

Who Benefits From Attending Effective High Schools? *

C. Kirabo Jackson
Northwestern University

Sebastián Kiguel
Northwestern University

Shanette C. Porter
Student Experience Research Network

John Q. Easton
UChicago Consortium on School Research

August 11, 2022

Abstract

We estimate the longer-run effects of attending an effective high school (one that improves a combination of test scores, survey measures of socio-emotional development, and behaviors in 9th grade) for students who are more versus less educationally advantaged (i.e., likely to attain more years of education based on 8th-grade characteristics). All students benefit from attending effective schools, but the least advantaged students experience larger improvements in high-school graduation, college going, and school-based arrests. This heterogeneity is not solely due to less-advantaged groups being marginal for particular outcomes. Commonly used test-score value-added understates the long-run importance of effective schools, particularly for less-advantaged populations. Patterns suggest this *partly* reflects less-advantaged students being relatively more responsive to non-test-score dimensions of school quality.

*Jackson: kirabo-jackson@northwestern.edu. Kiguel: skiguel@u.northwestern.edu. Porter: shanette@studentexperiencenetwork.org. Easton: jqeaston@uchicago.edu. The authors thank the staff at Chicago Public Schools, particularly the Office of Social and Emotional Learning and the University of Chicago Consortium on School Research, for providing access to, and information about, the Chicago Public Schools data. This paper benefited from discussion with seminar participants at the UChicago Consortium, and data management was facilitated by their archivist, Todd Rosenkranz. The authors acknowledge funding for this research from the Bill Melinda Gates Foundation. The content is solely the responsibility of the authors.

I Introduction

A growing body of research in the social sciences finds that schools have important causal effects on students' longer-term outcomes. For example, some charter schools increase college-going (e.g., [Angrist et al. 2016](#); [Sass et al. 2016](#)), some selective schools improve educational attainment, wages, and health (e.g., [Jackson 2010](#); [Beuermann and Jackson 2020](#)), and winning a school-choice lottery may increase college-going for girls and reduce interaction with law enforcement among certain boys (e.g., [Deming 2011](#) [Deming et al. 2014](#)). However, questions remain about whether the benefits of attending better schools differ for better or worse-prepared students. We seek to understand if “effective schools” (i.e., those that improve both test scores and Socioemotional Development (SED)) confer similar longer-run impacts on more and less *educationally advantaged* (i.e., likely to attain more years of education based on 8th-grade characteristics) students.

In principle, the least educationally advantaged students may benefit most from effective schools because they may have more room for improvement. Conversely, if “*skills beget skills*” ([Cunha et al., 2010](#)), effective schools, on average, may have small impacts on the least advantaged. The existing empirical evidence on this topic is mixed and (to overcome selection issues) has been focused on small numbers of oversubscribed charter schools or elite schools that rely on admission lotteries or admission tests.¹ Moreover, because these studies rely on comparisons among *applicants* to these special oversubscribed schools (who may differ from typical students), patterns in these studies may differ from those in the broader student population ([Bruhn, 2020](#)). That is, disadvantaged students who apply to elite schools or charter schools are unlikely to be representative of the typical disadvantaged student ([Hoxby and Murarka, 2009](#)). As such, whether the causal impacts of attending a better school differ by academic advantage across a representative sample of schools and students is unknown. By exploring differences in the effect of attending more effective schools across all schools and all students in a large public school district, we shed light on this issue.

A second motivation for our work is that both economists and social psychologists have found that differences in SED or (or non-cognitive skills) may explain attainment gaps by gender ([Jacob, 2002](#)) and socio-economic status ([Claro et al. 2016](#); [Liu 2020](#)). Moreover, experimental studies in psychology find that (a) students from low-income families or who are academically lower-achieving may benefit from mindset interventions ([Sisk et al., 2018](#)), and (b) interventions that promote a sense of belonging are beneficial for the educational outcomes of minoritized (including Black and Latinx) youth ([Walton and Cohen 2007](#); [Walton and Cohen 2011](#); [Gray et al. 2018](#);

¹[Angrist et al. \(2012\)](#), [Walters \(2018\)](#), and [Cohodes et al. \(2020\)](#) finds that less-advantaged Boston area charter applicants benefit more from attending oversubscribed charter schools. Conversely, looking at Charter-like schools in India [Kumar \(2020\)](#) finds little difference between more and less advantaged students. Looking at elite school attendance, [Dustan et al. \(2017\)](#) and [Oosterbeek et al. \(2020\)](#) find negative effects in Mexico City and Amsterdam, while [Shi \(2020\)](#) finds the opposite in North Carolina, and [Barrow et al. \(2020\)](#) report mixed results in Chicago.

Brady et al. 2020; Murphy et al. 2020). If so, measures of school quality that exclude impacts on SED may miss important components of school quality for disadvantaged or minority populations. To shed light on this issue, we identify school effectiveness explained by test score value-added only versus value-added in other dimensions (i.e., socio-emotional development and behaviors) and examine these different effects by educational advantage.

We leverage detailed data from Chicago Public Schools obtained from the [UChicago Consortium on School Research](#). These data link K12 students to high schools and colleges along with test scores, administrative records, and self-reported survey measures of SED over time. Our project entails three key steps: (1) First, we categorize students. We use student behaviors, survey measures, and test scores in 8th grade to predict their educational outcomes years later. We then use this model to create a latent educational advantage index for each student. (2) Next, we identify schools' impacts on student outcomes (i.e., value-added) by comparing end-of-year outcomes across schools while conditioning on lagged outcomes and other covariates. Value-added models yield results similar to causal estimates from randomized school assignments in several U.S. settings (Deming et al. 2014; Angrist et al. 2020). As in Jackson et al. (2020) we estimate schools' impacts on test scores and socio-emotional survey measures in 9th grade. We build on Jackson et al. (2020) by *also* estimating school impacts on behaviors (Heckman and Kautz 2012) and combining effects across all outcomes to create an overall school effectiveness index. (3) Finally, we estimate the effect on educational attainment and school-based arrests of attending a more effective school for students with different levels of educational advantage. We also explore differences between schools that improve test scores versus other dimensions (i.e., behaviors and survey measures).

First, we describe the different indexes. (1) The educational advantage index differentiates between groups of students who are more or less likely to graduate high school, enroll in college, and attend a 4-year college. Students who are low in this index are more likely to have low 8th grade test scores, low socio-emotional measures, and more absences and disciplinary incidents than those who are high on the index. Students low on this index are also more likely to come from low-income homes and be male and Black – student populations that are hypothesized to benefit the most from socio-emotional interventions. (2) Turning to our measure of school effectiveness, our *overall* school effectiveness index is a better predictor of school effects on longer-run outcomes (i.e., educational attainment and crime) than test score value-added alone. The effectiveness index is weakly related to school demographics and student-teacher ratios and is strongly correlated with college-going rates. We describe the characteristics of effective schools, show the consistency of our estimates with existing studies, and validate our estimates as reflecting causal impacts.

Focusing on heterogeneous effects, *all* students benefit from attending more effective schools – rejecting a model in which only the most advantaged, or marginal, students benefit from better schools. Looking at the longer-run outcomes, less advantaged students experience larger marginal

effects. Specifically, for those in the bottom decile of the distribution, attending a school at the 85th percentile of the effectiveness distribution versus one at the median is associated with a 3.8pp increase in high school graduation, a 3.6pp increase in college-going, and a 2.1pp reduction in being arrested. The corresponding estimates for those in the top decile is a 0.5pp increase in high school graduation, a 1.5pp increase in college-going, and a 0.22pp reduction in being arrested. These patterns are similar within ethnic groups and by gender, so are not driven by comparisons across broad demographic groups. We test for whether the pattern of effects reflects disadvantaged students being more likely to be marginal for certain outcomes. This may explain the patterns for arrests and high-school graduation, but does *not* explain the patterns for college outcomes.

To help explain these differential impacts by educational advantage, we look at the different components of school effectiveness. While test score value-added and the effectiveness index have similar marginal effects for the most educationally advantaged, for the least advantaged, overall effectiveness predicts much larger effects on longer-run outcomes than test score value-added alone. Additional patterns suggest that this is due, *in part*, to less academically-advantaged students being relatively more responsive to school impacts on SED or soft skills.

Our estimates suggest that the least-advantaged students gain the most from attending effective schools. However, in our data, the least-effective schools are disproportionately attended by the least-advantaged students. We simulate a policy that shuts down the least effective schools and re-assigns students to other schools. This leads to sizable gains for reassigned students, small negative spillover effects for students at receiving schools, and positive benefits overall – underscoring the likely gains from improving the schools attended by the least-advantaged students.

We contribute to the school-quality literature in several ways. We demonstrate that across all public schools (as opposed to a handful of elite or charter schools) in a large district, all students benefit from attending more effective schools with particularly large benefits for disadvantaged students. We also contribute to this literature by moving beyond a test-score measure of effectiveness and incorporating both survey and behavior-based measures of soft skills. We show how students with varying levels of educational advantage benefit from schools that raise test scores versus socio-emotional skills and behaviors. Finally, we show how test-score measures of school quality may understate the benefits of effective schools – particularly for disadvantaged students.

The remainder of the paper proceeds as follows: Section II describes the data used, and Section III details the methods used. Section IV validates our methodology as representing causal impacts. The results are presented in Section V, and Section VI concludes.

II Data

As in Jackson et al. (2020), we use administrative data from Chicago Public Schools (CPS) obtained from the UChicago Consortium on School Research. CPS is a large urban school district. Our data covers on six cohorts of 9th-graders between the spring of 2012 and 2017 during which there were 150 general education public (neighborhood /charter/ vocational/ magnet) high schools.² We can estimate effectiveness for 143 of these schools³ that serve largely ethnic minority students (40% Black and 46% Latinx) and economically disadvantaged students (85%).⁴ We only include the first observations for each 9th-grader ($n=160,148$) to remove sample selection biases due to grade repetition. For high school graduation and school-based arrests, we focus on three cohorts of 9th graders between 2012 and 2014 ($n=82,146$), and for college outcomes, we focus on two cohorts of 9th graders in 2012 and 2013 ($n=55,560$) because these students are old enough to have attended college.⁵ The data and sample are summarized in Table 1, and discussed below.

Survey Measures: Some of our key variables are survey measures of social-emotional development (SED). The SED constructs captured by these surveys are hypothesized to be particularly important for the success of disadvantaged youth. Responses are collected by CPS on a survey administered to students in 2008-09, and then every year from 2010-11 onward. These survey items are not part of Chicago’s accountability system and response rates were high (75%). However, nonresponse was higher for low-achievers (Appendix Table A2). Note that our analysis of impacts on longer-run outcomes is based on all students irrespective of survey completion. Each survey measure was comprised of several items and students responded to each item using point scales to indicate agreement (e.g., 1=Strongly disagree, to 4=Strongly agree). Rasch analysis was used to model responses and calculate a score for each student on each construct (for measure properties, see Appendix Table A3). Following Jackson et al. (2020), we combine the interpersonal-related questions into a **Social Index**⁶ and the work-related questions into a **Work Hard Index**.⁷ To

²CPS offers more than 160 high school programs (link). Our count only includes the general education programs.

³Of the 7 excluded schools, 5 are only observed in the data for a single year – precluding an out-of-sample school estimate – and 2 do not have enough data to estimate value-added (i.e., fewer than 10 students with survey data and test score). The included schools cover over 95 percent of all 9th-graders in these cohorts.

⁴Note that tests were not administered in 2017 and are therefore missing in that year.

⁵See Appendix Table A1 for a breakdown of the number of observations for each outcome by cohort.

⁶Two of the SED survey measures relate to one’s relationship with others in the school. The first is Interpersonal Skills, and the second is a measure of Belonging. **Interpersonal Skills** includes: I can always find a way to help people end arguments. I listen carefully to what other people say to me. I’m good at working with other students. I’m good at helping other people. **Belonging** includes: I feel like a real part of my school. People here notice when I’m good at something. Other students in my school take my opinions seriously. People at this school are friendly to me. I’m included in lots of activities at school.

⁷Three survey measures capture students’ orientation toward hard work. These are Academic Effort, the perseverance facet of Grit, and Academic Engagement. **Academic Effort** includes: I always study for tests. I set aside time to do my homework and study. I try to do well on my schoolwork even when it isn’t interesting to me. If I need to study, I don’t go out with my friends. **Grit** includes: I finish whatever I begin. I am a hard worker. I continue steadily toward

create each index, we take the first predicted principal factor among the survey items. This is a weighted average of the individual items that represents the maximum variance direction in the data (Jolliffe, 2002). We then standardize the predicted factors to create indexes that have mean zero and unit variance. All indexes used in this paper adopt this approach and have correlation over 0.95 with the simple mean of standardized items.

Behavior Measures: Motivated by work showing that impacts on behaviors measure skills not well captured by test score impacts (e.g., Heckman et al. (2013); Jackson (2018); Liu and Loeb (2019), Petek and Pope (2020)), we augment the survey-based measures used in Jackson et al. (2020) and *also* include student behaviors from CPS administrative data. These include the number of excused and unexcused absences, severe disciplinary incidents (eligible for suspension), and days a student is suspended, in each grade. In the analytic sample with non-missing behaviors ($n=157,628$), the average 9th grader is absent 15 days and suspended 0.8 days. Approximately 7.7% of these are involved in a severe disciplinary incident. We summarize these three 9th-grade measures by forming a **behaviors Index** which is the standardized first predicted principal factor.

Test Score Measures: The “hard” skills measure in our data is standardized test scores. To allow for comparability across grades, test scores were standardized to be mean zero unit variance within grade and year among all CPS test takers. Tests were not administered to 9th graders in Spring 2017, so there is nontrivial missingness for this 9th-grade outcome. However, Appendix Table A2 shows that the sample of test takers is similar to the full sample. We form a **Test Score Index** using the standardized first predicted principal factor among standardized math and English scores.

Long-Run Outcomes: A key longer-run outcome is having a school-related arrest (among those old enough to have graduated high school). These are arrests for activities conducted on school grounds, during off-campus school activities, or due to a referral by a school official. During our sample period, 3.7 percent of first-time 9th graders had a school-based arrest, 5.2 percent of males, and 7.7 percent of Black males. Roughly 20 percent of juvenile arrests in 2010 were school-based arrests (Kaba and Edwards, 2012), so these have important long-term implications. Our other longer-term outcomes include high school graduation and enrollment and persistence in college. High school completion is obtained from school leaving files from the years 2010 through 2018. We define a student as having graduated high school if they are marked as leaving high school because they graduated. Note that this definition does not include equivalency credentials such as a GED. About 74 percent of first-time 9th graders in CPS graduate high school. Our second key long-run outcome is enrollment in college. Our college data come from the National Student Clearinghouse (NSC) and are merged with all CPS graduates. We code a student as enrolling in

my goals. I don't give up easily. **Academic Engagement** includes: The topics we are studying are interesting and challenging. I usually look forward to this class. I work hard to do my best in this class. Sometimes I get so interested in my work I don't want to stop.

college if they are observed in the NSC data within two years of expected high school graduation (2012 and 2013 cohorts only). Students who did not graduate high school are coded as not enrolling in college, so our outcomes indicate “*graduating high school and enrolling in college.*” Because graduating high school is a prerequisite for college enrollment, these two definitions are largely the same.⁸ About 53 percent of first-time 9th graders enrolled in college. We further divide college enrollment into 2-year and 4-year college. In our sample, 34 and 28 percent of students enroll in a 4-year or 2-year college within 2 years of expected graduation, respectively.

III Methods

Our analysis involves three main steps: (1) First, we calculate an educational advantage score for each student by estimating their predicted educational attainment based on a rich set of covariates using an ordered probit. We place students into deciles from least to most likely to attain more years of education. (2) Second, following Jackson et al. (2020), we identify schools that improve students’ SED and test scores in 9th grade. In addition, we estimate school value-added on student behaviors using the same method. We combine school effects on the different 9th-grade measures - which are predictive of students’ long-term outcomes - into an index of school effectiveness. (3) Finally, we estimate the effect of attending a more effective school among students of differing educational advantage to assess who benefits from attending better schools. We also explore the effects of different value-added dimensions to shed light on whether schools that are better in some dimensions (cognitive, socio-emotional, or behaviors) are better for some students than for others.

III.1 Classifying Students

To classify students along a single dimension, we rank students by their likelihood to attain more years of education. We refer to students who are more likely to attain more years of education (based on observed characteristics *before* entering high school) as more educationally advantaged. We exploit the fact that we have a rich set of observable characteristics that may predict educational attainment and also multiple measures of educational attainment. In principle, with a single measure of educational attainment (say high school graduation) one could predict high-school completion based on observed covariates in 8th grade. However, because some characteristics may matter more for higher levels of education (such as 4-year college attendance), it is helpful to model the relationship between these covariates and 4-year college going also. If the underlying educational advantage predicts both high-school completion and college-going (or any other educational attainment level), one can model a student’s underlying educational advantage (in a way that will predict multiple educational attainment margins) using a rank-ordered probit.

⁸A limitation of our data is that it will miss students who enroll as 9th graders but who graduate from a high school outside of CPS. This is relatively rare, particularly for the populations for which we find the most pronounced effects.

The basic idea is that some underlying educational advantage, y^* , is a linear function of observable characteristics X so that $y^* = X\pi + \varepsilon$. Individuals with higher levels of educational advantage attain higher levels of education, where there are some unobserved thresholds between education levels. That is, for all individual i

$$y_i = \begin{cases} \text{No High School} & y_i^* \leq \chi_1 \\ \text{Graduate High School} & \chi_1 > y_i^* \leq \chi_2 \\ \text{Attend a 2-Year College} & \chi_2 > y_i^* \leq \chi_3 \\ \text{Attend a 4-Year College} & y_i^* > \chi_3 \end{cases}$$

The probability of observing outcome $y_i = k$ is then $Pr(y_i = k) = Pr(\chi_{k-1} < X\pi \leq \chi_k)$. The probability of observing the data is the product of these probabilities across all individuals i . Assuming a normally distributed error term, we solve for the set of estimates $(\hat{\pi}, \hat{\chi}_{k-1}, \hat{\chi}_{k-1}, \hat{\chi}_{k-1})$ that are most consistent with the observed data by estimating an ordered probit model by maximum likelihood.

Our predictors of the education outcomes include measures of lagged test scores (quadratics of 8th grade math and ELA), 8th grade survey measures, and lagged behaviors.⁹ We also include demographics (lunch status, race, gender, and interactions between race and gender). Once the parameter estimates have been estimated, we take the fitted values of the latent variable, $\hat{L}\hat{A}_i = X\hat{\pi}$, as our estimated latent educational advantage. *Note that, we use leave-year-out models to avoid mechanical correlation between our predicted and actual education levels for each student i .* As such, each student's predicted educational advantage index is based on the relationship between covariates and educational attainment in *other* cohorts. However, to show the relationship between the advantage index and the observable covariates we present the coefficient estimates from the ordered probit model for all students in the 2012 and 2013 cohorts in Appendix Table A4.

Differences in Incoming Attributes by Educational Advantage

We present summary statistics for the top and bottom deciles of the educational advantage distribution in the middle and right panels of Table 1. This categorization captures important differences between students, both in terms of demographics and achievement. The top decile contains almost three times more females than the bottom decile (71% versus 23%), about 8 times fewer students in special education (5.5% versus 46%), and less than half the share of students who qualify for free lunch (43% versus 95%). The top decile has more white students than the bottom decile (23% versus 3.8%), more Asian students (19% versus 0.14%), but with lower shares of Latinx students (34.1% versus 40.7%) and Black students (22.5% versus 55%). Regarding academic achievement, students in the lowest decile have 8th and 9th-grade test scores more than two standard deviations below those in the top decile. Students in the top decile also have fewer absences (5.6 compared

⁹Because these variables do not have a lot of variation in early grades, we include an indicator for being in the top quartile of absences in 8th grade and an indicator for having any severe disciplinary incidents in 7th or 8th grade.

to 33.7 days) and days suspended (0.065 vs. 2.9 days), and are involved in fewer severe incidents (0.007 vs 0.29), relative to the lowest decile in 9th-grade.

Differences in Outcomes by Educational Advantage

To illustrate the differences in our main longer-run outcomes by the latent educational advantage index, we compute the average of our key outcomes for each percentile of the index. This is presented graphically in [Figure 1](#). At the bottom of the index (the bottom 20 percent), even though about 40 percent of students graduate from high school, few (about 17 percent) go to any college, and even fewer (10 percent) attend a 4-year college. In the middle of the distribution (between the 40th and 60th percentiles), the high school graduation rate is about 75%, the college-going rate is about 50% and both the 4-year and 2-year college-going rates are around 25%. As one looks to the top of the distribution (the top 20%), the high school graduation rate is above 90%. Interestingly, the 4-year college-going rate increases to about 70%, while the 2-year college rate remains at 25%. That is, as one goes up the educational advantage distribution, 4-year college going increases but 2-year college going does not. Indeed at the very top of the educational advantage distribution, the 2-year college rate declines with educational advantage. Looking at arrests, school-based arrests are largely concentrated among students with very low educational advantage. For the bottom 20% the arrest rate is roughly 8 percent, while for those above the median it is almost zero (0.02%). Indeed, in the very bottom decile, the arrest rate is a sizable 12.5 percent (see [Table 1](#)). It is important to note that even though our educational advantage varies by ethnicity and gender, the patterns also hold within groups. [Appendix Figure A1](#) shows that the average outcomes by educational advantage within groups are similar to those overall so that the heterogeneity analysis is not merely based on comparisons across the broad demographic groups.

III.2 Classifying Schools

We use value-added models to estimate schools' impacts on 9th-grade SED, behaviors, and test scores. We then combine the value-added estimates across outcomes to form an overall school effectiveness index. We present evidence that these impacts can be interpreted causally.

