcomes increased for all outcomes except for the hourly wage in log terms (column B). The within-program

correlations also increased by similar magnitudes (column C), indicating that the “match quality” under the high-stakes tests is unlikely to be the reason of improved informativeness.

Our point estimates imply economically significant improvements of the informativeness of the high-stakes exam. For example, a 100 percentage point increase in the use of the ENEM by federal universities is associated with a 0.03 increase in the correlation between exam scores and college persistence after 3 years, which is 23 percent of the correlation found in the low-stakes exam. Similarly, we find a 62 percent increase in the correlation between scores and program completion, a 13 percent increase in the correlation between scores and federal university enrollment, and 23 percent in the correlation between scores and hourly wage.

Figure 4 shows that the increase of informativeness was relatively uniform across four subjects. The translucent area represents the low-stakes mean correlations (similar to column A of Table 6) and the opaque area represents DD coefficients of the correlation between subject scores and outcomes (similar to column B of Table 6). All four subject scores experienced similar improvements in the informativeness for all three outcomes. The informativeness of writing scores increased moderately for earning a college degree and college persistence but not for hourly wage.

Our finding that test prep increases the informativeness of scores can explain why many colleges around the world use high-stakes tests for admissions. Similarly, our finding sheds light on why many U.S. colleges prefer to use SAT or ACT “superscores” (Goodman et al., 2020), which are based on the maximum of each subject score across all of the student’s test attempts. Relative to average scores or scores from the first attempt, superscores are more correlated with the amount of test prep, which may increase the informativeness of scores for outcomes that colleges care about.

While the results suggest that either the correlation or the skill accumulation channels must be at work, the score-level results do not distinguish between these channels. In results yet to be written up, we find evidence that the results in Table 6 and Figure 4 are primarily driven by a correlational channel, i.e., the high-stakes ENEM scores were more correlated with abilities or socioeconomic characteristics that help students succeed in college.

6. CONCLUSION

This paper exploited a natural experiment in Brazil to examine how the use of high-stakes standardized admission exams affects inequality in exam scores and their information content. From 2009–2017, Brazil’s system of highly-selective federal universities transitioned from institution-specific admission exams to a national standardized test called the ENEM. Since the ENEM exam was also used for high school accountability, many students took the exam regardless of its role in college admissions. This setting allowed us to focus on

a comparable population of students and ask how the distribution and informativeness of exam scores changed as the ENEM's role in college admission grew.

We found that socioeconomic test score gaps increased by roughly 10 percent when federal universities began to use the ENEM for college admissions. High-income students earned better scores on the higher-stakes exam, suggesting that they exerted more effort or engaged in test prep to boost their performance. This shows that high-stakes exams can give wealthy students a leg up in college admissions, consistent with a common criticism of these exams.

Yet we also found that the adoption of the ENEM exam by federal universities made the exam scores *more* informative for students' college outcomes. Consistent with this increase in informativeness, we found that the performance of high-income students improved across a wide range of subjects and question types, suggesting that any test prep was not narrowly-targeted. This shows that high-stakes exams provide information on individual characteristics that help students succeed academically, such as the capacity to learn new material.

Increasingly, U.S. colleges are reducing their reliance on high-stakes exams in favor of lower-stakes admission signals like high school grades. Our findings show that this change will help them to diversify their student bodies. But if these schools wish to maintain their average graduation rates, our results show that they must also find other ways of identifying students who are likely to succeed.

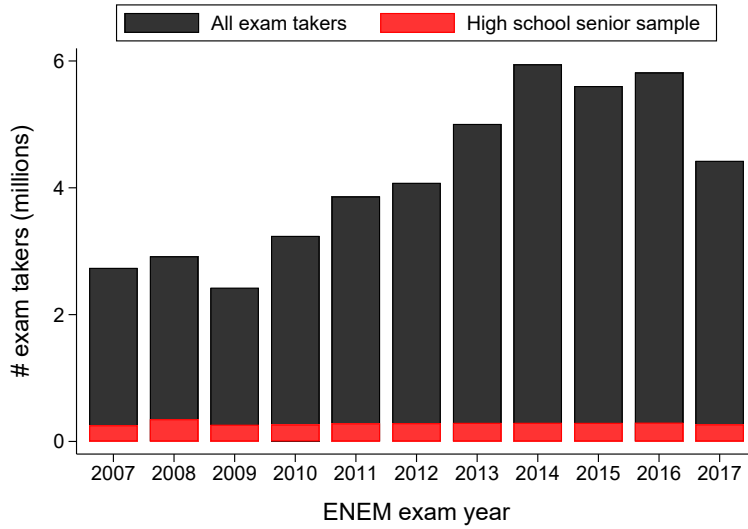
REFERENCES

- Abdulkadiroğlu, A., P. A. Pathak, and A. E. Roth (2005). The New York City high school match. *American Economic Review* 95(2), 364–367.
- Akhtari, M., N. Bau, and J.-W. P. Laliberté (2020). Affirmative action and pre-college human capital. NBER Working Paper No. 27779.
- Amrein, A. L. and D. C. Berliner (2002). High-stakes testing & student learning. *Education Policy Analysis Archives* 10, 18.
- Azmat, G., C. Calsamiglia, and N. Iriberry (2016). Gender differences in response to big stakes. *Journal of the European Economic Association* 14(6), 1372–1400.
- Bagde, S., D. Epple, and L. Taylor (2016). Does affirmative action work? caste, gender, college quality, and academic success in India. *American Economic Review* 106(6), 1495–1521.
- Bertrand, M., R. Hanna, and S. Mullainathan (2010). Affirmative action in education: Evidence from engineering college admissions in India. *Journal of Public Economics* 94(1), 16–29.
- Bettinger, E. P., B. J. Evans, and D. G. Pope (2013). Improving college performance and retention the easy way: Unpacking the ACT exam. *American Economic Journal: Economic Policy* 5(2), 26–52.
- Bhattacharya, D., S. Kanaya, and M. Stevens (2017). Are university admissions academically fair? *Review of Economics and Statistics* 99(3), 449–464.
- Buchmann, C., D. J. Condrón, and V. J. Roscigno (2010). Shadow education, american style: Test preparation, the sat and college enrollment. *Social Forces* 89(2), 435–461.
- Bulman, G. (2015). The effect of access to college assessments on enrollment and attainment. *American Economic Journal: Applied Economics* 7(4), 1–36.
- Cai, X., Y. Lu, J. Pan, and S. Zhong (2019). Gender gap under pressure: Evidence from China’s national college entrance examination. *Review of Economics and Statistics* 101(2), 249–263.
- Callaway, B. and P. H. Sant’Anna (2021). Difference-in-differences with multiple time periods. *Journal of Econometrics* 225(2), 200–230.
- Card, D., A. Mas, and J. Rothstein (2008). Tipping and the dynamics of segregation. *The Quarterly Journal of Economics* 123(1), 177–218.
- Chetty, R., J. N. Friedman, E. Saez, N. Turner, and D. Yagan (2020). Income segregation and intergenerational mobility across colleges in the United States. *The Quarterly Journal of Economics* 135(3), 1567–1633.
- College Board (2017). SAT Suite of Assessments Annual Report. <https://reports.collegeboard.org/media/pdf/2017-total-group-sat-suite-assessments-annual-report.pdf> (accessed July 2023).
- Cullen, J. B., M. C. Long, and R. Reback (2013). Jockeying for position: Strategic high school choice under Texas’ Top Ten Percent plan. *Journal of Public Economics* 97, 32–48.
- Daugherty, L., P. Martorell, and I. McFarlin (2014). Percent plans, automatic admissions, and college outcomes. *IZA Journal of Labor Economics* 3(1), 1.
- De Chaisemartin, C. and X. d’Haultfoeuille (2020). Two-way fixed effects estimators with heterogeneous treatment effects. *American Economic Review* 110(9), 2964–96.
- Duckworth, A. L., P. D. Quinn, D. R. Lynam, R. Loeber, and M. Stouthamer-Loeber (2011). Role of test motivation in intelligence testing. *Proceedings of the National Academy of*

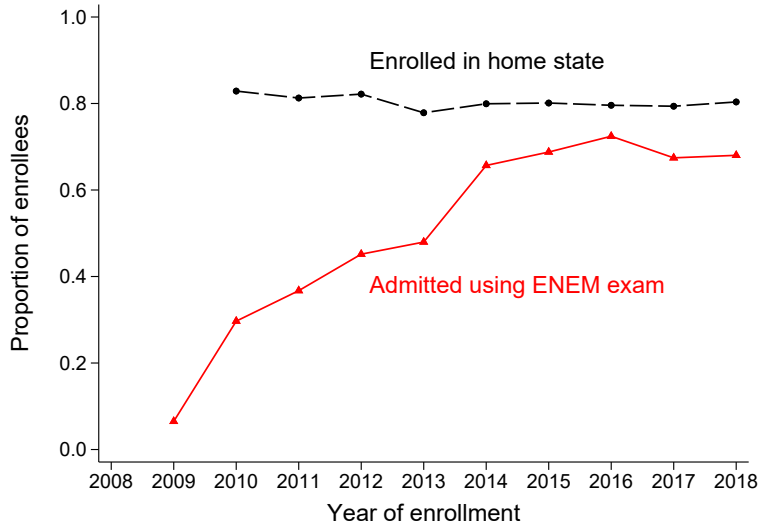
- Sciences* 108(19), 7716–7720.
- Durlauf, S. N. (2008). Affirmative action, meritocracy, and efficiency. *Politics, Philosophy & Economics* 7(2), 131–158.
- Fernandes, S. (2015). Neoliberalization of education in Brazil: the impact of cursinhos and the private sector on pedagogical practices and access to university. *Canadian Journal of Latin American and Caribbean Studies/Revue canadienne des études latino-américaines et caraïbes* 40(3), 351–367.
- Finn, B. (2015). Measuring motivation in low-stakes assessments. *ETS Research Report Series* 2015(2), 1–17.
- Fischer, P. E. and R. E. Verrecchia (2000). Reporting bias. *The Accounting Review* 75(2), 229–245.
- Frankel, A. and N. Kartik (2019). Muddled information. *Journal of Political Economy* 127(4), 1739–1776.
- Frankel, A. and N. Kartik (2022). Improving information from manipulable data. *Journal of the European Economic Association* 20(1), 79–115.
- Gneezy, U., J. A. List, J. A. Livingston, X. Qin, S. Sadoff, and Y. Xu (2019). Measuring success in education: The role of effort on the test itself. *American Economic Review: Insights* 1(3), 291–308.
- Goodman, J., O. Gurantz, and J. Smith (2020). Take two! SAT retaking and college enrollment gaps. *American Economic Journal: Economic Policy* 12(2), 115–58.
- Goodman, S. (2016). Learning from the test: Raising selective college enrollment by providing information. *Review of Economics and Statistics* 98(4), 671–684.
- Hoxby, C. and S. Turner (2013). Expanding college opportunities for high-achieving, low income students. Stanford Institute for Economic Policy Research Discussion Paper.
- Hoxby, C. M. (2009). The changing selectivity of American colleges. *The Journal of Economic Perspectives* 23(4), 95–118.
- INEP (2019a). Microdata for ENEM exam (Exame Nacional do Ensino Médio). Instituto Nacional de Estudos e Pesquisas Educacionais Anísio Teixeira, Brasília, Brazil. <https://www.gov.br/inep/pt-br/aceso-a-informacao/dados-abertos/microdados/enem> (accessed September 2019).
- INEP (2019b). School-level ENEM scores (Exame Nacional do Ensino Médio). Instituto Nacional de Estudos e Pesquisas Educacionais Anísio Teixeira, Brasília, Brazil. <https://www.gov.br/inep/pt-br/aceso-a-informacao/dados-abertos/microdados/enem-por-escola> (accessed September 2019).
- INEP (2022). Microdata for Higher Education Census (Censo da Educação Superior). Instituto Nacional de Estudos e Pesquisas Educacionais Anísio Teixeira, Brasília, Brazil. <https://www.gov.br/inep/pt-br/aceso-a-informacao/dados-abertos/microdados/censo-da-educacao-superior> (accessed April 2022).
- Jacob, B. A. (2005). Accountability, incentives and behavior: The impact of high-stakes testing in the chicago public schools. *Journal of public Economics* 89(5-6), 761–796.
- Jalava, N., J. S. Joensen, and E. Pellas (2015). Grades and rank: Impacts of non-financial incentives on test performance. *Journal of Economic Behavior & Organization* 115, 161–196.
- Kain, J. F., D. M. O’Brien, and P. A. Jargowsky (2005). *Hopwood and the Top 10 Percent Law: How they have Affected the College Enrollment Decisions of Texas High School*

