

IZA DP No. 7694

**Peer Effects in Disadvantaged Primary Schools:
Evidence from a Randomized Experiment**

Heather Antecol
Ozkan Eren
Serkan Ozbeklik

October 2013

Peer Effects in Disadvantaged Primary Schools: Evidence from a Randomized Experiment

Heather Antecol

*Claremont McKenna College
and IZA*

Ozkan Eren

Louisiana State University

Serkan Ozbeklik

Claremont McKenna College

Discussion Paper No. 7694
October 2013

IZA

P.O. Box 7240
53072 Bonn
Germany

Phone: +49-228-3894-0
Fax: +49-228-3894-180
E-mail: iza@iza.org

Any opinions expressed here are those of the author(s) and not those of IZA. Research published in this series may include views on policy, but the institute itself takes no institutional policy positions. The IZA research network is committed to the IZA Guiding Principles of Research Integrity.

The Institute for the Study of Labor (IZA) in Bonn is a local and virtual international research center and a place of communication between science, politics and business. IZA is an independent nonprofit organization supported by Deutsche Post Foundation. The center is associated with the University of Bonn and offers a stimulating research environment through its international network, workshops and conferences, data service, project support, research visits and doctoral program. IZA engages in (i) original and internationally competitive research in all fields of labor economics, (ii) development of policy concepts, and (iii) dissemination of research results and concepts to the interested public.

IZA Discussion Papers often represent preliminary work and are circulated to encourage discussion. Citation of such a paper should account for its provisional character. A revised version may be available directly from the author.

ABSTRACT

Peer Effects in Disadvantaged Primary Schools: Evidence from a Randomized Experiment

We examine the effect of peer achievement on students' own achievement and teacher performance in primary schools in disadvantaged neighborhoods using data from a well-executed randomized experiment in seven states. Contrary to the existing literature, we find that the average classroom peer achievement *adversely* influences own student achievement in math and reading in linear-in-means models. Extending our analysis to take into account the potential non-linearity in the peer effects leads to non-negligible differences along the achievement distribution. We test several models of peer effects to further understand their underlying mechanisms. While we find no evidence to support the monotonicity model and little evidence in favor of the ability grouping model, we find stronger evidence to support the frame of reference and the invidious comparison models. Moreover, we also find that higher achieving classes improve teaching performance in math. Finally, using a simple policy experiment we find suggestive evidence that tracking students by ability potentially benefits students who end up in a low achievement class while hurting students in a high achievement class.

JEL Classification: I21, J24

Keywords: peer effects, student achievement, random assignment

Corresponding author:

Heather Antecol
The Robert Day School of Economics and Finance
Claremont McKenna College
500 E. Ninth St.
Claremont, CA 91711
USA
E-mail: hantecol@cmc.edu

1. Introduction

Throughout the past several decades there has been a nation-wide debate focusing on how to improve student achievement in the United States. The debate was fueled by the influential Coleman Report of 1966 which questioned the long-standing belief that school funding was a key determinant of student achievement.¹ The report instead highlighted the importance of alternative determinants—e.g., family background and socio-economic status, teacher quality, and peer quality—which could have differential effects on students in schools in economically disadvantaged neighborhoods relative to students in schools in more economically advantaged neighborhoods (Coleman 1966). Not surprisingly, the report spawned a flurry of new research among social scientists, as well as a shift in policy-makers' education goals. However, it continues to be the case that there is little, if any, agreement about which specific education policies are more effective in improving student achievement (for reviews of the literature see Hanushek 2006; Hanushek and Rivkin 2006).