Identifying School Impacts on SED, Behaviors, and Test Scores

We seek to isolate the causal effects of individual schools on student measure $q \in Q = \{\text{test scores, work hard, social, behaviors}\}$ by comparing measure q at the end of 9th grade to those of similar students (with the same survey measures, course grades, incoming test scores, discipline, attendance, and demographics, all at the end of 8th grade) at other schools. School j 's value-added on measure q reflects how much school j increases measure q between 8th and 9th grade relative to the changes observed for similar students (based on all the attributes above) who attended different schools. We model the 9th grade measure q of student i who attends school j with observable

characteristics Z_{ijt} in year t as (1) below.

$$q_{ijt} = \underbrace{\beta_q Z_{ijt}}_{\text{Effect of Observables}} + \underbrace{\tau_{t,q}}_{\text{Cohort Fixed Effect}} + \underbrace{\alpha_{j,q} + \varepsilon_{ijt,q}}_{\text{Combined error } v_{ijt,q}} \quad (1)$$

Z_{ijt} includes lagged measures (i.e., 8th grade test scores, surveys, behaviors), gender, ethnicity, free-lunch status, and the socio-economic status of the student's census block.¹⁰ $\tau_{t,q}$ is a cohort fixed effect. To account for correlation between schools and covariates, we follow Chetty et al. (2014) and include a school-specific intercept, $\alpha_{j,q}$, so that β_q is estimated using variation across students at the same school, and $\varepsilon_{ijt,q}$ is a within-school student-level error.¹¹ Because β_q is estimated using time-varying within-school variation in individual student attributes, $\alpha_{j,q}$ may reflect the effect of persistent differences in peer demographics across schools. Where $v_{ijt,q} = \alpha_{j,q} + \varepsilon_{ijt,q}$, is the true combined error, $u_{ijt,q}$ is the empirical combined residual. The school-year average combined residuals from this regression is our estimated impact on measure q of attending a school in a given year, $\hat{\theta}_{jt,q}^{VA}$. Formally, where N_{jt} is the number of students attending school j in year t ,

$$\hat{\theta}_{jt,q}^{VA} = \sum_{i \in jt}^{jt} (u_{ijt,q}) / N_{jt} \quad (2)$$

When using value-added to *predict* outcomes for a particular cohort, we exclude data for that cohort when estimating value-added to avoid mechanical correlation. To aid precision, we follow Chetty et al. (2014) and use value-added with drift which places more weight on value-added for *other* years that are more highly correlated with the prediction year.¹² Our leave-year-out predictor for measure q in year t is (3) where the vector of weights $\hat{\psi}_q = (\hat{\psi}_{t-l,q}, \dots, \hat{\psi}_{t-1,q}, \hat{\psi}_{t+1,q}, \dots, \hat{\psi}_{t+l,q})'$ are selected to minimize mean squared forecast errors.

$$\hat{\mu}_{jt,q} = \sum_{m=t-l}^{t-1} \hat{\psi}_{m,q} [\hat{\theta}_{jm,q}^{VA}] \quad (3)$$

A school's predicted value-added on measure q is our best prediction *based on other years* of how much that school will increase measure q between 8th and 9th grade relative to the improvements of similar students at other schools. We use leave-year-out predictions for all analyses, but for brevity, refer to them simply as value-added.

¹⁰The census block SES measure is the average of occupation status and education levels in the block.

¹¹Results are similar using models that exclude $\alpha_{j,q}$ and include school-level averages of the individual covariates.

¹²If all years value-added were equally predictive of outcomes in year t , then the best leave-year-out predictor for a school would be the average value-added for that school *in all other years*. However, adjacent years tend to be more highly correlated with one another than less temporally proximate years (see the top panel of Table 2).

Correlations Across Effects on Different Measures

Each value-added measure may represent impacts on a different dimension or may reflect the same underlying school quality. To assess this, we correlate school impacts across these four measures. **For these cross-sectional correlations only** we only need one observation per school. Accordingly, our value-added measure for each school in all years, $\hat{\theta}_{j,q}^{VA}$, uses the full sample and is given by (4) below, where N_j is the number of students assigned to school j across all years.

$$\hat{\theta}_{j,q}^{VA} = \sum_{it \in j}^j (u_{ijt,q}) / N_j \equiv \underbrace{\theta_{j,q}^{VA}}_{\text{Real value-added on outcome } q} + \underbrace{\zeta_{j,q}}_{\text{Estimation error}} \quad (4)$$

Where q_1 and q_2 connote different outcomes, we report the raw pairwise correlations between the value-added for different outcomes (i.e., $\text{Corr}(\hat{\theta}_{j,q_1}^{VA}, \hat{\theta}_{j,q_2}^{VA})$ in the middle panel of Table 2). Because each of these value-added is measured with error, the true correlations could be higher than this (due to attenuation bias from random estimation errors), or lower than this (if the measurement errors are correlated across outcomes in the same year). Accordingly, we follow [Beuermann et al. \(2022\)](#) and implement a split-sample approach that uncovers the real correlation between effect across outcomes under the assumption that estimation errors are unrelated over time both within and across outcomes (as in [Kane and Staiger \(2008\)](#) and [Jackson \(2013\)](#)).¹³

We report these clean dissattenuated correlations in the bottom panel of Table 2. The general patterns of these clean dissattenuated correlations are broadly similar to those of the raw correlations. Test score value-added is strongly related to the social dimension of SED value-added ($\rho = 0.85$), and has moderate correlations with the work hard dimension of SED value-added ($\rho = 0.32$), and the behaviors value-added ($\rho = 0.44$). This indicates that there may be some underlying dimension of school quality that is associated with higher value-added in all these dimensions. Another interesting pattern is that behaviors value-added are more strongly related to test-score value-added than the survey-based SED value-added – suggesting that surveys and be-

¹³If measurement errors are uncorrelated over time, then the correlation between the value-added estimated using data only during the even years for one outcome ($\hat{\theta}_{j,\text{even},q_1}^{VA} = \theta_{j,q_1}^{VA} + \zeta_{j,\text{even},q_1}$) and those only using the odd years for the other outcome ($\hat{\theta}_{j,\text{odd},q_2}^{VA} = \theta_{j,q_2}^{VA} + \zeta_{j,\text{odd},q_2}$) would not be biased by correlated errors across outcomes because $\text{corr}(\zeta_{j,\text{even},q_1}, \zeta_{j,\text{odd},q_2}) = 0$. However, this even-odd cross outcome correlation ($\hat{\rho}_{12}^{\text{even-odd}}$) will be attenuated by random estimation errors. When two variables are measured with random errors, the raw correlations between the two variables (in this case, $\hat{\rho}_{12}^{\text{even-odd}}$) reflect the true correlation times the square root of the product of the reliability of each outcome ([Spearman, 1904](#)). In this case, that is $\hat{\rho}_{12}^{\text{even-odd}} = \rho_{12}^{\text{even-odd}} \sqrt{(R_{\text{even},q_1} R_{\text{odd},q_2})}$, where R_{even,q_1} is the reliability of $\hat{\theta}_{j,\text{even},q_1}^{VA}$ and R_{odd,q_2} is the reliability of $\hat{\theta}_{j,\text{odd},q_2}^{VA}$. As such, one can disattenuate the raw correlation by dividing by the square root of the product of the reliability ratios for each measure ([Spearman, 1987](#)). The reliability of each measure can be obtained using the correlation between even and odd year estimates for the same outcome (i.e., $\hat{\rho}_{11}^{\text{even-odd}}$ and $\hat{\rho}_{22}^{\text{even-odd}}$). With the clean raw correlations ($\hat{\rho}_{12}^{\text{even-odd}}$) and the estimated reliability ratios ($\hat{\rho}_{11}^{\text{even-odd}}$ and $\hat{\rho}_{22}^{\text{even-odd}}$), we compute clean dissattenuated correlation estimates $r_{12} = [\hat{\rho}_{12}^{\text{even-odd}}] / (\sqrt{(\hat{\rho}_{11}^{\text{even-odd}} \hat{\rho}_{22}^{\text{even-odd}})})$.

aviors may measure different dimensions of socio-emotions skills. Consistent with this notion, the behaviors value-added is moderately correlated with the social dimension of SED value-added ($\rho = 0.26$), and weakly related to the work hard dimension of SED value-added ($\rho = 0.084$). While the relatively low correlations among the survey- and behavior-based measures of socio-emotional skills is an important finding, we leave exploration into *why* to future work.

Creating an Overall School Effectiveness Index

To further understand correlation patterns in the data, we conduct factor analysis of the school effects (Appendix Table A5). The model finds that a single underlying factor explains almost all the common variation in these value-added. This single factor is positively related to all the value-added indicating that it is related to the schools’ quality across all dimensions. As such, we combine our value-added (work hard, social, test-scores, and behaviors) into a single index of school effectiveness. Our overall index is the predicted first principal factor of these four variables. The overall index, $\hat{\omega}_{jt}$, is a weighted average of the different value-added estimates given by (5) and represents a measure of school impacts on 9th grade measures that is shared across the SED (work hard and social), test score, and behavior dimensions.¹⁴

$$\hat{\omega}_{jt} = (0.3)\hat{\mu}_{jt,testscores} + (0.34)\hat{\mu}_{jt,workhard} + (0.26)\hat{\mu}_{jt,social} + (0.1)\hat{\mu}_{jt,behaviors} \quad (5)$$

The weightings indicate that the overall index value-added is positively related to value-added in all dimensions – capturing the best single summary measure of school quality based on these different value-added measures. We standardize the quality index to be mean zero, unit variance. As we show in Section IV, the index is generally a better predictor of school impacts on longer-run outcomes than the value-added on the individual measures.

III.3 Some Measurable Differences Between More and Less Effective Schools

To provide some sense of these schools that are effective based on this metric, we present binned scatterplots of observable school attributes by 20 ventiles of school effectiveness in Figure 2. We also report the correlation between the school attribute and school effectiveness. A few patterns emerge. First, despite a strong correlation between average outcomes and student demographics, among the least and most effective schools there are small differences in the proportion of non-white students (0.95 vs 0.93 percent in the top and bottom ventiles, respectively), and those on free and reduced-priced lunch (0.93 vs. 0.91 percent in the top and bottom ventiles, respectively). This echoes results from randomized lotteries in New York and Denver (Angrist et al., 2022). The most effective schools tend to be larger than the least effective schools and have slightly larger class sizes. Noble charter schools are over-represented among the most effective schools (consistent with

¹⁴The weights have been normalized to sum to 1 for ease of interpretation.

lottery-based studies (e.g., [Davis and Heller \(2019\)](#)), and selective-enrolment schools are slightly more likely to be among the most effective schools (consistent with [Barrow et al. \(2020\)](#) who find positive effects on college-going using a Regression Discontinuity design).¹⁵ As one might expect, given the weak correlations with incoming demographics, more effective schools are also those that have better high school graduation and college-going rates on average.

III.4 Estimating School Effectiveness Impacts by Educational Advantage

To quantify the effect of attending a school with one standard deviation higher predicted overall effectiveness, we regress each outcome on the standardized school effectiveness index (plus controls). Specifically, where Y_{ijt} is an outcome, and $\hat{\omega}_{jt}$ is the standardized out-of-sample predicted effectiveness, we estimate the following model by OLS.

$$Y_{ijt} = \delta \hat{\omega}_{jt} + \beta_1 Z_{ijt} + \tau_t + \varepsilon_{ijt} \quad (6)$$

All variables are as defined above and τ_t is a year fixed-effect. In addition to the controls in (1), we include school-level averages of all individual attributes. Standard errors are clustered at the school level.¹⁶ To flexibly present differences in the marginal impacts by student type, we estimate (6) separately for each decile d of the estimated educational advantage index, and we plot the decile-specific marginal effects δ_d against the educational decile d (along with confidence intervals for each decile-specific effect). We formally test the hypothesis that the marginal effects are linearly related to educational advantage by pooling the data for all students and adding an interaction between the estimated advantage decile and school effectiveness ($L\hat{A}_i \times \hat{\omega}_{jt}$) and fixed effects for each educational advantage decile (α_d) as in equation (7).

$$Y_{ijt} = \delta_1 \hat{\omega}_{jt} + \delta_2 (L\hat{A}_i \times \hat{\omega}_{jt}) + \beta_1 Z_{ijt} + \tau_t + \alpha_d + \varepsilon_{ijt} \quad (7)$$

Under no linear relationship between the marginal effect of schools and educational advantage, the slope between effectiveness and the marginal effect (i.e., δ_2) is zero.¹⁷ To take the estimated impacts of effectiveness as reflecting schools' causal impacts requires that, on average, there are no unobserved differences in the determinants of outcomes between students that attend more and less effective schools. We provide several empirical tests suggesting causal effects in Section IV.

¹⁵See column 7 of Table 4 in [Barrow et al. \(2020\)](#). [Angrist et al. \(2019\)](#), show that their somewhat negative results on certain outcomes may reflect differences in the counterfactual schools student may attend.

¹⁶Individuals with missing 8th grade surveys or test scores are given imputed values. We regress each survey measure or test score on all observed pre-8th grade covariates. We then obtain predicted 8th grade values based on these regressions, and replace missing values with the predictions. Results are similar with and without imputation.

¹⁷To assuage concerns that the advantage is a generated regressor, as a robustness check, we estimate equation (7) but instrument for the interaction ($L\hat{A}_i \times \hat{\omega}_{jt}$) with 8th-grade covariates (a test-score index, a surveys index and a behaviors index) each interacted with estimated effectiveness. The results are very similar (see Appendix Table A10.)

IV Validating the Method

Test Score Value-Added Versus the School Effectiveness Index

Before exploring heterogeneous effects, we first present the average impacts. Table 3 reports the coefficient on the effectiveness index in a regression of various outcomes on the index and controls for the full sample. The point estimate is the difference in outcomes associated with attending a school with 1σ higher estimated effectiveness (i.e., going from a school at the median to one at the 85th percentile of the effectiveness distribution). As a basis for comparison, we also report the estimated effect of the value-added on the individual dimensions. However, we focus the discussion on the effect on the overall effectiveness index and test score value-added. We refer to schools with a higher estimated overall school effectiveness index as more effective schools.

The top row shows that more effective schools improve 9th-grade test scores, socio-emotional development in 9th grade (as measured by surveys), and behaviors in 9th grade. Specifically, *on average*, a 1σ increase in effectiveness increases test scores by 8.9 percent of a standard deviation, socio-emotional development by 10.2 percent of a standard deviation, and behaviors by 5.7 percent of a standard deviation. Note that social and work hard are very highly correlated, so we combine these two SED measures into a single survey measure.¹⁸ Not surprisingly, more effective schools also improve longer-run outcomes on average. A 1σ increase in effectiveness increases high school graduation by 2.4 percentage points, college going (within 2 years of high school completion) by 2.57 percentage points, and decreases the likelihood of having a school-based arrest by 0.8 percentage points. All of these estimates are significant at the 1 percent level.

Our use of the index (as opposed to using test score impacts only) is motivated by Jackson et al. (2020) showing that a combination of school impacts on test scores and surveys better predict both short and long-run outcomes than test scores alone, and Jackson (2018) showing that a combination of teacher impacts on test scores and behaviors better predict long-run outcomes than test scores alone. We show this to be the case here also. In the third row, we show the estimated impact of a one standard deviation increase in test-score value-added on these same outcomes. Test score value-added *does* predict impacts on both short- and long-run outcomes, but these impacts are smaller than those based on the effectiveness index *including on test scores themselves*. For all outcomes, the improvement associated with a 1σ increase in effectiveness is greater than that of a 1σ increase in test-score value-added. For the longer-run outcomes, the marginal impacts of the effectiveness index are between 50 and 250 percent larger than that for test score value-added alone (Table 3).¹⁹ We shed light on the extent to which test scores understate the benefits of attending

¹⁸We provide analogous entries of Table 3 in Appendix Table A6 where the two survey measures are separated.

¹⁹Table 3 also presents the estimated impacts of value-added on the surveys and behaviors. Remarkably, for surveys, test scores, high-school graduation, and college enrollment, the school effectiveness index is more predictive of impacts

effective schools varies by educational advantage in V.3 after presenting patterns for the overall effectiveness index.

IV.1 Testing For Selection

Because students are not randomly assigned to schools, one may worry that students could have better outcomes at more effective schools not because the schools *causally* improve outcomes but because students who are predisposed to success (in unaccounted-for ways) attend these schools. Moreover, if unaccounted-for selection is different across educational advantage groups, it could yield patterns that suggest heterogeneity even when there is none. While there is no way to prove that the effectiveness estimates are unrelated to unobserved determinants of outcomes, we present several tests to show that this is likely satisfied in our setting.

Control Function Approach

Altonji and Mansfield (2018) show that in settings where individuals or families sort into treatments, group-level averages of observed individual characteristics tend to be correlated with averages of unobserved student characteristics. As such, group-level averages of observed individual characteristics can potentially serve as a control function and remove all the across-group variation in both observable *and unobservable* individual characteristics. Moreover, Agrawal et al. (2019) show that this condition holds “*when observed and unobserved interactions are introduced to the production function.*” We show that this is likely satisfied in our setting by controlling for (1) school-level average math scores, (2) school-level average absences, and (3) the share of white students at the school. To assess the plausibility of this, we show that conditional on these three school-level averages, predicted outcomes based on a rich set of individual characteristics are unrelated to our effectiveness index. While our test is similar to that in Altonji et al. (2005) and Oster (2019), our control function method is justified by theory and includes a very small number of variables relative to the full set of individual predictors of student outcomes.

We predict each outcome based on a linear regression of that outcome on rich individual attributes available in our data (8th-grade test scores, surveys, and behaviors and Grade Point Average (GPA), gender, race, free-lunch status, special education status, and neighborhood socioeconomic status).²⁰ Appendix Figure A2 shows a binned scatterplot of the predicted outcome against the

than any of the individual measures. This indicates that the effectiveness index is a good summary measure of school “effectiveness” for these outcomes. However, the behaviors value-added appear to have more predictive power for behaviors and school-based arrests than overall effectiveness. This indicates that school impacts on behaviors capture some meaningful dimension of school quality that is not fully captured by the index and which is predictive of behaviors and school-based arrests.

²⁰Where $(\hat{Y}_{ijt}|Z)$ is the predicted outcome given all the observed covariates, we estimate the following model by Ordinary Least Squares (OLS).

$$(\hat{Y}_{ijt}|Z) = \delta_p \hat{\omega}_{jt} + \tau_t + v_{ijt} \quad (8)$$

actual outcomes used in this paper. The predicted outcomes track actual outcomes very well.²¹ We then examine if the effectiveness index is correlated with predicted outcomes (i.e., a weighted average of *all* observable student-level characteristics that best predicts the student-level outcomes) in a regression model with only the three aforementioned school-level controls.²² In all models (see columns 4, 8, and 12 in Table 4), school effectiveness is not significantly related to predicted outcomes (conditional on average 8th grade scores, average 8th absences, and percent white at the school) and the point estimates are small. While this shows no selection on *observables* conditional on these key school-level controls, it validates the use of the control function approach suggested in [Altonji and Mansfield \(2018\)](#), to remove selection on *unobservables*.

Further Evidence of No Selection on Unobservables

Here we present further evidence of no selection on unobservables. While most schools have residential attendance zones, in Chicago, almost two-thirds of children attend schools other than their zoned school ([Hing and Jenniver, 2019](#)). As such, selection on unobservables can occur due to (1) selection of entire families to neighborhoods and therefore zoned schools, and (2) selection of individual students (even within families) to particular schools outside their residential zoned school. We show that neither form of selection seems to be operative in our setting, which, taken together, suggests that our estimates are largely unbiased.

Using Variation Across Attendance Boundaries. If our results were driven by individual students selecting to schools outside their residentially-zoned schools, then the effectiveness of the residentially-zoned school would be unrelated to student outcomes. We assess this by constructing instruments that remove the sorting bias that may exist when individuals chose to attend a school outside their zoned area. We instrument for the effectiveness of the school attended with the effectiveness of the residentially-assigned school. Because families *almost* always have the same address in our database, this approach is largely based on comparisons across families. The first stage regression is strong – yielding first stage F-statistics above 300. The two-stage-least-squares (2SLS) regressions are reported columns 2,6 and 10 of Table 4. The OLS estimates are reported as a basis for comparison in columns 1, 5, and 9. For the short run outcomes (test scores, surveys, and behaviors), the 2SLS estimates are all positive, significant at the 1 percent level, and of the same order of magnitude as the OLS estimates. Similarly, for all the main long run outcomes (arrests, high-school graduation, college-going, and 4-year college going), the point estimates are positive, and of the same order of magnitudes as the OLS estimates. For 3 out of 4 outcomes, the 2SLS

²¹The R-squared is above 0.2 for surveys, behaviors, and test scores, high-school graduation, any college enrollment, and 4-year college enrollment. Those for 2-year college going and arrests are 0.04 and 0.082, respectively.

²²Recall that we use out-of-sample estimates and equation (3) controls for individual characteristics *within* schools. As such, the fact that we control for some of these observables *within* schools when estimating value-added out-of-sample, **does not** imply no correlation between value-added *across* schools and observables *in-sample*.

estimate is significant at the 5 percent level. These patterns rule out that our main results are driven by selection of individual students to schools outside their residential attendance zone. Because these 2SLS models rely primarily on comparisons across families in different attendance zones, these estimates will only be biased if those families that attend the zoned schools self-select into neighborhoods along unobserved dimensions that are correlated with school effectiveness. To rule out this possibility of bias, we also examine variation *within families*.

Using Variation Within Families: If the 2SLS results were driven by selection of families to school assignment zones, then there would be no difference in outcomes among members of the same family (in the same attendance zone) but who attend different schools. We find no evidence of this. To isolate within-family variation, we use a subset of the data in which we can identify siblings. We can identify 19,420 families after 2015 in which more than one sibling is observed in 9th grade.²³ There are 7786 families with multiple 9th graders with observed test score measures and 19,420 with observed behaviour measures. Among these families, roughly half have some variation in school attended. While the effective samples are much smaller than the OLS samples for the within-family models (about 14 percent for behaviors and only 9 percent for test scores), they are large enough for reliable within-family estimates of school impacts on the short-run outcomes. We remove the correlation with potentially confounding fixed family characteristics (i.e., the selection of families to neighborhoods) by comparing students from the same family who attended different schools. This is achieved by adding a family fixed effect to our main model in equation (6). The within-family estimates are presented in columns 3, 7 and 11 of Table 4. For the short-run outcomes (test scores, surveys, and behaviors), the within-family estimates are all positive, significant at the 1 percent level, and of the same order of magnitude as the OLS estimates – effectively ruling out that the 2SLS results were driven by selection of families to residential zones.

For longer-run outcomes measured at the end of high school, there is considerably less data – limiting our ability to reliably estimate within-family models. There are 3,902 families with multiple children old enough to have graduated high school. This yields effective samples much smaller than the OLS samples for the within-family models on these outcomes (about 5 percent for high school graduation and arrests), which may not be reliable. Consistent with this, the standard errors for the within-family estimate on these two outcomes are much larger than those for OLS. However, the within-family estimates are both similar to the OLS estimates and significant at the 5 percent level. Unfortunately, the effective sample for college outcomes is less than half that for high school outcomes, so due to lack of sufficient data, we cannot reliably estimate the within-family models for the college outcomes.²⁴ However, the consistent within-family patterns for the short

²³Because we cannot identify *all* siblings prior to 2015, these data are imperfect and incomplete. However, if we are able to find similar effects in this small sub-sample as in the broader sample, it would be compelling evidence that our estimates are not biased by family selection to neighborhoods.

²⁴There are only 1,634 families with multiple children old enough to have enrolled in college. Among these, very

and medium-term outcomes suggest that this would likely also hold true for the college outcomes.

Considering all the Tests Together

If our estimates were biased by selection, one would expect that strong predictors of outcomes would be related to our estimated value-added (conditional on a small number of controls)– but this is not the case. If our results were driven by selection to school assignment zones *across* families, it would bias our 2SLS results but not our sibling results. If our results were driven by selection to schools *within* families (and assignment zones), it would bias our sibling results but not our 2SLS results. The similarity between the 2SLS and within-family results (coupled with the lack of selection on observables conditional on the three school-level variables) suggests that neither is biased. While none of these tests is dispositive in isolation, in light of validation work showing that value-added estimates tend to be largely unbiased (e.g., Angrist et al. (2020); Deming et al. (2014)), together they show that our estimated school impacts likely reflect causal impacts.