- Graduates*. Texas School Project, University of Texas at Dallas.
- Kapor, A. (2015). Distributional effects of race-blind affirmative action. Working Paper.
- Long, M. C. (2004). Race and college admissions: An alternative to affirmative action? *Review of Economics and Statistics* 86(4), 1020–1033.
- Machado, C. and C. Szerman (2021). Centralized college admissions and student composition. *Economics of Education Review* 85, 102184.
- MacLeod, W. B., E. Riehl, J. E. Saavedra, and M. Urquiola (2017). The big sort: College reputation and labor market outcomes. *American Economic Journal: Applied Economics* 9(3), 223–61.
- Mello, U. (2022). Centralized admissions, affirmative action, and access of low-income students to higher education. *American Economic Journal: Economic Policy* 14(3), 166–97.
- Neal, D. and D. W. Schanzenbach (2010). Left behind by design: Proficiency counts and test-based accountability. *The Review of Economics and Statistics* 92(2), 263–283.
- Niu, S. X. and M. Tienda (2010). The impact of the Texas Top Ten Percent Law on college enrollment: A regression discontinuity approach. *Journal of Policy Analysis and Management* 29(1), 84–110.
- Ors, E., F. Palomino, and E. Peyrache (2013). Performance gender gap: does competition matter? *Journal of Labor Economics* 31(3), 443–499.
- Otero, S., N. Barahona, and C. Dobbin (2021). Affirmative action in centralized college admission systems: Evidence from Brazil. Working paper.
- Pallais, A. (2015). Small differences that matter: Mistakes in applying to college. *Journal of Labor Economics* 33(2), 493–520.
- RAIS (2022). Relação Anual de Informações Sociais. Ministério do Trabalho e Emprego, Brasília, Brazil. <http://www.rais.gov.br/> (accessed April 2022).
- Ramey, G. and V. A. Ramey (2010). The rug rat race. *Brookings Papers on Economic Activity* 41(1 (Spring)), 129–199.
- Reyes, G. (2023). Cognitive endurance, talent selection, and the labor market returns to human capital. arXiv preprint arXiv:2301.02575.
- Riehl, E. (2023). Do less informative college admission exams reduce earnings inequality? evidence from Colombia. *Journal of Labor Economics*. Forthcoming.
- Riehl, E. and M. Welch (2023). Accountability, test prep incentives, and the design of math and English exams. *Journal of Policy Analysis and Management* 42(1), 60–96.
- Rothstein, J. M. (2004). College performance predictions and the SAT. *Journal of Econometrics* 121(1), 297–317.
- Scott-Clayton, J. (2012). Do high-stakes placement exams predict college success? ccrc working paper no. 41. *Community College Research Center, Columbia University*.
- Segal, C. (2012). Working when no one is watching: Motivation, test scores, and economic success. *Management Science* 58(8), 1438–1457.
- Soares, J. A. (2015). *SAT wars: The case for test-optional college admissions*. Teachers College Press.
- Westrick, P. A., J. P. Marini, L. Young, H. Ng, D. Shmueli, and E. J. Shaw (2019). Validity of the SAT for predicting first-year grades and retention to the second year. Technical report, College Board.
- Wise, S. L. and C. E. DeMars (2005). Low examinee effort in low-stakes assessment: Problems and potential solutions. *Educational assessment* 10(1), 1–17.

FIGURES AND TABLES



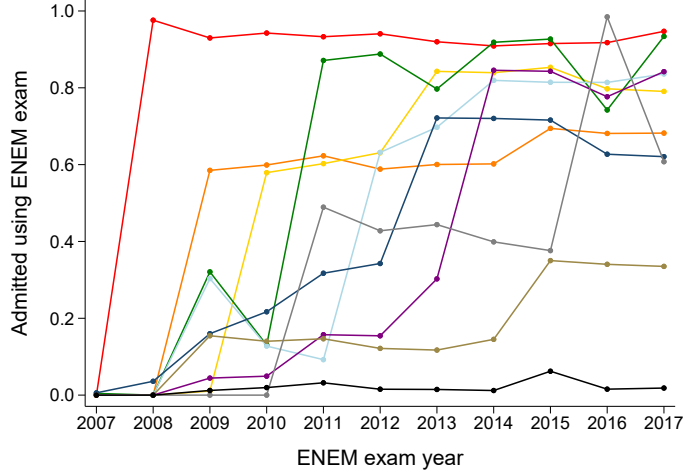
Panel A. Number of ENEM exam takers



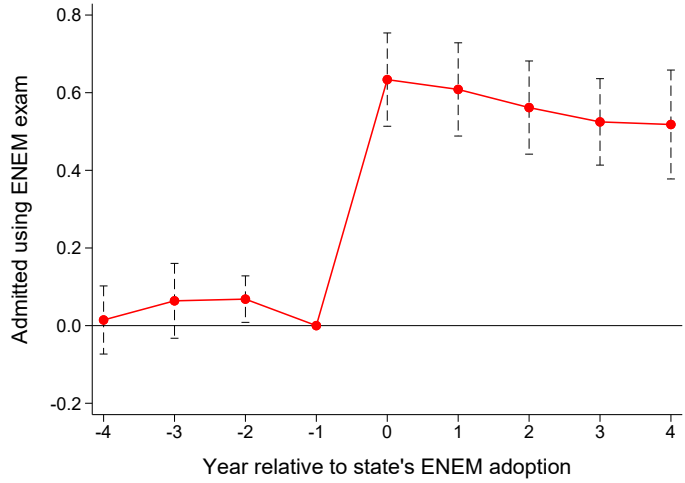
Panel B. Federal university enrollment

FIGURE 1. Adoption of ENEM exam by federal universities

Notes: This figure illustrates the impact of increased ENEM stakes on the number of exam takers and federal university enrollees through the ENEM. Panel A shows the total number of exam takers of the ENEM (including the pre-2009 version) over the 2007-2017 period. Each bar displays the overall number of ENEM exam takers (black area) and the number of exam takers in our analysis sample of high school graduates (red area). Panel B plots the fraction of new enrollees to federal universities admitted through the ENEM exam (red line) and that enrolled in the state where they were born (black dashed line) over the 2009-2018 period.



Panel A. Proportion of federal university enrollees admitted using ENEM exam by ENEM adoption year



Panel B. Event study for proportion of federal university enrollees admitted using ENEM exam

FIGURE 2. Variation in ENEM exam adoption by federal universities across states and years

Notes: This figure illustrates the staggered adoption of the ENEM exam by federal universities. The outcome in both panels is the proportion of new enrollees in federal universities in state s who were admitted using the ENEM exam administered in year t (the calendar year prior to enrollment), which we denote by $\text{ProportionENEM}_{st}$. Panel A plots the mean of $\text{ProportionENEM}_{st}$ in groups of state(s) based on their ENEM adoption year, τ_s^* , denoted by the legend. Panel B plots event-study coefficients, β_l , from the state (s) \times year (t) \times pairwise group (g) level regression

$$\text{ProportionENEM}_{stg} = \gamma_{sg} + \gamma_{tg} + \sum_{l=-7}^7 \beta_l \mathbb{1}\{t - \tau_s^* = l\} + \epsilon_{stg}$$

where l denotes years relative to τ_s^* . We include state \times group dummies, γ_{sg} , year \times group dummies, γ_{tg} , and dummies for years l , $\mathbb{1}\{t - \tau_s^* = l\}$, omitting $l = -1$. The graph plots the β_l coefficients from $l = -4$ to 4. Dashed lines depict 95% confidence intervals using standard errors clustered at the state level. See Appendix Tables A2–A3 for details on $\text{ProportionENEM}_{st}$ and each state's ENEM adoption year, τ_s^* . See the text in Section 3.4 and Appendix Table A4 for details on the dataset of pairwise ENEM adoption years that we use for our event studies.

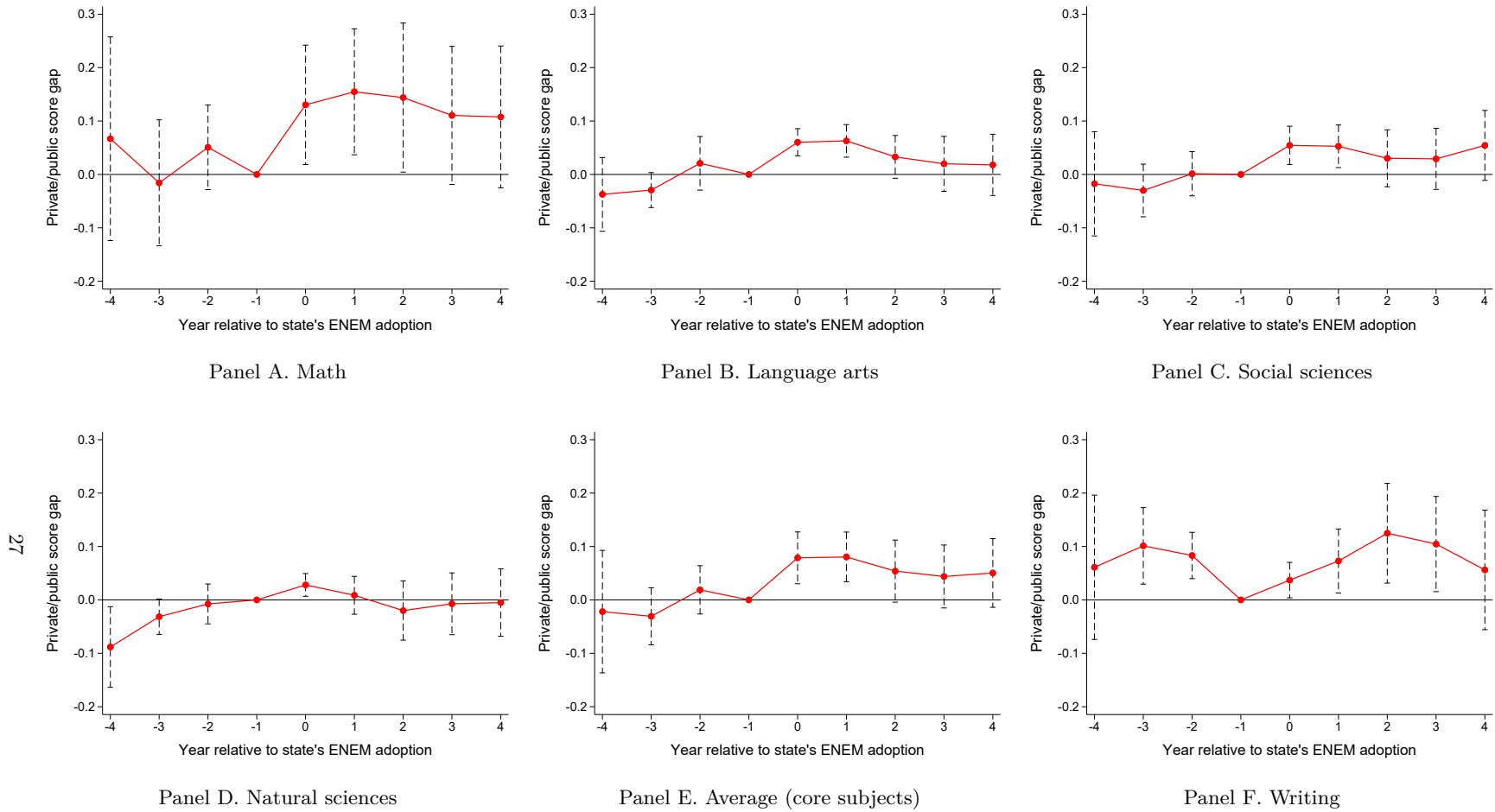


FIGURE 3. Event studies for effects of ENEM stakes on private/public test score gaps

Notes: This figure plots event study estimates of the effects of ENEM stakes on test score gaps between private and public school students. The sample includes all pairwise combinations of ENEM adoption years for which we can estimate treatment effects using 2009–2017 exam takers (the boxed cells in Appendix Table A4). The dependent variables are ENEM subject scores in SD units. “Average (core subjects)” is the average score across math, language arts, natural science, and social science. Each panel plots the β_l^{gap} coefficients (y -axis) from our event study regression (4) for years $l = -4$ to 4 relative to the state’s ENEM adoption year, $\tau_{s(i)}^*$ (x -axis). Dashed lines depict 95% confidence intervals using standard errors clustered at the state level.