One education policy that has received recent attention is improving teacher quality in K-12 education. Specifically, President Obama argued at the Urban League's 100th Anniversary Convention on July 29, 2010 that "...teachers are the single most important factor in a child's education from the moment they step in to the classroom."² As such, one of the main goals of the Obama Administration's education policy, Race to the Top (RTTT), is to improve student achievement by focusing on enhancing teacher quality and accountability, particularly in

¹ This debate has gained further momentum in the last decade as American students' test score outcomes continue to lag behind their counterparts from many other developed countries (Fleischman et al. 2010)

² <http://www.whitehouse.gov/the-press-office/remarks-president-education-reform-national-urban-league-centennial-conference>

disadvantaged neighborhoods with lower achieving students.³ How can teacher quality and performance be improved in primary schools in disadvantaged neighborhoods? We argue that one potential avenue that can help shed some light on this question is to further our understanding of how peers influence student achievement and teacher performance in these disadvantaged schools. This avenue of study should not only help pinpoint which types of classrooms enhance teacher performance in disadvantaged neighborhoods but also help policymakers assess the effectiveness of complementary policies such as tracking and/or racial/ethnic and economic desegregation.

Any study attempting to measure the causal effect of peer quality on student achievement has to deal with two important identification issues. First, it is a well-known fact that students are not randomly assigned to schools or classrooms largely because of families, school administrators, or teachers (see for example, Clotfelter, Ladd, and Vigdor 2006; and Kane et al. 2011). This identification issue is often referred to as the selection problem (Sacerdote 2001). Failure to account for non-random sorting of students in a regression framework would result in biased coefficient estimates of peer effects as there are likely to be observable and unobservable factors that affect both student achievement and peer quality. Second, it is often difficult to disentangle the impact that the peer group has on the student from the impact the student has on the peer group. A regression of, say, own achievement on contemporaneous average achievement of peers is problematic as these outcomes are jointly determined and peer achievement is likely to be endogenous in the model. This is usually referred to as the endogeneity or the reflection problem (Manski 1993; Moffitt 2001; Sacerdote 2001).

³ Additional goals of RTTT include “state success factors”, “standards and assessment”, “data systems to support instruction”, “general selection criteria”, and STEM (U.S. Department of Education, 2009).

The existing literature on peer effects and student outcomes in grades K through 12 (henceforth referred to as the existing literature) generally relies on panel or repeated cross-sectional data sets and uses within-school/grade variation in achievement or some other school/grade characteristics to measure peer effects (see for example, Hanushek et al. 2003; Vigdor and Nechyba 2007; Ammermueller and Pischke 2009; Lavy, Silva, and Weinhardt 2012; Lavy, Paserman, and Schlosser 2012; and Burke and Sass 2013) to overcome the threats to identification (i.e., under the assumption that the within-school [or grade] variation is random, the coefficient estimates on peer measures produce reliable estimates of peer effects). We are aware of only three studies (Hoxby and Weingarth 2006; Duflo, Dupas, and Kremer 2011; and Imberman, Kugler, and Sacerdote 2012) that overcome the threats to identification by exploiting random assignment of students to groups and generally use pre-determined outcomes as measures of peer effects.⁴ Moreover, the existing literature generally focuses on schools in one state and/or school district.⁵ Finally, with the exception of Hoxby and Weingarth (2006), Duflo, Dupas, and Kremer (2011), Imberman, Kugler, and Sacerdote (2012), and Lavy, Paserman, and Schlosser (2012), the existing literature has primarily focused on documenting the existence of peer effects as opposed to identifying the underlying mechanisms behind how peer effects might work.

The findings based on linear-in-means models in the existing literature are mixed at best. While some studies find positive and significant effects of average peer achievement on students' own achievement (see for example, Hoxby 2000; Boozer and Cacciola 2001; Hanushek

⁴ In a correctly specified natural/randomized experiment absent of sorting, the results from an OLS regression would yield causal peer effects.