V Results

V.1 Impacts of School Effectiveness on Short-Run Measures: By Advantage

We now consider how effects vary by students’ *ex-ante* educational advantage. We estimate these same regressions for each measure and for each decile. To summarize these ten regressions per outcome, we plot the point estimate and 95% confidence intervals for each estimate in Figure 3. Appendix Table A8 reports the point estimates for individual deciles, Appendix Table A9 reports the point estimates for the top and bottom three deciles, and Appendix Table A10 reports the point estimate on the interaction term of interest from Equation 7. We also plot the linear relationship between the estimated effect and the educational advantage decile along with the *p*-value associated with the null hypothesis of no linear relationship between educational advantage and the benefit of attending a more effective school. See Appendix D for more details.

The top left panel of Figure 3 shows the effect of attending a school one standard deviation higher in school effectiveness on students’ 9th grade test scores. All students benefit from attending a more effective school, with larger effects for the less advantaged. The *p*-value associated with the hypothesis that the effect varies linearly with educational advantage is 0.07 – weakly suggestive of larger marginal benefits for the least educationally advantaged. We now turn to the survey measures. The middle panel shows the effect of attending a school one standard deviation higher on the school effectiveness index on students’ socio-emotional measures in 9th grade by decile of

few families have variation in school attended and college status – so the within-family models are inconclusive on this sample. To show the implications of the lack of variation in this sub-sample, in Appendix Table A7 we show that using the sample of families for which there are multiple children old enough to have enrolled in college, none of the within-family relationships shown for test scores, surveys, and behaviors persist – indicating that there is not enough variation among this sub-sample for reliable inference for *any* outcome.

predicted educational attainment. As with test scores, all students benefit from attending a more effective school. However, for surveys, one cannot reject that the impacts of attending a more effective school on 9th grade survey measures are the same throughout the educational advantage distribution. We report the results for 9th grade behaviors in the right top panel of Figure 3. Unlike the socio-emotional and test score measures, one strongly rejects the hypothesis of the same marginal effect for all advantage groups at the 1 percent level. Effective schools have the strongest effect on the observed behaviors for students in the lower end of the advantage distribution. For a student in the lowest (first) decile, attending a school 1 standard deviation higher in school effectiveness improves the behavior index by 0.185 standard deviations. Meanwhile, for students in the top (tenth) decile, the behavior index only improves by 0.012 standard deviations. One interpretation of this pattern is that schools have heterogeneous effects on students across the distribution. However, the small impacts for students at the top of the distribution may be driven by a lack of variation among these students. Specifically, students in the top decile are very unlikely to be involved in a disciplinary incident (0.007 compared to 0.29 in the bottom decile) and have a low absence rate (5.6 days compared to 34 days in the bottom decile), so there is relatively little room for improvement. We show that this likely explains the seemingly heterogeneous impacts on the behaviors index in Section V.1.1.

V.1.1 Testing for Mechanical Heterogeneity

The basic idea of this test is that for most models of binary outcomes (such as a logit, or probit), with the same change in underlying skills the marginal effects will be largest for groups that are marginal (with probability of success close to 0.5), and smallest for groups with probabilities of success farthest from 0.5 (i.e., zero or one).²⁵ Under purely mechanical heterogeneity (i.e., the same underlying change in latent disposition across groups), there will be (a) a negative relationship between the absolute value of the marginal effect and the absolute difference between the group success rate and 0.5, and (b) a predicted marginal effect of zero for groups with an absolute difference of from 0.5 of 0.5. We assume a linear relationship, for simplicity, and test for this using the regression in (9) below.

$$|\delta_g| = \alpha + \pi \times (|p_g - 0.5|) + v_g \quad (9)$$

Where $|\delta_g|$ is the absolute value of the marginal effect for decile group g , and p_g is the average success rate for decile group g , then π represents the relationship between the marginal effect and the distance between the baseline success rate for a group and 0.5.

Conducting this test for some of the binary outcomes underlying the behaviors provides strong evidence of such “mechanical” effects (see Appendix Figure B1). For the likelihood of being chronically absent, having any suspensions, or having any disciplinary infections, one rejects that

²⁵We expand upon the logic of this test formally in Appendix B.

the slope is zero at the 1 percent level and one cannot reject that the effect is zero for those who are very likely or unlikely to have a success. To explore whether the heterogeneity is *all* mechanical or reflects some real heterogeneity, we simulate the distribution of linear slopes (as in Figure 3) under mechanical heterogeneity alone (See Appendix Figure B2).²⁶ The actual slopes observed are within the simulated range expected under pure mechanical heterogeneity for most of the binary outcomes comprising the behaviors index. These additional tests are consistent with the similar effect on the continuous measures (i.e., test scores and surveys) across the educational advantage distribution – suggesting little differential effect on underlying skills in the short run.

V.2 Impacts on Longer-Run Outcomes: By Advantage

We now examine similar figures for the longer-run outcomes (the middle and lower panels of Figure 3). Looking at high school graduation, the marginal impacts of school effectiveness are much larger for students at the bottom of the educational advantage index than those at the top. One rejects that the linear relationship between the marginal effect and educational advantage is zero at the 1 percent level. Indeed, for those in the bottom decile, a 1σ increase in effectiveness increases high school completion by 3.8 percentage points (p -value <0.01) compared to only 0.5 percentage points (p -value >0.10) in the top decile. Relative to each groups' baseline level, this is about a 9 percent increase for those at the bottom of the distribution compared to a .6 percent increase for the top. While the differences in the changes in graduation rates across groups are real and economically meaningful, one may wonder if this pattern is due to more students at the bottom being on the margin of high school graduation. We assess this using the test detailed in Section V.1.1. The results of this test are summarized in Appendix Figure B1, which plots the linear relationships between the absolute value of the marginal effects against the absolute deviation of the group's baseline success rate from 0.5. For high school graduation, the p -value on the slope is significant at the 1 percent level and the implied effect is near zero for groups that are very likely or unlikely to have a success. Also, in Appendix Figure B2, the actual slope is within the distribution of simulated slopes under mechanical heterogeneity – suggesting that the observed heterogeneity in high school graduation is largely mechanical.

Next, we examine enrolling in any college (2-year or 4-year) within two years of expected high-school completion. There are benefits for all groups. However, there are larger increases at the bottom of the distribution, and these differences are statistically significant (the p -value on the slope relating the marginal effect and educational advantage is 0.0007). More specifically, for the bottom decile, a 1σ increase in effectiveness increases college-going by 3.55 percentage points

²⁶That is, we implement a probit model using the actual data, estimate each individual's latent disposition, add a constant marginal treatment effect to each latent disposition, and then feed this back into the probit function to simulate "fake" binary treatments driven by mechanical heterogeneity. We do this 1000 times to simulate the distribution of slopes one would observe due to mechanical heterogeneity alone.

(p -value <0.01) compared to about 1.5 percentage points (p -value >0.1) in the top decile. Given the large differences in base rates, the differences in relative marginal impacts are sizable. For the bottom decile, a 1σ increase in effectiveness increases college-going by about 19 percent compared to under 2 percent in the top decile. Unlike for high school graduation, the heterogeneous effect on college-going is unlikely to be driven by differences in baseline success rate. The first evidence of this is that students in the middle third have college-going rates around 50 percent, so if all of the differences are due to differences in the proportion of marginal students, one might expect the largest college-going impacts for this group. The results are inconsistent with this idea because the largest increases are among the bottom third. The tests in Appendix B support these observations. That is, the linear relationships between the absolute value of the marginal effects and the absolute deviation of the group's baseline success rate from 0.5 cannot be distinguished from zero (the p -value is 0.58), and the estimated slope is well outside the range of simulated slopes under mechanical heterogeneity. This suggests that the observed heterogeneity reflects heterogeneous effects on the underlying disposition to attend college.

Looking at college type reveals some interesting patterns. The average results in Table 3 suggested no impacts on 2-year college going. However, the heterogeneous impacts provide an explanation for the null result on average. Among the bottom three deciles of the educational advantage distribution (who are least likely to attend *any* college), attending a more effective school increases 2-year college going by 1.1 percentage points, but among the top three deciles (who are more likely to attend college), attending a more effective school *reduces* 2-year college going by 1.4pp. The increase in college-going overall indicates that this reduction in 2-year college going for the advantaged students reflects shifting from 2-year to 4-year programs. One can see this clearly when looking at 4-year college going. While all groups have increased 4-year college going, the groups with the largest increases in 4-year college going are those with the reductions in 2-year college going. Specifically, for the bottom three deciles, a 1σ increase in effectiveness increases 4-year college-going by 3.6pp compared to over 5pp for the middle three deciles and 3.3pp for the top three deciles. In sum, the increase in college-going overall is due to increased 2-year and 4-year college going among the bottom third of the educational advantage distribution, and an increase in 4-year college-going among the top two-thirds of the educational advantage distribution driven by both (a) increased college going among those who would not have attended college and (b) shifting from a 2-year college to a 4-year college.

Another economically important result is that the increase in 4-year college going is very similar for those in the top and bottom three deciles even though the base rates are very different (about 15 versus 66 percent). This indicates that the increases in college going, and those for 4-year institutions are not limited only to populations with students on the margin. Indeed, the formal tests fail to reject the null hypothesis of no “mechanical heterogeneity” at the 5 percent significance level, and

one rejects that the observed heterogeneity is all mechanical – indicating larger effects on underlying dispositions to attend 2-year and 4-year colleges among the least educationally advantaged. Remarkably, attending an effective high school can lead to sizable increases in college going, even among student populations for which that may seem unlikely. That is, among the bottom three deciles, a 1σ increase in effectiveness increases 4-year college-going by about 25 percent.

If the marginal college enrollees from the less advantaged groups are less well-prepared for college, they may be less likely to persist and no more likely to earn a college degree (Jackson, 2014). Since roughly half of college attrition occurs in the first year (NSCRS, 2012), persistence through the first year is a key predictor of college success. As such, we also examine impacts on college persistence beyond freshman year. There are positive impacts throughout the educational advantage distribution, and one cannot reject equality of impacts through the distribution (p -value=0.19). Also, there is no linear association between base rates and the marginal effects – indicating sizable benefits to attending more effective school, particularly for less educationally-advantaged students. From a policy perspective, the similar effects on college persistence across the distribution are important because they imply that the marginal college goers are equally likely to persist irrespective of educational advantage (suggestive of real long-term gains among all students).

Finally, we examine whether a student had ever had a school-based arrest. Because this is a relatively rare outcome among students at the top of the educational advantage distribution, one would not expect much effect at the top of the distribution. Indeed, this is what one observes. Among students in the bottom decile, a 1σ increase in effectiveness decreases in-school arrests by 2.1 percentage points (p -value<0.01) compared to only 0.22pp in the top decile (p -value<0.05). Even though there are statistically significant effects even among those at the top of the educational advantage distribution, the marginal effects are much more pronounced for those at the bottom. Given the long-term implications of these school-based arrests, this implies sizable long-term benefits to attending effective schools, particularly for those who are least likely to complete high school. Note that this likely represents a *lower* bound on the effect on arrests because students who may have dropped out of school will not receive a school-based arrest. While these sizable benefits for the least-advantaged students are economically important, one may wonder whether the heterogeneity observed reflects differences in the likelihood of being marginal for an arrest. Our formal tests of this suggest that this is the case (see Appendix B). For arrests, the underlying base rate strongly predicts the marginal effect ($p < 0.01$), the predicted effect is close to zero for very high and very low base rate groups, and the estimated slope is well within the range of simulated slopes. The data are consistent with effective schools having similar effects on disposition toward crime, which manifest most strongly among those who are marginal.

Despite theories suggesting that the gains should be largest for the least advantaged and others suggesting the opposite, the results show that all students benefit from attending more effective

schools. The results are inconsistent with the notion that the least advantaged are unable to benefit from attending better schools (or benefit less). On the contrary, the evidence is more supportive of the notion that they benefit more – consistent with findings from oversubscribed charter schools in Boston (e.g., Angrist et al. (2012), Walters (2018), and Cohodes et al. (2020)). To assuage concerns that these patterns reflect comparisons across school type or broad demographic characteristics, Appendix C shows that these patterns hold within gender categories, within ethnic/racial groups, and only among traditional public schools.

V.3 Heterogeneity in the Importance of Dimensions Missed by Test Scores

Our effectiveness measure reflects a combination of school impacts on test scores, surveys, and behaviors. One might wonder if measures of school quality that exclude impacts on SED miss important components of school quality particularly for disadvantaged populations. To shed light on this, in Figure 4 we plot the *difference* between the marginal effect for the overall effectiveness and that for test score value-added (along with the confidence intervals for the difference in marginal effects, and the linear relationship between the *difference* and educational advantage).²⁷ Note that this *is not* a plot of the marginal effect of non-test-score dimensions of skills, but a plot of the marginal effect of the dimensions of school quality that are unrelated to test score effects. This analysis explores whether test score value-added misses key dimensions of school quality in ways that vary by educational advantage.²⁸

Looking at Short Run 9th Grade Skill Measures

Looking at the short-run outcomes, the non-test-score dimensions of school quality have more explanatory power (above and beyond test score value-added) for the less educationally advantaged. We first examine 9th grade test scores in the top left panel of Figure 4. The plot of overall school effectiveness unexplained by test score value-added against educational advantage has a negative slope (p -value <0.10). That is, the test scores of the less advantaged are *relatively* more responsive to improvements in the dimensions of school quality unrelated to test score value-added than those of more advantaged students. Indeed, for the most advantaged, the difference between the marginal effect of test scores value-added and the overall index are near zero and statistically indistinguishable from each other, while the marginal effect of the overall index is clearly larger than that of test score value-added for the least advantaged. Looking to the effect on 9th grade survey-based SED measures, the difference between the marginal effect of the overall index and test-score value-added is relatively similar (and positive) through the distribution of advantage. While the slope between the differences in marginal impacts and academic advantage is negative,

²⁷Practically, we do this by stacking the two regressions into a single model.

²⁸We test this formally using a stacked linear regression that assesses whether the difference between the marginal effect for the overall effectiveness and that for test-score value-added decreases for higher deciles. This is tested using an interaction between the educational advantage and the difference between the two school quality measures.

it is not statistically significantly different from zero (p -value=0.79). For behaviors, the school effects unexplained by test-score value-added is appreciably larger for the less-advantaged. The slope between the impact unexplained by test score value-added and educational advantage is negative (p -value<0.01)- suggesting that the behaviors of the less advantaged are *relatively* more responsive to improvements in the non-test-score dimensions of school quality than those of more advantaged students. We caution that because there is less variation in behaviors for the most advantaged, this particular result for behaviors may be somewhat mechanical.

In sum, for all three 9th grade outcomes, the impacts of school effectiveness unexplained by test score value-added is larger for the less advantaged, and the slope is at least marginally statistically significant for two of them (including test scores for which there is similar variation among more and less advantaged groups). Taken together, this suggests that the extent to which test score value-added misses key components of school quality *tends to be* greatest for less-advantaged students. Put differently, in terms of the short-run outcomes, the least-advantaged students appear to benefit the most from those components of school quality unmeasured by test score value-added. We now show that this pattern also holds for the longer-run outcomes.

Looking at Longer-Run Outcomes

The results for high school graduation are in the middle left panel. The difference between the overall effectiveness effect and that for test-score value-added is largest among the least advantaged (about 1pp) and near zero for the most advantaged. Indeed, the slope between the marginal unexplained effects and educational advantage is negative (p -value<0.001). As with the short-run measures, while test score value-added is a reasonable predictor of the benefits of attending a more effective school (that is, effective in multiple dimensions) on high-school graduation for educationally advantaged students, test score value-added may be a particularly poor predictor of overall school effects on high school completion for less-advantaged students. To see the importance of this, consider the common approach of using average test score value-added to predict effects. Note that while using *any* average measure will lead to inaccurate benefits for the most and least advantaged students, this is particularly so for average test score value-added that is commonly used. Ignoring heterogeneity and focusing on test score value-added, the average effect on high-school completion of attending a school with 1σ higher test value-added is 1.06pp (see Table 3). Indeed, for students in the top three deciles of the educational advantage distribution, the effect of attending a more effective school overall (using effects on multiple dimensions) is similar (0.78pp), but for students in the bottom three deciles, the overall school effectiveness effect is 3.76pp – more than three times larger than the effect implied by average test score value-added alone. These differences are economically meaningful and would affect any cost-benefit calculations regarding the benefits of improving schools, or any calculation of the distributional effects of improving schools.

The pattern is similar for college-going (middle panel). The marginal school-effectiveness effect unexplained by test score value-added is larger for less-advantaged students. For the bottom three deciles, the marginal effect of the index is about 1pp larger than that for test score value-added, while this difference is near zero for the top two deciles. Consistent with this, one rejects that the unexplained gap is the same for all groups at the 0.01 significance level. Importantly, the potentially heterogeneous effects on college going cannot be due to differences in the likelihood of being marginal across groups (since we find little evidence of this for the overall index). As such, this pattern is likely driven by the fact that the college-going outcomes of the least-advantaged students are relatively more responsive to school impacts on non-test-score dimensions of skills. As with high-school graduation, the commonly used average test score value-added is a much worse predictor of school effects on college-going for less-advantaged students than the more-advantaged. Specifically, the average effect on college going of attending a school with 1σ higher test value-added is 1.68 percentage points. For students in the top three deciles of the educational advantage distribution, the effect of the overall index is similar (1.48pp), but for students in the bottom three deciles of the educational advantage distribution, it is about 3.7pp – more than twice as large as that implied by average test score value-added alone.

We now turn to school-based arrests (lower right panel). A plot of the difference between the estimated effects based on overall school effectiveness and test-score value-added by educational advantage shows a clear negative relationship ($p - value < 0.01$). We caution that this pattern *may* be an artifact of different groups being differentially marginal for arrests. Irrespective of the reasons, given that having an arrest has real economically meaningful implications, the documented heterogeneity has real-world and policy implications. To see the importance of accounting for heterogeneity and also school effects on non-test-score dimensions, consider the following calculations. The average effect on the likelihood of arrest of attending a school with 1σ higher test value-added is 0.48pp. For students in the top three deciles of the educational advantage distribution, the effect of the overall index is similar (0.26 percentage points), but for students in the bottom third of the educational advantage distribution, the effect of attending a 1σ more effective school is about 1.49pp – about three times as large. While any average measure of school quality will be somewhat inaccurate for the least advantaged, the extent to which average test-score value-added (the most commonly used metric) understates the benefits (as measured by arrests) to attending a more effective school for less-advantaged populations is considerable.

In sum, we document a consistent pattern across several outcomes where test value-added misses important dimensions of school quality, particularly for the least advantaged. We are careful to note that this pattern may be due to (a) general differences in marginal effects across groups due to different groups being differentially marginal for particular outcomes, and/or (b) less educationally advantaged populations being particularly sensitive to improvements in the non-test-score

dimension of school quality. We cannot *rule out* the first explanation. However, the fact that we find similar patterns across all outcomes (including those where we find no evidence that differences in marginal effect are related to differences in being marginal such as test scores and college going) suggests that the second explanation is partly operative. As such, the patterns are broadly consistent with work in psychology suggesting that less advantaged students may enjoy particularly large benefits from interventions that promote socio-emotional development (e.g., [Walton and Cohen 2007](#); [Walton and Cohen 2011](#); [Gray et al. 2018](#); [Sisk et al. 2018](#)).

V.4 Policy Implications

V.4.1 Distribution of Effectiveness by Advantage

Our results indicate that the least advantaged students may benefit the most from attending more effective schools. An important policy implication is that investments in high school effectiveness should enhance outcomes for all students while closing gaps in outcomes between more and less educationally advantaged students. However, to assess the potential for a more targeted approach, we explore whether school effectiveness is evenly distributed by educational advantage. At the lowest-performing schools (bottom 10%) about 15 percent come from the lowest decile of advantage while about 3 percent come from the top decile of advantage. Similarly, at the highest performing schools (top 10%) about 3.7 percent come from the lowest decile of advantage, while as much as 23 percent come from the top decile of advantage. Consistent with this, on average, those in the top decile attend schools with 1σ higher effectiveness than those in the bottom decile. However, there is considerable overlap in the quality of schools attended by students across educational advantage groups (see Appendix D). This indicates that the potential gains to a more equitable distribution of students across schools are economically significant. The patterns also indicate that improving the outcomes of the least educationally advantaged will require targeting schools that disproportionately enroll these students. We consider such a targeted policy below.

V.4.2 School Closure Policy

One policy that has been employed is the shutting down of ineffective schools ([Brummet 2014](#)) and (importantly) the relocation of affected students to high-performing schools. We run a simple simulation where we remove the lowest performing ten schools (8 percent) and reallocate those students to one of the ten closest non-low-performing high schools. The static change in value-added is merely the change in the estimated value-added of the attended school (which will mechanically increase for those displaced and be unchanged for other students). As a rough estimate of any *indirect* effects associated with the reshuffling of peers, we (a) estimate the cross-sectional relationship between peer advantage at a school and school value-added and then (b) multiply this by the change in peer advantage at the school induced by this policy. See appendix D for additional details.

For the 8 percent of students who would have attended the least effective schools, the increase in total value-added is 1.78σ , while the decrease for other students is roughly 0.014σ (attributed to the inflow of less-advantaged peers). This implies sizable benefits for those displaced students. Given the distribution of advantage across schools, for the relocated students, high-school graduation would increase about 4.6pp (7.5pp for the least advantaged movers), college-going would increase by 4.7pp (7pp for the least-advantaged), and school-based arrests would fall by roughly 1.83pp (4.2pp for the least advantaged). For the receiving students, because the spillover effects are small, the changes are very small. We now consider the change in effectiveness for decile groups irrespective of mover and receiver status. The simulated change in effectiveness is very small and positive for the top decile (0.0044σ), around 0.1σ for deciles in the middle, and 0.23σ for the least-advantaged decile. That is, even for the most advantaged decile, which has the largest negative peer effect, *on average*, the total effect is non-negative – indicating efficiency gains. These gains are even more pronounced because the marginal effects are largest among the most disadvantaged (who experience the largest increases in effectiveness), on average. Because the policy only moves a small fraction of students, the overall effects are modest *overall*. For the bottom decile, this policy could increase graduation rates by 0.8pp, increase college going by 0.82pp, and reduce arrests by 0.5pp, with smaller effects for the more advantaged groups.

We caution readers that our back-of-the-envelope calculations do not capture any disruption effects or ill-effects associated with possibly having to cross gang lines to go to a different school ([link](#)). However, they do serve as illustrations of plausible benefits associated with selectively closing ineffective schools that are disproportionately attended by the least-advantaged students.

V.4.3 Without Surveys

Our findings are more broadly applicable if one constructs school effectiveness measures based on readily available data (such as test scores and behaviors) as opposed to surveys (which must be administered). To examine how different our results would be if one excluded school effects on SED (i.e., using test scores and behaviors only), we recreate Figure 3 but exclude the school impacts on survey measures (Appendix Figure D2). While the pattern of results is very similar, the magnitude (and precision) of the heterogeneity is less than when using all available information (i.e., including the survey measures). For example, while the marginal effects on high-school graduation and college-going at and above the median advantage are similar with and without SED, for the less advantaged, the marginal effects are almost twice as large with SED than without (about 0.02 pp on both outcomes without SED, compared to 0.04pp with SED). At the same time, the effects on arrests are quite similar with and without SED. While including the surveys adds additional explanatory power, the results indicate that one can identify schools that will meaningfully impact less-advantaged students using readily available administrative data.