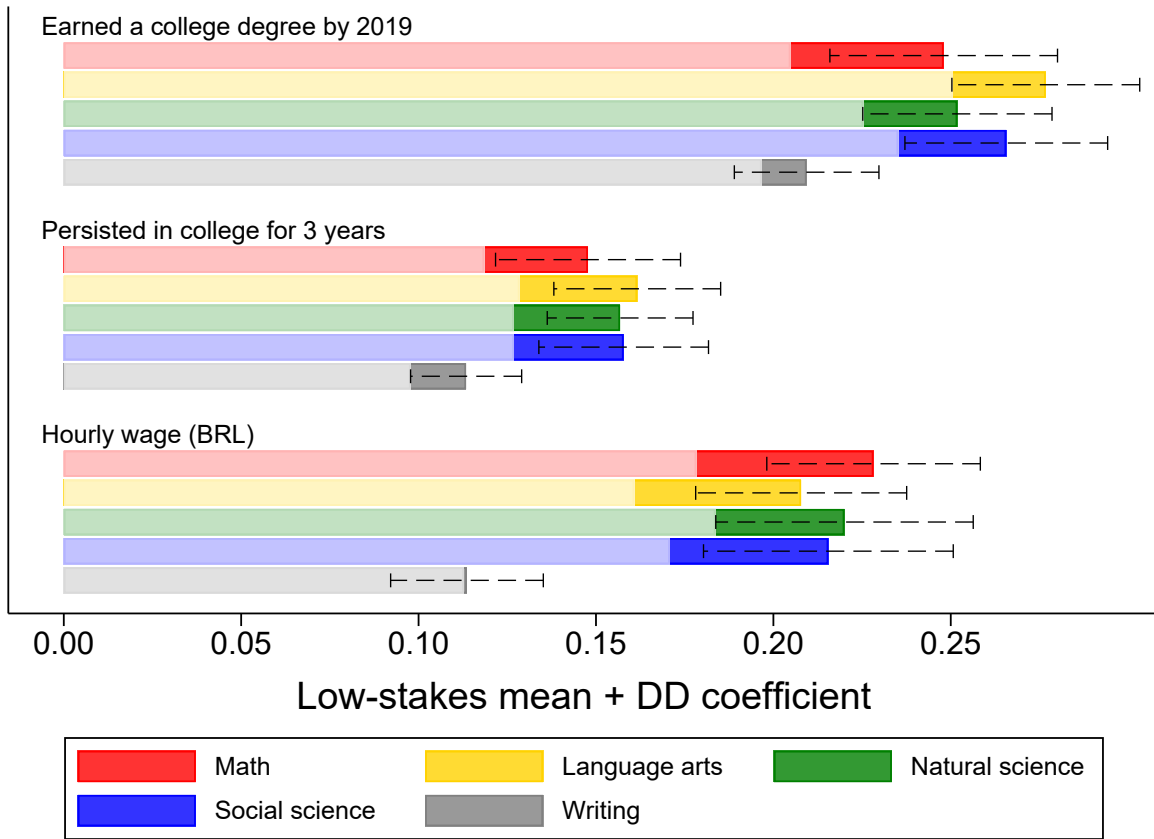


FIGURE 4. Effects on the correlation between outcomes and test scores by subject

Notes: This figure illustrates the effects of the ENEM stakes on the correlation between student outcomes and test scores by subjects. The sample in which we calculate the correlations includes 2009–2017 ENEM exam takers in our high school senior sample who have valid outcomes.

The translucent area represents the low-stakes mean correlations between the subject scores and the outcomes indicated by the heading. The opaque area represents β coefficients from regression (5) for the subject scores. The dashed lines represent the 95% confidence intervals of the estimates.

TABLE 1. Summary statistics for 2009–2017 ENEM exam takers

	(A)	(B)	(C)	(D)	(E)	(F)
			Analysis sample (high school seniors)			
	All exam takers	All HS seniors	All schools	Private schools	Public schools	Private/public gap
Panel A. Exam taker characteristics						
Age at exam	22.14	18.55	17.91	17.44	18.13	−0.70
Female	0.58	0.59	0.58	0.55	0.60	−0.05
White	0.40	0.44	0.51	0.69	0.43	0.26
Black	0.12	0.11	0.08	0.04	0.10	−0.06
Brown	0.44	0.42	0.37	0.23	0.44	−0.21
Mother attended college	0.15	0.18	0.27	0.56	0.13	0.44
Father attended college	0.11	0.13	0.21	0.49	0.07	0.41
Family income > 2x min. wage	0.35	0.38	0.49	0.85	0.32	0.52
Private high school	0.24	0.24	0.32	1.00	0.00	1.00
Panel B. ENEM scores						
Math score	−0.03	−0.01	0.32	1.28	−0.13	1.42
Language arts score	0.08	0.04	0.24	0.78	−0.01	0.79
Natural science score	−0.17	−0.18	0.05	0.75	−0.28	1.03
Social science score	0.30	0.22	0.43	1.07	0.14	0.93
Average score (core subjects)	0.05	0.02	0.30	1.12	−0.09	1.20
Writing score	−0.41	−0.38	−0.13	0.50	−0.43	0.93
Panel C. College and labor market outcomes						
Ever enrolled in college			0.76	0.95	0.67	0.27
Enrolled in a federal university			0.16	0.26	0.11	0.15
Graduated college within 5 years			0.17	0.23	0.15	0.08
Ever graduated college			0.31	0.43	0.25	0.18
Persisted in college for 3 years			0.66	0.73	0.61	0.12
Fraction of college credits completed			0.69	0.75	0.64	0.11
Appears in RAIS			0.26	0.20	0.29	−0.08
Hourly wage (BRL)			48.89	70.03	41.59	28.44
Number of exam takers	40,391,604	11,626,416	2,512,214	807,293	1,704,921	2,512,214
Number of high schools	46,584	45,867	3,276	1,437	1,839	3,276

Notes: This table reports summary statistics on the ENEM exam takers. Column (A) includes all the ENEM exam takers with a valid score (i.e., non-zero, non-missing) in all four subjects of the ENEM. Column (B) includes all the exam takers who were high school seniors. Column (C) includes all the exam takers in our analysis sample. Columns (D) and (E) include the exam takers in our analysis sample that attended a private and public high school, respectively. Column (F) displays the difference between columns (D) and (E).

Panel (A) describes demographic characteristics of the exam takers, including age, gender, race, parental education, family income, and whether they attended a private high school. Panel (B) reports the average ENEM scores (in SD units) in the respective samples. “Average score (core subjects)” is the average score across math, language arts, natural science, and social science. Panel (C) displays the college and labor market outcomes for the exam takers in our analysis sample. The last two rows reports the number of exam takers and high schools in the respective samples. See Appendix C.1 for detailed variable definitions.

TABLE 2. Balance tests for exam taker characteristics

Dependent variable	(A)	(B)	(C)	(D)	(E)
	Low-stakes mean	DD coefficients			
	All schools	All schools	Private schools	Public schools	Private/ public gap
Panel A. Number of exam takers per school \times year					
Number of exam takers	151.713	14.225 (8.985)	1.308 (18.395)	14.966* (8.073)	-13.658 (18.169)
Log number of exam takers	4.713	0.076 (0.065)	-0.004 (0.118)	0.089 (0.056)	-0.093 (0.107)
Panel B. Demographic characteristics of exam takers					
Age at exam	18.190	0.030 (0.054)	0.010 (0.015)	0.073 (0.083)	-0.063 (0.075)
Female	0.599	-0.014** (0.006)	-0.013** (0.005)	-0.012* (0.006)	-0.000 (0.006)
White	0.469	-0.007 (0.010)	-0.009 (0.008)	0.001 (0.010)	-0.010 (0.012)
Mother attended college	0.258	0.001 (0.008)	0.009 (0.010)	0.000 (0.005)	0.009 (0.010)
Father attended college	0.195	0.006 (0.006)	0.013 (0.011)	0.005 (0.004)	0.008 (0.009)
Family income > 2x min. wage	0.476	0.003 (0.022)	0.016 (0.010)	-0.002 (0.028)	0.018 (0.031)
Joint balance test (p value)		0.159	0.308	0.206	0.708
Panel C. Predicted score based on demographics					
Predicted ENEM score	0.181	0.004 (0.011)	0.017 (0.012)	-0.001 (0.011)	0.018 (0.014)
N (# exam takers)	492,436	2,512,214	807,293	1,704,921	2,512,214

Notes: This table displays balance tests for the ENEM participants in our analysis sample. The sample includes 2009–2017 ENEM exam takers in our high school senior sample (column C of Table 1). The dependent variables are: the number of exam takers in levels and logs (Panel A); exam taker demographic characteristics (Panel B); and the predicted value from a regression of ENEM scores (averaged across math, language arts, natural science, and social science) on age, gender, and dummies for race, mother’s education, father’s education, and family income bins (Panel C). Our dependent variables are high school \times year totals (Panel A) and averages (Panel B–C) of these variables.

Column (A) shows the mean of each dependent variable in exam cohorts prior to each state’s ENEM adoption year, i.e., cohorts with $\text{HighStakes}_{st} = 0$. Columns (B)–(D) display β coefficients from equation (1) estimated using all students, private students, and public students, respectively. Column (E) displays β^{gap} coefficients from equation (2) estimated using all students. The last row of Panel B shows the p value from an F test that the coefficients in Panel B are jointly equal to zero.

Parentheses contain standard errors clustered at the state level.

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

TABLE 3. Effects of ENEM adoption on test scores in public and private high schools

Dependent variable	(A)	(B)	(C)	(D)	(E)
	Low-stakes mean	DD coefficients			
	Private/ public gap	All schools	Private schools	Public schools	Private/ public gap
Math score	1.358	0.022 (0.055)	0.143** (0.058)	-0.015 (0.070)	0.158* (0.079)
Language arts score	0.837	0.035 (0.035)	0.068*** (0.020)	-0.008 (0.034)	0.076*** (0.026)
Natural science score	1.059	0.026 (0.040)	0.062* (0.031)	-0.003 (0.042)	0.065* (0.034)
Social science score	1.010	0.019 (0.034)	0.056* (0.029)	-0.024 (0.029)	0.081*** (0.023)
Average score (core subjects)	1.229	0.029 (0.043)	0.095** (0.036)	-0.014 (0.045)	0.110** (0.040)
Writing score	0.784	0.049 (0.035)	0.165** (0.072)	0.063* (0.033)	0.102* (0.058)
<i>N</i> (# exam takers)	492,436	2,512,214	807,293	1,704,921	2,512,214

Notes: This table shows how the increase in the stakes of the ENEM exam impacted scores for private and public high school students. The sample includes 2009–2017 ENEM exam takers in our high school senior sample (column C of Table 1). The dependent variables are ENEM subject scores in SD units. “Average score (core subjects)” is the average score across math, language arts, natural science, and social science.

Column (A) shows the mean private/public score gap in exam cohorts prior to each state’s ENEM adoption year, i.e., cohorts with $\text{HighStakes}_{st} = 0$. Columns (B)–(D) display β coefficients from equation (1) estimated using all students, private students, and public students, respectively. Column (E) displays β^{gap} coefficients from equation (2) estimated using all students.

Parentheses contain standard errors clustered at the state level.

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

TABLE 4. Robustness checks on the effects of ENEM adoption on private/public test score gaps

	(A)	(B)	(C)	(D)	(E)	(F)
Dependent variable	Benchmark model	Demographic controls	Binary treatment	Stacked regression	2009 vs 2013 adopters (2007–2012)	Stacked regression (2007–2017)
Math score	0.158* (0.079)	0.128* (0.067)	0.105** (0.050)	0.116** (0.055)	0.086*** (0.019)	0.061** (0.024)
Language arts score	0.076*** (0.026)	0.072*** (0.023)	0.042** (0.017)	0.055*** (0.014)	0.106*** (0.031)	0.062*** (0.022)
Natural science score	0.065* (0.034)	0.074** (0.036)	0.032 (0.030)	0.025 (0.016)	0.059*** (0.017)	0.022 (0.014)
Social science score	0.081*** (0.023)	0.053* (0.028)	0.042** (0.020)	0.053** (0.020)	0.046 (0.035)	0.042 (0.025)
Average score (core subjects)	0.110** (0.040)	0.094** (0.040)	0.064** (0.030)	0.072*** (0.023)	0.088** (0.029)	0.049** (0.024)
Writing score	0.102* (0.058)	0.144** (0.058)	0.035 (0.044)	0.023 (0.028)	0.058 (0.047)	0.064*** (0.021)
<i>N</i> (# exam takers)	2,512,214	2,512,214	2,512,214	5,858,862	1,099,500	15,738,474
Treatment variable:	Continuous	Continuous	Binary	Binary	Binary	Binary
Demographic controls:		Yes				
Level of dataset:	HS × year	HS × year	HS × year	Stacked	HS × year	Stacked
Included exam cohorts:	2009–2017	2009–2017	2009–2017	2009–2017	2007–2012	2007–2017

32

Notes: This table examines the robustness of our estimates for the effects of ENEM stakes on the private/public school gap in test scores. Our main sample includes 2009–2017 ENEM exam takers in our analysis sample (column C of Table 1). In columns (E)–(F), we add in 2007–2008 ENEM exam takers from the same set of high schools. The dependent variables are ENEM subject scores in SD units. “Average score (core subjects)” is the average score across math, language arts, natural science, and social science. Columns (E)–(F) include scores from the 2007–2008 ENEM tests; in these columns, we standardize scores to have mean 0 and SD 1 within each year of our sample. For the 2007–2008 exams, “average score” is the reported core-component score, and we compute math, language arts, natural science, and social science scores by categorizing the multiple choice questions into these subjects and then estimating scores using the IRT parameters. See Appendix C.1 for details.