⁵ For example, Burke and Sass (2013) examine public school students in grades 3-10 in Florida; Hoxby and Weingarth (2005) examine the Wake County Public School district in North Carolina; Angrist and Lang (2004) examine primary school students in Boston; Betts and Zau (2004) examine the San Diego Unified School District in California. In contrast, Imberman, Kugler and Sacerdote (2012) examine primary schools in two states: the Houston Independent School District in Texas and Louisiana.

et al. 2003; Betts and Zau 2004; Hoxby and Weingarth 2006; Vigdor 2006; Vigdor and Nechyba 2007; Ammermueller and Pischke 2009; Carman and Zhang 2012; and Lavy, Paserman, and Schlosser 2012), others find small to no effects (see for example, Angrist and Lang 2004; Lefgren 2004; Lavy, Silva, and Weinhardt 2012; Imberman, Kugler, Sacerdote 2012; and Burke and Sass 2013). The common perception from several of these studies is that it is not only the high ability students but also those at the bottom of the achievement distribution who seem to benefit from higher achieving peers. With that said, the peer effects estimates are not identical across different achievement groups and the impacts generally exhibit nonlinearities with no consensus on who benefit the most from better peers (see for example, Imberman, Kugler, and Sacerdote 2012; Lavy, Silva, and Weinhardt 2012; and Burke and Saas 2013).^{6,7}

The purpose of this paper is to contribute to the existing literature by examining the effect of peer achievement on students' own achievement in primary schools in *disadvantaged neighborhoods* in seven states. We use data from a well-executed randomized experiment which helps us avoid the aforementioned selection problem, and thus allows us to measure the causal effect of peer quality on student achievement. Moreover, we analyze non-linearities in peer effects, as well as perform detailed sub-group analyses (i.e., student gender, race/ethnicity, and free-lunch eligibility status). Following Hoxby and Weingarth (2006) and Imberman, Kugler, and Sacerdote (2012), we also explicitly examine how peer effects might work. Specifically, we focus on four potential models of peer effects: the monotonicity model (i.e., the effects of peers

⁶ There is also a large literature examining the effect of peers on student outcomes in college. The results from these studies are again mixed. Studies either find small positive effects (Sacerdote 2001; Zimmerman 2003), large positive effects (Stinebrickner and Stinebrickner 2006; Carrell, Fullerton, and West 2009), or no effects (Foster 2006; Lyle 2007). Moreover, there are a number of recent studies that examine peer effects in labor markets (see for example, Arcidiacono and Nicholson 2005; Bandiera, Barankay and Rasul 2005; Falk and Ichino 2006; Mas and Moretti 2009; Guryan, Kroft, and Notowidigdo 2009; and Brown 2011) and on social and behavioral outcomes (see for example Case and Katz 1991; Gaviria and Raphael 2001; Ludwig, Duncan and Hirsfield 2001; and Kling, Ludwig, and Katz 2005).

⁷ For a detailed review of the empirical peer effects literature see Sacerdote (2011) and of the theoretical peer effects literature see Epple and Romano (2011).

on student achievement is increasing in peer quality); the invidious comparison model (i.e., higher ability peers adversely influence the outcomes of students who are moved to a lower position in the local achievement distribution), the ability grouping (boutique) model (i.e., student performance is highest when their peers are similar to themselves), and the frame of reference model or the reverse big fish in a little pond model (i.e., higher ability peers adversely influence the outcomes of students due to a lower academic self-concept).⁸ Finally, we examine the influence of peers on teachers' performance at the classroom-level which to date has received almost no attention in the existing economics literature largely due to data constraints.⁹

Contrary to the existing literature, we find that the average classroom peer achievement *adversely* influences own student achievement irrespective of subject or group in linear-in-means models, although the effect is imprecisely estimated for certain subgroups. Extending our analysis to take into account the potential non-linearity in the peer effects leads to non-negligible differences along the achievement distribution. Focusing first on reading test scores, we find that an improvement in average peer quality for the full sample substantially hurts students both at the bottom and top of the achievement distribution but does not seem to affect middle ability students. The subgroup patterns for students in all achievement groups essentially mirror those found for the full sample, however the effects are estimated more precisely for certain subgroups (particularly at the top of the distribution). Turning to math test scores, we find negative effects of average peer quality for the full sample over the entire achievement distribution, although the coefficient estimates are imprecisely estimated. The patterns for the full sample appear to extend

⁸ The frame of reference model has received little attention in the existing economics literature. The main exceptions are Pop-Eleches and Urquiola (2013) and Johnson and Mood (2008) who examine this effect among secondary students in Romania and Sweden, respectively. Both studies find evidence supporting the predictions of the model. Using data from gifted and talented classrooms, Bui, Craig and Imberman (2012) also find some evidence supporting the invidious comparison and/or the frame of reference models.