VI Conclusions

It is known that schools can have meaningful impacts on both short- and longer-run outcomes. However, whether all students benefit similarly from attending better schools is not well understood. Moreover, the extent to which more or less advantaged students benefit differently from school quality in different dimensions (cognitive versus socio-emotional and behaviors value-added) is unknown. We speak to these issues by examining the effect of attending a more effective school (one that improves a combination of test scores, survey-based SED measures, and behaviors) for more- and less-advantaged students. Importantly, we do this for a representative set of schools and students – so that our results are more generalizable than existing work.

We show that all students benefit from attending effective schools, and that the marginal effects are larger for less-advantaged students. While some of the effect heterogeneity is due to less-advantaged groups being marginal for some outcomes, this “mechanical heterogeneity” *does not* explain larger college-going effects for the least advantaged student (among which college-going rates are very low). We show that dimensions of school quality unexplained by test-score value-added have the largest impacts for the less-advantaged students– which is, in part, due to less-advantaged students benefiting more from non-test-score dimensions of school quality. It is worth noting that our results may be caused by (a) more effective schools focusing inputs on disadvantaged student populations or (b) different students responding differently to the same inputs. While we focus on documenting the existence of this heterogeneity, further work is needed to uncover why– which will, in turn, better inform policy. Our findings reinforce the importance of accounting for soft skills while also accounting for effect heterogeneity. The patterns we uncover suggest that if one were to use test-based measures of school quality alone and ignore effect heterogeneity, one would dramatically understate the benefits to attending better schools for those students who may need access to better schools the most.

References

- Mohit Agrawal, Joseph G. Altonji, and Richard K. Mansfield. Quantifying family, school, and location effects in the presence of complementarities and sorting. *Journal of Labor Economics*, 37:S11–S83, 1 2019. ISSN 0734306X. doi: 10.1086/701012. URL <https://www.journals.uchicago.edu/doi/10.1086/701012>.
- Joseph G. Altonji and Richard K. Mansfield. Estimating group effects using averages of observables to control for sorting on unobservables: School and neighborhood effects. *American Economic Review*, 108: 2902–46, 10 2018. ISSN 0002-8282. doi: 10.1257/AER.20141708.
- Joseph G. Altonji, Todd E. Elder, and Christopher R. Taber. Selection on observed and unobserved variables: Assessing the effectiveness of catholic schools. *Journal of Political Economy*, 113:151–184, 2 2005. ISSN 00223808. doi: 10.1086/426036/0.
- Joshua Angrist, Peter Hull, Parag A. Pathak, and Christopher R. Walters. Simple and credible value-added estimation using centralized school assignment. 12 2020. doi: 10.3386/W28241. URL <https://www.nber.org/papers/w28241>.
- Joshua A Angrist, Peter Hull, Parag Pathak, and Christopher R Walters. Race and the mismeasure of school quality race and the mismeasure of school quality *. *Working Paper*, 2022. URL <https://www.usnews.com/education/k12/>.
- Joshua D. Angrist, Susan M. Dynarski, Thomas J. Kane, Parag A. Pathak, and Christopher R. Walters. Who benefits from kipp? *Journal of Policy Analysis and Management*, 31:837–860, 9 2012. ISSN 02768739. doi: 10.1002/pam.21647. URL <http://doi.wiley.com/10.1002/pam.21647>.
- Joshua D. Angrist, Sarah R. Cohodes, Susan M. Dynarski, Parag A. Pathak, and Christopher R. Walters. Stand and deliver: Effects of boston’s charter high schools on college preparation, entry, and choice. *Journal of Labor Economics*, 34:275–318, 4 2016. ISSN 0734306X. doi: 10.1086/683665. URL <https://www.journals.uchicago.edu/doi/abs/10.1086/683665>.
- Joshua D Angrist, Parag A Pathak, and Román Andrés Zárate. Choice and consequence: Assessing mismatch at chicago exam schools choice and consequence: Assessing mismatch at chicago exam schools. 2019. URL <http://www.nber.org/papers/w26137>.
- Lisa Barrow, Lauren Sartain, and Marisa de la Torre. Increasing access to selective high schools through place-based affirmative action: Unintended consequences. *American Economic Journal: Applied Economics*, 2020. ISSN 1945-7782. doi: 10.1257/APP.20170599.
- Diether Beuermann, C Kirabo Jackson, Laia Navarro-Sola, Francisco Pardo, and Inter-American Development Bank. What is a good school, and can parents tell? evidence on the multidimensionality of school output. *Review of Economic Studies*, 12 2022. ISSN 1556-5068. doi: 10.3386/W25342. URL <https://www.nber.org/papers/w25342>.
- Diether W. Beuermann and C. Kirabo Jackson. The Short and Long-Run Effects of Attending The Schools that Parents Prefer. *Journal of Human Resources*, pages 1019–1053R1, apr 2020. ISSN 0022-166X.

- doi: 10.3368/jhr.57.3.1019-10535r1. URL <http://jhr.uwpress.org/content/early/2020/04/03/jhr.57.3.1019-10535R1><http://jhr.uwpress.org/content/early/2020/04/03/jhr.57.3.1019-10535R1.abstract>.
- Shannon T. Brady, Geoffrey L. Cohen, Shoshana N. Jarvis, and Gregory M. Walton. A brief social-belonging intervention in college improves adult outcomes for black americans. *Science Advances*, 6:eayy3689, 4 2020. ISSN 23752548. doi: 10.1126/sciadv.aay3689. URL <http://advances.sciencemag.org/>.
- Jesse Bruhn. The consequences of sorting for understanding school quality, 2020.
- Quentin Brummet. The effect of school closings on student achievement. *Journal of Public Economics*, 119: 108–124, 11 2014. ISSN 0047-2727. doi: 10.1016/J.JPUBECO.2014.06.010.
- Raj Chetty, John N. Friedman, and Jonah E. Rockoff. Measuring the Impacts of Teachers I: Evaluating Bias in Teacher Value-Added Estimates. *American Economic Review*, 104(9):2593–2632, sep 2014.
- Susana Claro, David Paunesku, and Carol S. Dweck. Growth mindset tempers the effects of poverty on academic achievement. *Proceedings of the National Academy of Sciences of the United States of America*, 113:8664–8668, 8 2016. ISSN 10916490. doi: 10.1073/pnas.1608207113. URL www.pnas.org/cgi/doi/10.1073/pnas.1608207113.
- Sarah R. Cohodes, Elizabeth M. Setren, and Christopher R. Walters. Can successful schools replicate? scaling up boston's charter school sector. *American Economic Journal: Economic Policy*, 2020. ISSN 1945-7731. doi: 10.1257/POL.20190259.
- Flavio Cunha, James J. Heckman, and Susanne M. Schennach. Estimating the technology of cognitive and noncognitive skill formation. *Econometrica*, 78:883–931, 5 2010. ISSN 0012-9682. doi: 10.3982/ecta6551. URL <https://onlinelibrary.wiley.com/doi/full/10.3982/ECTA6551><https://onlinelibrary.wiley.com/doi/abs/10.3982/ECTA6551><https://onlinelibrary.wiley.com/doi/10.3982/ECTA6551>.
- Matthew Davis and Blake Heller. No excuses charter schools and college enrollment: New evidence from a high school network in chicago. *Education Finance and Policy*, 14:414–440, 7 2019. ISSN 1557-3060. doi: 10.1162/EDFP_A.00244.
- David J. Deming. Better Schools, Less Crime? *. *The Quarterly Journal of Economics*, 126(4):2063–2115, nov 2011.
- David J. Deming, Justine S. Hastings, Thomas J. Kane, and Douglas O. Staiger. School Choice, School Quality, and Postsecondary Attainment, 2014.
- Andrew Dustan, Alain de Janvry, and Elisabeth Sadoulet. Flourish or fail: The risky reward of elite high school admission in mexico city. *Journal of Human Resources*, 52:756–799, 6 2017. ISSN 15488004. doi: 10.3368/jhr.52.3.0215-6974R1. URL <http://jhr.uwpress.org/content/52/3/756><http://jhr.uwpress.org/content/52/3/756.abstract>.
- De Leon L. Gray, Elan C. Hope, and Jamaal S. Matthews. Black and belonging at school: A case for interpersonal, instructional, and institutional opportunity structures. *Educational Psychologist*, 53:97–113,

- 4 2018. ISSN 00461520. doi: 10.1080/00461520.2017.1421466. URL <https://www.tandfonline.com/doi/abs/10.1080/00461520.2017.1421466>.
- James Heckman, Rodrigo Pinto, and Peter Savelyev. Understanding the mechanisms through which an influential early childhood program boosted adult outcomes. *American Economic Review*, 103:2052–2086, 10 2013. ISSN 00028282. doi: 10.1257/aer.103.6.2052.
- James J. Heckman and Tim Kautz. Hard evidence on soft skills. *Labour Economics*, 19:451–464, 8 2012. ISSN 0927-5371. doi: 10.1016/J.LABECO.2012.05.014.
- Geoff Hing and Smith Richards Jenniver. Chicago school choice in charts - Chicago Tribune, 2019. URL <https://www.chicagotribune.com/ct-chicago-school-neighborhood-enrollment-charts-20160106-htmlstory.html>.
- Caroline M Hoxby and Sonali Murarka. Charter schools in new york city: Who enrolls and how they affect their students' achievement. 4 2009. doi: 10.3386/W14852. URL <https://www.nber.org/papers/w14852>.
- C. Kirabo Jackson. Do Students Benefit from Attending Better Schools? Evidence from Rule-based Student Assignments in Trinidad and Tobago*. *The Economic Journal*, 120(549):1399–1429, dec 2010. ISSN 00130133. doi: 10.1111/j.1468-0297.2010.02371.x.
- C. Kirabo Jackson. Match quality, worker productivity, and worker mobility: Direct evidence from teachers. *Review of Economics and Statistics*, 95:1096–1116, 10 2013.
- C. Kirabo Jackson. Do college-preparatory programs improve long-term outcomes? *Economic Inquiry*, 52: 72–99, 1 2014. ISSN 1465-7295. doi: 10.1111/ECIN.12040. URL <https://onlinelibrary.wiley.com/doi/full/10.1111/ecin.12040><https://onlinelibrary.wiley.com/doi/abs/10.1111/ecin.12040><https://onlinelibrary.wiley.com/doi/10.1111/ecin.12040>.
- C. Kirabo Jackson. What do test scores miss? the importance of teacher effects on non-test score outcomes. *Journal of Political Economy*, 5 2018. doi: 10.3386/w22226.
- C. Kirabo Jackson, Shanette C. Porter, John Q. Easton, Alyssa Blanchard, and Sebastián Kiguel. School effects on socio-emotional development, school-based arrests, and educational attainment. *American Economic Review: Insights*, 2020. doi: 10.1257/AERI.20200029.
- Brian A. Jacob. Where the boys aren't: Non-cognitive skills, returns to school and the gender gap in higher education. *Economics of Education Review*, 21:589–598, 12 2002. ISSN 02727757. doi: 10.1016/S0272-7757(01)00051-6.
- I T Jolliffe. *Principal Component Analysis, Second Edition*. Springer, 2002.
- Mariame Kaba and Frank Edwards. Policing Chicago Public Schools:A Gateway to the School-to-Prison Pipeline. Technical report, 2012. URL <http://cpdincps.com/fullreport>.
- Thomas Kane and Douglas Staiger. Estimating teacher impacts on student achievement: An experimental evaluation, 12 2008.

- Naveen Kumar. Public school quality and student outcomes: Evidence from model schools in india, 2020. URL <https://drive.google.com/file/d/1-Va6T6WLZzIi1H0UpTBxR0WUo04takXz/view>.
- Airan Liu. Non-Cognitive skills and the growing achievement Gap. *Research in Social Stratification and Mobility*, page 100546, sep 2020. ISSN 02765624. doi: 10.1016/j.rssm.2020.100546. URL <https://linkinghub.elsevier.com/retrieve/pii/S0276562420300822>.
- Jing Liu and Susanna Loeb. Engaging teachers: Measuring the impact of teachers on student attendance in secondary school. *Journal of Human Resources*, pages 1216–8430R3, 7 2019. ISSN 0022-166X. doi: 10.3368/jhr.56.2.1216-8430r3. URL <http://jhr.uwpress.org/content/early/2019/07/02/jhr.56.2.1216-8430R3><http://jhr.uwpress.org/content/early/2019/07/02/jhr.56.2.1216-8430R3.abstract>.
- Mary C. Murphy, Maithreyi Gopalan, Evelyn R. Carter, Katherine T.U. Emerson, Bette L. Bottoms, and Gregory M. Walton. A customized belonging intervention improves retention of socially disadvantaged students at a broad-access university. *Science Advances*, 6:eaba4677, 7 2020. ISSN 23752548. doi: 10.1126/sciadv.aba4677. URL <http://advances.sciencemag.org/>.
- NSCRS. Fall 2012 yearly success and progress rates - national student clearinghouse research center, 2012. URL <https://nscresearchcenter.org/snapshot-report-yearly-success-and-progress-rates-2019/>.
- Hessel Oosterbeek, Nienke Ruijs, and Inge de Wolf. Using admission lotteries to estimate heterogeneous effects of elite schools, 2020. URL <https://www.dropbox.com/s/1p4t9cxhosctyrs/eliteschools2020march.pdf?dl=0>.
- Emily Oster. Unobservable selection and coefficient stability: Theory and evidence. *Journal of Business and Economic Statistics*, 37:187–204, 4 2019. ISSN 15372707. doi: 10.1080/07350015.2016.1227711.
- Nathan Petek and Nolan G Pope. The multidimensional impact of teachers on students, 2020.
- Tim R. Sass, Ron W. Zimmer, Brian P. Gill, and T. Kevin Booker. Charter High Schools' Effects on Long-Term Attainment and Earnings. *Journal of Policy Analysis and Management*, 35(3):683–706, jun 2016.
- Ying Shi. Who benefits from selective education? evidence from elite boarding school admissions. *Economics of Education Review*, 74:101907, 2 2020. ISSN 02727757. doi: 10.1016/j.econedurev.2019.07.001.
- Victoria F. Sisk, Alexander P. Burgoyne, Jingze Sun, Jennifer L. Butler, and Brooke N. Macnamara. To what extent and under which circumstances are growth mind-sets important to academic achievement? two meta-analyses. *Psychological Science*, 29:549–571, 4 2018. ISSN 0956-7976. doi: 10.1177/0956797617739704. URL <http://journals.sagepub.com/doi/10.1177/0956797617739704>.
- C Spearman. The proof and measurement of association between two things. *Source: The American Journal of Psychology*, 15:72–101, 1904.
- C. Spearman. The proof and measurement of association between two things. by c. spearman, 1904. *The American journal of psychology*, 100:441–471, 1987. doi: 10.2307/1422689.

Christopher R. Walters. The demand for effective charter schools. *Journal of Political Economy*, 126:2179–2223, 12 2018. ISSN 1537534X. doi: 10.1086/699980. URL <https://www.journals.uchicago.edu/doi/abs/10.1086/699980>.

Gregory M. Walton and Geoffrey L. Cohen. A question of belonging: Race, social fit, and achievement. *Journal of Personality and Social Psychology*, 92:82–96, 1 2007. ISSN 00223514. doi: 10.1037/0022-3514.92.1.82.

Gregory M. Walton and Geoffrey L. Cohen. A brief social-belonging intervention improves academic and health outcomes of minority students. *Science*, 331:1447–1451, 3 2011. ISSN 00368075. doi: 10.1126/science.1198364. URL <https://pubmed.ncbi.nlm.nih.gov/21415354/>.

Tables and Figures

Table 1: Summary Statistics

	Analytic Sample		Bottom Decile of Educational Advantage		Top Decile of Educational Advantage	
	Mean	SD	Mean	SD	Mean	SD
<i>Demographics</i>						
Female	0.5013	0.5000	0.2295	0.4205	0.7090	0.4543
Special education (IEP)	0.1816	0.3855	0.4611	0.4985	0.0550	0.2280
Free lunch	0.7748	0.4177	0.9520	0.2138	0.4311	0.4952
Reduced-price lunch	0.0748	0.2630	0.0168	0.1285	0.1661	0.3721
Census Block SES	-0.4485	0.8732	-0.5744	0.8141	-0.1166	0.9166
White	0.0896	0.2856	0.0378	0.1908	0.2306	0.4212
Black	0.4043	0.4908	0.5495	0.4976	0.2254	0.4179
Native American	0.0017	0.0413	0.0024	0.0493	0.0024	0.0493
Asian/Pacific Islander	0.0323	0.1768	0.0014	0.0370	0.1899	0.3922
Latino	0.4584	0.4983	0.4069	0.4913	0.3412	0.4741
<i>9th grade Intermediate Outcomes</i>						
Test Scores in 9th Grade	-0.0095	0.9913	-1.0053	0.6532	1.3844	0.7314
Work Hard in 9th Grade	0.0162	1.0014	-0.2696	1.0380	0.3825	0.9819
Social in 9th Grade	0.0025	0.9989	-0.2425	1.0456	0.3425	0.9792
Surveys in 9th Grade	0.0090	1.0013	-0.2993	1.0460	0.4080	0.9723
Behavior in 9th Grade	0.1081	0.7869	-0.5155	1.5421	0.3472	0.1838
Days Absent in 9th Grade	14.9988	18.6334	33.6946	27.7080	5.5881	7.7442
Days Suspended in 9th grade	0.8042	3.2886	2.9187	6.5937	0.0653	0.6887
Diciplinary Incidents in 9th Grade	0.0769	0.4183	0.2910	0.8709	0.0069	0.0953
<i>8th Grade Measures</i>						
Math in 8th Grade	0.1992	0.9330	-0.8528	0.6151	1.7506	0.7655
ELA in 8th Grade	0.2038	0.9312	-0.8792	0.8332	1.5784	0.7948
Emotional Health in 8th Grade	0.0703	0.8896	-0.1918	0.8831	0.3175	0.9211
Academic Engagement in 8th Grade	0.2649	0.9067	0.1446	0.8476	0.3423	1.0038
Grit in 8th Grade	0.0434	0.8303	-0.3444	0.8699	0.4596	0.7952
School Connectedness in 8th Grade	0.1400	0.8937	-0.0228	0.8594	0.4301	0.9710
Study Habits in 8th Grade	0.1522	0.8833	-0.2094	0.8455	0.6695	0.9512
Absences in 8th Grade	8.7060	8.5615	19.2280	12.4523	4.6015	3.8142
GPA in 8th Grade	2.8017	0.7776	2.0370	0.7739	3.6012	0.4825
Days Suspended in 8th Grade	0.4366	1.8074	2.2078	4.3782	0.0177	0.2514
Incidents in 8th Grade	0.0639	0.3329	0.3702	0.8326	0.0005	0.0328
<i>Long Term Outcomes</i>						
Any school-Based arrest	0.0372	0.1891	0.1236	0.3291	0.0047	0.0684
Graduation	0.7406	0.4383	0.4333	0.4956	0.9334	0.2493
Enrolled in any college within 2 years	0.5322	0.4990	0.1842	0.3877	0.8615	0.3455
Enrolled in a 4 year college within 2 years	0.3432	0.4748	0.0722	0.2588	0.7638	0.4248
Enrolled in a 2 year college within 2 years	0.2761	0.4471	0.1280	0.3342	0.2514	0.4338
N	160148		16015		15995	

Notes: Number of observations may vary by variable due to missingness and variation in cohorts for which a variable was collected. For more information see Appendix Table A1. Note that because the sample size is not perfectly divisible by ten and because observations with the same value of educational advantage are placed in the same decile group, the decile groups are slightly unequal in size.

Table 2: Temporal Stability of Value-Added and Correlations Across Value-Added

Correlations of Value-Added Within Outcomes Across Time					
	Test-Score value-added	Social value- added	Work value-added	Hard Behavior value-added	
t+1	0.4522	0.5586	0.3847		0.6022
t+2	0.1884	0.3780	0.2158		0.5338
t+3	0.4012	0.3425	0.2303		0.4801
t+4	0.4674	0.3308	0.2937		0.3475

Correlations of Average School-Level Value-Added Across Outcomes (143 Schools)					
Test Scores value-added	1				
Social value-added	0.551	1			
Hard Work value-added	0.190	0.397	1		
Behavior value-added	0.362	0.184	0.055	1	

Disattenuated Correlations of Average School-Level Value-Added Across Outcomes (143 Schools)					
Test Scores value-added	1.000				
Social value-added	0.852	1.000			
Hard Work value-added	0.322	0.780	1.000		
Behavior value-added	0.438	0.259	0.084	1.000	

Notes: All reported results are restricted to school-year cells with at least 10 respondents. The **top panel** reports, for each 9th grade measure, the correlations between a schools value-added in year t and value-added for years t+1, t+2, t+3, and t+4. The **bottom panel** reports the correlations between the value-added (estimated across all years) for the 9th grade measures.

Table 3: Average Impacts of School Effectiveness and Value-Added

	1	2	3	4	5	6	7	8	9
	Test scores 9th Grade	Surveys 9th Grade	Behaviors 9th Grade	HS Gradu- ation	School-Based Arrests	Enrolled in Any College Within 2 Years	Enrolled in 4-Year College Within 2 Years	Enrolled in 2-Year College Within 2 Years	Persists in College After 1 Year
School Effectiveness Index	0.0890*** (0.0133)	0.102*** (0.00755)	0.0571*** (0.0127)	0.0241*** (0.00430)	-0.00804*** (0.00271)	0.0257*** (0.00654)	0.0416*** (0.00903)	-0.00570 (0.00456)	0.0203*** (0.00550)
Socioemotional Value-Added	0.0623*** (0.0126)	0.0840*** (0.00818)	0.0276** (0.0106)	0.0197*** (0.00397)	-0.00564** (0.00235)	0.0192*** (0.00588)	0.0321*** (0.00866)	-0.00541 (0.00396)	0.0153*** (0.00481)
Test-Score Value-Added	0.0682*** (0.0129)	0.0445*** (0.00849)	0.0267*** (0.00830)	0.0106*** (0.00340)	-0.00479*** (0.00177)	0.0168*** (0.00564)	0.0252*** (0.00713)	-0.00200 (0.00430)	0.0137*** (0.00498)
Behavior Value-added	0.0268** (0.0122)	0.0475*** (0.0153)	0.228** (0.0164)	0.0116** (0.00485)	-0.0114*** (0.00407)	0.0167*** (0.00558)	0.0263*** (0.00654)	-0.00527 (0.00404)	0.0151*** (0.00396)
Observations	102,235	120,129	157,628	82,146	82,146	55,564	55,564	55,564	55,564

Robust standard errors in parentheses

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

Note: Each point estimate comes from a separate regression

Results are based on regression of outcomes on a single measure of out-of-sample school impacts (overall effectiveness, test score value-added, socio-emotional value-added, or behaviour value-added). All models include individual demographic controls (race / ethnicity, free and reduced price lunch, and gender), 8th grade lags (math and ELA test scores, survey measures, absences, and discipline), and school-level averages for all the demographics and lagged measures, as well as year fixed effects. We also include the socio-economic status of the student census block proxied by average occupation status and education levels. Missing 8th grade measures were imputed using 7th grade measures and demographic characteristics.

