

Racial Discrimination in Grading: Evidence from Brazil

FERNANDO BOTELHO
RICARDO MADEIRA
MARCOS A. RANGEL



Racial Discrimination in Grading: Evidence from Brazil

Fernando Botelho (fbotelho@usp.br)

Ricardo Madeira (rmadeira@usp.br)

Marcos A. Rangel (marcos.rangel@duke.edu)

Abstract:

We investigate whether racial discrimination taking the form of the biased assessment of students is prevalent within Brazilian schools. Robust evidence is drawn from unique data pertaining to middle-school students and educators. After holding constant performance in blindly scored tests of proficiency and behavioral traits, we find that teacher-assigned Mathematics grades suffer from cardinal and ordinal biases. We unveil strong indications that these effects result from incomplete information issues highlighted in models of statistical discrimination which are made particularly salient by social promotion schemes currently operational in our context.

Keywords: race; schooling; grading; standardized tests; statistical discrimination. .

JEL Codes: I21, J15, I24.

Racial Discrimination in Grading: Evidence from Brazil

By FERNANDO BOTELHO, RICARDO MADEIRA AND MARCOS A. RANGEL*

We investigate whether racial discrimination taking the form of the biased assessment of students is prevalent within Brazilian schools. Robust evidence is drawn from unique data pertaining to middle-school students and educators. After holding constant performance in blindly scored tests of proficiency and behavioral traits, we find that teacher-assigned Mathematics grades suffer from cardinal and ordinal biases. We unveil strong indications that these effects result from incomplete information issues highlighted in models of statistical discrimination which are made particularly salient by social promotion schemes currently operational in our context.

JEL: I21, J15, I24.

Keywords: race, schooling, grading, standardized tests, statistical discrimination.

Evidence of a negative association between individual characteristics used to infer African ancestry and educational attainment abounds.¹ Equally notorious is the resilience of achievement gaps across cohorts of Black and White children

* Botelho: University of Sao Paulo, Av. Prof. Luciano Gualberto 908, fbotelho@usp.br. Madeira: University of Sao Paulo, Av. Prof. Luciano Gualberto 908, rmadeira@usp.br. Rangel: Duke University, 194 Rubenstein Hall, 302 Towerview Drive, Durham, NC 27708-0239, USA, marcos.rangel@duke.edu. Funding from CAPES-Brazilian Ministry of Education (*Observatorio da Educacao FEA-USP, Project 3313*) and INEP-Brazilian Ministry of Education (*Nucleo de Estudos da Educacao - Fipe*) are gratefully acknowledged. We benefited from comments received at the NBER Education Meeting, the PAA Meeting, the NEUDC Conference, the Economic Demography Workshop at PAA, and from seminar participants at Duke, Paris School of Economics, Vanderbilt, Princeton, PUC-Rio, EESP-FGV, Stanford, UC-Riverside, U. of Delaware, UIC, and USP. Part of the analysis presented here was performed while Rangel was a visiting scholar at Princeton University. He is thankful for the hospitality of the Research Program in Development Studies and of the Program in Latin American Studies. Errors, omissions, and ideas expressed in the text are the sole responsibility of the authors. These do not represent the views of funding agencies. Previous versions of this article have circulated under the title: "Discrimination Goes to School? Racial Differences in Performance Assessments by Teachers."

¹Data portraying such historically-rooted patterns have been drawn from different countries and under a variety of institutional settings. For comparative international studies see Alexander et al. (2001); Herring et al. (2004); Telles (2004); and Telles and Steele (2012).

(Neal, 2006). These are further emphasized by longitudinal studies showing that Black disadvantages emerge during infancy and remain pretty much intact while children attend school.² Because evidence regarding racial differences in *expected* returns to human capital accumulation is scant, a better understanding of obstacles to the acquisition of skills and educational credentials by Blacks seems warranted.

Here, we propose discrimination within racially-integrated schools as a candidate explanation for the patterns described above and subsequently examine its prevalence in Brazil. We recognize that such a phenomenon may manifest itself in many different ways within a classroom. Yet, we focus on a very specific one: a teacher's biased evaluation of students with respect to their scholastic proficiency and aptitude (i.e.: grading). We employ uniquely detailed administrative data from the state of Sao Paulo covering approximately 277 thousand eighth-graders spread across 10.6 thousand public-school classrooms in 2010. Our inference is based on the careful contrasting of teachers' subject-specific grades and scores from end-of-year standardized (and blindly marked) proficiency tests covering the same official curriculum delivered in regular classes.

The analyses show that portions of teachers' assessments in Mathematics not explained by proficiency scores are associated with pupils' racial background. Our most conservative estimates indicate that there are statistically significant under-scoring and under-ranking of Blacks relative to Whites. The measured racial gap in promotion rates between equivalently proficient and well-behaved students corresponds to a 4.1% increase in the retention probability for the average Black. Focusing exclusively on the ordinality aspect we also uncover a gap that translates into a 4.5% reduction on the probability of Blacks being graded above the classroom median. In practice, these work as if teachers were "taxing" the average Black student's performance in proficiency tests by 0.03 to 0.04 of one standard

²See Phillips et al. (1998); Hedges and Nowell (1999); Reardon (2008); and Madeira and Rangel (2013). Cautionary notes on these findings can be found in Bond and Lang (2013).

deviation at the time her competence is being assessed. These results are shown robust to possible omissions of behavioral attributes and to the likely incidence of measurement error on scores from standardized tests used as covariates in our estimations. They are also very much in line with the expected subtlety of this particular form of discrimination.

Once the existence of racial gaps in teacher assessments is established, we rely on economic theory to examine its likely source in our context. We draw from a rich literature on statistical and screening discrimination.³ We map our setting into these studies by focusing on two institutional aspects. First, teachers are limited by imperfect screening technology in the process of scholastic competence's measurement and, once assigned to students of a given level (whose admission is decided by a third party), are solely responsible for promotion and ranking decisions. Second, due to a number of policies implemented since the late 1990's, a dramatic increase in access to public education has been observed. We highlight in particular the adoption of *social promotion* schemes between the fifth and seventh grades, which eliminated proficiency-based retention in those grades. In practice, by establishing lenient standards for the admission of students into eighth grade, such policy has disproportionately benefited Blacks (who are over-represented among pupils with lower proficiency). In other words, social promotion has operated as affirmative action within the Brazilian school system we study. Eighth-grade teachers were well aware of the implications of such policy, and priors regarding students' proficiency may have been downgraded as a result. Therefore, we hypothesize that when teachers issue report cards assessing the competence of their students, subtle biases may be generated by the weighted combination of noisy information extracted from their own screening exams and stereotyped priors.

We then present evidence on the validity of such theoretical reasoning. Em-

³Aigner and Cain (1977); Borjas and Goldberg (1978); Lundberg and Startz (1983); Coate and Loury (1993); Cornell and Welch (1996); Altonji and Pierret (2001); Blume (2006); Bjerk (2008); and Lehmann (2011).

ploying a strategy similar in spirit to the one in the seminal work of Altonji and Pierret (2001), we examine whether the duration of interaction between teachers and students produces different assessment patterns. The basic idea is that the longer pupil and teacher interact, the smaller is the role of biased priors that emphasize racial identity and the larger is the role of hard-to-measure signals of proficiency.⁴ In this regard, our empirical exercises unveil that while racial gaps in promotion rates and ranking are salient for Black and White students attending classes with a teacher for the first time, no significant disparities are found among those that have already had classroom interactions with that instructor before eighth grade. Teachers seem to learn about a student's true "type" over academic years and once they are fully aware of grading standards previously used (i.e.: their own).⁵ These findings not only are compatible with a learning model of statistical discrimination but also lend further support to our baseline results' robustness with regards to potential omission of student-specific characteristics.

Despite the subtlety of the effects, there are a number of reasons to believe they are quite relevant. The gaps we uncover are equivalent to 40% of the raw (within-classroom) grade differential associated with having a mother with a college degree or more versus a mother with a high-school degree only. Alternatively, they correspond to 30% of the effect of being a child of parents that own their home. Ultimately, following the insightful discussion in Bond and Lang (2013), our measure of racial discrimination depends on the translation of intrinsically ordinal scales into meaningful quantities, something that could only be achieved in case they were to be directly associated with longer-term outcomes such as employment and wages.

With that argument in mind, we believe that the implication of our findings can be far reaching, and certainly go beyond differentials on level promotion and on ranking of students. There is an enormous potential for feedback effects in

⁴The same notion of interactions and learning is also central in Lundberg and Startz (2007).

⁵Tests of learning in the context statistical discrimination can also be seen in Autor and Scarborough (2008); Lange (2007); List (2004); and Farber and Gibbons (1998).

our context. This is the case because we detect discrimination in grading during the transition between middle and high-schools, at a time when Brazilian parents invariably find themselves in the position of investors relying on the asset-return evaluations of more informed experts. For our purposes, the key element of this reasoning is that teacher communications may steer investment decisions in one way or the other.⁶ That is to say; parents (and children themselves) likely update investment (and effort) decisions after extracting information from report cards issued by teachers. Therefore, if children's perceived competence increases the returns or reduces the costs of investments, as in the traditional *Beckerian* human-capital framework, this mechanism can reinforce racial gaps in the accumulation of human capital. In this case, intra-classroom evaluation biases may very well lead to gaps in attainment, school choice, future scholastic performance and, ultimately, labor market outcomes.⁷

Considering the role played by misinformation in the results presented here, and beyond its scientific interest, we draw three lessons for education policy from our analysis. First, curbing teacher rotation can be particularly important for Black students (over and beyond any effect on learning *per se*) because increasing interactions between a group of students and a given teacher diminishes the influence of noise on the evaluation of scholastic proficiency. The more a teacher gets acquainted with a given student, the less relevant the pupil's race becomes for evaluation purposes.⁸ Second, direct investment in teacher training with regard to the design of exams and tests may be warranted. Well-designed questions are easier to grade and more likely to differentiate students on the most relevant dimensions of proficiency. Finally, because blindly graded proficiency tests are regularly taken by students under standard school accountability' systems, and

⁶Lam et al. (2006) examines the significant effect of performance measurement's precision over high-school dropout behavior in South Africa, for example.

⁷See Mechtenberg (2009) for a formalization of an argument like this. See also Lundberg and Startz (1983), who are explicit in modeling human capital investments' response to the presence of discrimination.

⁸This adds interesting elements to the beneficial effects of grade/subject-specific teacher experience discussed by Ost (2014).

despite the intrinsic noisy nature of such scores, the generation of individual report cards could aid teachers in their competence evaluations. Particularly under social promotion schemes like the one we study, this additional information should make teachers better able to evaluate their students without resorting to racially biased priors. Above all, public schools and respective education authorities could do a better job on their use of performance information in order to maximize efficiency. The reduction of grading discrimination of the sort we uncover would be an added bonus.

We also believe our results can be used to conjecture on the effects of affirmative action over access to education and accumulation of human capital, a theme of prime importance as Brazil adopts racial quotas in access to college and in allocation of publicly-funded scholarships. While the profession has focused on behavioral responses among those targeted by such policies (the effort choices of high-schoolers that are granted easier access to college), we advocate that the role of instructors within colleges be considered.⁹ We believe that the latter's priors could be affected by the enactment of quota policies (stereotyping). Because GPA and course performance are intimately connected with drop-out and graduation rates, ceilings may be imposed on the progress of the population it was designed to help by distorting subjectively assigned course grades.¹⁰ This is an indirect policy implication we draw from our exercises.

The remainder of this article is organized as follows. Section 1 briefly reviews the literature on teacher perceptions and grading discrimination. Section 2 discusses the institutional background and describes the data we employ. Section 3 outlines a conceptual framework that guides the empirical analysis we perform. Section 4 presents our empirical strategy and the econometric identification strategy. Results and discussions are presented in Section 5. Section 6 concludes.

⁹See Assuncao and Ferman (2013) on the early Brazilian experience with quotas, Cortes and Zhang (2012) for a discussion in the context of the Top 10% Program in Texas, and Cotton et al. (2014) for a controlled experiment.

¹⁰In many ways this is similar to the original stereotyping-affirmative action nexus proposed by Coate and Loury (1993) in the context of labor markets.

I. Related literature

Despite being the first study to examine racial bias using Brazilian student-level data in this degree of detail, we are well aware that the question of whether teachers treat Black and White children differently is not new. In fact, there is a tradition within the sociology literature of directly examining whether teacher bias is a factor in course-grade assignment in the United States (Bowles and Gintis, 1976; Farkas et al., 1990; Rist, 1973; Rosenthal and Jacobson, 1968; Sexton, 1961). Both large- (Sewell and Hauser, 1980; Williams, 1976) and small- (Leiter and Brown, 1985; Natriello and Dornbusch, 1984) scale empirical studies do not detect significant biases. There is also a considerable number of contributions from the social psychology literature focusing on teacher's perceptions of Black and White children (see Ferguson, 1998, 2003 and references therein), which again only unveils weak relationships between Black stereotypes and measures of discriminatory actions.¹¹

Our work complements more recent studies from the education and economics literatures. Shay and Jones (2006) and Dorsey and Colliver (1995), examine quasi-experimental variation provided by institution-level policy changes regarding anonymity in the grading processes applied to college/graduate students and do not detect significant racial differentials. Figlio (2005) examines whether teachers' overall perception of a given student is affected by the "Blackness" of her first name, even after controlling for performance in standardized examinations. Using data from one school district in Florida, he uncovers evidence of lower teacher expectations for those perceived to have African American ancestry. Burgess and Greaves (2013) investigate differences in teacher grading according to ethnic background using observational data from England, finding significant under-assessment of Black Caribbean and Black African pupils. Finally, Hinnerich et al (2011a, 2011b) conduct audit-like studies by transcribing and blindly re-grading

¹¹See review of studies in Dovidio et al (1996). Demeis and Turner (1978), unlike most of this literature, find significant discrimination against Blacks in an experimental setting.

tests assessed by teachers in Sweden and estimate gender (insignificant) and nationality (significant) gaps. Similar exercise conducted in Germany by Sprietsma (2013) also uncovers biases against exam solutions which had Turkish-sounding names randomly allocated to them (relative German-sounding names).