⁹ We are aware of one Israeli study, by Lavy, Paserman, and Schlosser (2012), that focuses on the association between peer effects and teachers' pedagogical practices in the classroom.

only to male students at all achievement levels and Hispanic students for the middle achievement group for whom we observe negative (and significant) peer effects. Viewing the complete set of results, we find no evidence to support the monotonicity model and little evidence in favor of the ability grouping model, while we find stronger evidence to support the frame of reference and the invidious comparison models. We also find that higher achieving classes improve teaching performance (i.e., lesson implementation and classroom culture) in math. Finally, using a simple policy experiment we find suggestive evidence that tracking students by ability potentially benefits students who end up in a low achievement class while hurting students who end up in a high achievement class.

2. Data and Tests for Random Assignment

2.1 Data

We use data from the Mathematic Policy Research, Inc. (MPR) Teachers Trained Through Different Routes to Certification (TTTDR) Private Use File. TTTDR is a randomized study of primary school students, which was conducted to assess the effectiveness of different teacher certifications on student outcomes. MPR began in 2003 by identifying as many schools with alternatively certified (AC) teachers as possible where AC teachers are those who become a classroom teacher prior to completing all required coursework and without having to complete a period of student teaching.¹⁰ In order to be eligible for the study, (i) schools had to have had at least one alternative certification (AC) and one traditional certification (TC) teacher in the same grade (i.e. kindergarten through grade 5); (ii) both AC and TC teachers had to have had five or

¹⁰ The AC programs differ on the selectivity criteria of their admission requirements. For instance, AC programs such as the Teach for America require a minimum GPA of 3.0 from the applicants. The AC teachers in the TTTDR sample come from programs with less selective entrance requirements by design as this maintained a fairer comparison between AC and TC teachers. We further note that the TTTDR study did not find any difference in the end of the academic year test scores between students taught by AC and TC teachers.

fewer years of experience, and (iii) both AC and TC teacher must have taught in regular classes and must have delivered both math and reading instruction to all their own students.¹¹ Among a compiled list of 170 eligible schools, a random sample of 60 schools, which included 90 AC and 90 TC teachers and more than 2,800 students, were selected in seven states between 2004 and 2006.^{12,13}

This data is ideal for our purposes because, within each school, all students in the same grade were randomly assigned to either an AC or a TC teacher before the start of the academic year. Therefore, the randomization is done at the block level such that each block represents classrooms in the same grade level in any given school. This process not only ensured that those students in AC and TC classrooms are comparable but also that the pre-treatment achievement of own students and the average pre-treatment achievement of their peers in each classroom are not correlated (this is discussed in further detail in Section 2.2).

After the random assignment and before the start of the academic year, the students were given math and reading tests based on the grade they completed in the previous year (which we call *baseline* outcome variables); then at the end of the academic year in which the study was conducted the students re-took math and reading tests based on the grade they just completed (which we call *endline* outcome variables). We use Normal Curve Equivalent (NCE) points in math and reading as our measures of baseline and endline test scores.¹⁴ The NCE scale has a mean of 50 and standard deviation of 21 nationally.

¹¹ Even though the requirements for teachers who pursue alternative routes to certification vary by state and district, the AC programs, on average, require significantly less education coursework than TC programs (see Constantine et al. 2009 for more details on AC and TC teachers).

¹² Due to the confidential nature of the data agreement, the sample sizes are rounded to the nearest tenth.

¹³ The states included in the TTTDR sample are California, Georgia, Illinois, Louisiana, New Jersey, Texas and Wisconsin. There were 20 school districts in the effective sample; 5 districts from California, 7 districts in total together from Georgia, Illinois, Louisiana and Wisconsin, 3 districts from New Jersey and 5 districts from Texas.