For the longer-run college outcomes, the sample includes two cohorts of first-time 9th graders in Spring 2013 and 2014. For the longer-run high-school outcomes, the sample includes three cohorts of first-time 9th graders in Spring 2012, 2013, and 2014. For the measures, the sample includes six cohorts of first-time 9th graders between Spring 2012 and 2017 **Note:** Sample sizes may differ across outcomes due to some missingness in 9th grade test scores and surveys.

Table 4: Testing for Selection

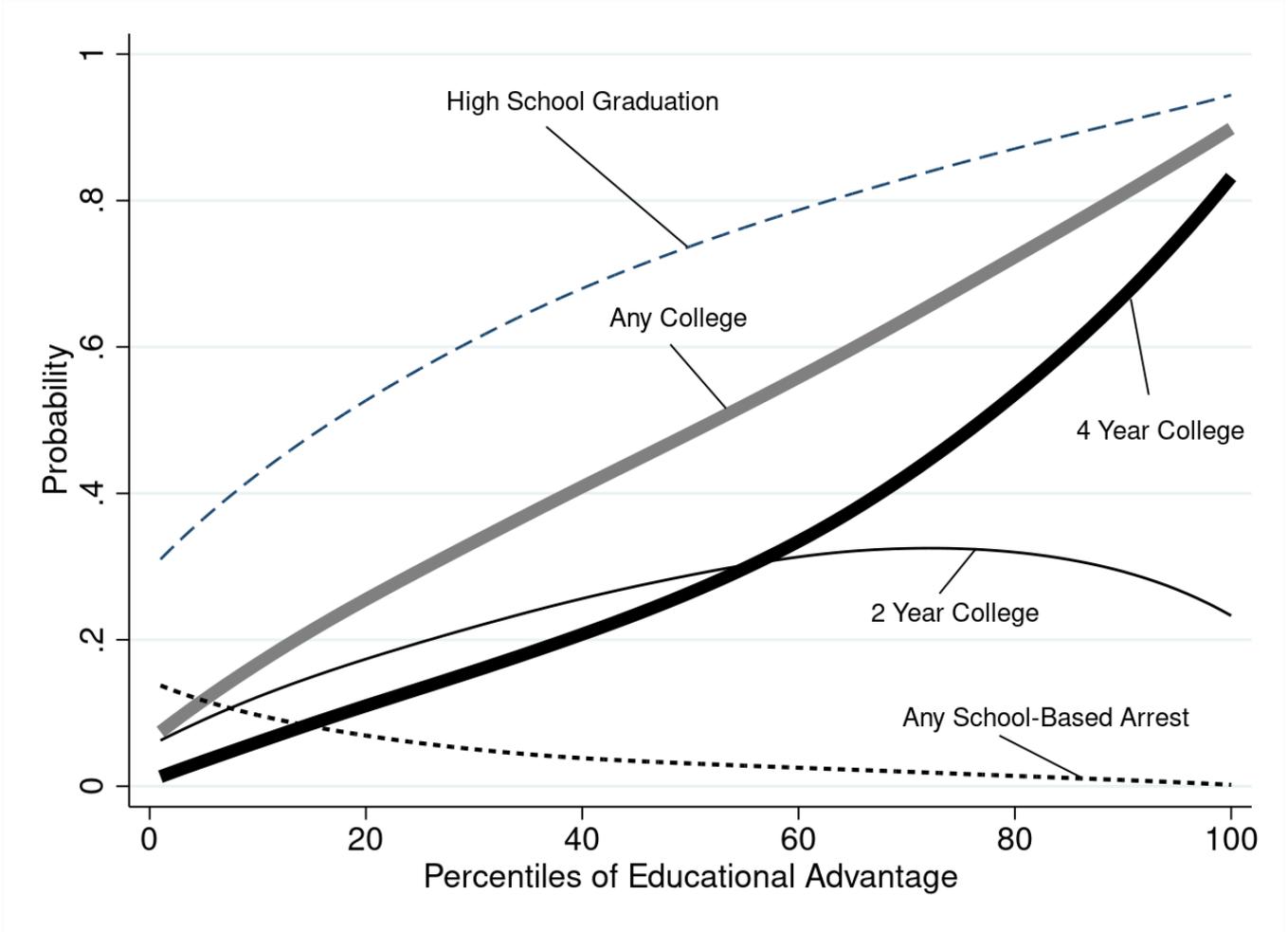
	Intermediate Outcomes				Main Longer-Run Outcomes				College Outcomes			
	1	2	3	4	5	6	7	8	9	10	11	12
	9th Grade Test Scores			Predicted	HS Graduation			Predicted	Enrolled in 4 Year College within 2 Years			Predicted
School Effectiveness Index	0.0890*** (0.0133)	0.0560*** (0.0177)	0.0527*** (0.0137)	-0.000275 (0.00430)	0.0241*** (0.00430)	0.0362*** (0.0121)	0.0171** (0.00845)	0.00266 (0.00183)	0.0416*** (0.00903)	0.0355** (0.0143)	- -	0.00347 (0.00467)
Observations	102,235	102,235	16,386	102,235	82,146	82,146	8,188	82,146	55,564	55,564	3,399	55,564
F-statistic on First Stage		321				393				374.9		
Number of Families			7786				3902				1634	
	9th Grade Survey Measures			Predicted	In-school Arrests			Predicted	Enrolled in 2 Year College within 2 Years			Predicted
School Effectiveness Index	0.102*** (0.00755)	0.105*** (0.0186)	0.0586*** (0.0173)	0.00260 (0.00852)	-0.00804*** (0.00271)	-0.0132** (0.00622)	-0.00957* (0.00501)	-0.000713 (0.000993)	-0.00570 (0.00456)	-0.0211* (0.0120)	- -	-0.000159 (0.00213)
Observations	120,129	120,129	27,103	120,129	82,146	82,146	8,188	82,146	55,564	55,564	3,399	55,564
F-statistic on First Stage		294.7				393				374.9		
Number of Families			13584				3902				1634	
	9th Grade Behaviors			Predicted	Enrolled in Any College within 2 Years			Predicted	Persist in College After 1 Year			Predicted
School Effectiveness Index	0.0571*** (0.0127)	0.0702*** (0.0207)	0.0317*** (0.0112)	0.00552 (0.00867)	0.0257*** (0.00654)	0.00650 (0.0156)	- -	0.00259 (0.00247)	0.0203*** (0.00550)	0.0153 (0.0139)	- -	0.00235 (0.00230)
Observations	157,628	157,628	41,711	157,628	55,564	55,564	3,399	55,564	55,564	55,564	3,399	55,564
F-statistic on First Stage		319.5				374.9				374.9		
Number of Families			19420				1634				1634	
Sibling FE			X				X				X	
School Assignment IV		X				X				X		

Robust standard errors adjusted for clustering at the school level.

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

Results are based on regression of outcomes on out-of-sample school effectiveness. All models include individual demographic controls (race/ethnicity, free and reduced-price lunch, and gender), 8th-grade lags (math and ELA test scores, survey measures, absences, and discipline), and school-level averages for all the demographics and lagged measures, as well as year fixed effects. We also include the socio-economic status of the student census block proxied by average occupation status and education levels. Missing 8th grade measures were imputed using 7th grade measures and demographic characteristics. For the longer-run college outcomes, the sample includes two cohorts of first-time 9th graders in Spring 2013 and 2014. For the longer-run high-school outcomes, the sample includes three cohorts of first-time 9th graders in Spring 2012, 2013, and 2014. For the measures, the sample includes six cohorts of first-time 9th graders between Spring 2012 and 2017. **Columns 4, 8, and 12:** Predicted outcomes are fitted values from a linear regression of said outcome on *all* observed controls. The predictors include lagged measures (i.e., 8th grade test scores, surveys, behaviors), gender, ethnicity, free-lunch status, and the socio-economic status of the student's census block. To avoid mechanical correlation, we use leave-year out predicted outcomes (i.e., predicted outcomes based on the relationship between the outcome and covariates in *other* years). Whether the predictions use relationships in-sample or out-of-sample, the results are the same. The reported point estimates are those on predicted outcomes on the value-added **while controlling** only for school-average math scores, school average absences, and the percent white at the school. **Note:** Sample sizes may differ across outcomes due to some missingness in 9th-grade test scores and surveys.

Figure 1. Average Outcomes: By Estimated Educational Advantage



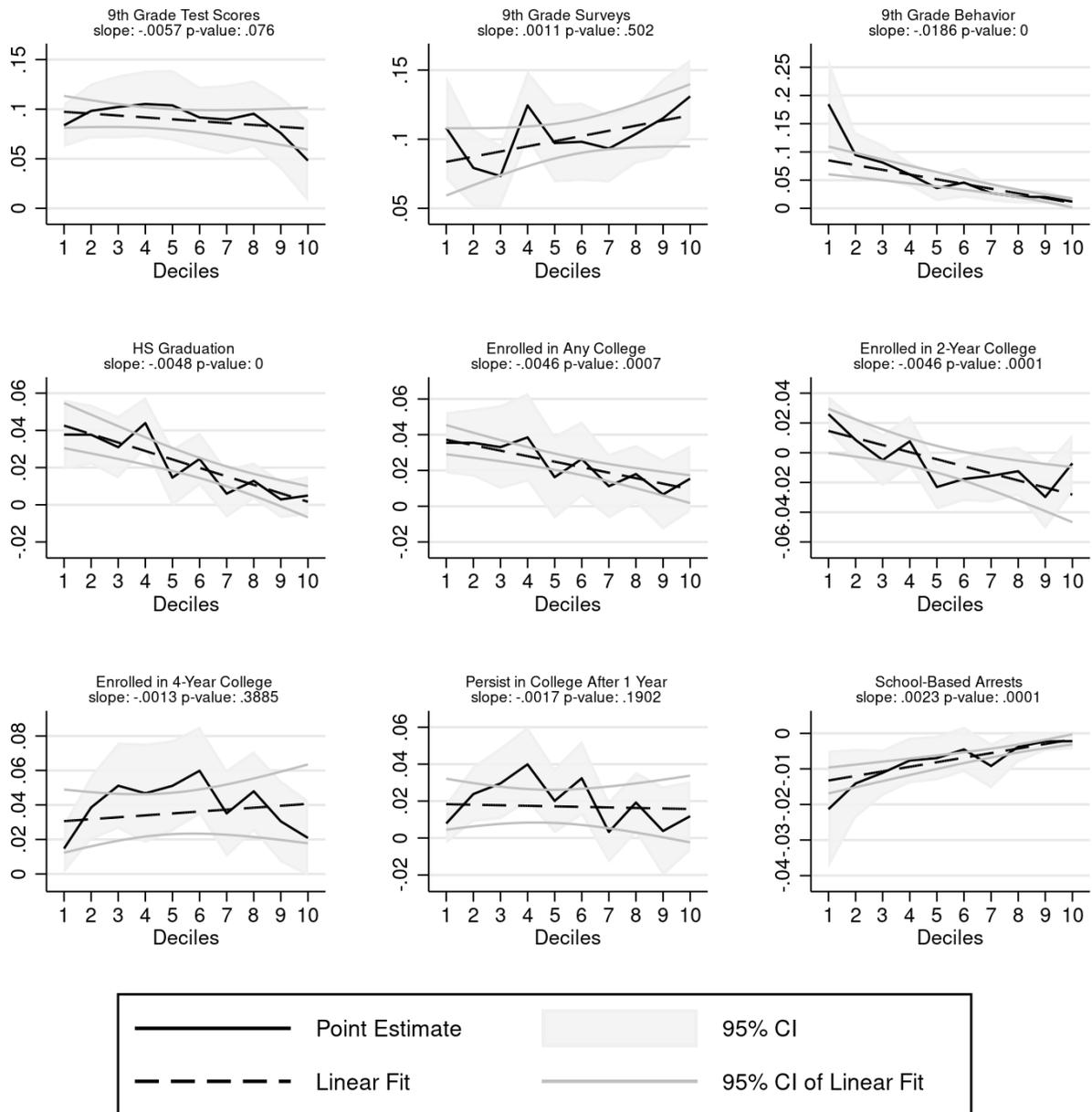
Notes: This figure plots the average of each outcome for different percentiles of the estimated educational advantage distribution. The predicted educational advantage is the fitted value from an ordered probit model predicting the level of education attained based on all 8th grade measures and demographics (*in all other years*). We present the coefficient estimates from the ordered probit model for the full sample in Appendix Table A4. We also present plots of outcome educational advantage within race and gender groups in Appendix Figure A2.

Figure 2. *Difference in Characteristics Among Most and Least Effective Schools*



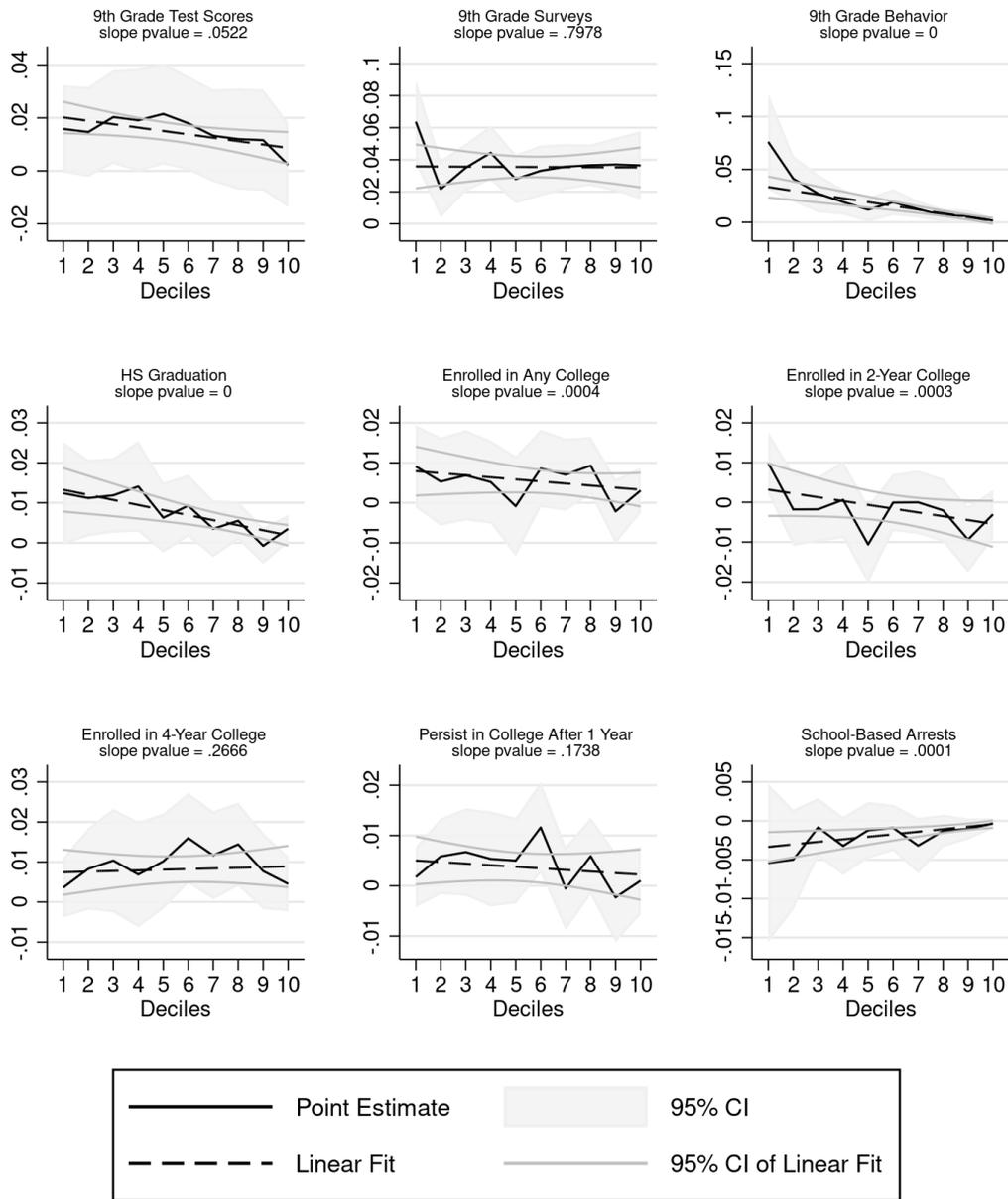
Notes: Each panel is a binned scatter-plot (20 bins) of the school characteristic against the standardized estimated school effectiveness. We report the raw correlation between the school characteristic and effectiveness below each plot for a sample of 131 schools. Share FRPL is defined as the share of students eligible for free or reduced-price lunch at a school. Note that class size and enrollment data come from the Illinois report card for the 2016-17 school year. Data are available at <https://www.illinoisreportcard.com>.

Figure 3. Impacts on Outcomes: By Estimated Educational Advantage



Notes: Each graph represents the marginal impact of a 1 standard deviation increase in overall school effectiveness for different deciles of the educational advantage distribution for a single outcome. Each panel presents the results of 10 separate regressions, each defined as in Equation (6). The 95 percent confidence interval for each point estimate is depicted by the grey shaded area. The dashed black line in each panel depicts the line of best fit for the relationship between deciles of educational advantage and the marginal effect, including a 95 confidence interval. The reported slope is the coefficient on the interaction between educational advantage and overall effectiveness, δ_2 , from equation 7. We also report the p-value associated with the null hypothesis that this slope is zero (i.e., that there is no heterogeneity).

Figure 4. *Impacts on Outcomes Unexplained by Test Score: By Educational Advantage*



Notes: Each graph depicts the *difference* between the marginal impact of a 1 standard deviation increase in the school effectiveness index and in test-score value-added. That is, it shows the impact of school effectiveness on each outcome that cannot be explained by test scores. Each regression model controls for the same covariates as in Equation (6). Each panel presents the results of 10 separate regressions. The 95 percent confidence interval for each point estimate is depicted by the grey shaded area. The dashed black line in each panel depicts the line of best fit for the relationship between deciles of educational advantage and the difference in the marginal effect for the index and test score value-added, including a 95 confidence interval. The reported p -value is for the coefficient on the difference between the interaction between educational advantage and overall effectiveness versus test-score value-added in a stacked model. The p -value is associated with the null hypothesis that the difference in effectiveness unexplained by test scores is the same by educational advantage.

VII Appendix

Table A1: Observations By Outcome and Cohort

	Cohort						Total
	2012	2013	2014	2015	2016	2017	
All	28,141	27,423	26,582	26,376	26,718	24,908	160,148
With Behaviors	27,540	26,991	26,229	26,009	26,358	24,501	157,628
With Test scores	24,785	20,885	24,343	15,252	16,970	0	102,235
With Surveys	20,259	20,540	19,279	19,706	21,174	19,171	120,129
With Graduation	28,141	27,423	26,582	0	0	0	82,146
With Arrests	28,141	27,423	26,582	0	0	0	82,146
With College	28,141	27,423	0	0	0	0	55,564

Notes: The table displays the number of observations with non-missing values for each outcome for each spring cohort of first time 9th-graders. The full sample includes all first-time 9th-graders attending a school for which we can calculate effectiveness. Students with behaviors have 9th-grade data for absences, incidents, and suspensions. Students with surveys have 9th-grade data for emotional health, academic engagement, grit, school connectedness, and study habits. Students with test scores have 9th-grade test data in math and English. Graduation and arrests are available for first-time 9th-graders in the Spring of 2012, 2013, and 2014. College outcomes are available for the Spring 2012 and 2013 9th-grade cohorts.

Table A2: Summary Statistics for Survey Completers and Non-Completers

VARIABLES	Full Analytic Sample		Completed Surveys in 9th Grade		Did Not Complete in 9th Grade		Completed Tests in 9th Grade	
	(1) mean	(2) sd	(3) mean	(4) sd	(5) mean	(6) sd	(7) mean	(8) sd
<i>Demographics</i>								
Female	0.501	0.500	0.512	0.500	0.469	0.499	0.506	0.500
Special education (IEP)	0.182	0.385	0.157	0.364	0.255	0.436	0.167	0.373
Free lunch	0.775	0.418	0.768	0.422	0.796	0.403	0.787	0.409
Reduced-price lunch	0.0748	0.263	0.0783	0.269	0.0641	0.245	0.0799	0.271
Census Block SES	-0.449	0.873	-0.454	0.881	-0.433	0.848	-0.474	0.865
White	0.0896	0.286	0.0949	0.293	0.0735	0.261	0.0795	0.270
Black	0.404	0.491	0.375	0.484	0.492	0.500	0.416	0.493
Native American	0.00171	0.0413	0.00164	0.0405	0.00192	0.0438	0.00165	0.0406
Asian/Pacific Islander	0.0323	0.177	0.0361	0.187	0.0208	0.143	0.0304	0.172
Latino	0.458	0.498	0.479	0.500	0.397	0.489	0.461	0.498
<i>9th grade Intermediate Outcomes</i>								
Test Scores in 9th Grade	-0.00948	0.991	0.0483	0.977	-0.196	1.015	-0.00948	0.991
Work Hard in 9th Grade	0.0162	1.001	0.0165	1.001			-0.000766	0.984
Social in 9th Grade	0.00246	0.999	0.00822	0.996			-0.00867	0.989
Surveys in 9th Grade	0.00899	1.001	0.00899	1.001			-0.00655	0.983
Behavior in 9th Grade	0.108	0.787	0.163	0.665			0.120	0.724
Days Absent in 9th Grade	15.00	18.63	12.90	15.54	21.49	24.87	13.18	14.91
Days Suspended in 9th grade	0.804	3.289	0.634	2.772	1.314	4.457	0.859	3.174
Diciplinary Incidents in 9th Grade	0.0769	0.418	0.0609	0.357	0.125	0.561	0.0737	0.376
<i>8th Grade Measures</i>								
Math in 8th Grade	0.199	0.933	0.258	0.928	0.0214	0.925	0.163	0.907
ELA in 8th Grade	0.204	0.931	0.265	0.914	0.0206	0.959	0.182	0.910
Emotional Health in 8th Grade	0.0703	0.890	0.0826	0.897	0.0334	0.867	0.0827	0.904
Academic Engagement in 8th Grade	0.265	0.907	0.271	0.918	0.246	0.872	0.253	0.902
Grit in 8th Grade	0.0434	0.830	0.0521	0.839	0.0174	0.803	0.0499	0.828
School Connectedness in 8th Grade	0.140	0.894	0.145	0.903	0.126	0.866	0.120	0.892
Study Habits in 8th Grade	0.152	0.883	0.165	0.897	0.114	0.839	0.140	0.882
Absences in 8th Grade	8.706	8.561	8.102	7.695	10.52	10.54	8.270	7.685
GPA in 8th Grade	2.802	0.778	2.848	0.771	2.661	0.781	2.779	0.762
Days Suspended in 8th Grade	0.437	1.807	0.351	1.547	0.693	2.409	0.447	1.724
Incidents in 8th Grade	0.0639	0.333	0.0520	0.285	0.0998	0.444	0.0651	0.324
<i>Long-term Outcomes</i>								
Any school-Based arrest	0.0372	0.189	0.0313	0.174	0.0531	0.224	0.0330	0.179
Graduation	0.741	0.438	0.778	0.415	0.638	0.481	0.785	0.411
Enrolled in any college within 2 years	0.532	0.499	0.577	0.494	0.408	0.492	0.572	0.495
Enrolled in a 4 year college within 2 years	0.343	0.475	0.378	0.485	0.247	0.432	0.366	0.482
Enrolled in a 2 year college within 2 years	0.276	0.447	0.296	0.456	0.222	0.416	0.298	0.457
N	160,148		120,129		40,019		102,235	

Notes: Survey completers are students who have 9th-grade data for emotional health, academic engagement, grit, school connectedness, and study habits. As such, we report averages for some measures even among non-completers because many non-completers are missing some data but not others. Those with complete test data have 9th-grade test data in math and English.