Column (A) replicates the estimates from column (E) of Table 3, which are the β^{gap} coefficients from equation (2). Column (B) estimates equation (2) including high school × year averages of age, gender, and dummies for race, mother’s education, father’s education, and family income bins. Column (C) estimates equation (2) replacing our continuous treatment variable, $\text{ProportionENEM}_{st}$, with our binary treatment variable, HighStakes_{st} . Columns (D)–(F) display estimates of β^{gap} from equation (3) using our stacked dataset, which contains pairwise combinations of ENEM adoption years (as described in Section 3.4). Column (D) includes all pairwise combinations for which we can estimate treatment effects using 2009–2017 exam takers (the boxed cells in Appendix Table A4). Column (E) includes 2007–2012 exam takers and a single pair of ENEM adoptions years: 2009 and 2013 (the bolded cells in Appendix Table A4). Column (F) includes all 2007–2017 exam takers and all pairwise combinations (all cells in Appendix Table A4).

Parentheses contain standard errors clustered at the state level. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

TABLE 5. Heterogeneity by topic area and competency — Math exam
 Dependent variable: Proportion correct answers

(A)	(B)	(C)	(D)	(E)	(F)
Question group	N_q	Public students		Private/public gap	
		Mean	β/mean	Mean	$\beta^{\text{gap}}/\text{mean}$
Panel A. All questions					
All questions	405	0.291	-0.018 (0.033)	0.176	0.136 (0.052)**
Panel B. Topic area (and competency reference numbers)					
Numbers (1–5)	67	0.307	-0.032 (0.035)	0.159	0.135 (0.069)*
Geometry (6–9)	57	0.317	0.009 (0.021)	0.160	0.099 (0.051)*
Measurements (10–14)	62	0.257	-0.022 (0.039)	0.193	0.150 (0.048)***
Proportions (15–18)	51	0.336	-0.024 (0.055)	0.225	0.159 (0.066)**
Algebra (19–23)	66	0.264	-0.011 (0.020)	0.172	0.076 (0.041)*
Interpreting data (24–26)	47	0.325	-0.043 (0.053)	0.193	0.177 (0.077)**
Statistics (27–30)	55	0.241	-0.001 (0.029)	0.137	0.158 (0.066)**
All coefficients equal (p value)			0.269		0.060
Panel C. Competencies (top 5 and bottom 5 by $\beta^{\text{gap}}/\text{mean}$)					
Calculate statistical quantities from data (27)	15	0.220	0.025 (0.047)	0.140	0.275 (0.156)*
Evaluate interventions using proportions (18)	12	0.293	-0.040 (0.080)	0.219	0.272 (0.078)***
Make inferences using data in tables/graphs (24)	14	0.261	-0.040 (0.027)	0.122	0.230 (0.102)**
Identify numerical patterns (2)	12	0.274	-0.050 (0.063)	0.155	0.198 (0.083)**
Solve problems using statistics (28)	16	0.225	-0.018 (0.026)	0.154	0.192 (0.079)**
...					
Use numbers to construct arguments (4)	15	0.266	-0.027 (0.036)	0.161	0.069 (0.076)
Use proportions to construct arguments (17)	12	0.218	0.024 (0.027)	0.186	0.069 (0.071)
Solve problems using geometry (8)	18	0.236	0.054 (0.015)***	0.147	0.063 (0.069)
Interpret Cartesian graphs (20)	11	0.541	-0.033 (0.050)	0.209	0.037 (0.187)
Evaluate interventions using statistics (30)	10	0.253	0.004 (0.043)	0.072	-0.149 (0.216)
21 coefficients equal (p value)			0.000		0.000

Notes: This table shows how the increase in ENEM stakes impacted students' performance on different topic areas and competencies of the math exam. The sample includes 2009–2017 ENEM exam takers in our high school senior sample (column C of Table 1). Regressions are at the high school (h) \times year (t) \times exam question (q) level. The dependent variable is the proportion of correct answers in each htq cell. We estimate regressions pooling across all math questions (Panel A) and separately for math questions in 7 topic areas (Panel B) and 30 competencies (Panel C) defined by ENEM test designers. Panel C reports only the top 5/bottom 5 competencies by the values in column (F). See Appendix C.4 for details on the categorization of ENEM math questions.

Column (A) defines the group of questions for each regression. Column (B) shows the number of questions in each group. Column (C) shows the mean proportion of correct answers for public school students in cohorts prior to each state's ENEM adoption year (i.e., cohorts with $\text{HighStakes}_{st} = 0$). Column (E) shows the mean private/public gap in the proportion of correct answers in those cohorts. Columns (D) and (F) display the β and β^{gap} coefficients from equation (2) estimated for each group of questions. We normalize β and β^{gap} so that they represent a percentage change from the means in columns (C) and (E). In Panel B, the last row reports p values from F tests that the 7 topic area coefficients in columns (D) or (F) are equal. In Panel C, the last row reports p values from F tests that 21 competency coefficients (the first 3 in each topic area) are jointly equal.

Parentheses contain standard errors clustered at the state level.

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

TABLE 6. Effects of an increase in ENEM stakes on the correlation of test scores and outcomes

Dependent variable	(A)	(B)	(C)	(D)	(E)
	Low-stakes mean	Benchmark model DD coefficients		Binary treatment DD coefficients	
	Raw corr.	Raw corr.	Within-program	Raw corr.	Within program
Panel A. Outcomes for all exam takers					
Enrolled in any college by 2019	0.372	0.036*** (0.006)		0.022*** (0.004)	
Enrolled in a federal university	0.400	0.053*** (0.015)		0.028*** (0.009)	
Finished college within 5 years of ENEM	0.121	0.014* (0.008)	0.026*** (0.007)	0.007 (0.004)	0.015*** (0.005)
Earned a college degree by 2019	0.257	0.033** (0.015)	0.038*** (0.010)	0.016 (0.011)	0.023*** (0.007)
Appears in RAIS in 2016–2018	−0.112	0.056 (0.044)	0.020* (0.011)	0.014 (0.025)	0.008 (0.007)
<i>N</i> (# exam takers)	336,175	1,266,412	1,266,412	1,266,412	1,266,412
Panel B. Outcomes for college enrollees					
Persisted in college for 1 year	0.064	0.008 (0.014)	0.024*** (0.007)	0.007 (0.011)	0.011* (0.006)
Persisted in college for 3 years	0.142	0.033** (0.013)	0.043*** (0.008)	0.018** (0.009)	0.024*** (0.006)
Completed program within 5 years	0.071	0.044** (0.018)	0.035*** (0.011)	0.022* (0.011)	0.018** (0.007)
Fraction of college credits completed	0.214	0.003 (0.014)	0.014 (0.013)	−0.011 (0.011)	0.013* (0.007)
<i>N</i> (# in higher ed.)	274,022	966,649	966,649	966,649	966,649
Panel C. Outcome for individuals in RAIS					
Hourly wage (BRL)	0.200	0.046** (0.018)	0.027** (0.010)	0.027*** (0.007)	0.012* (0.006)
Log hourly wage	0.362	−0.029** (0.014)	−0.001 (0.010)	−0.017* (0.008)	−0.001 (0.006)
<i>N</i> (# in RAIS)	80,382	328,773	328,773	328,773	328,773

Notes: This table shows how the increase in the stakes of the ENEM exam impacted the correlation between test scores and student outcomes. The sample in which we calculate the correlations includes 2009–2017 ENEM exam takers in our high school senior sample who have valid outcomes. The dependent variables are correlation coefficients between the specified outcome and the average score across math, language arts, natural science, and social science.

Column (A) shows the mean correlation coefficients in exam cohorts prior to each state’s ENEM adoption year, i.e., cohorts with $\text{HighStakes}_{st} = 0$. Columns (B)–(C) display β coefficients from equation (5) using raw and within-program correlation coefficients, respectively. Columns (B)–(C) display the same coefficients estimated with our binary treatment variable, HighStakes_{st} .

Parentheses contain standard errors clustered at the state level. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

Appendix — For Online Publication

A. APPENDIX FIGURES AND TABLES

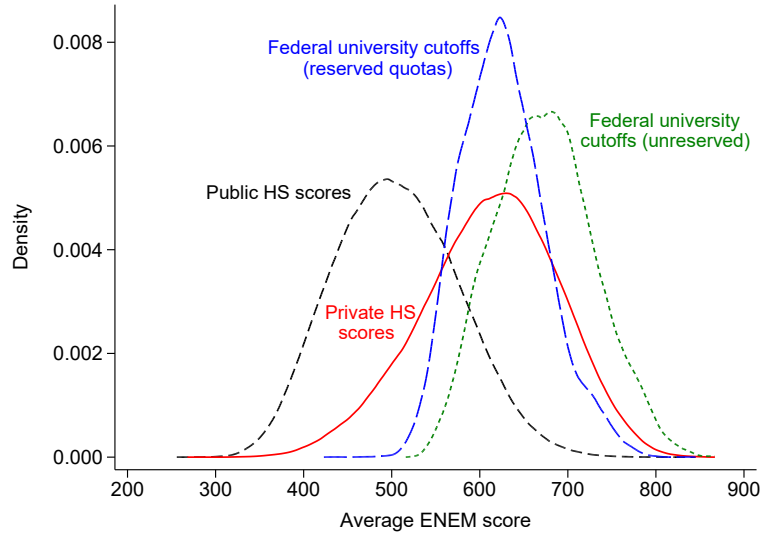
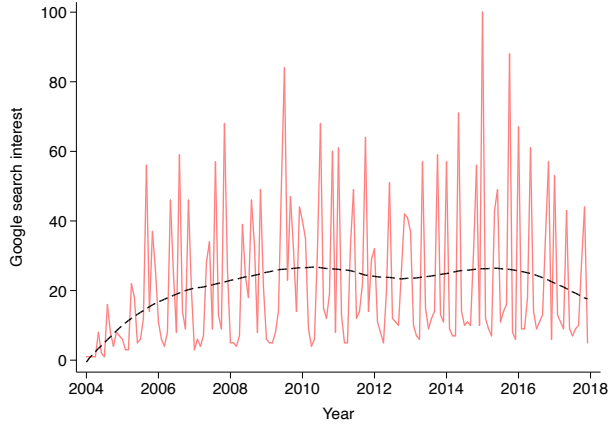


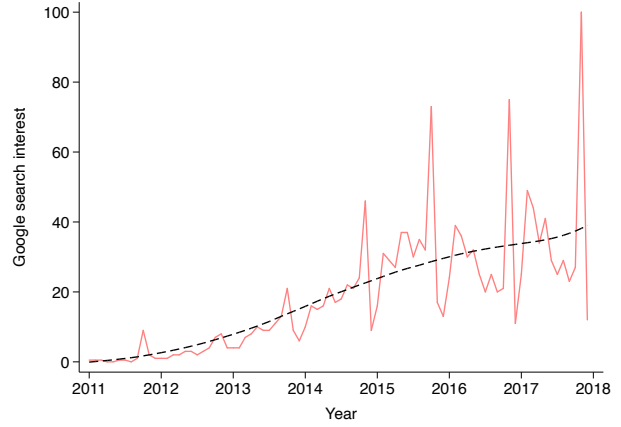
FIGURE A1. Distributions of average ENEM scores and federal university admission cutoffs

Notes: This figure compares the distributions of actual ENEM scores and federal university cutoff scores. The solid red line shows the distribution of average ENEM scores for private school students in our sample who took the ENEM exam in exam cohorts prior to each state’s ENEM adoption year, i.e., cohorts with $\text{HighStakes}_{st} = 0$. The dashed black line shows the same distribution for public school students. In both distributions, average ENEM scores are the average score across math, language arts, natural science, and social science.

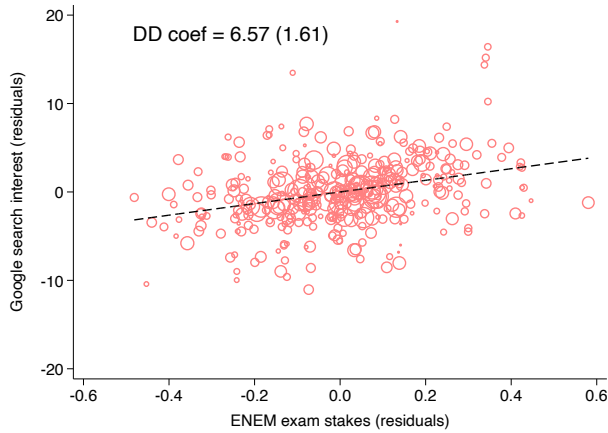
The green short-dashed line plots the distribution of cutoff scores for unreserved admissions to all federal university programs in 2016. The blue long-dashed line plots the same distribution for reserved quotas at federal university programs, which include quotas for public high students, low-SES students, and/or underrepresent minority students.