¹⁴ The students were administered two reading tests (reading comprehension and vocabulary). The sum of scores from these two tests establishes total reading score and our measure of student achievement in reading. There

The sample attrition in TTTDR data set is relatively small, but we still lose roughly 7 (8) percent of the initial reading (math) sample because of missing test scores.¹⁵ After dropping these observations, our estimation sample consists of 2,610 (2,580) students for the reading (math) test score sample from classes taught by 180 teachers. To ensure the student composition was unaffected by the sample attrition, Constantine et al. (2009) show the attrition rates in the AC and TC samples were almost identical and did not differ significantly between the two types of classrooms (Appendix A, pp. A13, Table A3 in Constantine et al. 2009).¹⁶

Besides test scores and the type of classroom (AC or TC classroom) the data set also contains information on the student's gender, race/ethnicity, and eligibility for free lunch. Table 1 present some features of the TTTDR student sample. Specifically, 34.5 (47.0) percent of the student body is black (Hispanic), while 9.2 percent is white.¹⁷ Moreover, students tend to come from low income families; roughly 75 percent of the effective sample is eligible for free lunch as opposed to 40 percent nationwide. Finally, the average baseline test scores for reading and math are roughly 39 and 42 NCE points for the full sample, respectively. Compared to the national average, the reading (math) scores are roughly 0.5 (0.4) of a standard deviation lower in the TTTDR sample. Overall, it is evident that the TTTDR sample consists of lower achieving students from disadvantaged neighborhoods.

The second and third columns of Table 1 report the average baseline characteristics of students in AC and TC classrooms, respectively. Under the assumption that the random

were also two different tests in math (math concepts and applications and math computation). Unlike reading, however, students in kindergarten and grade 1 were not administered math computation test. Thus our measure in math achievement is scores from math concepts and applications only.

¹⁵ Students test scores are missing either because they moved out of school district or they did not take endline tests.

¹⁶Ideally, we would like to run a regression of the non-response indicator on an AC classroom dummy along with the baseline characteristics. Even the restricted version of the data set, however, does not include any information on those moving out of the school district and on students not taking the test.

¹⁷ The remaining 9.3 percent of the student body indicated "other" race. The survey instrument does not provide details on what this category includes.

assignment is implemented correctly, baseline characteristics of students in AC and TC classrooms must be similar. To test this, as in Krueger and Whitmore (2001), we run a regression of the AC indicator variable on each baseline characteristic conditional on block fixed effects (the dependent variable taking the value of one if the student is in a AC classroom and zero otherwise). The fourth column of Table 1 displays the coefficient estimates from this exercise. None of the coefficient estimates are statistically significant at conventional levels. By the nature of randomization in TTTDR, it is important to note that we include the block fixed effects in all of our specifications throughout the paper. The use of conditional randomization is a very common practice in the education literature (see for example, Sacerdote 2001; Carrell, Fullerton, and West 2009; and Duflo, Dupas and Kremer 2011).

Finally, TTTDR includes information on teacher characteristics including gender, race/ethnicity, teaching experience, hours of instruction for certification, and SAT Composite Score. Not surprisingly given our sample is comprised of primary schools, roughly 90 percent of teachers are female (see Column 1 of Appendix Table A1), however AC teachers are less likely to be female relative to their TC counterparts (see Columns 2 and 3 of Appendix Table A1). TC teachers are less racially/ethnically diverse than their AC counterparts. Specifically, roughly 72 (45) percent of TC (AC) teachers are white. By construction AC and TC teachers have similar levels of teaching experience, roughly 3 years. Finally, TC teachers have roughly 2 times more teaching training than their AC counterparts, although this difference is somewhat less pronounced for math.