Table A3: Psychometric Properties of SED measures (as reported by the University of Chicago Consortium on School Research): 2011 through 2013

Measure	School Year	Separation	Reliability	Item Infits	Item Outfits
Grit	2010-11	1.68	0.74	0.84, 0.76, 0.71, 1.24	0.85, 0.76, 0.71, 1.19
Social Skills	2010-11	1.69	0.74	1.08, 1.36, 1.41, 1.11	1.05, 1.33, 1.44, 1.15
Academic Effort	2010-11	1.74	0.75	0.85, 1.22, 1.1, 0.91	0.82, 1.17, 1.12, 0.94
Academic Engagement	2010-11	1.59	0.7	0.49, 0.56, 0.71, 0.56	0.49, 0.57, 0.72, 0.58
Belonging	2010-11	2.07	0.81	0.93, 1.02, 0.99, 0.96, 1.29	0.91, 0.97, 0.99, 0.93, 1.33
Grit	2011-12	1.54	0.7	0.8, 0.73, 0.68, 1.19	0.81, 0.57, 0.6, 0.42
Social Skills	2011-12	1.68	0.74	1.37, 1.36, 1.28, 1.06	1.68, 1.24, 1.18, 0.95
Academic Effort	2011-12	1.75	0.75	0.85, 1.22, 1.08, 0.92	0.82, 1.17, 1.1, 0.96
Academic Engagement	2011-12	1.56	0.71	0.54, 0.53, 0.47, 0.69	0.56, 0.55, 0.48, 0.71
Belonging	2011-12	2.13	0.82	0.98, 1.28, 0.91, 1.02, 0.97	0.97, 1.32, 0.89, 0.97, 0.94
Grit	2012-13	1.55	0.71	0.77, 0.69, 0.63, 1.13	0.79, 0.7, 0.63, 1.1
Social Skills	2012-13	1.67	0.74	1.3, 1.37, 1.23, 1.04	1.55, 1.25, 1.12, 0.94
Academic Effort	2012-13	1.77	0.76	0.86, 1.2, 1.13, 0.94	0.83, 1.15, 1.15, 0.97
Academic Engagement	2012-13	1.57	0.71	0.55, 0.54, 0.47, 0.69	0.57, 0.56, 0.48, 0.70
Belonging	2012-13	2.14	0.82	0.95, 1.28, 0.90, 1.03, 0.96	0.95, 1.31, 0.87, 0.98, 0.93

Notes. The reported statistics are from internal documentation at the University of Chicago Consortium on School Research where Rasch analysis was performed on individual survey items. All measures are anchored to 2010-11 step and item difficulties. Infit and outfit measures greater than 1 indicate underfit to the Rasch model and values lower than 1 indicate overfit. Generally, infit and outfit values in the range of 0.6-1.4 are considered reasonable for survey measures. Reliability represents individual reliability and includes extreme people. The patterns are very similar for years 2013 through 2018.

Table A4: Ordered Probit Parameter Estimates

	Educational Advantage		cont'd
8th Grade Math	0.270*** (0.0111)	Native	-0.414 (0.280)
8th Grade Math Squared	0.00434 (0.00830)	Asian	0.0820 (0.166)
8th Grade ELA	0.177*** (0.0102)	Latinx	-0.295* (0.155)
8th Grade ELA Squared	0.00845* (0.00462)	Other Race	0.162 (0.259)
Emotional Health in 8th Grade	-0.0185** (0.00743)	Female	0.0274 (0.124)
Academic Engagement in 8th Grade	-0.0154** (0.00694)	Female * White	0.0684 (0.113)
Grit in 8th Grade	0.0768*** (0.00803)	Female * Black	0.322** (0.131)
School Connectedness in 8th Grade	-0.0173** (0.00699)	Female * Native	0.250 (0.259)
Study Habits in 8th Grade	0.0910*** (0.0102)	Female * Asian	0.104 (0.156)
8th Grade Top 25% of Absences	-0.564*** (0.0155)	Female * Latinx	0.221* (0.123)
Serious Incidents in 7th or 8th Grade	-0.370*** (0.0292)	Female * Other Race	-0.473 (0.317)
Receive Free Lunch	-0.162** (0.0658)	/cut1	-1.148*** (0.174)
Receive Reduced Price Lunch	0.0889 (0.0661)	/cut2	-0.485*** (0.179)
White	-0.203 (0.137)	/cut3	0.0876 (0.185)
Black	-0.343** (0.169)	Observations	57,093

Robust standard errors in parentheses

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

This reports the parameter estimates from an ordered probit model predicting the level of educational attainment (dropout, graduate high school, attend 2-year college, attend 4-year college) as a function of observable 8th-grade characteristics. Note that the sample size is larger than the analytic sample used for the main outcome analysis because this model uses all observations with college data, which includes 1529 observations for individuals who attend schools that do not have valid value-added estimates. The results are virtually identical if we restrict the prediction model only to use those individuals in the main analytic long-term sample.

Table A5: Factor Analysis

	Eigenvalue	Difference	Proportion
Factor 1	1.27178	1.05385	1.1393
Factor 2	0.21793	.	0.1952

Rotated factor loadings (pattern matrix) and unique variances

	Factor1	Factor2	Uniqueness
Workhard value-added	0.6287	-0.0316	0.6038
Social value-added	0.7328	0.179	0.4309
Test Score value-added	0.5443	0.2787	0.626
Behavior value-added	0.2079	0.3274	0.8496

Method: principal factors

Rotation: orthogonal varimax (Kaiser off).

Note: The proportion explained by this factor is greater than one because the model also includes factors with negative eigenvalues.

Table A6: Effect of SED Value-Added on Average Intermediate and Long-Term Student Outcomes

	(1) Test scores 9th Grade	(2) Surveys 9th Grade	(3) Behaviors 9th Grade	(4) HS Gradua- tion	(5) School- Based Arrests	(6) Enrolled in Any College Within 2 Years	(7) Enrolled in 4-Year College Within 2 Years	(8) Enrolled in 2-Year College Within 2 Years	(9) Persists in College After 1 Year
Work hard Value-Added	0.0533*** (0.0138)	0.0727*** (0.00907)	0.0215** (0.00991)	0.0180*** (0.00415)	-0.00522** (0.00225)	0.0175*** (0.00584)	0.0285*** (0.00964)	-0.00448 (0.00407)	0.0145*** (0.00469)
Social Value-Added	0.0761*** (0.0121)	0.102*** (0.00887)	0.0394*** (0.0139)	0.0219*** (0.00471)	-0.00565* (0.00297)	0.0206*** (0.00652)	0.0352*** (0.00875)	-0.00555 (0.00449)	0.0152*** (0.00563)
Observations	102,235	120,129	157,628	82,146	82,146	55,564	55,564	55,564	55,564

Robust standard errors in parentheses

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

Note: Results are based on separate regressions of outcomes on out-of-sample socioemotional value-added, disaggregated into the social well-being and work hard constructs. All models include individual demographic controls (race / ethnicity, free and reduced price lunch, and gender), 8th grade lags (math and ELA test scores, survey measures, absences, and discipline), and school-level averages for all the demographics and lagged measures, as well as year fixed effects. We also include the socio-economic status of the student census block proxied by average occupation status and education levels. Missing 8th grade measures were imputed using 7th grade measures and demographic characteristics. For the longer-run college outcomes, the sample includes first-time 9th-grade students in the Spring of 2012 and 2013. For the longer-run high-school outcomes, the sample includes first-time 9th grade students in 2012, 2013, and 2014. For the measures, the sample includes six cohorts of first-time 9th grade students between 2012 and 2017. Sample sizes may differ across outcomes due to some missingness in 9th-grade test scores and surveys. Each point estimate is based on a separate regression.

Table A7: Results Using Small Within-Family College Sample

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
	Test scores 9th Grade	Surveys 9th Grade	Behaviors 9th Grade	HS Gradua- tion	School- Based Arrests	Enrolled in Any College Within 2 Years	Enrolled in 4- Year College Within 2 Years	Enrolled in 2-Year College Within 2 Years	Persists in College After 1 Year
School Effectiveness Index	-0.0175 (0.0200)	0.00638 (0.0272)	0.0152 (0.0299)	-0.00979 (0.0124)	-0.00410 (0.00739)	0.00108 (0.0132)	-0.00618 (0.00917)	0.00500 (0.0124)	0.00328 (0.00774)
Observations	1,943	2,439	3,357	3,399	3,399	3,399	3,399	3,399	3,399
Sibling FE	X	X	X	X	X	X	X	X	X
Number of Families	940	1178	1614	1634	1634	1634	1634	1634	1634

Robust standard errors in parentheses

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

Note: Results are based on regressions of outcomes on out-of-sample school-effectiveness using only within-family variation and limiting the sample to those old enough to have college outcomes. All models include individual demographic controls (race / ethnicity, free and reduced price lunch, and gender), 8th grade lags (math and ELA test scores, survey measures, absences, and discipline), and school-level averages for all the demographics and lagged measures, as well as year fixed effects. We also include the socio-economic status of the student census block proxied by average occupation status and education levels. Missing 8th grade measures were imputed using 7th grade measures and demographic characteristics.

Table A8: Marginal Effects For Individual Deciles

Deciles	(1) Test scores 9th Grade	(2) Surveys 9th Grade	(3) Behaviors 9th Grade	(4) HS Gradu- ation	(5) School- Based Arrests	(6) Enrolled in Any College Within 2 Years	(7) Enrolled in 4-Year College Within 2 Years	(8) Enrolled in 2-Year Col- lege Within 2 Years	(9) Persists in College After 1 Year
1	0.0838*** (0.0111)	0.108*** (0.0190)	0.185*** (0.0403)	0.0378*** (0.00951)	-0.0213** (0.00836)	0.0355*** (0.00874)	0.0147** (0.00719)	0.0259*** (0.00621)	0.00789 (0.00551)
2	0.0983*** (0.0141)	0.0793*** (0.0143)	0.0945*** (0.0204)	0.0377*** (0.00823)	-0.0141*** (0.00497)	0.0355*** (0.00962)	0.0386*** (0.00959)	0.00857 (0.00740)	0.0238*** (0.00768)
3	0.102*** (0.0162)	0.0734*** (0.0115)	0.0813*** (0.0163)	0.0310*** (0.00857)	-0.0111*** (0.00333)	0.0331*** (0.0120)	0.0512*** (0.0128)	-0.00493 (0.00896)	0.0296*** (0.00990)
4	0.105*** (0.0171)	0.124*** (0.0127)	0.0601*** (0.0111)	0.0439*** (0.00713)	-0.00763** (0.00327)	0.0385*** (0.0126)	0.0468*** (0.0147)	0.00761 (0.00886)	0.0399*** (0.0104)
5	0.104*** (0.0182)	0.0972*** (0.0145)	0.0359*** (0.0118)	0.0146* (0.00792)	-0.00697** (0.00319)	0.0162 (0.0118)	0.0511*** (0.0135)	-0.0231*** (0.00770)	0.0201** (0.00898)
6	0.0917*** (0.0158)	0.0982*** (0.0144)	0.0458*** (0.0135)	0.0246*** (0.00728)	-0.00453 (0.00324)	0.0265** (0.0107)	0.0598*** (0.0130)	-0.0177** (0.00752)	0.0324*** (0.0103)
7	0.0895*** (0.0178)	0.0933*** (0.0125)	0.0270*** (0.00720)	0.00600 (0.00661)	-0.00918*** (0.00323)	0.0112 (0.00919)	0.0351*** (0.0129)	-0.0155* (0.00920)	0.00330 (0.00849)
8	0.0954*** (0.0171)	0.104*** (0.0110)	0.0204*** (0.00607)	0.0129** (0.00501)	-0.00380* (0.00221)	0.0180** (0.00830)	0.0479*** (0.0119)	-0.0124 (0.00838)	0.0192** (0.00843)
9	0.0759*** (0.0187)	0.115*** (0.0146)	0.0202*** (0.00607)	0.00290 (0.00511)	-0.00229* (0.00138)	0.00661 (0.01000)	0.0305** (0.0121)	-0.0298*** (0.0110)	0.00381 (0.0123)
10	0.0482** (0.0211)	0.131*** (0.0136)	0.0121*** (0.00331)	0.00496 (0.00534)	-0.00221* (0.00113)	0.0154 (0.00945)	0.0209* (0.0112)	-0.00714 (0.00990)	0.0118 (0.00963)

Robust standard errors in parentheses adjusted for clustering at the school level.

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

We report the main result for regression models estimated on different sub-samples of the data. Each point estimate comes from a separate regression. Results are based on regression of outcomes on the out-of-sample overall school effectiveness index. All models include individual demographic controls (race/ethnicity, free and reduced-price lunch, and gender), 8th-grade lags (math and ELA test scores, survey measures, absences, and discipline), and school-level averages for all the demographics and lagged measures, as well as year fixed effects. We also include the socio-economic status of the student census block proxied by average occupation status and education levels. Missing 8th grade measures were imputed using 7th grade measures and demographic characteristics. For the longer-run college outcomes, the sample includes two cohorts of first-time 9th graders in Spring 2013 and 2014. For the longer-run high-school outcomes, the sample includes three cohorts of first-time 9th graders in Spring 2012, 2013, and 2014. For the measures, the sample includes six cohorts of first-time 9th graders between Spring 2012 and 2017.

Table A9: Marginal Effects For Select Decile Groups

	Test Scores						Surveys					
	Bottom Decile	Top Decile	Bottom Deciles	3	Top Deciles	3	Bottom Decile	Top Decile	Bottom Deciles	3	Top Deciles	3
School Effectiveness Index	0.0838*** (0.0111)	0.0482** (0.0211)	0.0975*** (0.0124)		0.0692*** (0.0165)		0.108*** (0.0190)	0.131*** (0.0136)	0.0868*** (0.0101)		0.115*** (0.00914)	
	Behaviors						Grad					
	Bottom Decile	Top Decile	Bottom Deciles	3	Top Deciles	3	Bottom Decile	Top Decile	Bottom Deciles	3	Top Deciles	3
School Effectiveness Index	0.185*** (0.0403)	0.0121*** (0.00331)	0.119*** (0.0247)		0.0184*** (0.00441)		0.0378*** (0.00951)	0.00496 (0.00534)	0.0376*** (0.00703)		0.00780** (0.00302)	
	College						Arrested					
	Bottom Decile	Top Decile	Bottom Deciles	3	Top Deciles	3	Bottom Decile	Top Decile	Bottom Deciles	3	Top Deciles	3
School Effectiveness Index	0.0355*** (0.00874)	0.0154 (0.00945)	0.0370*** (0.00856)		0.0148** (0.00623)		-0.0213** (0.00836)	-0.00221* (0.00113)	-0.0149*** (0.00495)		-0.00262** (0.00122)	
	4-Year College						2-Year College					
	Bottom Decile	Top Decile	Bottom Deciles	3	Top Deciles	3	Bottom Decile	Top Decile	Bottom Deciles	3	Top Deciles	3
School Effectiveness Index	0.0147** (0.00719)	0.0209* (0.0112)	0.0364*** (0.00908)		0.0329*** (0.00898)		0.0259*** (0.00621)	-0.00714 (0.00990)	0.0110* (0.00587)		-0.0143** (0.00666)	

Robust standard errors in parentheses

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

We report the main result for regression models estimated on different sub-samples of the data. Results are based on regression of outcomes on the out-of-sample overall school effectiveness index. All models include individual demographic controls (race/ethnicity, free and reduced-price lunch, and gender), 8th-grade lags (math and ELA test scores, survey measures, absences, and discipline), and school-level averages for all the demographics and lagged measures, as well as year fixed effects. We also include the socio-economic status of the student census block proxied by average occupation status and education levels. Missing 8th grade measures were imputed using 7th grade measures and demographic characteristics. For the longer-run college outcomes, the sample includes two cohorts of first-time 9th graders in Spring 2013 and 2014. For the longer-run high-school outcomes, the sample includes three cohorts of first-time 9th graders in Spring 2012, 2013, and 2014. For the measures, the sample includes six cohorts of first-time 9th graders between Spring 2012 and 2017.

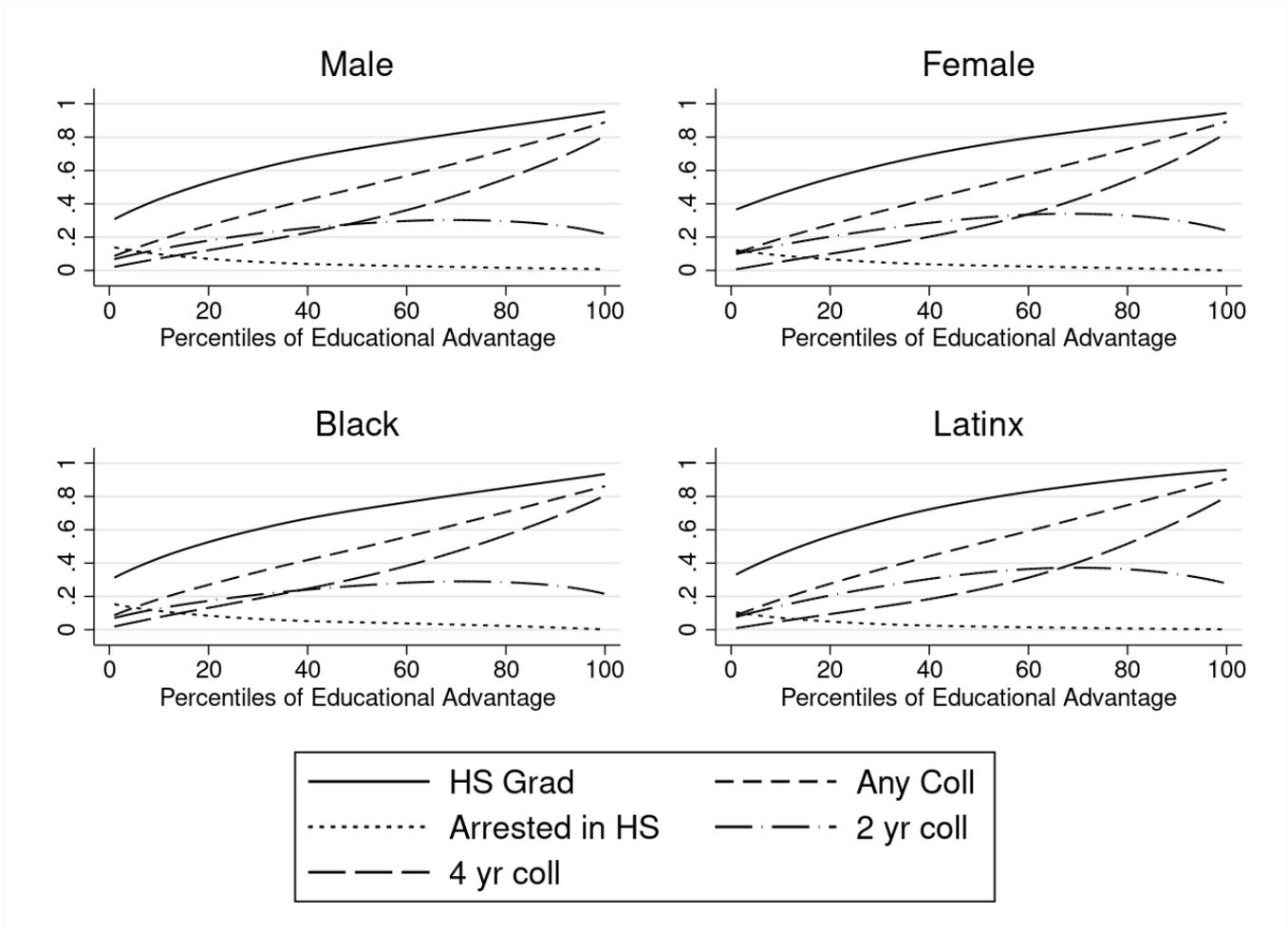
Table A10: The Interaction Effect

	OLS								
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
	Test Scores	Surveys	Behaviors	Graduate	Arrested	Any College	4-Year College	2-Year College	Persists
Effectiveness*Advantage	-0.00566* (0.00317)	0.00114 (0.00170)	-0.0186*** (0.00293)	-0.00481*** (0.000886)	0.00233*** (0.000566)	-0.00463*** (0.00134)	-0.00135 (0.00156)	-0.00460*** (0.00116)	-0.00174 (0.00132)
	IV								
	(10)	(11)	(12)	(13)	(14)	(15)	(16)	(17)	(18)
Effectiveness*Advantage	-0.0107* (0.00637)	0.00275 (0.00245)	-0.0196*** (0.00395)	-0.00320** (0.00129)	0.00226*** (0.000817)	-0.00392** (0.00190)	-0.000283 (0.00226)	-0.00581*** (0.00125)	-0.00245 (0.00174)
Observations	102,235	120,129	157,628	82,146	82,146	55,564	55,564	55,564	55,564

Robust standard errors in parentheses. *** p<0.01, ** p<0.05, * p<0.1

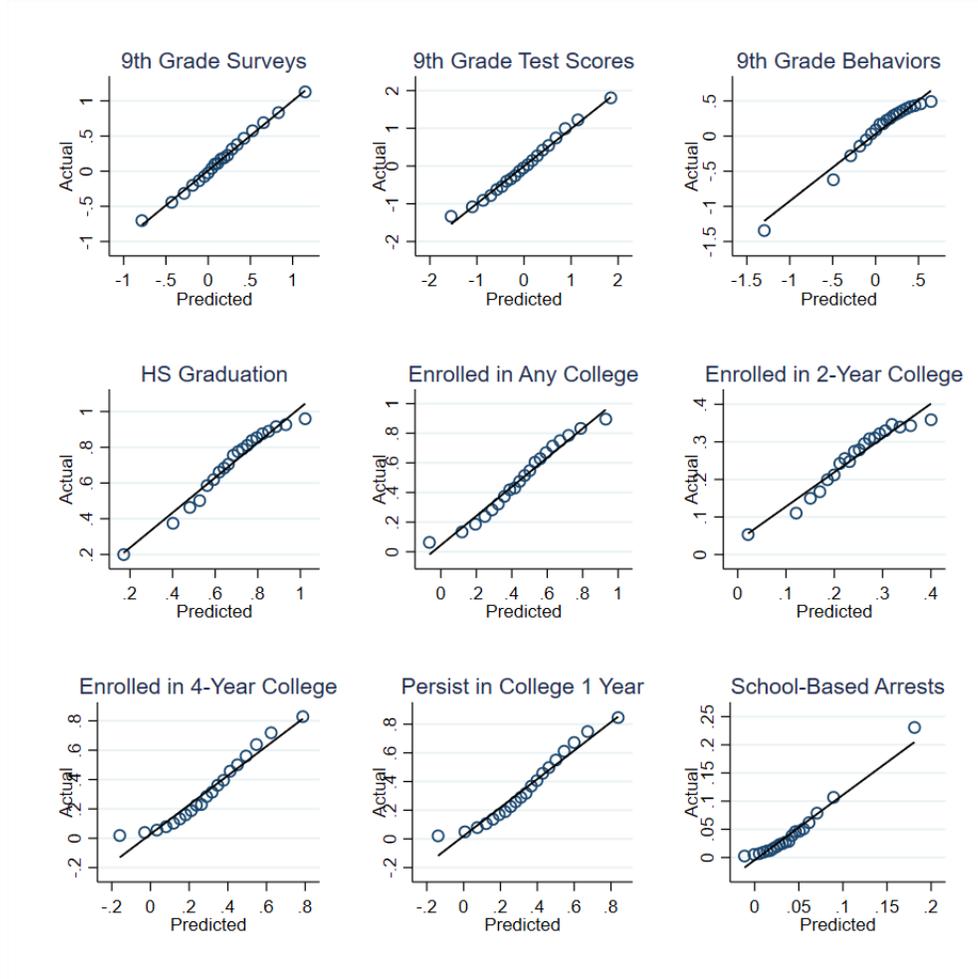
Notes: The OLS coefficient on the interaction between overall effectiveness and the decile of education advantage is presented in the top row. The bottom row presents Instrumental Variables (IV) estimate of this same parameter while instrumenting for the interaction with the interaction between 8th-grade measures (test scores, surveys, and behaviors) and effectiveness. The IV models all have first-stage F-statistics above 100. All models include individual demographic controls (race/ethnicity, free and reduced-price lunch, and gender), 8th grade lags (math and ELA test scores, survey measures, absences, and discipline), and school-level averages for all the demographics and lagged measures, as well as year fixed effects. We also include the socio-economic status of the student census block proxied by average occupation status and education levels. Missing 8th grade measures were imputed using 7th grade measures and demographic characteristics. For the longer-run college outcomes, the sample includes first-time 9th-grade students in the Spring of 2012 and 2013. For the longer-run high-school outcomes, the sample includes first-time 9th grade students in 2012, 2013, and 2014. For the measures, the sample includes six cohorts of first-time 9th grade students between 2012 and 2017. Sample sizes may differ across outcomes due to some missingness in 9th-grade test scores and surveys. Each point estimate is based on a separate regression.