Particularly important and convincing studies detecting discrimination in grading are the ones reported by Lavy (2008) and Hanna and Linden (2012). Lavy (2008) capitalizes on a natural experiment in Israeli high-schools. The author cleverly explores the fact that students take two examinations that cover the same material and have the same format during their senior year, and that the grading of each exam happens under different anonymity regimes. Focusing on gender differentials, his findings indicate that male students receive lower marks in the non-blindly graded exams (relative to those blindly scored), and that these differences are larger (in absolute value) than among girls. Blind/non-blind contrasts are also skillfully explored in a randomized control trial designed and implemented by Hanna and Linden (2012). The authors identify statistically significant positive differences between blind and non-blind scores for members of lower castes in India (relative to upper castes), which is clear evidence of discrimination.

The discussion presented here plays on three advantages of our context with respect to other studies in the literature. First, the sheer size of and level of detail in our data base allows us to convey a complete portrait of teacher and student-body characteristics associated with discrimination in actual classroom environments. Teachers grading in experimental settings may very well reveal different discriminatory behavior due to the one-shot nature of the event (even when hypothetical biases are curbed by incentivizing schemes). Regular teachers are gatekeepers, being (and feeling) responsible for assigning credentials that will follow a child for life. Second, in our context there are both weak regulation of grading and non-disclosure of information regarding standardized test performance to acting parties (teachers or students) before pupils' final assessments are processed. In this way, the present paper explores an environment in which:

i) subtle discriminatory behavior is hardly detected by school authorities or students themselves, and ii) last minute reactions to performance information are not sought by evaluators or by those being evaluated. It is also true that given the nature of desegregation within schools and classrooms in Brazil we are nonetheless unlikely to observe racial disparities that are not subtle. After all, students do interact at a level that gross misrepresentation of relative grades could easily be observed (students do see each other's graded exams) and lead to re-grading petitions. Finally, here we can explore detailed information on the longitudinal relation between teachers and pupils in order to closely examine theories of learning under statistical discrimination, and to shed light on the mechanism behind grading discrimination in our context.

II. Data and institutional background

A. Student-level data

The Sao Paulo's Secretary of Education collects detailed information on the universe of students and teachers in the state's educational system. Considering only regular primary and secondary schools, official records indicate that enrollment corresponded to approximately 6 million primary, middle and high-school students in 2010. Among eighth-graders, 67% were served by schools directly administered by the state authority, with the remaining share being evenly split between municipal and private institutions. Using confidential individual identifiers we merged information from four distinct sections of the Secretary's data bank: matriculation information, teachers' allocation to classrooms, transcript records and standardized tests of proficiency.¹² We turn to the description of each one of these.

Matriculation information covers all schools in the state of Sao Paulo, be they

¹²The Secretary has never attempted to combine these data. There are different departments in charge of each of these sections, and communication between them is scant. This is the first time these data have been used in an integrated format.

private or public. These records are centralized by the Secretary of Education through its role as a regulating agency for private and municipal schools. Matriculation within the public system is defined in terms of a school's catchment area (districting). Parents apply for a slot and pupils are assigned to the school serving the requested level closest to their residence. The centralization of information allows tracking students within the school system, and across classroom over the years. Our working data set covers the 2007-2012 period.

Records of teacher allocations to classrooms for the years 2007 to 2011 were also obtained. These files contain basic demographics (race, age, gender) for all the teachers in the system, and can be linked longitudinally. Combined with the matriculation records, we are able to map all teachers with whom each student had classes in the three years prior to eighth grade.

We also take advantage of the administrative data set on teachers' assessments of individual students between 2007 and 2011. This data set contains detailed information regarding scores and attendance records for all students in schools directly administered by the state's school authority. The complete set of report cards available to us includes information on every school subject. In eighth grade, in which teachers are fully specialized by subject, these correspond to Language (Portuguese), Mathematics, History, Geography, Sciences, Physical Education, and the Arts.

Assessment data started being centrally recorded after the adoption of a uniform criterion-referenced rule in 2007. According to official guidelines, all teachers assign numeric integer grades ranging from 0 to 10, with a passing grade set at 5 points for all disciplines. Attendance is recorded in percentage points (0-100 interval). Teachers and school administrators are not given instructions on how to attribute grades as a function of a student's observed proficiency level beyond the guidelines imposed by their uniform school curriculum. The state administration provides pedagogical material and teachers are supposed to evaluate students according to proficiency in its content. Nonetheless, no explicit guidance regarding

the design of evaluations is given (except for questions included in the back of teachers' booklets), and teachers still have great autonomy to define evaluation technology/methods and to allocate students across the 11 grading categories.

The final data set employed in our analysis provides results from standardized scores in the context of Sao Paulo's Performance Evaluation System - (SARESP- *Sistema de Avaliacao de Rendimento do Estado de Sao Paulo*). The system consists of an annual statewide exam taken by public school students in grades 2 and 4 (elementary school), 6 and 8 (middle school), and 11 (high school). Here we employ data from the 10th to the 13th editions (2007 to 2010), with over 1.5 million test-takers in approximately 5,050 schools covered in the latter year. Of this total, 420 thousand were eighth-graders (87.4% attendance rate in this particular level). As an integral part of the testing procedures, parents, students, and teachers also answer a survey that covers socioeconomic status, demographics (including race), study habits, teaching and pedagogical practices, and perceptions about the school environment, among other issues.

The main purpose of the SARESP exam is to measure the students' proficiency on the subjects assigned to each specific grade according to a predetermined curriculum, which is imposed on schools by the state authority. The exam in 2010 had questions covering Math and Portuguese language. For students in eighth grade, each exam contained 30 multiple choice questions. The exams were taken in late November (Spring), right before the end of the academic year, during regular-class meeting times, and in the same classrooms in which students sit for lectures. Teachers from different schools and levels were mobilized to supervise students during the test and grading was electronically conducted. Microdata on these tests' results were made available in the form of proficiency scores in each subject. These scores were computed using Item Response Theory (IRT) methods. Importantly, individual-level results from SARESP, past or current, are *never* made available to children, parents, teachers or schools.¹³

¹³For years prior to 2010 we were also granted access to IRT-scores from proficiency tests in Science

B. Racial gaps in Brazil

The discussion of racial differentials in Brazil is somewhat paradoxical. On the one hand, widespread racial mixing in marriage and the desegregation of housing markets have helped spread the view of a Brazilian “haven of racial reconciliation and affinity” (see Richman, 1999). On the other hand, there is overwhelming evidence that such racial tolerance indicators coexist with pertinent differences between Whites and non-Whites in terms of wages and other measures of economic well-being (see Arias et al., 2004 and Perry et al., 2006). In fact, the 2005 Human Development Report (United Nations) states that racial difference in economic achievement is one of the main social challenges facing Brazil. The report goes on to suggest that anti-discrimination policies should be central to any poverty reduction program implemented in the country. According to the 2010 Brazilian population census, adult male Whites have 8.4 years of completed education while the corresponding quantity for Blacks is 6.4 years. This lower educational attainment goes hand in hand with log-wage gaps of approximately 0.40 points. These gaps are of equal size when we restrict the sample to the state of Sao Paulo, which is the geographic area of focus for our analysis.

There have been important recent advancements towards the (potential) closing of racial gaps coming about as a result of *colorblind* social policies, however. Starting in the mid-1990s demand- and supply-side initiatives began to be undertaken, including the early steps and expansion of *Bolsa Familia*'s conditional cash-transfer program, and innovations in the allocation of federal budget toward school maintenance and teacher salaries under the *Fundo de Manutencao e Desenvolvimento do Ensino Fundamental (FUNDEF)*. Under this new institutional setting, standard educational policy targets rapidly improved: an unprecedented and significant increase in the rates of enrollment of school-aged children all over

and evaluations of an essay-based portion of the Language exam. The latter covered four different dimensions of writing ability: theme (ability to keep the text within the proposed theme); vocabulary and pronoun-noun concordance; cohesion and coherence (text organization); and syntax and subject-verb/time concordance.

the country. Based on household survey data, Madeira and Rangel (2013) show trends in enrollment for children aged 6 or 7 in the state of Sao Paulo between 1989 and 2009 by race. Aggregate enrollment figures went from somewhere around 75% in 1990 to more than 95% (or nearly universal coverage) by 2010. Importantly, from a racial perspective this increased access to schooling had a major influence on the composition of the student body, increasing the participation of a deprived portion of the population (among which non-Whites were overrepresented). In essence, Black-White gaps in enrollment among young children have virtually been eliminated in the state by the end of the period studied.

The absence of racial gaps in initial enrollment does not imply a closing of attainment gaps, however. For the country as a whole, the evidence on this dimension is mixed, while in Sao Paulo the patterns seem more favorable. We conjecture (but do not directly examine) that the adoption of a *social promotion* scheme in Sao Paulo is at least in part responsible for a faster convergence in education attainment between Blacks and Whites. Starting in 1998 such policy grouped contiguous primary school grades into two cycles, with retention only occurring at the end of each of them. Cycle 1 encompasses grades 1 to 4 (elementary) and cycle 2 covers grades 5 to 8 (middle school). Under this regulation, a student is promoted to the next level within a cycle if she attends more than 75% of the classes (and has no record of extreme disciplinary problems), irrespective of her mastery of the material covered during the academic year. Insufficient performance can only result in retention at the end of each cycle.¹⁴

[Figure 1 here]

In fact, trends are more pronounced in Sao Paulo than in other parts of the country, and the timing of convergence coincides with the policy's adoption. Yet what most substantiates this argument is the comparison of year-to-year attrition probabilities between middle-schools directly managed by Sao Paulo's school au-

¹⁴Several international organizations, including the World Bank, support this policy as an effective way to curb low grade completion and to decrease drop-out rates. The general lines of the argument are that grade retention could adversely affect some of the students' "non-cognitive" skills (like confidence and self-esteem), increasing anxiety levels and hampering their learning process. See King et al. (2008).

thority, and those run by municipal authorities. The former were all under social promotion during the 2006-2010 period we examine. Meanwhile, among the municipality schools only a small minority were under the same promotion scheme during that time. Figure 1 reproduces a simple computation of Black/White relative survival probabilities in both school systems. Assuming parity at 5th grade (one Black student per White student) we see that within the system adopting social promotion Blacks' relative attrition is lower than within municipal systems in every single year examined.¹⁵ We also observed in auxiliary exercises (Figure 2) that weaker students were relatively more likely to benefit from social promotion between 6th and 7th grades in particular. Here, as in the case of increased access, even if not aimed directly at racial issues, by benefiting students at the bottom of the skill distribution, social promotion had a disproportional effect on primary-school re-enrollment (higher) and retention (lower) rates among Blacks. Importantly, the increase in access to higher levels of basic education was gender-neutral. Both Black girls and Black boys were both more likely to reach eighth grade as a result of this policy.

We keep these recent trends in racial inclusion in perspective. The analyses that follow focus on how they are likely to affect the experiences of Black and White children that reach the final grade of middle school, right before racial differentials in enrollment rates and attrition dramatically increase at the high-school level.

C. Descriptive statistics

Our working data set was obtained after imposing restrictions based on the availability of both transcripts and (concurrent and past) test scores data for at least 75% of the students in a given 8th-grade classroom at the end of 2010.¹⁶ We also restricted our analysis to classrooms with non-homogeneous racial com-

¹⁵We use longitudinal matriculation records to compute these, and a description of the data is provided above. It is important to note that they are *unconditional* average transition rates.

¹⁶Since we cannot be sure that those not taking the test are a random subsample of students, our option to restrict the sample in this way was to make sure that relative rankings were closer to be representative of what happens in the actual classroom.

position (at least one Black and one White student) and fifteen or more students. We were left with observations on 277,444 students in 10,614 classrooms across 3,511 schools. Students that self-declared as Black or White are the main focus of the analysis, but our models are estimated including (and identifying) individuals classified under other races. As in most exercises in the social sciences that consider race, we implicitly assume that those that discriminate (teachers/employers) and those that are discriminated against (pupils/workers) agree on the racial classification captured in the records. We identify as Black all students that have been declared as such in any survey or enrollment documentation between 2005 and 2012. Columns 1 to 3 of Table A1, in the Appendix, presents descriptive statistics for our working data set. It is easy to see that, as expected, in pretty much every dimension in which we compare Blacks and Whites (and that are later used as control variables in our analysis), the former are unfavorably compared to the latter.