2.2 Are Peers Randomly Assigned?

Although we have shown some preliminary evidence on the random assignment of students within two types of classrooms, it is imperative for the purpose of our study to validate the random assignment of peers (absence of sorting) within blocks. A typical test for this is to run an OLS regression of student i 's pre-determined achievement on the pre-determined average achievement of i 's peers, controlling for any variable on which randomization was conditioned on and is given by

$$TS_{icb}^{base} = \pi_0 + \pi_1 \overline{TS}_{-i,cb}^{base} + \eta_b + u_{icb} \quad (1)$$

where TS_{icb}^{base} is the subject-specific baseline test score for student i in classroom c and block b , $\overline{TS}_{-i,cb}^{base}$ is the average peer baseline subject-specific test score in classroom c and block b excluding student i , η_b is a set of block fixed effects (i.e., classrooms in the same grade level in any given school), and u_{icb} is the error term. Under the assumption that peers are randomly assigned, one would expect the estimate of π_1 to be equal to zero. This approach is the common practice in the peer effects literature to test for randomization (see for example, Sacerdote 2001; Foster 2006; and Carrell, Fullerton, and West 2009).

Guryan, Kroft, and Notowidigdo (2009) however recently showed that the mechanical relationship between own ability and average ability of peers (i.e., peers of high achieving students are chosen from a block with a slightly lower mean achievement than peers of low achieving students) may cause the aforementioned falsification exercise to produce negative and

statistically significant coefficient estimates for π_1 . Random assignment may not appear random, while positive sorting of students to classrooms may appear random.¹⁸

Given the bias is a by-product of the differences in the average achievement level of the group once the student i is withdrawn, the proposed solution in Guryan, Kroft, and Notowidigdo (2009) is to control for this relevant group mean in the falsification regressions. Specifically, the revised falsification test equation is given by

$$TS_{icb}^{base} = \pi_0 + \pi_1 \overline{TS}_{-i,cb}^{base} + \pi_2 \overline{TS}_{-i,b}^{base} + \eta_b + u_{icb} \quad (2)$$

where $\overline{TS}_{-i,b}^{base}$ is the mean achievement of students in block b and all other variables are as previously defined. Using simulations, Guryan, Kroft, and Notowidigdo (2009) show that equation (2) is a well-behaved randomization test and if the student assignment to classrooms is truly random, we would expect the coefficient estimate $\hat{\pi}_1$ to be equal to zero.

Table 2 presents our results from various falsification tests. The first and second columns of Table 2 report the results from estimating equation (1) for baseline reading and math test scores, respectively, while the third and fourth columns report the results from estimating equation (2). In the absence of correction, the correlation between own and peers' baseline achievement is negative and statistically significant for both the reading and math test score samples. As previously noted, however, ignoring the bias in the randomization tests leads to the erroneous conclusion that students are negatively sorted within each block. Once we do the correction the coefficient estimates on average peer baseline test scores are insignificant and almost equal to zero in magnitude irrespective of subject.¹⁹

¹⁸ It is also important to note that as the size of the randomization group (block in our case) grows the contribution of each student to the average ability goes down and the magnitude of the bias from the falsification exercise is also reduced. The average class and block sizes in our study are 15.1 and 32, respectively.

¹⁹ The coefficients from the falsification test remain intact when we control for classroom type (AC or TC classroom).

3. Empirical Methodology and Results

3.1 Empirical Methodology

Having shown that students are randomly assigned to classrooms, we now turn to the estimation of peer effects on student achievement. To begin with, we first analyze the peer effects using linear-in-means models, where we regress endline test scores on average peer baseline test scores along with students' own baseline scores and block fixed effects. In a randomized experiment setting, it is a well-known fact that controlling for the baseline characteristics does not affect the consistency of the estimates; however, it helps increase efficiency (Frölich and Melly, 2013). To this end, we estimate the following equation for the full sample and by subgroups (i.e., student gender, student race/ethnicity, and student free lunch eligibility status):

$$TS_{icb}^{end} = \beta_0 + \beta_1 \overline{TS}_{-i,cb}^{base} + SC'_{icb} \delta + TC'_{cb} \gamma + \eta_b + e_{icb} \quad (3)$$