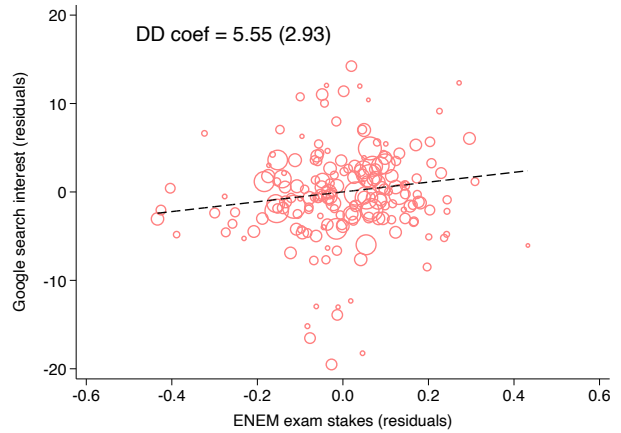
Panel A. National trend in Google search for “ENEM”



Panel B. National trend in Google search for “Descomplica”



Panel C. Effect of ENEM stakes on Google search for “ENEM”



Panel D. Effect of ENEM stakes on Google search for “Descomplica”

FIGURE A2. Google search trends for “ENEM” and the online test prep service “Descomplica”

Notes: This figure shows how Google search trends for “ENEM” and the online test prep service “Descomplica” varied over time at the national and state levels. Panels A and C show results for the search “ENEM” using data from 2004–2017. Panels B and D show results for the search “Descomplica” using data from 2011–2017; we do not include years prior to 2011 because search volume for “Descomplica” was low and state-level data are noisy.

Panels A–B plot monthly Google search interest for the entire country of Brazil (solid red lines) plus non-parametric predicted values (black dashed lines). A value of 100 is the peak popularity for the term. A value of 50 means that the term is half as popular. A score of 0 means there was not enough data for this term.

Panels C–D plot state \times year level means of Google search interest (y -axis) against our continuous treatment variable, $\text{ProportionENEM}_{st}$ (x -axis). We residualize both variables on state and year dummies, and plot the residuals along with a linear regression line (black dashed line). We also display DD coefficients (and standard errors) from estimating equation (1) using Google search interest for each term as the dependent variable. These DD coefficients are equivalent to the slopes of the black dashed lines in Panels C–D.

TABLE A1. Summary of Brazilian high school and college markets

(A)	(B)	(C)	(D)	(E)	(F)	(G)
Institution type	# of schools	Prop. of schools	# of students	Prop. of students	# students per school	Attended a private HS
Panel A. High school seniors in 2009						
Federal high schools	100	0.004	9,772	0.005	98	1.000
State high schools	16,583	0.702	1,823,524	0.849	110	0.000
Municipal high schools	373	0.016	23,156	0.011	62	0.000
Private high schools	6,567	0.278	290,366	0.135	44	1.000
All high schools	23,623	1.000	2,146,818	1.000	91	0.140
Panel B. New college enrollees in 2009						
Federal universities	59	0.025	225,112	0.108	3,815	0.471
State universities	40	0.017	119,489	0.057	2,987	0.370
Municipal universities	9	0.004	22,453	0.011	2,495	0.319
Private universities	225	0.094	1,018,698	0.489	4,528	0.458
Public technical colleges	168	0.070	55,609	0.027	331	0.259
Private technical colleges	1,888	0.790	640,021	0.307	339	0.331
All colleges	2,389	1.000	2,081,382	1.000	871	0.401

Notes: This table presents summary statistics on the Brazilian high school and college systems. Panel A presents statistics for students who were high school seniors in 2009 using data from a national primary and secondary school census (*Censo Escolar*). Panel B presents statistics for students who were new college enrollees in 2009 using data from a national higher education census (*Censo da Educação Superior*).

Column (A) categorizes high schools by ownership (federal, state, municipal, or private), and it categorizes colleges by both ownership and institution type (university or technical college). University includes both *Universidade* and *Centro Universitário* institutions. Technical colleges include *Faculdade*, *Instituto Federal de Educação Ciência e Tecnologia*, and *Centro Federal de Educação Tecnológica* institutions. Column (B) shows the number of schools in each category, and column (C) shows the proportion of schools. Column (D) shows the number of students who attended schools in each category, and column (E) shows the proportion of students. Column (F) shows the number of students per school (column D divided by column B). Column (G) shows the proportion of students at each school type who attended a private high school. Throughout the paper, we include the small number of federal high schools in the group of private high schools since both tend to enroll wealthier and higher-achieving students.

TABLE A2. Proportion of federal university enrollees admitted using the ENEM by state and year

(A)	(B)	(C)	(D)	(E)	(F)	(G)	(H)	(I)	(J)	(K)	(L)	(M)
Proportion admitted using ENEM by exam year ($\text{ProportionENEM}_{st}$)												
State	# 2009 enrollees	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017
Pernambuco (PE)	7,375	0.00	0.98	0.93	0.94	0.93	0.94	0.92	0.91	0.92	0.92	0.95
Amazonas (AM)	2,821	0.00	0.00	0.32	0.36	0.38	0.42	0.43	0.42	0.50	0.47	0.47
Espirito Santo (ES)	3,302	0.00	0.00	0.99	0.92	1.00	0.90	0.92	0.89	0.87	0.94	0.92
Maranhão (MA)	2,359	0.00	0.00	0.98	0.88	0.95	0.94	0.91	0.87	0.89	0.78	0.92
Mato Grosso (MT)	3,582	0.00	0.00	0.64	0.93	0.90	0.86	0.83	0.94	0.95	0.89	0.97
Paraná (PR)	6,820	0.00	0.00	0.53	0.51	0.51	0.47	0.58	0.60	0.63	0.57	0.57
Rio Grande do Sul (RS)	12,723	0.00	0.00	0.42	0.37	0.42	0.41	0.39	0.41	0.62	0.63	0.62
São Paulo (SP)	5,774	0.00	0.00	0.87	0.95	0.89	0.80	0.74	0.70	0.74	0.77	0.71
Ceará (CE)	3,704	0.00	0.00	0.00	0.85	0.85	0.82	0.85	0.83	0.87	0.77	0.81
Mato Grosso do Sul (MS)	2,710	0.00	0.00	0.00	0.66	0.63	0.54	0.65	0.72	0.71	0.62	0.56
Paraíba (PB)	7,160	0.00	0.00	0.02	0.40	0.45	0.55	0.94	0.90	0.92	0.93	0.90
Acre (AC)	955	0.00	0.00	0.00	0.00	0.96	0.76	0.94	0.89	0.82	0.86	0.83
Alagoas (AL)	3,008	0.00	0.00	0.00	0.00	0.90	0.83	0.87	0.88	0.94	0.95	0.93
Piauí (PI)	3,592	0.01	0.00	0.55	0.29	0.82	1.00	0.71	0.97	0.96	0.58	0.98
Rio Grande do Norte (RN)	6,443	0.00	0.00	0.34	0.15	0.06	0.61	0.70	0.82	0.82	0.79	0.82
Rondônia (RO)	1,106	0.00	0.00	0.08	0.02	0.31	0.77	0.69	0.79	0.81	0.95	0.96
Bahia (BA)	6,251	0.00	0.00	0.16	0.16	0.31	0.36	0.90	0.83	0.82	0.81	0.77
Distrito Federal (DF)	5,296	0.00	0.00	0.00	0.01	0.10	0.04	0.23	0.24	0.20	0.21	0.21
Minas Gerais (MG)	20,918	0.00	0.08	0.19	0.20	0.29	0.36	0.72	0.74	0.75	0.72	0.73
Pará (PA)	2,937	0.00	0.00	0.00	0.00	0.03	0.09	0.89	0.82	0.78	0.76	0.81
Rio de Janeiro (RJ)	16,871	0.02	0.02	0.22	0.37	0.49	0.45	0.77	0.77	0.78	0.50	0.48
Roraima (RR)	842	0.00	0.00	0.00	0.07	0.11	0.17	0.48	0.46	0.38	0.41	0.39
Goiás (GO)	3,636	0.00	0.00	0.00	0.01	0.14	0.19	0.36	0.85	0.83	0.77	0.84
Tocantins (TO)	1,602	0.00	0.00	0.14	0.15	0.20	0.08	0.18	0.83	0.88	0.78	0.84
Santa Catarina (SC)	4,632	0.00	0.00	0.15	0.14	0.15	0.12	0.12	0.15	0.35	0.34	0.34
Amapá (AP)	309	0.00	0.00	0.00	0.00	0.49	0.43	0.44	0.40	0.38	0.98	0.61
Sergipe (SE)	3,309	0.00	0.00	0.01	0.02	0.03	0.02	0.01	0.01	0.06	0.02	0.02
All states	140,037	0.00	0.07	0.30	0.37	0.45	0.48	0.66	0.69	0.72	0.67	0.68

Notes: This table shows the proportion of federal university enrollees who were admitted using the ENEM by state and year. Column (A) lists the 27 states of Brazil. Column (B) shows the number of new federal university enrollees in each state in the 2009 calendar year. Columns (C)–(M) show the proportion of new federal university enrollees who were admitted using the ENEM based on the year students took the ENEM exam (the calendar year prior to enrollment). The sample for these statistics is new enrollees in bachelor’s programs at federal universities using data from the Brazilian higher education census.

The numbers in columns (C)–(M) are the values we use for our continuous treatment variable, $\text{ProportionENEM}_{st}$, where s denotes states and t denotes years. Bolded numbers represent state \times years that we classify as high stakes using our binary treatment variable, HighStakes_{st} . See Section 3.3 for details on the definition of $\text{ProportionENEM}_{st}$ and HighStakes_{st} .

TABLE A3. Characteristics of federal universities by state’s ENEM adoption year

(A)	(B)	(C)	(D)	(E)	(F)	(G)	(H)	(I)	(J)	(K)	(L)
	Year of state’s ENEM adoption (τ_s^*)										
Characteristic	2008	2009	2010	2011	2012	2013	2014	2015	2016	Never	All
Panel A. University characteristics											
# states	1	7	3	3	2	6	2	1	1	1	27
# federal universities	3	15	6	3	3	24	2	1	1	1	59
# 2009 enrollees	7,375	37,381	13,574	7,555	7,549	53,115	5,238	4,632	309	3,309	140,037
Mean university size	2,458	2,492	2,262	2,518	2,516	2,213	2,619	4,632	309	3,309	2,373
Mean cutoff score (2016)	670	682	669	656	657	711	660	707	670	655	688
Panel B. Characteristics of 2009 enrollees											
Age at enrollment	23.64	24.53	23.80	24.31	24.17	24.15	23.60	24.38	26.81	24.76	24.21
Female	0.49	0.46	0.48	0.50	0.45	0.49	0.48	0.44	0.51	0.48	0.48
White	0.42	0.67	0.33	0.46	0.40	0.58	0.46	0.85	0.40	0.30	0.56
Black	0.07	0.07	0.40	0.08	0.07	0.09	0.08	0.04	0.15	0.10	0.12
Brown	0.46	0.19	0.24	0.44	0.51	0.31	0.43	0.09	0.41	0.57	0.29
Private high school	0.40	0.43	0.60	0.26	0.74	0.52	0.58	0.79	0.87	0.62	0.52

Notes: This table display characteristics of federal universities and their student bodies by their state’s ENEM adoption year. Columns (B)–(K) categorize the federal universities by the year in which their state adopted the ENEM exam, τ_s^* , as defined in Section 3.3. Column (L) includes all federal universities. The ENEM adoption years for each state are:

- 2008: Pernambuco.
- 2009: Amazonas, Espirito Santo, Maranhão, Mato Grosso, Paraná, Rio Grande do Sul, São Paulo.
- 2010: Ceará, Mato Grosso do Sul, Paraíba.
- 2011: Acre, Alagoas, Piauí.
- 2012: Rio Grande do Norte, Rondônia.
- 2013: Bahia, Distrito Federal, Minas Gerais, Pará, Rio de Janeiro, Roraima.
- 2014: Goiás, Tocantins.
- 2015: Santa Catarina.
- 2016: Amapá.
- Never: Sergipe.

Data on enrollment size and student characteristics are from the Brazilian higher education census. In Panel A, the number of universities, the number of enrollees, and the mean university size are defined using new 2009 enrollees in bachelor’s programs at federal universities. In Panel B, some demographic variables are missing in the 2009 census year, so we compute student characteristics using students who enrolled in 2009 but appear in any census year between 2009–2018.

The mean cutoff score (2016) is from a public data request from the centralized admission platform SISU (*Sistema de Seleção Unificada*). These averages correspond to non-reserved quotas for bachelor’s degree programs at federal universities in the year 2016. The cutoff scores are typically weighted averages of ENEM scores in up to five subjects (math, language arts, natural science, social science, and writing) and are presented in ENEM scale score units. We obtained the SISU data in March 2020 at:

<http://www.consultaesic.cgu.gov.br/busca/dados/Lists/Pedido/Item/displayifs.aspx?List=0c839f31-47d7-4485-ab65-ab0cee9cf8fe&ID=518622&Web=88cc5f44-8cfe-4964-8ff4-376b5ebb3bef>.