Focusing more specifically on scholastic performance, Figure 3 plots the cumulative distribution function of test scores (left) and teacher-assigned grades (right). These represent the main control and the main dependent variables in the econometric exercises that follow, respectively. Even with all of the observed progress in attainment, we can still find sizable differences in achievement between Blacks and Whites in Sao Paulo. For the students in our sample, differentials amount to 0.25 of one standard deviation (holding constant classroom fixed-effects). A similar pattern is observed in the distribution of teacher-assigned grades, with a disproportionate concentration of Blacks among those obtaining lower marks. Average differences in grades are approximately 5.6 points in a 0-100 scale.

[Figure 3 here]

Finally, in Figure 4 we plot the (Lowess) smoothed raw relationship between teacher-assigned grades and test scores in our data. This figure summarizes the main exercise of this article. For every level of test performance, Blacks receive lower grades from their teachers. The econometric strategy described below and

all our empirical estimations are in essence an attempt to verify whether these gaps are indeed there even after we both hold constant other productive attributes that make Black and White students different in the eyes of their teachers and address measurement error challenges. However, before examining the data in more detail, we turn to a simple conceptual framework that guides our estimations and orients the interpretation of results.

[Figure 4 here]

III. Conceptual framework

We focus our attention on a stylized description of grading that leads directly into our empirical specifications. The model is by no means general, but rather is used as a rhetorical device to emphasize a particular source of racial differentiation in teachers' assessments. In principle, there are two basic reasons for teachers to systematically mis-evaluate the competence of students with certain characteristics. First, teachers may merely like/dislike people with those traits, imposing rewards/punishments that can take both cardinal and ordinal forms. Second, teachers may attempt to be more sophisticated, evaluating (hard to measure) competence by also using observed characteristics perceived to be correlated with the former. In this case, the characteristics themselves convey information, and can "help" teachers generate better assessments. These alternative sources of discrimination are well known in the economics literature. The first is a loose representation of taste discrimination (Becker, 1957), whereas the second falls under the realm of statistical discrimination (Arrow, 1971; Phelps, 1972; Aigner and Cain, 1977). In our model we highlight the operation of the second, concentrating sole attention on the screening role of eighth-grade instructors.

The basic intuition is that teachers have access to noisy signals of the students' proficiency in Math, and observe both their behavior in class and their racial identities. We define an objective function for graders of school work by assuming they operate as statisticians compelled to maximize the power of the

hypothesis test embedded in the evaluation of a student's competence. In addition, we impose that teachers weight Type I and Type II errors symmetrically (i.e.: excessive lenience and excessive rigor are equally unwelcome). Evaluation errors can be reduced by exerting more screening effort, something we implicitly assume teachers either dislike (utility costs) or have limited access to due to high monetary/opportunity costs, or even that school authorities set the number of tests that can be applied to students in a given year (costs of effort could then be modeled as a function of distance to the "norm", such as in Holmstrom and Milgrom, 1991).¹⁷

Schematically, teacher r inelastically employs a grading/evaluation effort level T_r and at the end of the school year assigns to each student i (in a group of size n_r) a grade g_{ir} taking into consideration i 's unobservable true competence (g_{ir}^*) in order to solve on expectation the following optimization problem:

$$(1) \quad \min_{g_i} E \left[\sum_{i=1}^n \frac{1}{2} (g_i - g_i^*)^2 \right],$$

where we omit teacher-level subscripts for clarity of exposition and impose symmetry and tractability by adopting a simple quadratic function for the disutility generated by evaluation errors.

Importantly, we allow teachers to broadly define competence. As in Mechtenberg (2009), they acknowledge true proficiency (p_i^*) and other directly observed scholastic attributes (\vec{a}_i) as elements to be rewarded. Mechtenberg (2009) refers to the latter as *attitudes*, which we envision as a broad concept that includes habits, styles, behavior, and any other personality trait deemed *productive* by

¹⁷One could also conceive a technological constraint that limits the choices of teaching and testing effort.

teachers.¹⁸ That is to say:

$$(2) \quad g_i^* = \alpha_1 p_i^* + \vec{a}_i' \vec{\alpha}_2$$

Teachers do not observe true proficiency directly, so we further assume that they collect a sequence of noisy (yet unbiased) signals $s_i^t = p_i^* + u_i^t$. Signals result from formulating and grading tests/exams, and hence we associate them with evaluation effort ($t = 1, 2, \dots, T$).¹⁹ The higher the effort, the more signals will be gathered about each student's proficiency. Teachers' estimator of proficiency can then be described as a combination of those signals and a prior for mean proficiency:

$$(3) \quad \hat{p}_i^* = \frac{\sigma_{p^*}}{\sigma_{p^*} + \sigma_{\bar{u}}} \bar{s}_i + \frac{\sigma_{\bar{u}}}{\sigma_{p^*} + \sigma_{\bar{u}}} \beta_1,$$

where $\bar{s}_i = \frac{\sum s_i^t}{T}$, $\sigma_{\bar{u}} = \frac{\text{var}(u_i^t)}{T}$ and σ_{p^*} represents the variance of actual proficiency within the student population, while β_1 indicates the average student's proficiency (prior).

Combining all the elements in the model, and defining $\theta = \frac{\sigma_{\bar{u}}}{\sigma_{p^*} + \sigma_{\bar{u}}}$, we reach the following optimal rule for grading:

$$(4) \quad g_i = \theta \alpha_1 \beta_1 + (1 - \theta) \alpha_1 \bar{s}_i + \vec{a}_i' \vec{\alpha}_2.$$

From this formulation there are two ways in which statistical racial differentiation can be depicted. The first, rational stereotyping, is based on the idea

¹⁸Our formulation could also allow for racial bias operating directly via teachers' definition of competence (which we would recognize as taste-based discrimination, nonetheless). There is an interesting parallel between this variation and racial bias in the perception of others' pain discussed in Trawalter et al. (2012).

¹⁹For clarity of exposition, measurement error in teacher's tests is considered classical. We acknowledge that, due the bounded nature of grading scales in most of these classroom tests, errors would be negatively correlated with the true proficiency level. As long as the absolute value of the covariance between the error and the true proficiency is smaller than the noise variance (Black et al., 2000), introducing non-classical measurement error does not alter in any way the main messages of the model.

that attributes including race (\vec{b}_i) can be informative in the computation of proficiency's best linear projection $E \left[p_i^* | s_i^1, \dots, s_i^T, \vec{b}_i, \vec{a}_i \right]$.²⁰ In other words, the formulation of priors regarding group's average proficiency encompasses the use of other individual characteristics.²¹

The case of racial discrimination at hand can be illustrated within our context. Due to social promotion in earlier grades, eighth-grade teachers know that a particularly lenient rule for promoting male and female students was used. They likely assume that such scheme disproportionately affected promotion rates among Blacks. In the absence of any other information teachers will therefore have lower expectations regarding the latter's proficiency levels. If we let \vec{b}_i be a scalar corresponding to an indicator $Black_i$ not included in \vec{a}_i , we can amend the optimal grading equation to:

$$(5) \quad g_i = \theta\alpha_1\beta_1 + (1 - \theta)\alpha_1\bar{s}_i + \vec{a}_i' \vec{\alpha}_2 + \theta\alpha_1\beta_2 Black_i.$$

The second (and not mutually exclusive) possibility is that racial biases materialize as screening discrimination. This is the case when the reliability of proficiency signals collected by teachers is a function of race. Lang (1986) raised this as possible result of communication difficulties between Whites (teachers) and Blacks (students), while Lundberg and Startz (2007) suggest that they are the outcome of differential rates of social interaction. In our model screening discrimination would be embedded on race-specific signal-to-noise ratios: θ_1 and $\theta_1 + \theta_2 Black_i$. Under these circumstances, the practical distinction with respect to Equation (5) would solely come from the inclusion of race-specific effects of average proficiency signals (slopes).

Notice that in any of these representations, racial bias is derived from the im-

²⁰At this point we do not take a stand on the elements shared by \vec{a}_i and \vec{b}_i , but elaborate on it in the empirical section below.

²¹Ben-Zeev et al (2014) provides interesting laboratory-based experimental evidence of racialized recall biases. In particular, Black man are remembered as lighter when subjects are offered a counter-stereotypic stimulus (regarding educational attainment).

precision on the information about proficiency contained in the signals. It follows that improvements in the signal-extraction technology should make race a less relevant element of the grade assignment process. At the same time, the relationship between grades and individual test scores should be strengthened. This would be the case if teachers were to (exogenously) increase grading effort, if new information were distributed, or if tests were made less noisy. We take this simple model to the data, emphasizing its prediction regarding learning a child's true proficiency. Further discussions on alternative specifications and identification challenges are presented in the empirical section below.

IV. Empirical strategy

A. Practical issues

The first practical challenge we face in our empirical strategy comes from the way grades are reported. A conceptual issue arises from the heterogeneity in different teachers' application of the grade scale. As in the case of comparing responses using a Likert scale, contrasting grades assigned by different teachers is not clear cut. While a classroom fixed-effect added to the regression accounts for different mean scores across classes, an issue of dispersion remains; that is, even after factoring out the class average, a one point gain in class *A* can hardly be compared to the same absolute gain in class *B* if they have different grading standards in the spread of grades.²² At first we simply put aside this concern and use grades as our dependent variable, but we do so recognizing that (within this scale) measured gaps have both cardinal and ordinal meanings. Nonetheless, we also focus solely on ordinal aspects by present results based on the converting of grades assigned by teachers into classroom-specific percentile rankings.

In order to facilitate the interpretation of the practical impacts of our main results we also present two alternative binary dependent variables. The first is

²²In other words, the non-additive nature of this grading heterogeneity implies that linear fixed-effects will not wash them out.

the only really cardinal measure available in our data: an indicator of minimum competence. This was made common across teachers by the central authority's establishment of a passing grade (set at 5). So, independently of a teacher's choices regarding dispersion of grades within a classroom (or her subjective understanding of one additional point in the scale), it will always be the case that those above or at grade 5 are deemed competent while those below are not. This cardinal notion ought to be common across all classrooms, even if in different levels of stringency (captured by a class fixed-effect). As an additional ordinal measure we consider an indicator for grades above the classroom's median grade.

A second practical concern is the different natures of the exams applied within the school context by teachers and the standardized tests adopted for external monitoring of learning. Since teachers receive a uniform curriculum, textbooks and practice exercises from the external examiner, their evaluations regarding proficiency should reflect the same skills and cognitive abilities as the standardized exam. Yet, it is plausible that proficiency in a given content can be measured by examining performance using different tasks (format). Take the case of Language evaluations, for example. Teachers most likely combine observations regarding reading, writing, and speaking abilities when assessing a student's language competence. Paper-and-pencil standardized tests implemented in our context can only capture reading skills using a multiple choice exam. This is a reason for restricting our analysis to Mathematics: we expect the objectivity inherent in the material to translate itself into skills more easily measured in a test-like format. It is still possible that teacher-designed Math exams also reward writing skills so that we flexibly include scores from past essays and concurrent Language tests as controls in our empirical model.

B. Econometric issues

In essence, we explore our information regarding scores in standardized Math and Language exams as a proxy for the average level of proficiency measured by

teachers in their own classroom examinations. Meanwhile, other skills also considered relevant by teachers are factored into the *productive attributes* term (\vec{a}_i). Therefore, we propose the following empirical representation that incorporates teacher/classroom fixed-effects (η_r) and a pupil-level disturbance term (ϵ_{ir}):

$$(6) \quad g_{ir} = \delta_1 f(\text{scores}_{ir}) + \vec{x}_{ir}' \vec{\delta}_{21} + \vec{z}_{ir}' \vec{\delta}_{22} + \vec{b}_{ir}' \vec{\delta}_3 + \eta_r + \epsilon_{ir},$$

where $f(\text{scores}_{ir})$ is a function test performance available in our data that replaces the “theoretical” average level of proficiency captured in teacher-designed examinations (\bar{s}_{ir}), and once again \vec{b}_{ir} lists elements affecting teachers’ priors with regard to proficiency. Meanwhile, in order to make explicit further challenges to our empirical exercise, the elements in the vector of scholastic attributes (\vec{a}_i) were also be decomposed into observed and unobserved components, with \vec{x}_{ir} representing elements observed both by teachers and the econometrician and \vec{z}_{ir} standing for those only observed by the former.

Given that our central objective is to consistently estimate δ_1 and $\vec{\delta}_3$, this simple empirical representation highlights the two main econometric problems we face: a) measurement error in proficiency scores, and b) unobserved heterogeneity.²³

Measurement error biases result from the fact that despite being associated to the average proficiency measured by teachers, our measure is necessarily noisier. An easy way to understand the discrepancy between the two is to consider that while teachers “draw” observations from multiple and heterogeneous tests, the econometrician only observes results from one of them. Those biases directly limit our ability to test the predictions from the aforementioned conceptual framework. Despite being graded in a scale, the measurement errors inherent to IRT scores we use in our analysis can be treated as classical. There are two basic reasons for that to be the case. First, the test scale in Sao Paulo is calibrated to cover proficiency

²³For a discussion of the effects of measurement errors and omission biases when using test scores as covariates, see Andrabi et al. (2011).

comparisons between 4th and 11th grades, so that its boundaries are unlikely to be reached within the population of 8th graders we study. Second, despite being based on exams with 30 items, the IRT estimation of proficiency is conducted in such a mechanism that the number of points in the scale is proportional to the number of possible combinations of the response-vector as well as on the proficiency level of those correctly responding each item (and is, therefore, not close to being a discrete scale). IRT scores are indeed less precise in the tails of the distribution, due to the scarcer information available for maximum-likelihood estimations embedded in those psychometric methods (see Samejima, 1994). This, however, does not mean that the errors are non-classical, yet it necessarily makes them heteroskedastic.²⁴

In the empirical exercises below we explore the fact that the individual results of standardized tests in Math and Language taken in previous years by each student are available in our data and, similarly to Andrabi et al. (2011), employ a fixed-effects instrumental variables estimation that should directly deal with the measurement error problem we face. Since we also have access to past proficiency tests covering Natural Sciences' material, we are in addition able to perform overidentification tests. Despite not being a panacea, such tests can be used to inform if there is any clear evidence against the validity of our instruments.