Figure A1. *Average Outcomes by Educational Advantage: By Race and Gender*



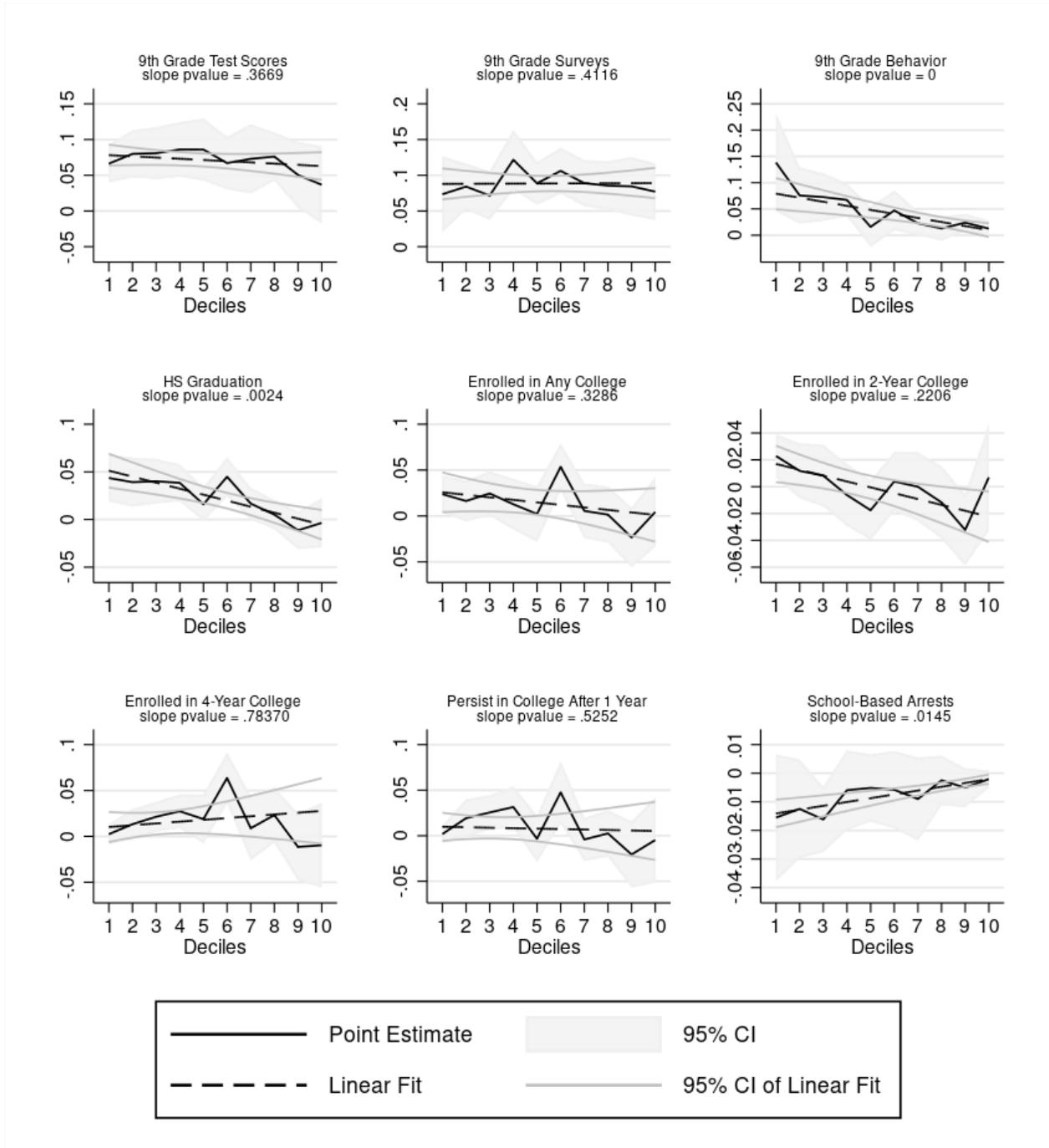
Notes: This figure plots the average of each outcome for different percentiles of the estimated educational advantage distribution by race and gender. The predicted educational advantage is the fitted value from an ordered probit model predicting the level of education attained based on all 8th grade measures and demographics (*in all other years*). We present the coefficient estimates from the ordered probit model for the full sample in Appendix Table A4.

Figure A2. *Actual Outcome by Predicted Outcome*



Notes: Each graph presents the average of the actual outcome for different groups of students by predicted outcome. The predicted outcomes are the fitted values from a regression of each outcome on all observed demographics and 8th grade measures based on students in *other* years. The predictors include lagged measures (i.e., 8th grade test scores, surveys, behaviors), gender, ethnicity, free-lunch status, and the socio-economic status of the student's census block.

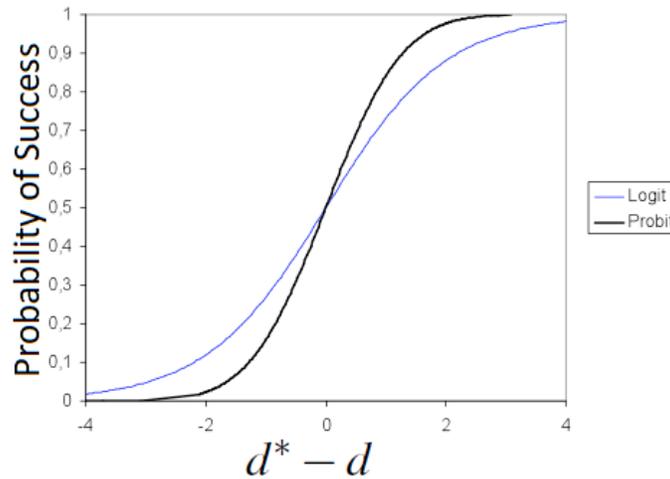
Figure A3. Impacts on Outcomes: By Educational Advantage (neighborhood schools only)



Notes: Each graph represents the marginal impact of a 1 standard deviation increase in overall school effectiveness for different deciles of the educational advantage distribution for a single outcome. Each panel presents the results of 10 separate regressions each defined as in Equation (6) but only on the sample of traditional public school students. The 95 percent confidence interval for each point estimate is depicted by the grey shaded area. The dashed black line in each panel depicts the line of best fit for the relationship between deciles of educational advantage and the marginal effect, including a 95 confidence interval.

Appendix B: The Test For Mechanical Heterogeneity

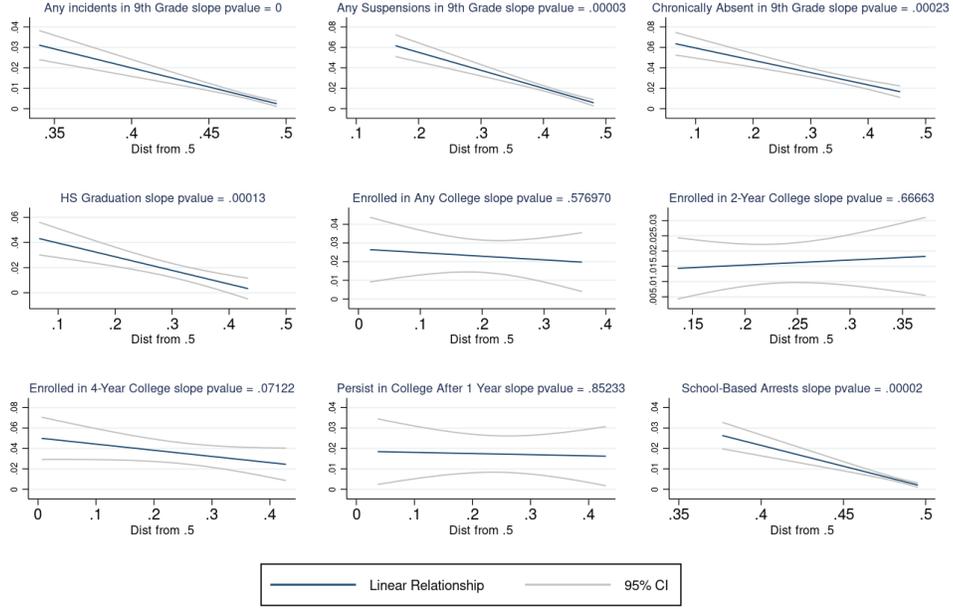
In models such as the logit or probit models, there is some continuous latent variable (d) that summarizes a predisposition to “success” (i.e., a positive outcome). The realized outcome ($D \in 0, 1$) is a function of this latent disposition plus some random error, ε , such that $D = 1$ if $d + \varepsilon > d^*$ and $D = 0$ otherwise, where d^* is some unobserved fixed threshold. The probability of success for an individual with disposition d is therefore $Pr(\varepsilon > d^* - d)$. This is a cumulative probability. See the logit and probit models depicted below. The change in the probability due to a marginal increase in d is therefore the probability density of ε at $d^* - d$. Under a symmetric single-peaked bell-shaped distribution of ε (as in a logit or probit model), for the same change in d , the observed change in probability will be largest (in magnitude) for individuals with d close to d^* and success probability close to 0.5, declining in magnitude for those with d farther from d^* and success probability farther from 0.5, and smallest for those with d very far from d^* and success probability farthest from 0.5 (i.e., close to 1 or 0).



For example, consider an outcome such as being suspended. The underlying disposition toward success (in this case being suspended) for any group is denoted d_g . Those in the top decile of educational advantage have baseline suspension rates below 1 percent and therefore have large $|d^* - d_g|$ compared to those in the bottom decile who have baseline suspension rates around 30 percent and therefore much smaller $|d^* - d_g|$. The above framework indicates that **with the same change in d_g** , the change in suspension rates will be larger for the bottom decile than the top.

This logic forms the basis for our test. We propose that if the differences in the marginal effect on binary outcomes across groups can be explained by differences in the baseline probabilities across groups, it would be indicative “mechanical heterogeneity”. In contrast, if differences in baseline probabilities across groups do not explain differences in marginal effects on these binary outcomes, it would imply that any observed heterogeneity is not mechanical and therefore reflects heterogeneous effects on skills and latent predispositions.

Figure B1. Relationship Between Probability of Success and Marginal Effects



Notes: Each graphs plots the linear relationship between the the absolute value of the marginal effect for each decile group and the difference between the average success rate for that same decile group and 0.5. This comes from the regression model laid out and detailed in equation (9). That is, for each binary outcome, we run the regression below

$$|\delta_g| = \alpha + \pi \times (|p_g - 0.5|) + v_g \quad (10)$$

where $|\delta_g|$ is the absolute value of the marginal effect for decile group g , and p_g is the average success rate for decile group g . The slope of the plotted lined in each graph is π , which represents the relationship between the absolute value of the marginal effect and the distance between the baseline success rate and 0.5. For each outcome, we report the p -value on the hypotheses that $\pi = 0$. If $\pi = 0$, it would imply that the differences in the marginal effect on binary outcomes can be explained by differences in the baseline probabilities across groups – which would be indicative “mechanical heterogeneity”. In contrast, if differences in baseline probabilities across groups do not explain differences in marginal effects on these binary outcomes, it would imply that any observed heterogeneity is not mechanical and therefore reflects heterogeneous effects on skills and latent predispositions.

VII.1 Simulation Evidence

Because there could be both mechanical heterogeneity and treatment heterogeneity, it is helpful to disentangle the two. To this aim, we simulate how much heterogeneity one would observe due to the “mechanical effect” alone. We then compare this to the heterogeneity observed. We conduct this simulation as below:

1. For each binary outcome, we use a probit model to estimate the latent variable for each observation based on observed covariates (excluding the school effect). Note that these are simply the fitted values of the linear portion of the probit model. The basic idea is that some underlying latent outcome, y^* , is a linear function of observable characteristics X plus a normally distributed mean-zero error term ε , so that $y^* = X\beta + \varepsilon$. For each individual i , the observed outcome is given by

$$y_i = \begin{cases} 1 & X\beta + \varepsilon > 0 \\ 0 & \text{otherwise} \end{cases}$$

Under the assumption that the error term is normally distributed, it follows that $Pr(Y = 1|X) = \Phi(X\beta)$ which can be estimated using a probit regression. The fitted values from a probit regression of a binary outcome on all the observed covariates (i.e., $X\hat{\beta}$) is an estimate of each persons latent disposition as a function of the observed covariates. The model also provides the standard error of this estimate, $s_{X_i\hat{\beta}}$, which we use to account for noise in the estimation of the latent index.

2. To simulate constant effects on the latent variable (i.e., no true heterogeneity), we add a constant marginal treatment effect (π) to this underlying latent variable. To match the observed marginal effect on average, we set this simulated marginal effect equal to the average effect (estimated using a separate probit model). This simulated latent variable is $y_{i,sim}^* = X_i\hat{\beta} + \pi\hat{\omega}_{jt}$. Note that we impose no real treatment heterogeneity by using the same marginal effect (π) for all observations.
3. Imposing the assumptions of the probit model, we add standard normal randomly distributed error (v_i). To account for the impact of estimation errors in the latent index, we also add a random noise (v_i) with standard deviation $s_{X_i\hat{\beta}}$. We then form a “fake” binary outcome as

$$y_{i,sim} = \begin{cases} 1 & X_i\hat{\beta} + \pi\hat{\omega}_{jt} + v_i + v_i > 0 \\ 0 & otherwise \end{cases}$$

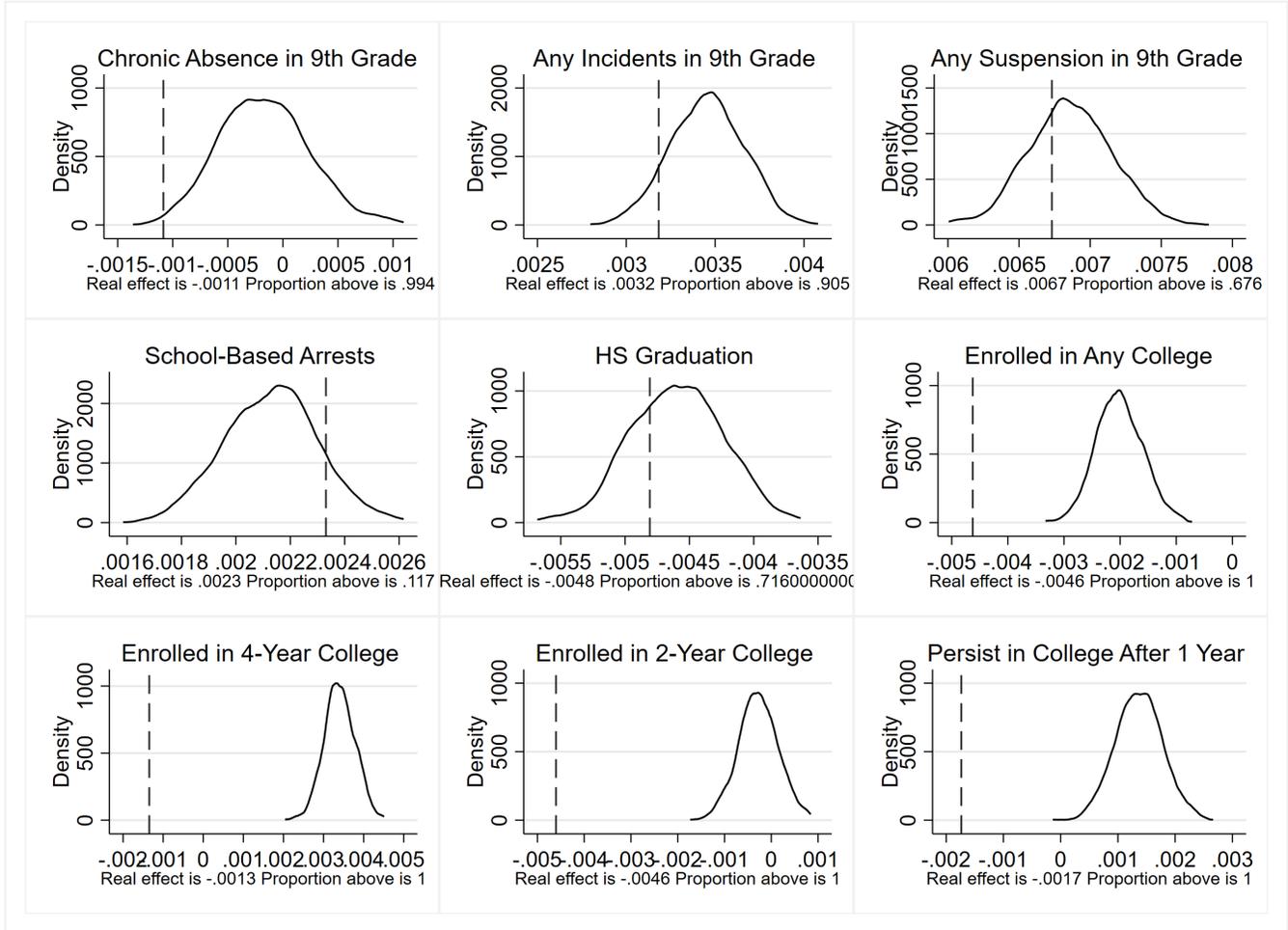
This fake binary outcome reflects what one would observe if the marginal treatment effects were constant on the latent variable.

4. We take δ_2 from estimation of Equation (7) on the simulated binary outcome.

We replicate these steps across 1000 random draws of both the estimation errors and the random errors to obtain a distribution of simulated coefficients on the interaction between the effectiveness and the decile of educational advantage under no treatment heterogeneity on the latent variable. If our estimated true heterogeneity (i.e., the slope of the interaction between school effectiveness and the educational advantage) is similar to the simulated distribution, it would imply that the observed heterogeneity is likely mechanical. Conversely, if the heterogeneity is greater than would be expected due to mechanical effects, it would be indicative of larger effects on the latent variable for less advantaged populations.

The plot of simulated effects is in Appendix Figure B2 below. As a frame of reference, we also plot the actual estimated slope and report the proportion of simulated slopes that lie above the real slope.

Figure B2. Simulated Distribution of Slopes Under Pure Mechanical Heterogeneity



Notes: This is the kernel density plot of the simulated slopes (i.e., the coefficient on the interaction between overall effectiveness and the decile of educational advantage) where only mechanical heterogeneity is imposed on the simulated data. We simulate “fake” binary outcomes based on 1000 random draws of the error term and estimation noise and then estimate equation 7. For each outcome, we plot the distribution of slopes under mechanical heterogeneity, and we report the proportion of simulated slopes that lie above the slope estimated on the real data.

Appendix C: Differences by Race and Gender and School Type

The summary statistics in Table 1 show that students in the bottom and top of the educational advantage distribution differ along both sex and ethnicity dimensions. As such, one may wonder if these patterns reflect gender or race differences, or if these are broad patterns that exist within demographic groups. To assess this, we implement analogous analyses using students from a particular group (males, females, black, Latinx). Appendix Figure A2 shows that the average outcomes by educational advantage within groups are similar to those overall so that the heterogeneity analysis within groups is meaningful. By and large, the patterns of results that we document across all groups exist within groups. As such, our results are not an artifact of making comparisons across sex or ethnic groups. There are, however, some differences we detail below.