TABLE A5. Effects of ENEM adoption on private-public test score gaps in alternative samples

	(A)	(B)	(C)	(D)	(E)
Dependent variable	Main sample	Appear any year	Appear pre-ENEM	Participation pre-ENEM	Participation all years
Math score	0.158* (0.079)	0.097 (0.058)	0.112* (0.060)	0.158** (0.070)	0.127* (0.069)
Language arts score	0.076*** (0.026)	0.048*** (0.017)	0.054*** (0.013)	0.088*** (0.028)	0.072*** (0.024)
Natural science score	0.065* (0.034)	0.047 (0.032)	0.046 (0.030)	0.102** (0.036)	0.080** (0.035)
Social science score	0.081*** (0.023)	0.066** (0.026)	0.069*** (0.025)	0.119** (0.044)	0.103** (0.043)
Average score (core subjects)	0.110** (0.040)	0.074** (0.035)	0.081** (0.032)	0.135*** (0.043)	0.110** (0.040)
Writing score	0.102* (0.058)	0.124** (0.048)	0.129** (0.049)	0.037 (0.105)	0.006 (0.119)
<i>N</i> (# exam takers)	2,512,214	10,991,098	6,774,892	856,165	718,748
Appear in INEP:	All years	Any year	At least 2005–2008	At least 2005–2008	All years
Participation rate:				≥ 50% in 2005–2008	≥ 50% in all years

Notes: This table examines the robustness of our results on private/public test score gaps in alternative high school samples. The dependent variables are ENEM subject scores in SD units. “Average score (core subjects)” is the average score across math, language arts, natural science, and social science. Each column estimates β^{gap} coefficients from equation (2) with a different underlying sample of high schools.

Column (A) replicates the estimates in the main sample, which consists of high schools that appeared in the INEP report in each year from 2005 to 2015. Column (B) relaxes the selection criterion and includes high schools that appeared at least once in the INEP report from 2005 to 2015. Column (C) includes high schools that appeared every year in the pre-ENEM period, i.e., from 2005 to 2008. In addition to the requirements in column (C), column (D) requires a participation rate of over 50% every year in the pre-ENEM period, meaning that at least 50% of the seniors in the high school took the ENEM in those years. Finally, column (E) requires that the high schools appeared in the INEP report each year from 2005 to 2015 and maintained a participation rate of over 50% in all years.

Note that since INEP requires a minimum participation rate for the high schools to be included in the report in certain years, columns (A)-(C) have this innate participation rate requirement. The requirements in columns (D) and (E) are in addition to the innate participation rate requirement of INEP. Details on the INEP requirement on participation rate can be found in Appendix C.3.

Parentheses contain standard errors clustered at the state level. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

TABLE A6. Effects of ENEM adoption on other test score gaps

Dependent variable	(A) Private/ public	(B) White/ Non-white	(C) Male/ Female	(D) Mom college/not	(E) Dad college/not	(F) High-/ low-income
Math score	0.158* (0.079)	0.069* (0.038)	-0.017 (0.016)	0.111* (0.058)	0.082 (0.060)	0.125* (0.067)
Language arts score	0.076*** (0.026)	0.027 (0.021)	-0.007 (0.012)	0.047 (0.034)	0.011 (0.034)	0.050* (0.025)
Natural science score	0.065* (0.034)	0.041* (0.023)	-0.019* (0.011)	0.039 (0.031)	-0.001 (0.033)	0.055** (0.025)
Social science score	0.081*** (0.023)	0.060*** (0.020)	-0.017 (0.011)	0.067** (0.030)	0.022 (0.028)	0.073*** (0.024)
Average score (core subjects)	0.110** (0.040)	0.057** (0.024)	-0.017 (0.012)	0.076* (0.041)	0.033 (0.039)	0.088** (0.038)
Writing score	0.102* (0.058)	0.005 (0.033)	-0.014 (0.014)	0.052 (0.039)	0.066 (0.045)	0.061* (0.032)
<i>N</i> (# exam takers)	2,512,214	2,387,052	2,512,214	2,489,743	2,489,191	2,487,270
Pre-ENEM average score gap	1.229	0.477	0.271	0.882	1.027	0.892

Notes: This table examines the effects of ENEM adoption on test score gaps between other historically advantaged and disadvantaged groups. The dependent variables are ENEM subject scores in SD units. “Average score (core subjects)” is the average score across math, language arts, natural science, and social science. The sample includes ENEM test takers in our high school senior sample (column C of Table 1). Some columns have slightly smaller samples due to some exam takers having missing demographic variables.

Coefficients in this table are estimated at the individual level instead of at the high school level. Column (A) displays β^{gap} coefficients from equation (2) at the individual level and replicates the estimates in column (E) of Table 3. Columns (B)-(F) estimates β^{gap} coefficients from analogous equations where the covariates are interacted with the indicator for being in the advantaged group, as indicated by the column headings. The bottom row shows the pre-existing test score gap in the average score (core subjects) in exam cohorts prior to each state’s ENEM adoption year, i.e., cohorts with $\text{HighStakes}_{st} = 0$.

Column (B) defines “non-white” to include black, brown, and indigenous students. Since yellow (Asian) students are not typically considered historically disadvantaged in educational contexts, we exclude them from this analysis. Column (D) and (E) defines “college” as having college or post-graduate degrees. Column (F) defines a student as “high-income” if his/her reported family income was greater or equal to two times the minimum wage in the year of the exam.

Parentheses contain standard errors clustered at the state level. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

TABLE A7. Robustness to controls for affirmative action adoption

	(A)	(B)	(C)	(D)
		Controls for affirmative action adoption		
Covariates	Benchmark model	Federal universities	Federal & state univ.	All universities
Panel A. Math score				
ProportionENEM _{s(h)t} × Private _h	0.158* (0.079)	0.151* (0.078)	0.157* (0.079)	0.158* (0.080)
ProportionAA _{s(h)t} × Private _h		0.005 (0.088)	0.039 (0.122)	0.029 (0.116)
Panel B. Language arts score				
ProportionENEM _{s(h)t} × Private _h	0.076*** (0.026)	0.080*** (0.023)	0.074*** (0.024)	0.084*** (0.026)
ProportionAA _{s(h)t} × Private _h		-0.109** (0.051)	-0.105 (0.071)	-0.126 (0.076)
Panel C. Natural science score				
ProportionENEM _{s(h)t} × Private _h	0.065* (0.034)	0.062* (0.033)	0.063* (0.034)	0.066* (0.035)
ProportionAA _{s(h)t} × Private _h		-0.024 (0.061)	-0.043 (0.085)	-0.016 (0.072)
Panel D. Social science arts score				
ProportionENEM _{s(h)t} × Private _h	0.081*** (0.023)	0.075*** (0.024)	0.080*** (0.022)	0.076*** (0.025)
ProportionAA _{s(h)t} × Private _h		0.029 (0.063)	0.036 (0.089)	0.091 (0.080)
Panel E. Average score (core subjects)				
ProportionENEM _{s(h)t} × Private _h	0.110** (0.040)	0.106** (0.040)	0.108** (0.040)	0.111** (0.042)
ProportionAA _{s(h)t} × Private _h		-0.029 (0.067)	-0.021 (0.096)	-0.006 (0.085)
Panel F. Writing score				
ProportionENEM _{s(h)t} × Private _h	0.102* (0.058)	0.083 (0.052)	0.103* (0.054)	0.086 (0.054)
ProportionAA _{s(h)t} × Private _h		0.199** (0.078)	0.206* (0.106)	0.285** (0.133)
<i>N</i> (# exam takers)	2,512,214	2,512,214	2,512,214	2,512,214

Notes: This table examines the robustness of our results on private/public school test score gaps to controls for affirmative action adoption. The sample includes ENEM test takers in our high school senior sample (column C of Table 1). The dependent variables are the ENEM subject scores listed in the panel titles (in SD units). Column (A) replicates our benchmark results from column (E) of Table 3, which are the β^{gap} coefficients on the interaction between ProportionENEM_{s(h)t} and a dummy for private high schools, Private_h from equation (2). In columns (B)–(D) we add in a measure of the adoption of affirmative action at the state × year level, ProportionAA_{s(h)t}, and its interaction with Private_h. We compute ProportionAA_{s(h)t} as the proportion of all new enrollees in state *s*(*h*) and year *t* who were admitting through reserved quotas using higher education census data. Columns (B)–(D) define ProportionAA_{s(h)t} using only federal universities, federal and state universities, and all universities, respectively.

Parenteses contain standard errors clustered at the state level.

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

TABLE A8. Heterogeneity by topic area — Language arts, natural science, and social science
 Dependent variable: Proportion correct answers

(A) Question group	(B) N_q	(C)		(D)	(E)	(F)
		Mean	β /mean	Public students	Private/public gap	β^{gap} /mean
Panel A. Language arts						
Communication (1–4)	54	0.468	−0.007 (0.020)	0.193	0.100 (0.047)**	
Foreign language (5–8)	80	0.424	−0.022 (0.031)	0.234	−0.022 (0.036)	
Body language (9–11)	29	0.521	−0.034 (0.019)*	0.139	0.114 (0.095)	
Art (12–14)	42	0.425	0.012 (0.028)	0.188	0.103 (0.038)**	
Literary text (15–17)	63	0.342	−0.012 (0.016)	0.157	0.143 (0.037)***	
Linguistics (18–20)	41	0.449	−0.007 (0.034)	0.193	0.104 (0.062)	
Argumentation (21–24)	63	0.450	0.000 (0.010)	0.171	0.073 (0.033)**	
Portuguese (25–27)	39	0.410	−0.027 (0.019)	0.205	0.082 (0.049)	
Social communication (28–30)	33	0.455	0.003 (0.023)	0.186	0.127 (0.035)***	
All coefficients equal (p value)			0.035		0.005	
Panel B. Natural science						
Human constructions (1–4)	57	0.274	0.007 (0.032)	0.166	0.077 (0.037)**	
Technology (5–7)	34	0.255	0.019 (0.020)	0.102	0.101 (0.038)**	
Environmental conservation (8–12)	58	0.360	0.013 (0.017)	0.178	0.020 (0.026)	
Ecosystems (13–16)	54	0.316	−0.012 (0.024)	0.204	0.091 (0.056)	
Scientific methods (17–19)	49	0.297	−0.019 (0.034)	0.189	0.032 (0.044)	
Physics (20–23)	57	0.269	−0.037 (0.031)	0.158	0.054 (0.027)*	
Chemistry (24–27)	65	0.246	−0.002 (0.010)	0.127	0.045 (0.054)	
Biology (28–30)	31	0.416	−0.029 (0.033)	0.192	0.093 (0.071)	
All coefficients equal (p value)			0.017		0.017	
Panel C. Social science						
Culture (1–5)	73	0.385	−0.027 (0.016)*	0.184	0.122 (0.040)***	
Geography (6–10)	66	0.353	0.011 (0.016)	0.214	0.030 (0.029)	
Social institutions (11–15)	74	0.376	−0.017 (0.020)	0.173	0.123 (0.029)***	
Technology (16–20)	61	0.376	−0.027 (0.014)*	0.186	0.049 (0.035)	
Citizenship (21–25)	62	0.432	0.004 (0.016)	0.192	0.078 (0.049)	
Society and nature (26–30)	69	0.405	−0.001 (0.019)	0.168	0.080 (0.043)*	
All coefficients equal (p value)			0.000		0.000	

Notes: This table shows how the increase in ENEM stakes impacted students’ performance on different topic areas of the language arts (Panel A), natural science (Panel B), and social science (Panel C) exams. This table is analogous to the math exam results in Panel B of Table 5. The sample includes 2009–2017 ENEM test takers in our high school senior sample (column C of Table 1). Regressions are at the high school (h) \times year (t) \times exam question (q) level. The dependent variable is the proportion of correct answers in each htq cell. We estimate regressions separately for questions in the topic areas of each subject as defined by ENEM test designers.

Column (A) defines the group of questions for each regression. Column (B) shows the number of questions in each group. Column (C) shows the mean proportion of correct answers for public school students in cohorts prior to each state’s ENEM adoption year (i.e., cohorts with $\text{HighStakes}_{st} = 0$). Column (E) shows the mean private/public gap in the proportion of correct answers in those cohorts. Columns (D) and (F) display the β and β^{gap} coefficients from equation (2) estimated for each group of questions. We normalize β and β^{gap} so that they represent a percentage change from the means in columns (C) and (E). The last row of each panel reports p values from F tests that the topic area coefficients in columns (D) or (F) are jointly equal.

Parentheses contain standard errors clustered at the state level.