Unobserved heterogeneity adds another layer of complications because even in the absence of measurement error in scores, elements of $b_{ir}^{\vec{}}$ may very well be related to elements of $z_{ir}^{\vec{}}$. In particular, we worry about behavioral indicators that are available to teachers during classroom interactions and are correlated with racial identity.²⁵ We take this very seriously and, in the exercises below, consider a number of proxies for behavior in an attempt to check the sensitivity of our results. We have explored information correlated with behavior from different sources such as: *i*) teacher attendance records, assuming the students

²⁴We have in fact experimented with measurement-error-corrected estimations employing Lewbel (2012)'s identification via heteroskedasticity (Rigobon, 2003). Since no additional insight relative to the estimations performed below was obtained, we have opted not to present it here.

²⁵Cornwell et al. (2013) face a similar issue in the case of gender differentials in grading.

that miss more classes are disengaged or poorly behaved even when attending (we used attendance to Language classes in the first half of the academic year to avoid confounding feedback effects); *ii*) parent-reported perceptions of student engagement, behavior, and effort in school-related activities; *iii*) student self-reported indicators of class absence and procrastination with homework; and *iv*) Physical Education (PE) grades (in the first half of the academic year). PE grades are under the responsibility of a different teacher. Athletic equipment and infrastructures, such as fields and tracks, are not available in most schools, and students usually perform simple calisthenics and routines during classes. In eighth grade, for instance, one can hardly argue that grades are assigned as a function of athletic skills. Instead, other traits often valued by teachers, such as obedience, respect for the other students, and the capacity to respond to simple commands, are likely more relevant.

Ultimately, our main empirical model consists of regressing grades on race, gender, age, essay scores, parental socio-demographics, and our proxies for behavior. These are all considered elements of the vector x_{ir}^{\rightarrow} while the remaining elements of z_{ir}^{\rightarrow} not observed by the econometrician are either absorbed by the classroom fixed-effects or by the disturbance term. $f(scores_{ir})$ is estimated as fourth-order polynomials of Math scores, a linear term for Language scores, and interactions between those.²⁶

C. Learning

We also extend the analysis to explore the heterogeneity of the parameters according to teacher and student-body characteristics. In particular we pay attention to the amount of knowledge a given teacher has about each of her pupils. Social interaction in the school neighborhood, tenure in a given school and duration of classroom-like interactions for a given student-teacher pair are our main

²⁶The use of either splines or indicator variables after discretizing the scales does not alter the inferences we perform. Moreover, whenever F-tests indicated that the fourth-order elements were not significant, we opted for presenting results based on a more parsimonious third-order polynomial.

candidates. In this way we examine the central prediction from our statistical discrimination conceptual framework: learning of students' true types should preclude the use of race as an indicator of scholastic competence.

In practice, and in the spirit of Altonji and Pierret (2001), we test whether racial differentials in teacher-assigned grades diminish as a teacher's information regarding students improves. By the same token we examine if such improved information also translates into increased weight given to proficiency signals when end-of-year Math evaluations are issued. If such coefficients are shown to conform with these predictions, we can be more confident that statistical discrimination is at play in our study's environment.

V. Results

A. General results

Table 1 presents results illustrating the impact of the addition of controls over racial differentials and over the marginal effect of proficiency scores (measured at the average performance level) in our two main dependent variables.²⁷ Panel A focuses on the Black-White gaps in final grades (0-100 scale). Group averages are presented in column 1. Considering all of the students in our sample, Whites are graded at 61.4 on average while grades among Blacks average 55.7. This difference is relatively unaffected by the inclusion of classroom fixed effects (column 3), indicating that racial segregation in assignment to classrooms or schools is unlikely to be behind the racial gaps. In column 4, individual demographic characteristics (gender and a second order polynomial on age) the polynomial for Math and Language contemporaneous standardized scores, and past performance in essays are included. Measured racial gaps are, not suprisingly, significantly reduced. Indeed, a large share of the competence differences seen by teachers is captured by performance in standardized tests of proficiency.

²⁷The sequential inclusion of controls should not be taken as representative of the influence they exert over the gaps we want to measure. See Gelbach (2009) for a methodological discussion.

[Table 1 here]

In column 5 we include family background and information on past year's Math grades as additional control variables, with the intention of capturing child-specific (time-invariant) abilities and competence aspects relevant to teachers that were not being previously controlled for. Proxies for a child's behavioral attributes (self-reported, parent-reported, school-reported), over and above those indirectly captured by family socio-economic background, are included in column 6. The inclusion of behavioral aspects has minimal impact on our estimates, suggesting that at this point very little is left out of the model. An inspection of the direct effects of these behavioral traits indicates significant results that go in the expected direction. Holding performance in tests and socio-demographics constant, Math grades improve (and significantly do so) when the child attends a higher proportion of classes, when she gets higher grades in physical education, when parents report her as dedicated to and motivated with school work and, ultimately, when she herself declares not to procrastinate on finishing homework.²⁸ Despite the reduction in size, estimated racial gaps are still statistically significant.

Finally, in columns 7 and 8 we tackle the robustness of our findings to the presence of measurement error on the proficiency score variables. As discussed above, because these are used as covariates in our analysis, biases on the estimation of *all* parameters are expected. We therefore employ polynomials of lagged test scores (resulting from tests taken in the most recent school year prior to the current one) as instrumental variables. Reflecting the cumulative nature of proficiency exams, past scores are very correlated with current ones (see first-stage summary statistics in the appendix Table A2). Moreover, over-identification tests suggest we have no obvious reason to distrust the validity of the sets of instruments employed (of course, Hansen's test remains only a necessary but not a sufficient condition for exogeneity) and, therefore, hints to the absence of unobserved het-

²⁸These coefficient estimates are not shown in Table 1 to preserve space, and are available upon request.

erogeneity issues. Once measurement error is accounted for, we encounter smaller racial differentials and at the same time larger slope parameters in the relation between Math grades and Math test scores (marginal effects at the average proficiency level). The racial gaps are still significant after this correction. Indeed, they are statistically significant even when we employ the more stringent Schwarz criterion.²⁹

We find that Blacks' average Math grades are 0.35 points below those of equally proficient and well-behaved Whites, or that the former are regularly ranked 0.7 percentiles behind the latter. These amount to 6 and 7% of the unconditional gaps, respectively. By taking the ratio of coefficients in the estimated model, we see that the Black-White differentials in teacher-assigned grades are equivalent to the marginal effect of a reduction of 0.03 to 0.04 of one standard deviation in proficiency scores. Incidentally, despite differences in methods, scales and context, our results are within range (in standardized-tests units) of those experimentally obtained by Hanna and Linden (2012).

Moreover, one may even argue that some of our control variables are the result of grading discrimination in their own right, inducing our models to underestimate the size of Black-White gaps. We see merit in such argumentation, since biased grading within the school year can indeed induce students to "misbehave", but prefer to be as conservative as possible in our empirical exercises. We restrict the analysis that follows to the use of a fully controlled model.

There are a number of reasons to believe these gaps are quite relevant. The gaps we uncover are equivalent to 40% of the raw (within-classroom) difference in terms of grades associated with having a mother with a college degree or more versus a mother with a high-school degree only. They equivalently correspond to 30% of the effect of being a child with parents that own their home. Ultimately, our measures of racial discrimination depends on the translation of intrinsically

²⁹Considering our very large sample, the Schwarz criterion, which sets critical values of significance as a function of sample sizes is indeed more appropriate to judge the statistical significance of results.

ordinal scales into meaningful quantities, something that could only be achieved in case they were to be directly associated with longer-term outcomes such as employment and wages (Bond and Lang, 2013). In the context of the present study there is an enormous potential for feedback effects, nonetheless. Teacher communications may steer educational investment decisions in one way or the other, with parents (and children themselves) updating decisions (and effort) after extracting information from report cards issued by teachers. In this case, intra-classroom evaluation biases may very well lead to gaps in attainment, school choice, future scholastic performance and, ultimately, labor market outcomes.

Column 9 in Table 1 reproduces these exercises with a focus on meaningful binary variables that summarize cardinal and ordinal gaps. According to these estimates, the measured racial gap in promotion rates between equivalently proficient and well-behaved students corresponds to a 4.1% increase in the 8th-grade retention probability for the average Black (Panel A). Focusing exclusively on the ordinality aspect (Panel B) we also estimate a gap that translates into a 4.5% reduction on the probability of Blacks being graded above the classroom median.³⁰ These meaningful effects are very much in line with the subtleties we expect to permeate racial discrimination in grading.

In order to gather a sense of the size of these effects and (possibly) the mechanics behind grading discrimination, we explore a simple simulation exercise. We start by converting the blindly graded proficiency scores into a classroom-specific 0-100 scale. Conversion is undertaken by: a) computing the difference between the score of student i and the minimum score in her classroom; b) dividing this quantity by the difference between the maximum and minimum score in that classroom, and; c) distributing this quantity in the range given by the teacher-assigned grades in that classroom. We then add a simulated discriminatory rounding routine while converting from a continuous into a discrete grade scale. This is done by

³⁰In principle, if Blacks take standardized tests more seriously than in-class examinations relative to their White counterparts we would expect to find results like these. This stereotype-threat-like argument is indeed valid, but one for which we do not have a direct empirical implication to be tested using our data. Such caveat applies to the whole literature on racial and gender gaps in test scores/grades.

assuming that every time a Black student's score lies in the $[h + 4.5, h + 5.4)$, where h is an integer in the 0-90 scale, the resulting grade is necessarily h . White students in the exact same situation have $h+10$ as their assigned grades, following unbiased rounding rules. We then compare the average racial gap in the biased and unbiased conditions. We find that this biased rounding generates racial gaps of 0.94 (in end-of-year grades) and 1.94 (in percentile rankings). Notice, therefore, that the differences between Blacks and Whites we estimate above are between a third and half the size of the ones in the simulated exercise. This is an interesting finding, as it gives us a notion that the results we find are in the ballpark of what happens when such a subtle discriminatory action is imposed over a share of the four times teachers edit report cards during the school year, for example. This exercise provides additional confidence that the results we uncover are not only important but also very plausible.

B. Robustness of main findings and modeling choices

We also explore expected heterogeneity in the size of racial differentials and its relation to some teacher characteristics (grading practices) to further examine the robustness of our findings to the omission of behavioral characteristics. In Table 2 (columns 1 and 2), before moving into the comparison across different data strata, we present a summary of the main effects under the full sample and under the subsample for which we have additional teacher characteristics (from survey questionnaires). The contrast between these indicate that we should not expect selection biases when dealing with the smaller sample.

In the first set of stratifications (columns 3 to 5) we examine if the gaps in evaluation we measure are not generated by unobserved heterogeneity biases. We explore a section of the questionnaire answered by teachers in the context of SARESP, in which opinions regarding the importance of objective instruments of evaluation (tests and exams) and also the importance of using more observational methods (classroom behavior, students' motivation, oral examinations, etc.) were

gathered. These questions were posed in an independent manner, so that there they are not excludable categories. We explore these responses to stratify teachers in three (not necessarily distinct) groups. Those that believe objective methods are very important, those to whom objective methods are not important, and those to whom subjective/observational methods are very important. We find no evidence that these groups discriminate against Blacks with different intensities. In fact, if anything, larger effects are found among those that believe in the objective evaluation of students. In our opinion, this is the first indication that imperfect information plays a central role in our findings: racial bias seems to occur among those that are trying to extract the most out of their noisy measures of proficiency.³¹

[Table 2 here]

Columns 6 and 7 are solely based on teacher demographics (obtained combining official assignment records and survey questionnaires). We re-estimate our model using fixed-effects instrumental variables techniques for different strata according to teacher's race, which is examined here to investigate in-group biases. We see that no clear pattern emerges from these. Despite significant results among Whites and not among Black and mixed-race teachers, we cannot rule that point estimates are the same. These findings seem incompatible with both the idea of taste discrimination (at least in its simplest format) against members of the out-group and the idea that teachers have more information about pupils of their own racial group. The very small number of Black teachers in our sample is surely a limitation when examining these hypotheses, but also give additional hints regarding overall racial discrimination in Sao Paulo.