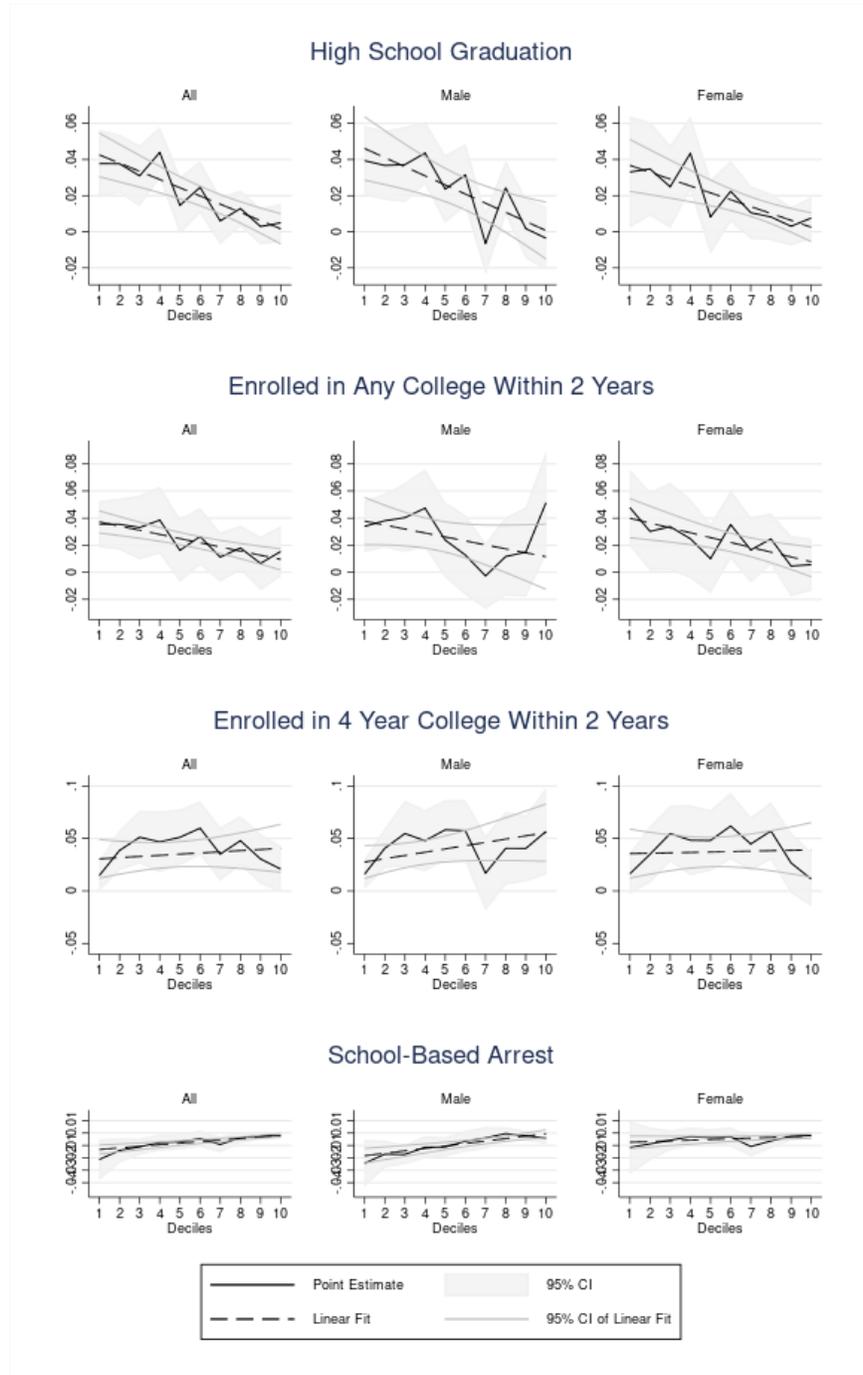
In Appendix Figure C1, the marginal effects are quite similar for both males and females; the average effects on high-school graduation are slightly larger for males, the effects on college going are slightly larger for females, and the effects on arrests rates of males somewhat larger for males than for females. Importantly, for both males and females, the less educationally-advantaged students experience larger marginal effects from attending a more effective school.

In Appendix Figure C2 we show effects for Black and Latinx students separately (other ethnic groups are too small to examine heterogeneous impacts). The arrest outcomes are much more sensitive to school effectiveness for Black students than Latinx students, while the educational attainment effects are particularly pronounced for Latinx students. In particular, among Black students in the bottom decile of the educational advantage distribution, a standard deviation increase in school effectiveness reduces the likelihood of a school-based arrest by about 3 percentage points (p -value <0.01), while that for Latinx students is about one percentage point. Looking at educational outcomes, for Latinx students in the bottom of the educational advantage distribution, a standard deviation increase in school effectiveness increases the likelihood of high-school graduation by over 5 percentage points (p -value <0.01), and in the middle of the distribution it increases the four-year college going rate by around 9 percentage-points.²⁹ The analogous numbers for Black students are around 3 percentage points for high school graduation and 1.5 percentage points for four-year college going. For both Black and Latinx students, the less educationally-advantaged students experience larger marginal effects from attending a more effective school. However, this heterogeneity by educational advantage is more pronounced for Latinx students – suggesting that less educationally-advantaged Latinx students may be particularly well-served by access to high-quality schools.

Since much evidence of differential school effectiveness is based on small samples of oversubscribed charter schools, one may wonder if our results hold only among traditional public schools. To assess this, we implement the entire analysis looking only at traditional public schools (Appendix Figure A3). The patterns we document are generally similar when restricted only to traditional public schools. This suggests that the patterns we document may generalize to other settings.

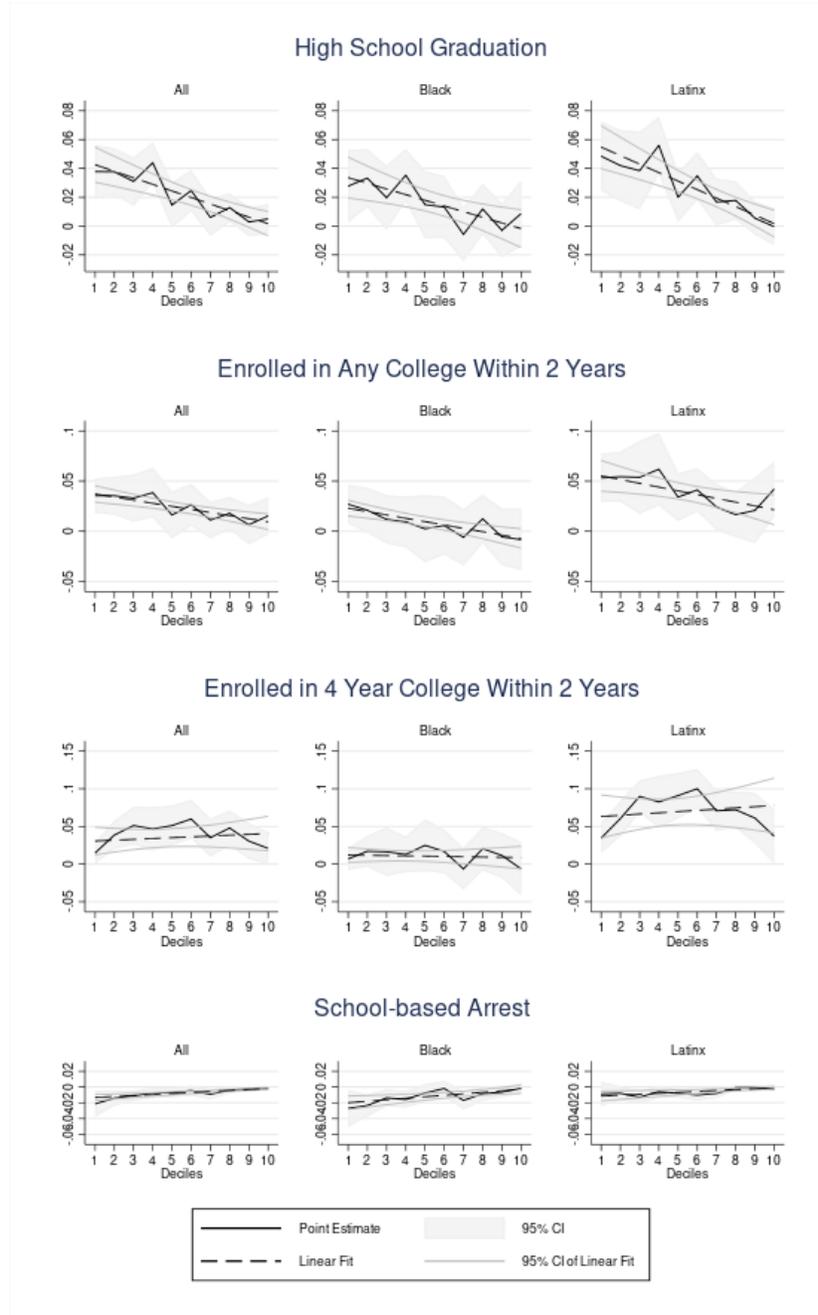
²⁹These relatively large college-going effects are consistent with Jackson (2014) finding particularly large college going responses among Latinx student to a college preparatory program in Texas.

Figure C1. Impacts on Outcomes: By Estimated Educational Advantage and Sex



Notes: Each graph represents the marginal impact of a 1 standard deviation increase in overall school effectiveness for different deciles of the educational advantage distribution for a single outcome and sub-population. Each panel presents the results of 10 separate regressions each defined as in Equation (6). The 95 percent confidence interval for each point estimate is depicted by the grey shaded area. The dashed black line in each panel depicts the line of best fit for the relationship between deciles of educational advantage and the marginal effect, including a 95 confidence interval.

Figure C2. Outcomes by Advantage: By Race / Ethnicity



Notes: Each graph represents the marginal impact of a 1 standard deviation increase in overall school effectiveness for different deciles of the educational advantage distribution for a single outcome and sub-population. Each panel presents the results of 10 separate regressions each defined as in Equation (6). The 95 percent confidence interval for each point estimate is depicted by the grey shaded area. The dashed black line in each panel depicts the line of best fit for the relationship between deciles of educational advantage and the marginal effect, including a 95 confidence interval.

Appendix D: Policy Implications Details

Distribution of Effectiveness by Advantage

Here, we explore the differences in school quality for those of different educational advantage. To this aim, we compute various percentiles of the school effectiveness index for students in each decile of the advantage index. This provides information about the extent of exposure to high-quality schools by educational advantage. We plot the percentiles for the deciles in Figure D1. One takeaway from this figure is that students of all educational advantage levels are exposed to schools that are both high and low on the effectiveness index. Indeed, the differences in school effectiveness within each decile (e.g., comparing the 5th to the 95th percentile of school effectiveness within a given educational advantage decile) are much larger than the differences in the same percentiles of effectiveness across educational advantage (e.g., comparing the 95th percentile of school effectiveness for the top and bottom deciles of educational advantage). However, there *are* economically significant differences across deciles. Looking across deciles of educational advantage, the most advantaged are exposed to more effective schools. Indeed, the 95th percentile of school effectiveness for the bottom and top deciles are about 1.4 and 2.4, respectively - a sizable 1σ difference.

Simulated School Closure Policy

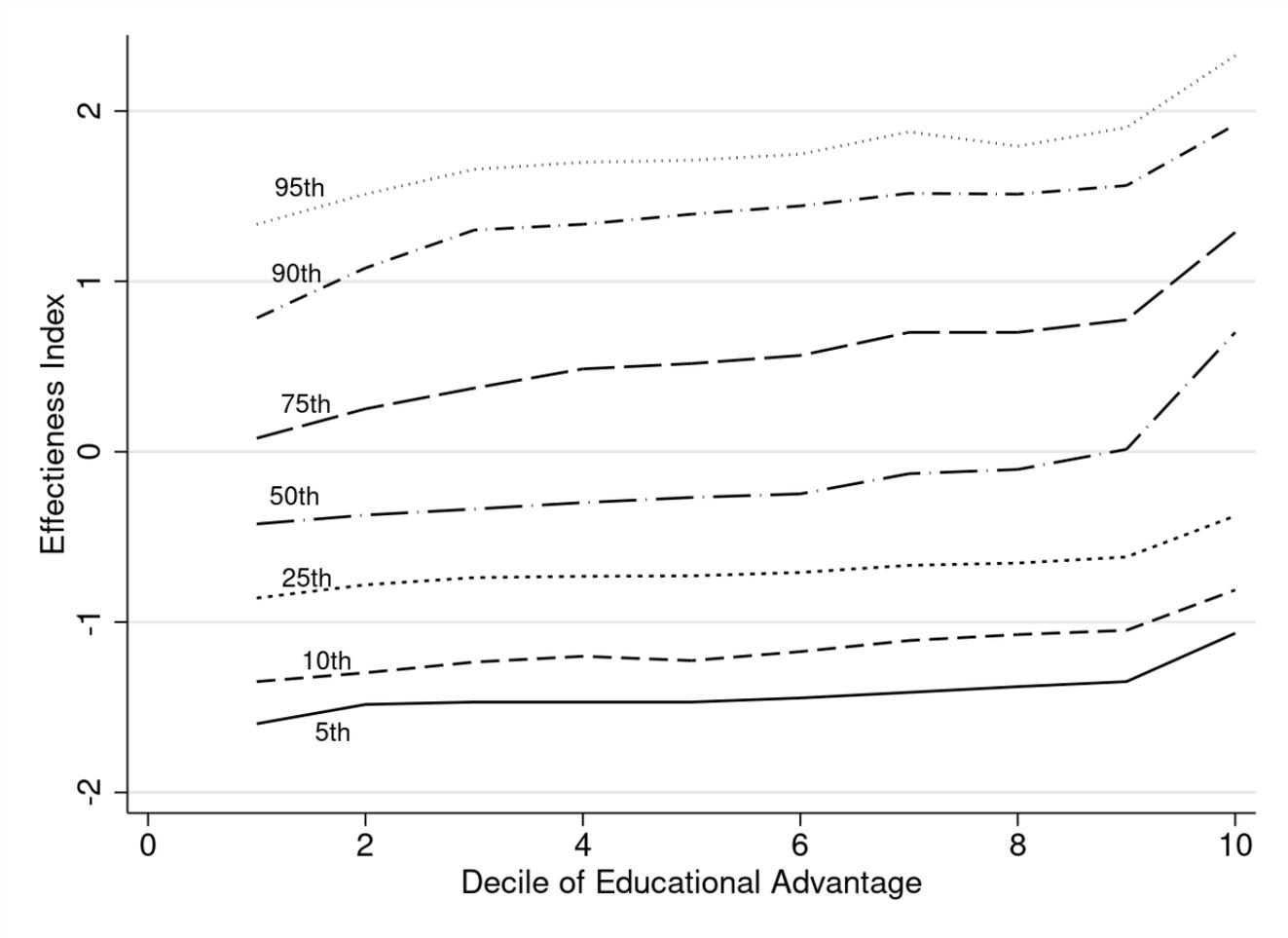
We simulate the effect of closing down the lowest performing eight schools (as measured by the effectiveness index) and then reassigning the previous enrollees to one of the ten nearest remaining schools. We define the direct effect as the change in estimated effectiveness associated with the move. For movers, this is simply the difference in estimated school effectiveness between the old low-performing school and the new higher-performing schools. For receiving students (i.e., those at other schools who are not reassigned), the direct effect is zero. We define the indirect effect of the policy as the change in school effectiveness that could be attributed to a change in peer composition. For this, we estimate the cross-sectional relationship between school effectiveness and the average educational advantage at the school. We find that when the average advantage increases by 1, school effectiveness declines by roughly 1.3σ . To define the indirect effect, we compute the change in average advantage due to the reassignment of students and then multiply this change by 1.3. The total effect reflects the sum of the direct and the indirect effects. We report the direct, indirect, and total effects for all students, movers, and receivers, by educational advantage decile in Appendix Table D1.

With the simulated change in effectiveness, we can calculate a simulated change in outcome by multiplying the changes in school effectiveness by the marginal effect of school effectiveness for the respective deciles from Appendix Table A8. The simulated policy impacts on high-school graduation, college going, and arrests are reported in Appendix Table D2.

Relationships without Surveys

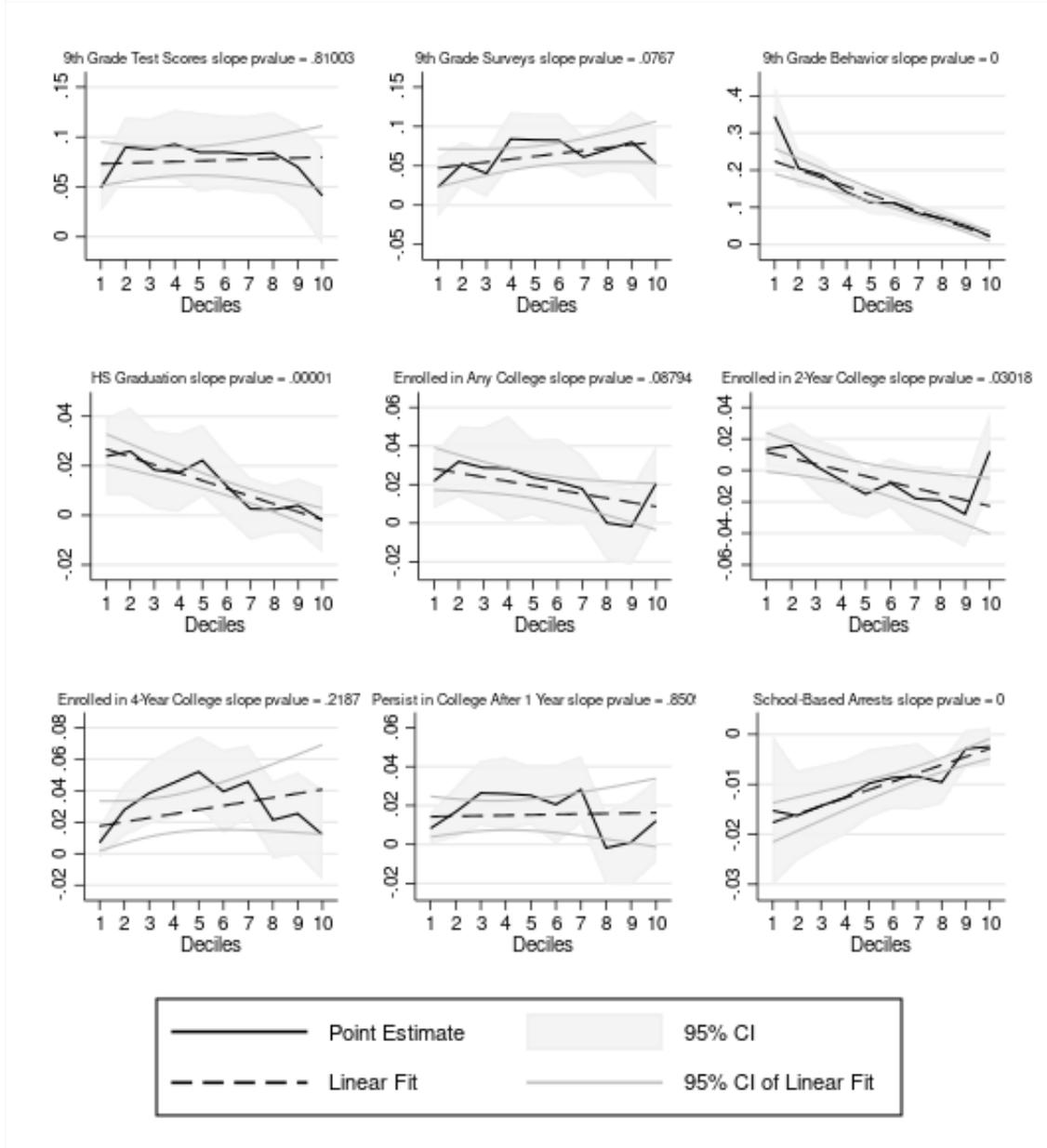
We replicate the analysis used to create Figure VI but where we estimate the overall effectiveness index only using test score and behaviors in Appendix Figure D2.

Figure D1. *Percentiles of Effectiveness Index: By Estimated Educational Advantage*



Notes: This figure plots the distribution of the overall effectiveness for students with different levels of educational advantage. The lines depict the 5th, 10th, 25th, 50th, 75th, 90th, and 95th percentile of school effectiveness that students in each decile are exposed to.

Figure D2. Main Results Without SED Effects



Notes: Each graph represents the marginal impact of a 1 standard deviation increase in overall school effectiveness based on school impacts on test scores and behaviors only (i.e., excluding survey-based measures of SED) for different deciles of the educational advantage distribution for a single outcome and sub-population. Each panel presents the results of 10 separate regressions, each defined as in Equation (6). The 95 percent confidence interval for each point estimate is depicted by the grey shaded area. The dashed black line in each panel depicts the line of best fit for the relationship between deciles of educational advantage and the marginal effect, including a 95 confidence interval.

Table D1: Simulated Change in School Effectiveness

Deciles	Total Effect			Direct Effect			Indirect Effect		
	All	Movers	Receivers	All	Movers	Receivers	All	Movers	Receivers
1	0.2332	1.9774	0.0051	0.1924	1.6638	0.0000	0.0407	0.3136	0.0051
2	0.1665	1.8291	-0.0012	0.1462	1.5955	0.0000	0.0203	0.2336	-0.0012
3	0.1489	1.8056	-0.0047	0.1356	1.5984	0.0000	0.0133	0.2072	-0.0047
4	0.1273	1.7451	-0.0070	0.1209	1.5781	0.0000	0.0063	0.1670	-0.0070
5	0.1262	1.8085	-0.0092	0.1205	1.6171	0.0000	0.0058	0.1913	-0.0092
6	0.1038	1.7526	-0.0123	0.1063	1.6172	0.0000	-0.0026	0.1353	-0.0123
7	0.1032	1.7593	-0.0148	0.1094	1.6446	0.0000	-0.0062	0.1147	-0.0148
8	0.0895	1.6602	-0.0185	0.1033	1.6067	0.0000	-0.0138	0.0535	-0.0185
9	0.0566	1.6024	-0.0255	0.0799	1.5835	0.0000	-0.0233	0.0190	-0.0255
10	0.0044	1.5611	-0.0402	0.0451	1.6192	0.0000	-0.0407	-0.0581	-0.0402

Notes: We obtain these effects by simulating a shutdown of the 10 worst-performing schools and reassigning students to one of the 10 nearest non-low performing schools. The direct effect is estimated by calculating the difference in school effectiveness of the mover's original school and their new school in the scenario of a simulated shutdown for each decile of educational advantage. The indirect effect is the effect of changes in peer composition. To calculate this, we regress schools' effectiveness on the average educational advantage of the students at the school. Then, we multiply the simulated change in peer composition by this coefficient. The total effect is the sum of the direct and indirect effects.

Table D2: Simulated Effect on Outcomes: By Mover Status and Advantage

Deciles	Graduation Effect			School-Based Arrests			Enrolled in Any College		
	All	Movers	Receivers	All	Movers	Receivers	All	Movers	Receivers
1	0.0088	0.0747	0.000191	-0.0050	-0.0422	-0.000108	0.0083	0.0702	0.000180
2	0.0063	0.0689	-0.000045	-0.0023	-0.0258	0.000017	0.0059	0.0649	-0.000043
3	0.0046	0.0559	-0.000146	-0.0017	-0.0200	0.000052	0.0049	0.0597	-0.000156
4	0.0056	0.0767	-0.000307	-0.0010	-0.0133	0.000053	0.0049	0.0672	-0.000269
5	0.0018	0.0264	-0.000134	-0.0009	-0.0126	0.000064	0.0021	0.0294	-0.000149
6	0.0026	0.0431	-0.000302	-0.0005	-0.0079	0.000055	0.0027	0.0464	-0.000324
7	0.0006	0.0106	-0.000089	-0.0009	-0.0162	0.000136	0.0012	0.0197	-0.000166
8	0.0012	0.0214	-0.000238	-0.0003	-0.0063	0.000070	0.0016	0.0299	-0.000333
9	0.0002	0.0046	-0.000074	-0.0001	-0.0037	0.000058	0.0004	0.0106	-0.000168
10	0.0000	0.0077	-0.000199	0.0000	-0.0034	0.000089	0.0001	0.0240	-0.000619

Notes: We obtain the above estimates by multiplying the total effect from Appendix Table D1 with the marginal effect of school effectiveness for each outcome for each decile (from Appendix Table A8). We calculate the effects for all students, movers, and receivers within each decile of educational advantage.