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

TABLE A9. Heterogeneity by topics in a *Me Salva!* study guide — Math exam
 Dependent variable: Proportion correct answers

(A) Question group	(B) N_q	(C) Public students		(E) Private/public gap	
		Mean	β/mean	Mean	$\beta^{\text{gap}}/\text{mean}$
Panel A. Study guide words matched to question text					
Not in study guide	262	0.306	-0.022 (0.036)	0.186	0.141 (0.054)**
In study guide	143	0.267	-0.010 (0.029)	0.161	0.132 (0.048)**
cone (<i>cone</i>)	8	0.471	-0.022 (0.056)	0.210	0.114 (0.143)
cube (<i>cubo</i>)	8	0.179	-0.002 (0.046)	0.182	0.161 (0.054)***
cylinder (<i>cilindro</i>)	10	0.374	-0.028 (0.034)	0.157	0.135 (0.140)
directly (<i>diretamente</i>)	8	0.284	0.033 (0.040)	0.218	0.138 (0.045)***
median (<i>mediana</i>)	12	0.226	0.038 (0.044)	0.130	0.422 (0.164)**
parallelepiped (<i>paralelepípedo</i>)	9	0.314	0.005 (0.110)	0.342	0.126 (0.081)
possibilities (<i>possibilidade</i>)	8	0.275	-0.053 (0.023)**	0.126	0.169 (0.134)
possible (<i>possíveis</i>)	10	0.310	-0.048 (0.074)	0.168	0.068 (0.133)
prism (<i>prisma</i>)	8	0.368	-0.005 (0.042)	0.136	0.104 (0.182)
probability (<i>probabilidade</i>)	25	0.234	-0.041 (0.035)	0.131	0.213 (0.078)**
pyramid (<i>pirâmide</i>)	10	0.337	-0.103 (0.053)*	0.138	0.303 (0.121)**
rectangle (<i>retângulo</i>)	10	0.277	0.020 (0.039)	0.192	0.126 (0.073)*
square (<i>quadrado</i>)	35	0.256	-0.003 (0.020)	0.197	0.046 (0.031)
triangle (<i>triângulo</i>)	8	0.231	-0.008 (0.027)	0.120	0.313 (0.079)***
Fewer than 8 occurrences	40	0.268	-0.020 (0.025)	0.146	0.107 (0.053)*
In vs. not in study guide (p value)			0.248	0.565	
Panel B. Study guide concepts matched to solutions					
Not in study guide	231	0.317	-0.016 (0.032)	0.193	0.117 (0.054)**
In study guide	173	0.259	-0.016 (0.035)	0.155	0.164 (0.050)***
Geometric formulas	64	0.261	-0.017 (0.035)	0.174	0.149 (0.043)***
Proportions (“Rule of 3”)	12	0.330	-0.034 (0.053)	0.205	0.280 (0.051)***
Manipulating fractions	58	0.248	-0.017 (0.033)	0.161	0.179 (0.057)***
Radicals	13	0.195	0.049 (0.028)*	0.139	0.172 (0.097)*
Combinatory and statistical analysis	25	0.224	-0.014 (0.019)	0.103	0.232 (0.141)
Probability	36	0.255	-0.041 (0.038)	0.138	0.183 (0.054)***
Trigonometric formulas	8	0.231	0.027 (0.028)	0.104	0.166 (0.064)**
In vs. not in study guide (p value)			0.964	0.075	

Notes: This table shows how the increase in ENEM stakes impacted students’ performance on math questions that are covered in a study guide by the test prep company *Me Salva!*. The sample includes 2009–2017 ENEM exam takers in our high school senior sample (column C of Table 1). Regressions are at the high school (h) \times year (t) \times exam question (q) level. The dependent variable is the proportion of correct answers in each htq cell. Panel A defines question groups based on whether the question text contains key words from the *Me Salva!* study guide, restricting to words that appear in 8+ questions. Panel B defines question groups based on the whether the text of the solution (prepared by another test prep company, *Descomplica*) requires concepts from the *Me Salva!* study guide. See Appendix C.4 for details on these data sources and our match between *Me Salva!* topics and ENEM questions.

Column (A) defines the group of questions for each regression. Column (B) shows the number of questions in each group. Column (C) shows the mean proportion of correct answers for public school students in cohorts prior to each state’s ENEM adoption year (i.e., cohorts with $\text{HighStakes}_{st} = 0$). Column (E) shows the mean private/public gap in the proportion of correct answers in those cohorts. Columns (D) and (F) display the β and β^{gap} coefficients from equation (2) estimated for each group of questions. We normalize β and β^{gap} so that they represent a percentage change from the means in columns (C) and (E). In both panels, the last row reports p values from an F test that the coefficients in the first and second rows are equal. Parentheses contain standard errors clustered at the state level.

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

B. THEORETICAL APPENDIX

This theoretical appendix presents a simple framework that helps to interpret our main empirical results and shed light on our analysis of the underlying mechanisms.

B.1. Framework setup. We consider a population of exam takers that are characterized by their socioeconomic status (SES) and their abilities. Let X_i denote an observable measure of individual i 's SES, e.g., an indicator for attending a private high school. We let a_i denote individual i 's ability for performing well on a college admission exam, which is not directly observable in data. We refer to a_i as *test ability* to emphasize that it may be distinct from other abilities that help the individual perform well in college and in the labor market.

We assume that the stakes of the college admission exam vary across cohorts, and that test takers in high-stakes cohorts engage in more test prep. For simplicity, we suppose that individuals are randomly assigned to either a low-stakes or a high-stakes exam cohort, and we let H_i be a binary indicator for the high-stakes cohort.²⁷ Individuals have a stronger incentive to perform well in the high-stakes exam cohort, and so we assume that these individuals engage in additional test prep. This additional prep increases their test ability by an amount that we denote by e_i . We interpret e_i broadly; for example, it may include studying test prep books, taking preparatory courses, focusing more intently during the school year, or exerting more effort on the exam.

Thus a test taker's preparedness for the exam, as a function of their cohort, is given by:

$$(B1) \quad \theta_i = a_i + H_i e_i.$$

where we refer to θ_i as *test skill*. In other words, individuals in low-stakes cohorts have test skill $\theta_i = a_i$, while individuals in high-stakes cohorts have test skill $\theta_i = a_i + e_i$.

An individual's score on the college admission exam is a noisy measure of their test skill. We let T_i denote individual i 's test score, which is observable in the data. We assume test scores are given by:

$$(B2) \quad T_i = \theta_i + \epsilon_i^T,$$

where ϵ_i^T is random noise that reflects variation in test performance due to factors like guessing and health on exam day.

Our interest is in the predictive power of the test score for measures of college success that matter to both individuals and colleges. We let Y_i denote an observable measure of college success, e.g., persisting in college after enrolling or completing a college degree. We assume Y_i is given by:

$$(B3) \quad Y_i = \alpha a_i + \beta H_i e_i + \gamma X_i + v_i + \epsilon_i^Y.$$

²⁷ Our empirical analysis relies on a parallel trends version of this assumption.

We allow college success to potentially depend on test ability, a_i , test prep, e_i , and SES, X_i . The parameter α represents the effect of test ability on college success; it is natural to assume $\alpha > 0$ since individuals with high test scores tend to perform better in college. Similarly, an individual’s SES may help them succeed in college above and beyond their test ability, and so we assume $\gamma > 0$. The parameter β allows for the possibility that test prep directly affects college performance, although, as we discuss below, the sign of this parameter is less clear. In addition to these three factors, we allow college success to depend on other abilities that are unrelated to test skill, which we denote by v_i , and a random noise term, ϵ_i^Y .

B.2. Effects of exam stakes on test score gaps. In Section 4, we find that increasing the stakes of a college admission exam increases test score gaps between high- and low-SES students. In our framework, this result can be written as:

$$(B4) \quad \text{cov}(T_i, X_i | H_i = 1) > \text{cov}(T_i, X_i | H_i = 0).$$

In other words, the covariance between test scores, T_i , and SES, X_i , is larger in the high-stakes cohort ($H_i = 1$) than in the low-stakes cohort ($H_i = 0$).

Plugging in equations (B1)–(B2) and simplifying, expression (B4) can also be written as:

$$(B5) \quad \text{cov}(X_i, e_i) > 0.$$

Expression (B5) states that high SES students engage in more test prep than low SES students when the stakes of the exam increase.

B.3. Effects of exam stakes on predictive validity. In Section 5, we examine how an increase in stakes impacts the predictive power of exam scores for college success, as well as the potential mechanisms for this effect.

In Table 6, we find that scores from higher stakes exams have more predictive power for various measures of college success. In our framework, this result can be written as:

$$(B6) \quad \text{cov}(T_i, Y_i | H_i = 1) > \text{cov}(T_i, Y_i | H_i = 0).$$

Using equations (B1)–(B3), expression (B6) simplifies to:

$$(B7) \quad \gamma \text{cov}(X_i, e_i) + (\alpha + \beta) \text{cov}(a_i, e_i) + \text{cov}(v_i, e_i) + \beta \text{var}(e_i) > 0.$$

Expression (B7) shows that there are three broad channels through which additional incentives for test prep could increase the predictive validity of exam scores:

- (1) **SES correlation:** $\gamma \text{cov}(X_i, e_i) > 0$. First, the high-stakes exam scores may be more predictive of college outcomes simply because they are more correlated with SES, which also tends to benefit students in college. For example, wealthy students may have greater access to test prep services, and family wealth also may help students

succeed in college. Since we find that $\text{cov}(X_i, e_i) > 0$, then all else equal, one would expect that higher-stakes exam scores would have more predictive validity.

- (2) **Ability correlation:** $(\alpha + \beta)\text{cov}(a_i, e_i) + \text{cov}(v_i, e_i) > 0$. Second, the predictive validity of scores may increase if the induced test prep is correlated with ability. Students who have higher test ability on the low-stakes exam may engage in more test prep when the stakes increase, i.e., $\text{cov}(a_i, e_i) > 0$. Additionally, students with higher ability to succeed in college ability may engage in more test prep for the higher-stakes exam, i.e., $\text{cov}(v_i, e_i) > 0$. The sign of these covariance terms is *a priori* less clear. On the one hand, students who know that they are likely to succeed in college may have the strongest incentives to prep for the higher-stakes exam. On the other hand, test prep may disproportionately benefit high-income but lower-ability students who are unlikely to be admitted to top colleges without prepping. Thus we do not have a strong prior on whether an ability correlation plays a role in the observed increase in exam score validity.
- (3) **Skill accumulation:** $\beta\text{var}(e_i) > 0$. Finally, test prep may be directly beneficial for college success. This channel is operative if the test skills that individuals accumulate from prepping, e_i , directly improve their college outcomes, Y_i . In our framework, this skill accumulation channel exists if $\beta > 0$. High-stakes exams are often criticized for creating incentives to engage in test-oriented learning that is not useful outside the exam. Thus it is possible that $\beta = 0$. It is also possible that $\beta < 0$ if test prep crowds out other useful learning.

There is a fourth potential channel that is outside the scope of our framework: high-stakes exam scores may have more predictive validity if they increase the “match quality” between individuals and college programs. This channel arises because an individual’s test score, T_i , may have a causal effect on their outcome, Y_i , through its influence on which college and/or major they attend. Thus high-stakes exam scores may be more informative for college success if the distribution of these scores leads to better student/college matches, e.g., on the basis of academic preparation. To distinguish between this channel and the three mentioned above, we follow the standard practice that testing agencies use to measure predictive validity. Specifically, in addition to estimating raw correlations between test scores, T_i , and outcomes, Y_i , we also estimate correlations after de-meaning each variable *within* college programs.

Distinguishing between the SES correlation, ability correlation, and skill accumulation channels is challenging because of unobserved abilities that impact both test scores and college success (i.e., a_i and v_i). The next subsection describes how we present some evidence on the skill accumulation channel. But regardless of which channel is at play, our finding that test prep increases the informativeness of scores can explain why many colleges around the world use high-stakes tests for admissions. Similarly, our findings shed light on why

many U.S. colleges prefer to use SAT or ACT “superscores” (Goodman et al., 2020), which are based on the maximum of each subject score across all of the student’s test attempts. Relative to average scores or scores from the first attempt, superscores are more correlated with the amount of test prep, which may increase the informativeness of scores for outcomes that colleges care about.

C. EMPIRICAL APPENDIX

C.1. **Variable definitions.** This section describes the main variables in our paper.

C.1.1. *Test scores.*

- **Subject scores.** The post-2009 ENEM scores, as reported to the public, are scaled to have a mean of 500 and an SD of 100 in the population of 2009 high school seniors who took the exam. Throughout the paper, we report ENEM scores in SD units relative to this population. For ENEM scores in math, language arts, natural science, and social science, our transformation is:

$$\text{Transformed subject score} = \frac{\text{Raw subject score} - 500}{100},$$

After transformation, a score of zero in our paper is equivalent to the performance of the average high school senior who took the ENEM in 2009, and a score of one is 1 SD higher within this population. These transformations preserve the comparability of test scores across cohorts.

The 2007–2008 ENEM reported only a single core-component score plus a writing score. To define scores for each subject, we first categorize the multiple choice questions into math, language arts, natural science, and social science, and then compute a separate score for each subject using the IRT parameters estimated from the response data.

Since the reference populations differ for the 2007–2008 and 2009–2017 exams, in regressions where we use scores from both periods, we standardize the scores to have mean 0 and SD 1 within each year of our sample.