In Table A3 in the Appendix we investigate the robustness of our formulation by examining if the marginal impact of proficiency tests over grades are different for Black and White students, as predicted by the screening discrimination version of

³¹In another exercise that examines unobserved heterogeneity biases we estimate if the difference in future drop-out rates between retained and non-retained Blacks were larger than among Whites. If they were it could mean that teachers observe other productive aspects on retained Blacks that they do not see in retained Whites. We find that this is not the case in our data covering 2011 and 2012, however.

the model above. From the estimates presented, we have no evidence to support the idea that slope coefficients should be student-race specific.³²

C. Learning-by-grading

In order to more directly examine the role of imperfect information we explore data on pupil-teacher matches by utilizing the longitudinal information on students' and teachers' assignment to classrooms. We actually map the individual-level acquaintance level between every student and their current teacher. In this case we emulate a student-specific change in grading-effort (T) exerted by her current teacher. Simply put, larger T 's should increase signal to noise ratios, increasing the marginal effect of (posterior) proficiency measures at the same time it reduces the one related to characteristics used to construct priors.

[Table 3 here]

It is clear from estimates in Table 3 that longer-term teacher-student interactions produce smaller grading gaps associated with racial identity. In other words, this empirical exercise reveals that while Black-White gaps in grades and rankings are salient for students attending classes with a teacher for the first time, no significant disparities are found among those that have already had classroom interactions with that instructor before eighth grade. It is also the case that acquaintance between teacher and students increase the weight given to proficiency scores on the determination of grades or rankings (steeper relation). Both these differences (in intercept and in slope) are shown statistically significant. In practice, Black students that have not interacted with their current teacher before eighth grade have their grades diminished by what is equivalent to a taxation of 0.06 of one standard-deviation over proficiency tests' performance. Those that are known to the teacher are not "taxed" at all. This is our main indication that imperfect information lies at the heart of the discrimination results we estimate.

Further experimenting with these ideas we examine stratifications based on

³²The conclusion remains unchanged if we restrict this analysis to classrooms with White teachers.

different levels of detail in the information teachers have about their students. In Table 4 we start by reproducing in column 1 the differences estimated in Table 3. Columns 2, 3 and 4 focus on the proportion of students in a classroom that are “known” to the current teacher. The idea here is that by knowing a sufficient number of students, teachers are able to employ relative references to grade their pupils. We do observed that to be the case, particularly regarding the ranking measures (which we would expect to be more prone to the use of relative referencing). In practice, racial gaps are not observed in classrooms where teachers have had past interactions with at least 50% of the students, while they are in case teachers know relatively less students. Once again conforming to the model, the former also give more weight to proficiency scores when assessing competence and defining end-of-year grades.

[Table 4 here]

In columns 5, 6 and 7 we turn to the idea that information flows result from teachers’ tenure in a given school (and with a given population of students or set of co-workers). When we estimate racial gaps employing this idea we find that indeed gaps are larger among teachers that have shorter tenure in the school. Differences in slopes are less precisely estimated but still support the role of information flows. We have also estimated models with teacher-student-specific acquaintance levels among longer tenure teachers and found room for “learning a student’s type” even among those (not shown). This most likely means that the level of detail regarding a student’s competence is finer when classroom interactions do occur than when information is provided via interactions with other teachers in the same school.

Finally, in columns 8 to 10 we examine if information is also spread via social interactions within the schools’ neighborhood. Interestingly, in this case those that would supposedly know more about the student population seem to discriminate more. This makes us believe that the information flows that translate into reduced discrimination need to be somewhat related to Math abilities, something

that teachers do capture within classroom/school settings, but that neighbors cannot easily infer.

The robustness of our learning argument can be further put to the test by examining an alternative explanation for the findings above. In particular, we investigate if the assignment of teachers to students captured in our proposed measure of knowledge above is not simply revealing that schools that assign teachers to the same students are culturally different (say, in terms of taste for discrimination) or even that students with specific behavioral characteristics are selected into longer-term interactions with teachers. In Table 5 (columns 1 to 3) we strongly reject both these threats by presenting evidence that neither *Math*-grade racial gap nor its relation to proficiency are a function of length of interaction time between a student and her *Language* teacher. It is important to emphasize that in our context the group of students is common for Math and Language teachers, so if there were a particular rule for allocating students to the same teacher year after year this should be true for both subjects.³³

[Table 5 here]

We reach the the exact same conclusion when we employ the identity of future Math teachers to measure acquaintance levels in columns 4 to 6. Teachers that *will* spend more time with a given student do not discriminate more or less today than those that will not. Alternatively, students that will spend more time interacting with their current Math teachers in the future do not have their racial identification playing a role on evaluations that is different than for those that that will not.

We conclude our analysis in Table 6 by verifying that our finding regarding learning and consequent reduction in racial gaps are not a result of the omission (at the level of the interaction effect) of other characteristics in our econometric specification. In particular, we examine if by including interactions between

³³Even though this is implicit in these robustness exercises, for completeness, we present descriptive statistics of students that do and do not spend more than one year interacting with a given teacher in columns 4 to 7 of Table A1 in the Appendix.

teacher-student relation indicators and other control variables we are able to eliminate the differences observed in the race coefficient. Strikingly, we see no reason to believe this is the case. From what we can tell from the estimates, a student racial identity is likely used in our context as an indicator of lower proficiency. Its impact over grades is remediated if teachers get to know (and test) students for longer than an academic year.

[Table 6 here]

From this exercise we incidentally uncover three other very interesting (yet less significant) patterns. First, social economic background variables seem to have a role similar to race in these grading decisions, with teachers looking at indicators of those (which are not necessarily directly observed) to draw their priors regarding a student's competence. This would be the case if clothing, mobile phone ownership, and their quality/brands were used to infer socio-economic status, for example. Of course, socio-economic background can also be translated into higher endowments of other productive skills, being legitimately evaluated as competence (e.g.: maternal education translating into better manners or better vocabulary) and considered even by teachers that know well the student. These interpretations are compatible with our findings. Second, behavioral traits have a role on the evaluation that is independent of how much a teacher knows about a specific student. That is to say: both those that are well acquainted with a student and those that are not weight behavioral traits in the same way when evaluating competence. Finally, our estimates indicate that the sizeable gender differentials in grading estimated by these models do not operate in the same way as racial gaps when it comes to information availability. Therefore we are inclined to believe that statistical discrimination of the form we modeled does not fit the case of gender discrimination in grading (see Botelho et al., 2013 for a more detailed discussion). We see this as reinforcing our interpretation of racial biases, however. Since we have no reason to believe that gender carries any information in terms of potential competence in our context (social promotion or other social

programs did not have any gendered effect on access to eight grade), a child's sex should not contribute to the formation of teachers priors in any particular way.

Taken together, main and auxiliary findings clearly substantiate the formulation of a conceptual framework like the one we proposed, particularly with regards to a well-defined channel of operation for information. Most importantly, however, our "dynamic" findings can be interpreted as reinforcing the robustness of our basic results to unobserved heterogeneity biases (in the form of omitted time-invariant student-specific characteristics). Since we see no indication that students interacting with a given teacher for more than an academic year have specific characteristics (among the many we indeed measure), there is no reason to believe omission biases would only be operating within the sample of first-time interactions between pupils and teachers.

Summing up, the econometric results presented here are not only robust to possible omissions of behavioral attributes and to the likely incidence of measurement error on scores from standardized tests, but are also very much in line with the expected subtlety of this particular form of discrimination within Brazilian schools.

VI. Conclusions

In this article, we empirically detect racial discrimination within racially integrated Brazilian eighth grade public-school classrooms. Math teachers' assessments of students with respect to scholastic proficiency and aptitude (grading) are found to be biased. White students are less likely to be deemed non-competent (below passing grade) than their equally proficient and equivalently well-behaved Black classmates. The former are also relatively more likely to be graded above their classroom median. Quantitatively, these correspond to a 4.1% increase in the retention probability and a 4.5% reduction in the probability of Blacks being at the top of their class grade distribution. Such effects are equivalent to "taxing" Blacks' performance in proficiency tests by 0.03 to 0.04 of one standard deviation.

These results are shown robust to possible omissions of a students' behavioral attributes and to the incidence of measurement error on scores from standardized tests. It turns out that well intentioned teachers issue report cards for their students with subtle biases (possibly incurred when rounding continuous marks into a discrete scale, for example) and in, this way, may end up adding obstacles to the acquisition educational credentials by Blacks. These are meaningful effects resulting from racial discrimination in grading. They are equivalent to 40% of the raw (within-classroom) difference in terms of grades associated with having a mother with a college degree or more versus a mother with a high-school degree only, or to 30% of the effect of being a child of parents that own their home.

We also find that these biases most likely result from imperfect information and statistical discrimination or, in other words, from the weighted combination of noisy proficiency signals extracted from teacher-designed exams and stereotyped priors. In the case explored here, stereotyping seems to have resulted from lenient standards for admission of students into eighth grade (which have disproportionately benefited Blacks). Improvements in the signal-extraction "technology" available to teachers make race a less relevant element of the grade assignment process and, at the same time, strengthen the relationship between grades and individual proficiency scores. This is clearly shown to be the case in our data, particularly when we use the length of classroom-interaction time between the teacher and a given student. In practice, Black students that have not interacted with their current teacher before eighth grade have their grades diminished by what is equivalent to a taxation of 0.06 of one standard-deviation over proficiency tests' performance. Those that are known to the teacher are not "taxed" at all.

Our findings lead to some education-policy implications. First, curbing teacher rotation (which is very high in our context) can be particularly important for Black students. Beyond their likely influence over learning, increased interactions between a group of students and a given teacher diminishes the influence of noise on the evaluation of scholastic proficiency. The more a teacher gets acquainted

with a given student, the less relevant the latter's race becomes for screening purposes. Second, direct investment in training of teachers with regards to the design of exams and tests may be warranted when attempting to curb discriminatory outcomes. If equipped with better signal extracting technologies, teachers would not need to resort to distorted priors. Third, educational governing bodies could promote the clear communication of (already collected) standardized test results at the individual level to teachers as a way of widening their information set about students' abilities. Finally, our results point to important nuances on the overall impact of affirmative action policies in admission to college (or social promotion schemes in basic education for that matter) in environments where the progress of those targeted by the policy depends on continued subjective evaluations of performance. Facilitating access can also impose lower ceilings.

In scientific terms, the results presented here indicate that well-designed randomized control trials focusing on the amount, type, and timing of information about individual students available to teachers can go a long distance in helping us understand the inner workings of grading discrimination within schools. They also pave the road to studies that focus on parental/student responses to school reports in terms of investments in the accumulation of human capital (or to parent-school communication mechanisms in general) as well as to more ambitious and controlled examinations of social promotion schemes effects over students' future outcomes. We leave these for our future research on such topics.

REFERENCES

- Aigner, D. and G. Cain (1977); "Statistical theories of discrimination in the labor market". *Industrial and Labor Relations Review*, 30, 175-187.
- Alexander, N.; A Guimaraes; C. Hamilton; L. Huntley and W. James (2001); "Beyond Racism: Race and Inequality in Brazil, South Africa, and the United States," L. Rienner Publisher, Boulder, CO.
- Altonji, J. and C. Pierret (2001); "Employer Learning and Statistical Discrimination". *The Quarterly Journal of Economics*, Vol.116. No. 1 (Feb.2001), pp. 313-350.

- Andrabi, T.; J. Das; A. Khwaja and T. Zajonc (2011); "Do Value-Added Estimates Add Value? Accounting for Learning Dynamics", *American Economic Journal: Applied Economics*, 3(3): 29-54.
- Arias, O.; G Yamada and L.Tejerina (2004); "Education, Family Background and Racial Earnings Inequality in Brazil", manuscript, Inter-American Development Bank, Washington, DC, USA.
- Arrow, K. (1971); "Some Mathematical Models of Race in the Labor Market" in Pascal, A. (ed.) *Racial Discrimination in Economic Life* (Lexington, MA: Lexington Books) 187-204.
- Assuncao, J. and B. Ferman (2013); "Does Affirmative Action Enhance or Undercut Investment Incentives? Evidence from Quotas in Brazilian Public Universities", unpublished manuscript, Sao Paulo School of Economics, Fundacao Getulio Vargas.
- Autor, D. and D. Scarborough (2008); "Does Job Testing Harm Minority Workers? Evidence from Retail Establishments", *The Quarterly Journal of Economics*, 123 (1): 219-277.
- Becker, G (1957); *The Economics of Discrimination*, The University of Chicago Press.
- Ben-Zeev, Avi; Tara C. Dennehy; Robin I. Goodrich; Branden S. Kolarik; and Mark W. Geisler (2014); "When an "Educated" Black Man Becomes Lighter in the Mind's Eye: Evidence for a Skin Tone Memory Bias", *SAGE Open*, 4(1), January.
- Bjerk, D. (2008); "Glass Ceilings or Sticky Floors? Statistical Discrimination in a Dynamic Model of Promotion and Hiring", *The Economic Journal*, 118 (July).
- Black, D.; M. Berger and F. Scott (2000); "Bounding Parameter Estimates with Nonclassical Measurement Error". *Journal of the American Statistical Association*, 95 (451), pp 739-748.
- Blume, L. (2006); "The Dynamics of Statistical Discrimination", *The Economic Journal*, Vol 116, November, pp. F480-F498.
- Bond, T. and K. Lang (2013); "The Evolution of the Black-White Test Score Gap in Grades K-3: The Fragility of Results". *The Review of Economics and Statistics*, Vol 95(5), pp 1468-1479.
- Borjas, G. and A. Goldberg (1978); "Biased Screening and Discrimination in the Labor Market". *The American Economic Review*, Vol. 68, No. 5 (Dec.,1978), pp. 918-922.
- Botelho, F.; R. Madeira and M. A. Rangel (2013); "Gender disparities in test scores and teacher assessments", Working paper, PLAS Princeton and University of Sao Paulo.
- Bowles, S. and H. Gintis (1976); *Schooling in Capitalist America*. New York: Basic Books.
- Burgess, S. and E. Greaves (2003); "Test Scores, Subjective Assessment, and Stereotyping of Ethnic Minorities". *Journal of Labor Economics*, Vol. 31(3), pp. 535-576.
- Coate, S. and G. Loury (1993); "Will Affirmative-Action Policies Eliminate Negative Stereotypes?", *The American Economic Review*, American Economic Association, vol. 83(5), pages 1220-40, December.

- Cornell, B. and I. Welch (1996); "Culture, Information, and Screening Discrimination". *Journal of Political Economy*, 1996, vol. 104, no. 3
- Cornwell, D.; B. Mustard and J. Van Parys (2013); "Noncognitive Skills and the Gender Disparities in Test Scores and Teacher Assessments: Evidence from Primary School." *Journal of Human Resources*, Vol. 48(1), pp. 236-264.
- Cortes, K. and L. Zhang (2012); "The Incentive Effects of the Top 10% Plan" unpublished manuscript, The Bush School of Government and Public Service.
- Cotton, C.; B. Hickman and J. Price (2014); "Affirmative Action and Human Capital Investment: Evidence from a Randomized Field Experiment". *NBER Working Paper Series*, WP No. 20397, August.
- Demeis, D. and R. Turner (1978); "Effects of Students' Race, Physical Attractiveness, and Dialect on Teachers' Evaluations". *Contemporary Educational Psychology* 3 (1):77-86.
- Dorsey, J. K. and J. Colliver (1995); "Effect of Anonymous Test Grading on Passing Rates as Related to Gender and Race". *Academic Medicine* 70 (4): 321-23.
- Dovidio, J. F.; J. Brigham; B. Johnson and S. Gaertner (1996); "Stereotyping, Prejudice, and Discrimination: Another Look" in N. Macrae; C. Stangor and M. Hewstone (eds.) *Stereotypes and Stereotyping*, New York, Guilford.
- Farber, H. and R. Gibbons (1996); "Learning and Wage Dynamics" *The Quarterly Journal of Economics*, 111 (4): 1007-1047.
- Farkas, G.; R. Grobe; D. Sheehan and Y. Shuan (1990); "Cultural Resources and School Success: Gender, Ethnicity, and Poverty Groups within an Urban School District" *The American Sociological Review*, 1990, Vol.55 (February: 127-142).
- Ferguson, R. (1998); "Can Schools Narrow the Black-White Score Gap?" in C. Jencks and M. Phillips (eds.) *The Black-white Test Score Gap*, The Brookings Institution.
- Ferguson, R. (2003); "Teachers' Perceptions and Expectations and the Black-White Test Score Gap". *Urban Education* 38 (4): 460-507.
- Figlio, D. (2005); "Names, Expectations and the Black-White Test Score Gap", NBER Working Paper No. 11195.
- Gelbach, J. (2009); "When Do Covariates Matter? And Which Ones, and How Much?"; Working Paper 09-07, University of Arizona 2009.
- Hanna, R. and L. Linden (2012); "Discrimination in Grading". *American Economic Journal: Economic Policy*, 4(4): 146-168.
- Hedges, L. and Nowell, A. (1999); "Changes in the Black-White Gap in Achievement Test Scores". *Sociology of Education*, 72:2, 111-135.
- Herring, C; V. Keith and K. Horton (2004); "Skin Deep: How Race and Complexion Matter in the Color-Blind Era" Institute for Research on Race and Public Policy.
- Hinnerich, B.; E. Hoglin and M. Johannesson (2011); "Are boys discriminated in Swedish high schools?" *Economics of Education Review*, Vol. 48(1), pp. 236-264.

- Hinnerich, B.; E. Hoglin and M. Johannesson (2011); "Ethnic Discrimination in High School Grading: Evidence from a Field Experiment." SSE/EFI Working Paper Series in Economics and Finance, 733.
- Holmstrom, B. and P. Milgrom (1991); "Multitask Principal-Agent Analysis: Incentive Contracts, Asset Ownership, and Job Design." *Journal of Law, Economics and Organizations*, Vol. 7, pp. 24-52.
- King, E.; P. Orazem and E. Paterno (2008); "Promotion With and Without Learning", World Bank Policy Research Working Paper, 10.1596/1813-9450-4722.
- Lam, D.; C. Ardington and M. Leibbrandt (2006); "Schooling as a Lottery: Racial Differences in School Advancement in Urban South Africa", University of Cape Town, Working Paper Number 56.
- Lang, K. (1986); "A Language Theory of Discrimination", *The Quarterly Journal of Economics* 101 (2): 363-382.
- Lavy, V. (2008); "Do Gender Stereotypes Reduce Girls' or Boys' Human Capital Outcomes? Evidence from a Natural Experiment". *Journal of Public Economics*, 92 (10-11): 2083-2105.
- Lehmann, J. K. (2011); "Job Assignment and Promotion Under Statistical Discrimination: Evidence from the Early Carrers of Lawyers". Unpublished manuscript, University of Houston.
- Leiter, J. and J. Brown (1985); "Determinants of Elementary School Grading". *Sociology of Education* 58: 166-80.
- Lewbel, A. (2012); "Using Heteroskedasticity to Identify and Estimate Mismeasured and Endogenous Regressor Models". *Journal of Business and Economic Statistics*, Volume 30, Issue 1, pp. 67-80.
- List, J. (2004); "The Nature and Extent of Discrimination in the Marketplace: Evidence from the Field". *The Quarterly Journal of Economics*, 119 (1): 49-89.
- Lundberg, S. and R. Startz (1983); "Private Discrimination and Social Intervention in Competitive Labor Markets". *The American Economic Review*, LXXIII, 340-347
- Lundberg, S. and R. Startz (2007); "Information and racial exclusion". *Journal of Population Economics*, vol. 20(3), pages 621-642, July.
- Madeira, R. and M. A. Rangel (2013); "Racial Achievement Gaps in Another America: Discussing Schooling Outcomes and Affirmative Action in Brazil" in J. Clarke (ed) *Closing the Achievement Gap from an International Perspective*, Springer Verlag.
- Metchemberg, L. (2009); "Cheap Talk in the Classroom: How Biased Grading at School Explains Gender Differences in Achievements, Career Choices and Wages", *The Review of Economic Studies*, 76, 1431-1459.
- Natriello, G. and S. Dornbusch (1984); "Bringing Behavior Back In: The Effects of Student Characteristics and Behavior on the Classroom Behavior of Teachers." *American Educational Research Journal* 20: 29-43.
- Neal, D. (2006); "Why has the Black-White Skill Convergence Stopped?" *Handbook of the Economics of Education*, Volume 1.

- Ost, Ben (2014); "How Do Teachers Improve? The Relative Importance of Specific and General Human Capital." *American Economic Journal: Applied Economics*, 6(2): 127-51.
- Perry, G.; O. Arias; J. H. López; W. Maloney and L. Servén (2006); "Poverty Reduction and Growth: Virtuous and Vicious Circles", The World Bank.
- Phelps, E. (1972); "The Statistical Theory of Racism and Sexism", *The American Economic Review*, American Economic Association, vol. 62(4), pages 659-61, September.
- Phillips M.; J. Crouse and J. Ralph (1998); "Does the Black-White Test Score Gap Widen After Children Enter School?" in C. Jencks and M. Phillips (eds.) *The Black-white Test Score Gap*, The Brookings Institution.
- Reardon, S. (2008); "Differential Growth in the Black-White Achievement Gap During Elementary School Among Initially High-and-Low-Scoring Students". Institute for Research on Education Policy and Practice, Working Paper: 2008-07.
- Rigobon, R. (2003); "Identification through Heteroskedasticity". *The Review of Economics and Statistics*, Vol. 85, No. 4, pp. 777-792.
- Rist, R. (1973); *The Urban School: A Factory for Failure*. Cambridge, MA: MIT Press.
- Rosenthal, L. and R. Jacobson (1968); *Pygmalion in the Classroom*. New York: Holt.
- Samejima, F. (1994); "Estimation of Reliability Coefficients Using Test Information Function and its Modifications" *Applied Psychological Measurement*, Vol. 18 (3): September, pp. 229-244.
- Sexton, P. (1961); *Education and Income*. New York: Viking.
- Sewell, W. and R. Hauser (1984); "The Wisconsin Longitudinal Study of Social and Psychological Factors in Aspiration and Achievements". Pp. 59-100 in *Research in Sociology of Education and Socialization*, vol.1, edited by Alan C. Kerchhoff. Greenwich, CT: JAI Press.
- Shay, S. and B. Jones (2006), "Anonymous Examination Marking at University of Cape Town: The Quest for an Agonising-Free Zone". *South Africa Journal of Higher Education* 20 (4): 528-46.
- Sprietsma, M. (2012); "Discrimination in grading: experimental evidence from primary school teachers." *Empirical Economics*, Vol. 45(1), pp. 523-538.
- Telles, E. (2004); "Race in Another America: The Significance of Skin Color in Brazil" Princeton University Press.
- Telles, E. and L. Steele (2012); "Pigmentocracy in the Americas: How is Educational Attainment Related to Skin Color?", *Americas Barometer Insights*, 2012, Number 73.
- Trawalter, S.; K. Hoffman and A. Waytz (2012); "Racial Bias in Perceptions of Others' Pain". *PLOS ONE*, November, Volume 7, Issue 11.
- United Nations (2005); "Brazil: National Human Development Report", New York, USA.
- Williams, T. (1976); "Teacher Prophecies and the Inheritance of Inequality". *Sociology of Education*, 49: 223-36.

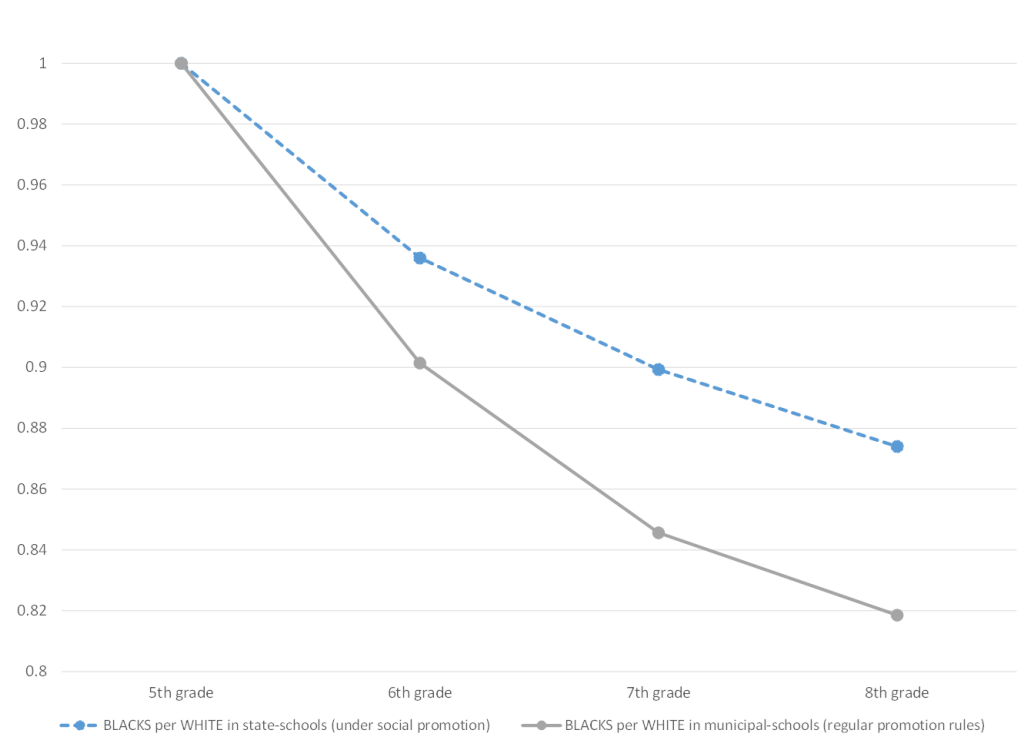


FIGURE 1. DIFFERENTIAL ATTRITION (5TH TO 8TH GRADE) BY RACE AND PROMOTION RULES.