- **Average scores (core subjects).** The post-2009 average scores are calculated by taking the average of four subject scores, and then standardize relative to the reference population. In practice, our transformation is:

$$\text{Transformed average score} = \frac{\text{Raw average score} - 500}{86.7},$$

where 86.7 is the SD of the average score in the reference population. After transformation, a score of zero in our paper is equivalent to the performance of the average high school senior who took the ENEM in 2009, and a score of one is 1 SD higher within this population.

For the 2007–2008 ENEM, the average score is defined as the single core-component score. In regressions where we use scores from both the 2007–2008 and 2009–2017 exams, we standardize the scores to have mean 0 and SD 1 within each year of our sample.

- **Writing scores.** The post-2009 writing score is also standardize relative to the reference population. In practice, our transformation is:

$$\text{Transformed writing score} = \frac{\text{Raw writing score} - 597}{137},$$

In regressions where we use scores from both the 2007–2008 and 2009–2017 exams, we standardize the writing scores to have mean 0 and SD 1 within each year of our sample.

C.1.2. *Exam-taker characteristics.* These variables were collected from a survey that applicants completed as part of the ENEM exam process.

- **Race.** In Brazil, race is commonly classified in five groups: *branco* (white), *pardo* (brown), *preto* (Black), *amarelo* (yellow), and indigenous. Since Asian and indigenous people represent a small proportion of the population in Brazil (less than 3 percent in our sample), we use indicators for three major racial groups: *branco* (white), *pardo* (brown), and *preto* (Black). We set the indicator variables to missing if the students declined to declare their racial identities (2 percent of the students in our data).
- **Parental education.** The measures for mother’s and father’s education consist of 8 categories from “none” to “post-graduate”. From these categorical variables, we derive two indicator variables “Mother attended college” and “Father attended college”, which equal 1 if the respective parent achieves an educational level of “college” or “post-graduate”.
- **Family income.** Family income is measured as multiples of the minimum wage in the year of the exam. We define an indicator variable “Family income > 2x min. wage” (or “High-income”), which equals 1 if the reported income is more than twice the minimum wage.
- **Private high school.** Throughout our analysis, we define “private high schools” to include both private and federal high schools (0.5 percent of students) since their students are comparable in terms of socioeconomic status and achievement. In contrast, “public high schools” include both state and municipality high schools.

C.1.3. *College and labor market outcomes.* These variables were collected from INEP’s higher education census (*Censo da Educação Superior*) for the years 2010–2019 (INEP, 2022) and Brazil’s employee-employer dataset, the RAIS (*Relação Anual de Informações Sociais*), for the years 2016–2018 (RAIS, 2022).

- **College enrollment.** We define an indicator variable “Ever enrolled in college” which equals 1 if the student has a record in the INEP’s higher education census for the years 2010–2019.

- **College persistence.** The higher education census contains information on a student’s enrollment year and last year on record. We define an indicator variable “Persisted in college for 3 years” which equals 1 if the student’s last year on record is greater or equal to 3 years after the enrollment year.
- **College graduation.** The higher education census contains information on a student’s enrollment and graduation year. We define an indicator variable “Ever graduated college” which equals 1 if the student has a non-missing graduation year. We define an indicator variable “Graduated college within 5 years” which equals 1 if the student’s graduation year is within 5 years of the enrollment year.
- **Fraction of college credits completed.** The higher education census contains information on a student’s credits completed in his/her program and the total credits required in the program curriculum. We define the variable “Fraction of college credits completed” as the credits completed on the student’s last record in the census divided by the total credits in the program.
- **Appears in RAIS.** “Appears in RAIS” is an indicator variable that equals 1 if a given student has a matching record in RAIS, which indicates that the student has engaged in formal employment.
- **Hourly wage (BRL).** We compute an individual’s hourly wage as average monthly earnings divided by average monthly contracted hours, both available in RAIS. The wages are expressed in terms of Brazilian Reals.

C.2. Data and merging. Our base dataset contains all individuals who took the ENEM exam in 2007–2017. This dataset includes student-level and question-level information. The student-level data includes self-reported demographic and socioeconomic status (SES) measures, such as sex, race, high-school type (public/private), parental education, and family income. The question-level data includes student responses to each exam question, the question subject, and skill tested. From this dataset, we only keep high-school students with a valid score (i.e., non-zero and non-missing) on each subject test. These restrictions exclude, for example, individuals who took the exam after graduating from high school or who missed one of the testing days.

To measure long-run outcomes, we combine the 2009–2014 ENEM records with two individual-level administrative datasets using individuals’ national ID numbers (*Cadastro de Pessoas Físicas*). The linkage was conducted in the secured data room at the INEP facilities in Brasilia, Brazil. We exclude students with missing national ID (0.04 percent) and those who took the exam more than three times in our data (0.10 percent).

We measure college outcomes using Brazil’s higher-education census from 2010–2019. This dataset offers comprehensive information about all college enrollees, including their university

of enrollment, major, the academic year when they enrolled, and their year of graduation. 65.6 percent of high-school seniors taking the ENEM during 2009–2014 appear in the census data.

We measure labor-market outcomes using an administrative employee-employer matched dataset called RAIS (*Relação Anual de Informações Sociais*) from 2016–2018. The RAIS contains data on workers employed in the formal sector. It does not include data on individuals working within the informal sector, those who are self-employed, or individuals who are currently unemployed. This dataset includes worker-level and firm-level information. The worker-level data includes educational attainment, occupation, and earnings. The firm-level data includes total employee count, the industry they operate within, and their geographical location. 32.9 percent of high-school seniors taking the ENEM during 2009–2014 appear in the RAIS. The relatively low match rate might be attributed to the fact that some individuals could still be enrolled in college. However, even when considering individuals who took the ENEM in 2009, the match rate remains comparably low with only 31.7 percent of them being matched to the RAIS.

C.3. Sample definition. This section describes criteria needed for high schools to be included in the INEP annual reports we leverage to create our sample.

As noted in the main text, our *high school graduate sample* consists of the set of high schools that were in all the yearly ENEM-performance reports created by the INEP during 2005-2015. Only high schools that meet two main conditions are included in the annual reports. First, the high school needs to have at least 10 test-takers who declared that they would graduate that year. Before 2011, all individuals taking the ENEM were recorded as test-takers for the calculation of the participation rate, regardless of whether they completed the test or not. Since 2011, to be considered a test-taker, an individual has to complete the four subject tests plus the writing essay, and obtain a non-zero score in all subjects.²⁸ Second, starting in 2009, the INEP required a minimum *participation rate* to be included in the report. The participation rate is the total number of ENEM test-takers in the high school divided by the number of enrolled students in the final years of high school, based on the records of the High School Census. Between 2009 and 2010, the minimum participation to be included in the report rate was 2 percent. In 2011, the minimum participation threshold increased to 50 percent.

²⁸ Only the following grades are considered in the report: 3rd and 4th grade of regular high school (*ensino médio regular 3a e 4a série*), 3rd and 4th grade of teaching track high school (*ensino médio magistério 3a e 4a série*), non-serialized regular and teaching track high school (*ensino médio não-seriado, regular e magistério*), vocational education and high school for youth and adult education (*educação profissionalizante e ensino médio para educação de jovens e adultos*).

TABLE C1. Topic areas and competencies for ENEM math exam

(A)	(B)	(C)	(D)
Topic area	Competency (and reference number)	Topic area	Competency (and reference number)
Numbers	Recognize numbers (1)	Proportions	Solve problems using proportions (16)
Numbers	Identify numerical patterns (2)	Proportions	Use proportions to construct arguments (17)
Numbers	Solve problems using numbers (3)	Proportions	Evaluate interventions using proportions (18)
Numbers	Use numbers to construct arguments (4)	Algebra	Identify algebraic relationships (19)
Numbers	Evaluate interventions using numbers (5)	Algebra	Interpret Cartesian graphs (20)
Geometry	Project 3D objects into 2D space (6)	Algebra	Solve problems using algebra (21)
Geometry	Identify geometric shapes (7)	Algebra	Use algebra to construct arguments (22)
Geometry	Solve problems using geometry (8)	Algebra	Evaluate interventions using algebra (23)
Geometry	Use geometry to construct arguments (9)	Interpreting data	Make inferences using data in tables/graphs (24)
Measurements	Identify units of measurement (10)	Interpreting data	Solve problems using data in tables/graphs (25)
Measurements	Use scales in everyday situations (11)	Interpreting data	Use tables/graphs to construct arguments (26)
Measurements	Solve problems using magnitudes (12)	Statistics	Calculate statistical quantities from data (27)
Measurements	Use measurements to construct arguments (13)	Statistics	Solve problems using statistics (28)
Measurements	Evaluate interventions using measurements (14)	Statistics	Use statistics to construct arguments (29)
Proportions	Identify proportional relationships (15)	Statistics	Evaluate interventions using statistics (30)

Notes: This table shows the 7 topic areas (columns A and C) and 30 competencies (columns B and D) for the ENEM math exam. Labels are translated and shortened by the authors from the descriptions in `Matriz_Referencia_Enem.pdf`, which is included with the microdata.

C.4. Categorization of math questions. This subsection provides details on the categorization of math questions that we use for the heterogeneity analyses in Table 5 and Appendix Table A9.

In Table 5, we categorize questions into 7 *topic areas* (Panel B) and 30 *competencies* (Panel C) defined by ENEM test designers. These topics areas and competencies are in the `ITENS_PROVA_****.TXT` files of the microdata. The labels for each topic area and competency are defined in `Matriz_Referencia_Enem.pdf`, which is included with the microdata and is also available online at the link in this footnote.²⁹ Table C1 shows the translated and shortened labels that we use for Table 5. We also present results by topic area for language arts, natural science, and social science in Appendix Table A8.

In Appendix Table A9, we define groups of questions based on whether the questions are related to topics covered in a study guide created by *Me Salva!*, which is a well-known Brazilian test prep company. The study guide is called *The Approved Book: One topic per day to pass the ENEM*.³⁰ In Panel A of Appendix Table A9, we match key words from the *Me Salva!* study guide to the text of each question. In Panel B of Appendix Table A9, we match concepts from the *Me Salva!* study guide to solutions to each question that were created by another well-known test prep company called *Descomplica*. *Descomplica* creates solutions to

²⁹ See: https://download.inep.gov.br/download/enem/matriz_referencia.pdf (accessed in June 2023).

³⁰ The Portuguese title is *O Livro do Aprovado: Um conteúdo por dia para passar no ENEM*. See: <https://cdn.mesalva.com/uploads/medium/attachment/MS2018-livro-do-aprovado.pdf> (accessed in June 2023).

each question from previous ENEM exams and makes them publicly-available online to help students prepare.³¹

We match key words and concepts in the *Me Salva!* study guide to ENEM questions using text analysis. The *Me Salva!* study guide contains tips for solving questions in seven different content areas, and in each content area there are key words that appear in bold in the text. For Panel A of Appendix Table A9, we search the question text for the key words from the study guide. For Panel B of Appendix Table A9, we search the *Descomplica* solutions for both the key words and for regular expressions that indicate questions in which the solution is likely to depend on concepts from the study guide. The content areas and search terms that we use are as follows; the key words that we use in both Panels A and B appear in italics, and the regular expression searches that we use in Panel B appear in plain text:

- **Geometric formulas.** *cilindro, cone, cubo, equilatero, esfera, hexagonal, hexagono, losango, paralelepipedo, piramide, prisma, quadrado, quadrangular, retangulo, trapezio, triangular, triangulo.*
- **Proportions (“Rule of 3”).** *diretamente, grandezas, inversamente, proporcionais, proporção, regra de três.*
- **Manipulating fractions.** Solutions that contain at least two fractions ($[\text{0-9}]/[\text{0-9}]$) and also an equals sign (=).
- **Radicals.** Solutions that contain a square or cube root sign ($\sqrt{\quad}$ or $\sqrt[3]{\quad}$).
- **Combinatory and statistical analysis.** *arranjo, combinação, mediana, moda, permutação, possibilidades.* Solutions that contain arrangement or combination notation ($A[\text{0-9}],[\text{0-9}]$ or $C[\text{0-9}],[\text{0-9}]$).
- **Probability.** *combinação, possíveis, probabilidade.* Solutions that contain combination notation ($C[\text{0-9}],[\text{0-9}]$).
- **Trigonometric formulas.** *cos, cossec, cosseno, cotg, sec, sen, seno, tangente, tg.*

In all cases we trim plural, adjective, and masculine/feminine endings to words before matching. In Panel A, we display results for key words that appear in eight or more questions, and we group all other words into the “Fewer than 8 occurrences” category. In Panel B, we display results separately for each of the seven *Me Salva!* content areas. In both panels, we also show results pooling across all questions that do/don’t match any search term in the study guide.

³¹ See: <https://descomplica.com.br/gabarito-enem/questoes/?cor=azul> (accessed in June 2023). One question is behind a paywall—Question 145 in the *Azul* book of the 2015 math exam—and thus we exclude it from our analysis in Panel B of Table A9).