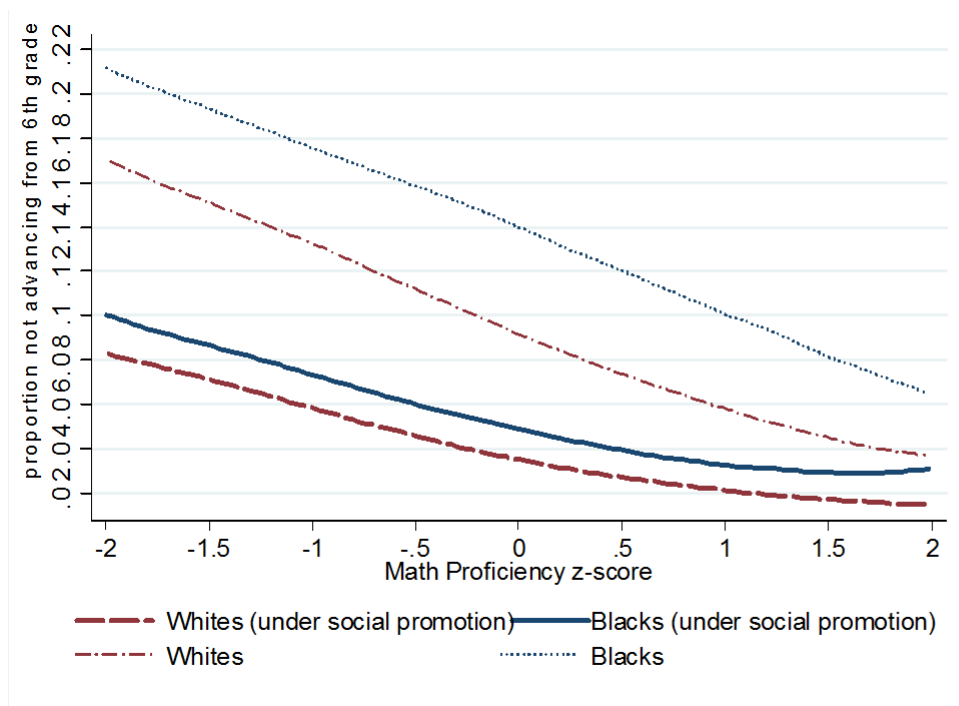


FIGURE 2. DIFFERENTIAL ATTRITION FROM 6TH GRADE BY RACE, PROMOTION RULES AND PROFICIENCY LEVEL.

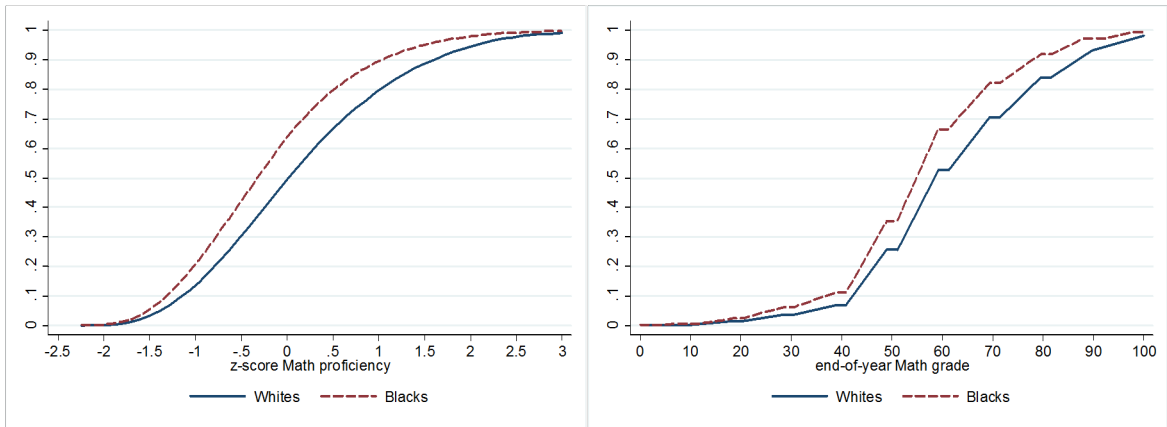


FIGURE 3. CUMULATIVE DISTRIBUTION FUNCTIONS FOR PROFICIENCY SCORES AND TEACHER-ASSIGNED GRADES FOR 8TH GRADERS.

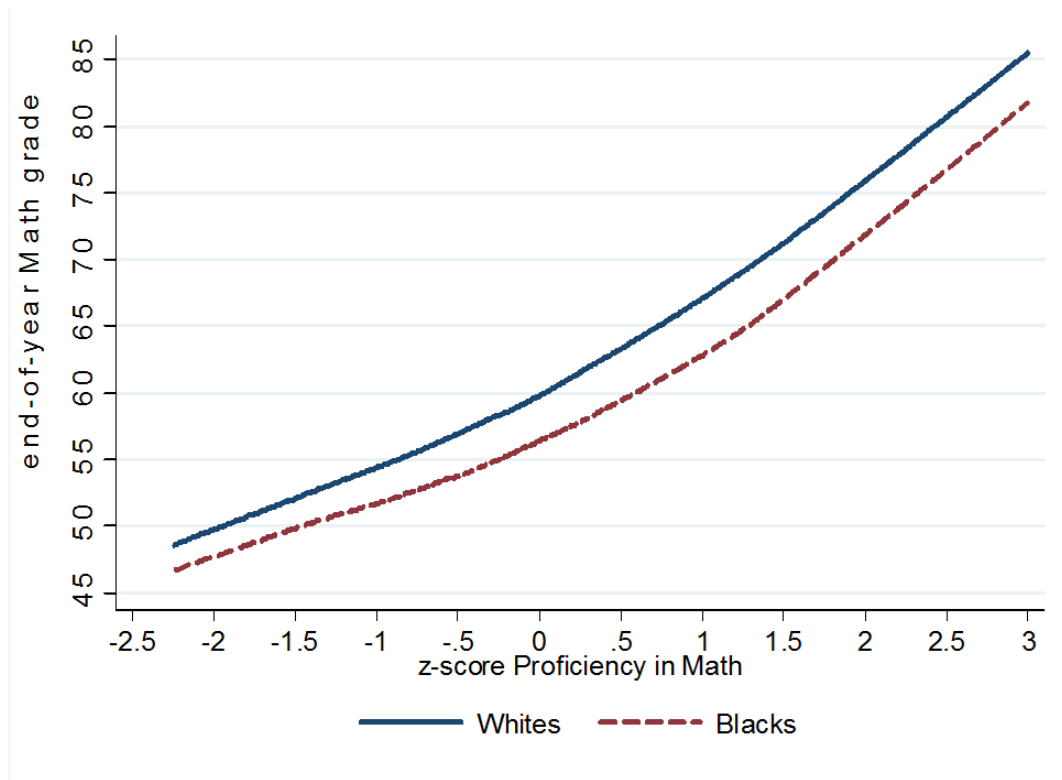


FIGURE 4. SMOOTHED RAW RELATION BETWEEN PROFICIENCY SCORES AND TEACHER-ASSIGNED GRADES FOR 8TH GRADERS.

Table 1
Unconditional and Conditional Racial Differentials in Grading - OLS and IV Estimations

		Black-White raw gaps - OLS			Black-White conditional gaps					
Averages		[1]	[2]	[3]	[4]	[5]	[6]	[7]	[8]	[9]
		Above or at passing grade (>=50)								
<i>Panel A: End-of-year assessment by teacher (0-100 scale)</i>										
White	61.4									
Black	55.7	-5.664*** (0.115)	-4.966*** (0.105)	-1.581*** (0.089)	-0.698*** (0.079)	-0.646*** (0.074)	-0.349*** (0.083)	-0.347*** (0.083)	-0.570*** (0.202)	
Proficiency in Math		3.806*** (0.058)	2.025*** (0.052)	1.823*** (0.048)	9.238*** (2.779)	8.475*** (0.782)	3.412* (1.890)			
Over-ID test (J-statistic [p-value])					1.488 [.2225]	1.478 [.2240]	0.495 [.4817]			
<i>Panel B: Intra-classroom percentile rank of end-of-year assessment by teacher (0-100)</i>										
White	41.9									Above classroom median grade
Black	32.1	-9.813*** (0.200)	-9.887*** (0.205)	-3.250*** (0.177)	-1.522*** (0.161)	-1.369*** (0.153)	-0.735*** (0.175)	-0.721*** (0.172)	-1.177*** (0.275)	
Proficiency in Math		8.124*** (0.117)	4.650*** (0.107)	4.268*** (0.102)	24.780*** (5.501)	20.442*** (1.540)	27.088*** (2.456)			
Over-ID test (J-statistic [p-value])					1.659 [.1977]	1.971 [.1604]	1.998 [.1575]			
<i>Controls</i>										
Classroom fixed-effects	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Child demographics	No	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Performance in standardized tests	No	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Family background + 2009 Math grade	No	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Behavioral traits	No	No	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Order of scores' polynomial	-	-	4th	4th	4th	4th	4th	3rd	3rd	3rd

Note: Standard-errors in parentheses are clustered at the classroom level. ***, **, * 1%, 5% and 10% significance levels. Sample consists of 277,444 students in 10,614 classrooms. Marginal effect of proficiency scores evaluated at the mean proficiency level (for the population) are presented. Controls consist of classroom fixed-effects, child's gender and age polynomial (second order), a 4th-order (3rd order) polynomial function of concurrent Math z-scores interacted with Language z-scores and past performance in essays. Family background includes maternal education, age, region of birth (in or out of state), home ownership, ownership of automobiles, and number of wc's in the household. Behavioral traits include reports of parents regarding child's interest for school work, effort regarding studies and overall behavior. They also include Physical Education grades and Language classes attendance rates for the first half of the school year. Finally, self-reported measures of behavior are included with indicators of procrastination with homework, class-skipping and interest in extra-curricular Math activities.

Table 2
Conditional Racial Differentials in Grading by Teacher's Evaluation Practices and Race -
IV Estimations

	Full sample [1]	Responding quests. [2]	Teacher's Grading Practices			Teacher's Race	
			Objective grader [3]	Subjective grader [4]	Non-Objective grader [5]	White grader [6]	Black + Mixed grader [7]
<i>Panel A: End-of-year assessment by teacher (0-100 scale)</i>							
Black	-0.347*** (0.083)	-0.340*** (0.089)	-0.426*** (0.145)	-0.380*** (0.105)	-0.283** (0.114)	-0.362*** (0.092)	-0.265 (0.194)
Proficiency in Math	8.475*** (0.782)	8.153*** (0.834)	8.687*** (1.560)	6.732*** (1.008)	7.987*** (0.994)	8.809*** (0.855)	7.021*** (1.948)
<i>Panel B: Intra-classroom percentile rank of end-of-year assessment by teacher (0-100)</i>							
Black	-0.721*** (0.172)	-0.743*** (0.187)	-0.957*** (0.305)	-0.884*** (0.220)	-0.627*** (0.236)	-0.741*** (0.191)	-0.582 (0.414)
Proficiency in Math	20.442*** (1.540)	20.412*** (1.642)	18.323*** (2.958)	17.748*** (1.967)	21.597*** (1.978)	21.277*** (1.666)	16.666*** (4.170)
Sample students	277,444	233,750	86,485	171,727	147,846	224,936	52,198
Sample teachers	10,614	8,925	3,305	6,548	5,641	8,596	2,006

Note: Standard-errors in parentheses are clustered at the classroom level. *** 1%, ** 5% and * 10% significance levels. See notes in Table 1.

Table 3
Conditional Racial Differentials in Grading and Learning Students'
Types - IV Estimations

	Math teacher knows student [1]	Math teacher does not know [2]	Difference [3] = [2] - [1]
<i>Panel A: End-of-year assessment by teacher (0-100 scale)</i>			
Black	-0.092 (0.172)	-0.427*** (0.095)	-0.335* (0.197)
Proficiency in Math	11.664*** (1.612)	7.517*** (0.925)	-4.147** (1.862)
Ratio of coefficients: χ -squared test [p-value]	0.29 [0.5911]	15.91 [0.0001]	5.85 [0.0136]
<i>Panel B: Intra-classroom percentile rank of end-of-year assessment by teacher (0-100)</i>			
Black	0.079 (0.364)	-0.960*** (0.197)	-1.038** (0.414)
Proficiency in Math	24.999*** (3.063)	19.119*** (1.826)	-5.880* (3.572)
Ratio of coefficients: χ -squared test [p-value]	0.05 [0.8292]	19.99 [0.0000]	8.55 [0.0035]

Note: Standard-errors in parentheses are clustered at the classroom level. *** 1%, ** 5% and * 10% significance levels. Sample consists of 277,444 students in 10,614 classrooms. Teachers are identified as knowing a given student if they have taught in classes to which the student was assigned between 2007 and 2009. Third column's ratio of coefficients reflect difference of ratios, not ratio of differences. See notes in Table 1.

Table 4
Conditional Racial Differentials in Grading and Learning Students' Types - IV Estimations by Information Level

	Classroom-level acquaintance rate			Tenure in school			Knowledge of neighborhood			
	Difference when knowledgeable specific to student [1]	Math teacher knows 50% or more of class [2]	Math teacher knows less than 50% of class [3]	Difference [4] = [3] - [2]	Math teacher in school 3 years or more [5]	Math teacher in school less than 3 years [6]	Difference [7] = [6] - [5]	Math teacher from school neighborhood [8]	Math teacher NOT from school neighborhood [9]	Difference [10] = [9] - [8]
Panel A: End-of-year assessment by teacher (0-100 scale)										
Black	-0.335* (0.197)	-0.113 (0.168)	-0.422**** (0.098)	-0.309 (0.194)	-0.220** (0.110)	-0.617**** (0.162)	-0.397** (0.195)	-0.390**** (0.101)	-0.198 (0.196)	0.193 (0.220)
Proficiency in Math	-4.147** (1.862)	10.736**** (1.251)	7.493**** (0.999)	-3.243** (1.601)	9.128**** (0.920)	5.053**** (1.935)	-4.075* (2.144)	8.105**** (0.907)	8.219**** (2.350)	0.114 (2.516)
Panel B: Intra-classroom percentile rank of end-of-year assessment by teacher (0-100)										
Black	-1.038** (0.414)	0.088 (0.353)	-1.006**** (0.203)	-1.093**** (0.408)	-0.507** (0.229)	-1.249**** (0.334)	-0.741* (0.404)	-0.926**** (0.213)	-0.169 (0.396)	0.757* (0.450)
Proficiency in Math	-5.880* (3.572)	23.823**** (2.461)	19.176**** (1.965)	-4.648 (3.144)	21.478**** (1.813)	17.378**** (3.728)	-4.100 (4.145)	19.759**** (1.778)	23.375**** (4.770)	3.617 (5.083)
Sample students	277,444	277,444	277,444	277,444	233,750	233,750	233,750	233,750	233,750	233,750
Sample teachers	10,614	10,614	10,614	10,614	8,925	8,925	8,925	8,925	8,925	8,925

Note: Standard-errors in parentheses are clustered at the classroom level. *** 1%, ** 5% and * 10% significance levels. Teachers are identified as knowing a given student if they have taught in classes to which the student was assigned between 2007 and 2009. This info is aggregated at the classroom level in columns (2) and (3). Tenured is defined from responses to teacher questionnaires. The same is the case for neighborhood, which is a function of how far the teacher has to travel to teach in a given school. See notes in Table 1.

Table 5
Conditional Racial Differentials in Grading and Learning Students' Types - IV Estimations for Falsification of Hypothesis

	Language teacher knowledge of student		Future Math teacher knowledge of student		
	Lang. teacher knows student [1]	Lang. teacher does not know [2]	Future Math teacher knows student [4]	Future Math teacher does not know [5]	
			Difference [3] = [2] - [1]	Difference [6] = [5] - [4]	
<i>Panel A: End-of-year assessment by teacher (0-100 scale)</i>					
Black	-0.347* (0.186)	-0.359*** (0.093)	-0.012 (0.208)	-0.463* (0.247)	0.137 (0.263)
Proficiency in Math	9.151*** (1.424)	8.358*** (0.939)	-0.793 (1.705)	9.512** (4.131)	-1.366 (4.219)
<i>Panel B: Intra-classroom percentile rank of end-of-year assessment by teacher (0-100)</i>					
Black	-0.267 (0.405)	-0.841*** (0.191)	-0.574 (0.446)	-0.754 (0.510)	-0.709*** (0.186)
Proficiency in Math	20.533*** (2.857)	20.645*** (1.842)	0.112 (3.402)	18.323** (7.918)	20.144*** (1.562)

Note: Standard-errors in parentheses are clustered at the classroom level. *** 1%, ** 5% and * 10% significance levels. Sample consists of 277,444 students in 10,614 classrooms. Teachers are identified as knowing a given student if they have taught in classes to which the student was assigned between 2007 and 2009 in columns (1) and (2). Future teachers are identified as knowing a given student if they are teaching and will teach in classes to which the student is going to be assigned in 2010 and 2011 in columns (4) and (5). See notes in Table 1.

Table 6
Conditional Racial Differentials in End-of-year assessment by teacher (0-100 scale) and Learning Students' Types -
IV Estimations for Signals Beyond Race and Interactions with Behavioral Traits

	Base Model		Interactions with SES added		Interactions with behavior added	
	Math teacher knows student [1]	Math teacher does not know [2]	Math teacher knows student [3]	Math teacher does not know [4]	Math teacher knows student [5]	Math teacher does not know [6]
Black	-0.092 (0.172)	-0.427*** (0.095)	-0.106 (0.186)	-0.423*** (0.093)	-0.109 (0.216)	-0.423*** (0.105)
Proficiency in Math	11.664*** (1.612)	7.517*** (0.925)	11.690*** (1.619)	7.516*** (0.925)	11.691*** (1.614)	7.519*** (0.925)
Male	-2.876*** (0.152)	-3.052*** (0.087)	-2.864*** (0.152)	-3.055*** (0.087)	-2.857*** (0.153)	-3.057*** (0.087)
<i>Family background (SES)</i>						
Mom HS grad.			0.274** (0.131)	0.309*** (0.074)	0.277** (0.132)	0.308*** (0.074)
Mom some college			0.120 (0.316)	0.385** (0.177)	0.131 (0.316)	0.383** (0.177)
Mom college grad.			0.254 (0.277)	0.427*** (0.150)	0.261 (0.277)	0.425*** (0.150)
Home ownership			0.074 (0.124)	0.160** (0.067)	0.073 (0.124)	0.161** (0.067)
<i>Behavioral traits</i>						
Well behaved (parental report)					1.164*** (0.121)	1.267*** (0.070)
Poorly behaved (parental report)					-0.813*** (0.218)	-0.553*** (0.125)
High-effort behaved (parental report)					0.953*** (0.142)	0.951*** (0.083)
Low effort (parental report)					-0.846*** (0.149)	-0.738*** (0.083)
Level of interest in school work (parental report)					0.490*** (0.032)	0.462*** (0.018)
PE grades					0.692*** (0.040)	0.763*** (0.023)
School attendance (Language classes)					0.118*** (0.008)	0.109*** (0.004)
School attendance (self-report)					1.184*** (0.185)	1.171*** (0.101)
Don't procrastinate (self-report)					1.741*** (0.128)	1.868*** (0.073)
Enrolled in Math extra curricular (self-report)					0.967*** (0.123)	0.972*** (0.067)

Note: Standard-errors in parentheses are clustered at the classroom level. *** 1%, ** 5% and * 10% significance levels. See notes in Table 1.

Table A1
Descriptive statistics - working sample and stratifications

	[1]	[2]	[3]	[4]	[5]	[6]	[7]
	Whites mean (se)	Blacks mean (se)	Black-White (classroom FE) mean (se)	Math Teacher Knows Student Black-White (classroom FE) mean (se)	Math Teacher DOES NOT Know Student Black-White (classroom FE) mean (se)	Lang. Teacher Knows Student Black-White (classroom FE) mean (se)	Lang Teacher DOES NOT Know Student Black-White (classroom FE) mean (se)
<i>Proficiency and school performance</i>							
8th grade Math scores (z-scores)	0.134 (0.0050)	-0.220 (0.006)	-0.249 (0.006)	-0.270 (0.014)	-0.241 (0.007)	-0.256 (0.014)	-0.246 (0.007)
8th grade Language scores (z-scores)	0.202 (0.0050)	-0.243 (0.006)	-0.343 (0.006)	-0.328 (0.014)	-0.344 (0.007)	-0.326 (0.014)	-0.344 (0.007)
Proficiency in past Math exam	217.755 (0.1800)	201.999 (0.242)	-11.521 (0.250)	-10.519 (0.552)	-11.801 (0.282)	-11.233 (0.580)	-11.527 (0.278)
Proficiency in past Language exam	215.870 (0.1920)	197.434 (0.262)	-14.018 (0.265)	-12.524 (0.578)	-14.324 (0.302)	-13.367 (0.608)	-14.061 (0.297)
Proficiency in past Sciences exam	238.236 (0.2330)	216.271 (0.321)	-16.501 (0.325)	-15.521 (0.718)	-16.633 (0.368)	-15.119 (0.723)	-16.761 (0.367)
<i>Family background</i>							
Mom HS grad (%)	22.109 (0.200)	16.976 (0.200)	-3.669 (0.300)	-3.483 (0.600)	-3.674 (0.300)	-3.116 (0.600)	-3.773 (0.300)
Mom some college (%)	3.008 (0.100)	2.425 (0.100)	-0.278 (0.100)	-0.379 (0.200)	-0.248 (0.100)	-0.344 (0.200)	-0.222 (0.100)
Mom college grad (%)	4.495 (0.100)	3.036 (0.100)	-0.751 (0.100)	-0.658 (0.300)	-0.778 (0.100)	-0.926 (0.300)	-0.696 (0.100)
Mom aged 16 to 24 (%)	0.657 (0.001)	1.076 (0.100)	0.283 (0.100)	0.263 (0.100)	0.283 (0.100)	0.533 (0.200)	0.214 (0.100)
Mom aged 25 to 34 (%)	16.863 (0.100)	17.282 (0.200)	0.577 (0.300)	0.280 (0.600)	0.653 (0.300)	-0.133 (0.600)	0.703 (0.300)
Mom aged 45 to 59 (%)	17.020 (0.100)	16.259 (0.200)	-0.005 (0.200)	-0.751 (0.500)	0.196 (0.300)	-0.180 (0.600)	0.020 (0.300)
Mom 60 or older (%)	0.733 (0.100)	1.040 (0.100)	0.320 (0.100)	0.578 (0.200)	0.246 (0.100)	0.348 (0.200)	0.313 (0.100)
Home owned (%)	54.130 (0.200)	50.666 (0.400)	-3.034 (0.300)	-3.355 (0.700)	-2.893 (0.400)	-2.479 (0.700)	-3.130 (0.400)
Number of cars owned	0.583 (0.003)	0.406 (0.004)	-0.119 (0.004)	-0.136 (0.010)	-0.115 (0.005)	-0.118 (0.011)	-0.119 (0.005)
Number of exclusive-use vc's in household	1.069 (0.004)	0.899 (0.006)	-0.111 (0.005)	-0.105 (0.012)	-0.112 (0.006)	-0.114 (0.011)	-0.109 (0.006)
<i>Behavioral traits</i>							
Well behaved (parental report)	0.426 (0.002)	0.316 (0.003)	-0.092 (0.003)	-0.094 (0.007)	-0.091 (0.004)	-0.087 (0.007)	-0.093 (0.004)
High effort (parental report)	0.161 (0.001)	0.135 (0.002)	-0.023 (0.002)	-0.024 (0.005)	-0.022 (0.003)	-0.013 (0.006)	-0.025 (0.003)
Interest level 0-10, (parental report)	5.892 (0.020)	5.315 (0.027)	-0.391 (0.022)	-0.346 (0.046)	-0.395 (0.025)	-0.263 (0.048)	-0.421 (0.025)
PE grade in first bi-monthly evaluation	7.015 (0.012)	6.686 (0.016)	-0.160 (0.011)	-0.168 (0.023)	-0.155 (0.012)	-0.111 (0.024)	-0.166 (0.012)
Language classes attendance in first bi-monthly evaluation	91.705 (0.052)	90.639 (0.077)	-0.528 (0.058)	-0.498 (0.123)	-0.536 (0.066)	-0.391 (0.125)	-0.542 (0.065)
Does not skip classes often (self report)	0.689 (0.002)	0.635 (0.003)	-0.031 (0.003)	-0.029 (0.006)	-0.031 (0.003)	-0.023 (0.007)	-0.033 (0.003)
Does not procrastinate with homework (self report)	0.250 (0.002)	0.183 (0.002)	-0.054 (0.003)	-0.059 (0.006)	-0.052 (0.003)	-0.054 (0.006)	-0.054 (0.003)
Enrolled in extr-curricular Math activities (self report)	0.446 (0.003)	0.371 (0.003)	-0.045 (0.003)	-0.047 (0.007)	-0.044 (0.003)	-0.048 (0.007)	-0.043 (0.003)

Note: Standard-errors in parentheses are clustered at the classroom level. Estimation of differences conducted including classroom fixed-effects. Samples consist of 277,444 students in 10,614 classrooms, of which 10.2

Table A2
First-Stage Regressions' Summary Statistics

	4th-order polynomial		3rd-order polynomial	
	F-test of instruments	P-values	F-test of instruments	P-values
<i>Endogenous variables</i>				
Proficiency score in Math (z-score)	734.89	0.000	863.64	0.000
Proficiency score in Math (z-score) squared	144.13	0.000	1645.96	0.000
Proficiency score in Math (z-score) to the third	348.09	0.000	666.91	0.000
Proficiency score in Math (z-score) to the fourth	152.27	0.000	-	
Language score x Proficiency score in Math (z-score)	756.6	0.000	3465.75	0.000
Language score x Proficiency score in Math (z-score) squared	745.81	0.000	1343.68	0.000
Language score x Proficiency score in Math (z-score) to the third	280.32	0.000	1032.86	0.000
Language score x Proficiency score in Math (z-score) to the fourth	251.02	0.000	-	
Proficiency score in Language (z-score)	3419.5	0.000	7994.92	0.000

Note: Samples consist of 277,444 students in 10,614 classrooms. Instruments are polynomials of past test scores in Math and Language and past test scores in Natural Sciences.

Table A3
Marginal Effects of Proficiency over Grades by Student's Race

	Base model [1]	Interacted model [2]	Differential [2]
<i>Panel A: End-of-year assessment by teacher (0-100 scale)</i>			
Proficiency in Math	8.475*** (0.782)		
Proficiency in Math x White		8.290*** (0.344)	
Proficiency in Math x Black		8.549*** (0.751)	0.259 (0.786)
<i>Panel B: Intra-classroom percentile rank of end-of-year assessment by teacher (0-100)</i>			
Proficiency in Math	20.442*** (1.540)		
Proficiency in Math x White		18.611*** (0.744)	
Proficiency in Math x Black		19.048*** (1.501)	0.437 (1.593)
<i>Sample students</i>	277,444	277,444	277,444
<i>Sample teachers</i>	10,614	10,614	10,614

Note: Standard-errors in parentheses are clustered at the classroom level. *** 1%, ** 5% and * 10% significance levels. Standard-errors in parentheses are clustered at the classroom level. Sample consists of 277,444 students in 10,614 classrooms. See notes in Table 1.