where TS_{icb}^{end} is the subject-specific endline test score for student i in classroom c and block b ; SC is a set of student characteristics (i.e., baseline subject-specific test scores, gender, race/ethnicity, and free lunch status), TC is a set of teacher characteristics (i.e., AC/TC status, gender, race/ethnicity and years of teaching experience), $\overline{TS}_{-i,cb}^{base}$ and η_b are as previously defined.²⁰ We also estimate two versions of equation (3) for the full sample and by subgroups to address the potential non-linearity in the peer effects. The first version is given by

$$E(TS_{icb}^{end} | Q_k^{base}) = \beta_0 + \beta_1 \overline{TS}_{-i,cb}^{base} + SC'_{icb} \delta + TC'_{cb} \gamma + \eta_b + e_{icb} \quad (4)$$

where Q_k^{base} is the student i 's grade and subject-specific baseline achievement quartile k ($k = \text{top } 25\%$; middle 25-75%; bottom 25%) and all remaining variables are as previously defined. We

²⁰ It is important to note that peer effects estimates from equation (3) are reduced form in the sense that equation (3) does not separately identify the effects of peers outcomes (endogenous effects) and peers' background characteristics (contextual effects).

estimate equation (4) separately for each quartile. The second version specifies a slightly different measure of peer quality and the estimation equation is given by

$$E(TS_{icb}^{end} | Q_k^{base}) = \beta_0 + \beta_{k,bottom} P_{-i,cb}^{bottom} + \beta_{k,top} P_{-i,cb}^{top} + SC_{icb}' \delta + TC_{cb}' \gamma + \eta_b + e_{icb} \quad (5)$$

where $P_{-i,cb}^{bottom}$ and $P_{-i,cb}^{top}$ represent the fraction of bottom 25% and top 25% of peers in classroom c and block b , respectively, based on the grade and subject-specific baseline test score distribution. Due to collinearity of the proportions, we omit the proportions of middle ability peers in each specification of equation (5). All other variables are defined as previously. Finally, we report the standard errors clustered at the block-level beneath each coefficient estimate.²¹

We examine four potential models through which peer effects might work. First, we examine the monotonicity model which implies that the effect of peers on student achievement is increasing in peer quality. Using equation (5) we test for two versions of the monotonicity model: weak monotonicity states $\beta_{k;top} > \beta_{k;bottom}$ and strong monotonicity states $\beta_{k;top} > \beta_{k;middle}$ and $\beta_{k;middle} > \beta_{k;bottom}$. for $k = top; middle; bottom$.

The second model we examine is the invidious comparison (proposed in Hoxby and Weingarth 2006) which states that higher ability peers adversely influence the outcomes of students who are moved to a lower position in the local achievement distribution, perhaps because of a fall in their self-esteem. Note that this model does not say anything about the impact of peers at the same ability level. Using equation (5) we test for the invidious comparison model as follows: for $k = top; middle; bottom$ and $j = top; middle; bottom$, the invidious comparison model states $\beta_{kj} < 0$ for $j > k$ and $\beta_{kj} > 0$ for $k > j$ where j denotes the grade and subject-specific baseline achievement quartile of peers.

²¹ The results are similar, although slightly more precisely estimated, if we instead cluster at the classroom-level and are available upon request.

The third model we examine is the ability grouping (boutique) (proposed in Hoxby and Weingarth 2006) which states that student performance is highest when their peers are similar to themselves. Using equation (5) we test for the ability grouping model as follows: for $k = top; middle; bottom$ and $j = top; middle; bottom$, the ability grouping model states $\beta_{kk} > \beta_{kj}$ for $j \neq k$.

The final model we examine relies on social comparison theory and frame of reference (Marsh and Parker 1984; Marsh 1987). In an educational setting, the theoretical model underlying the frame of reference mechanism states that students compare their own academic achievement with the achievement of peers and use this social comparison for forming their own academic self-concept where academic self-concept is defined as one's knowledge and perceptions about one's academic ability. In this context, academic self-concept depends not only on one's own achievement but also on the achievement of a reference group. Consider a high achieving student in a regular classroom is assigned to a gifted classroom; the student in this new environment may become an average student relative to his peers. According to Marsh and Hau (2003), this then can have adverse effects on the student's academic self-concept as he is no longer “a big fish in a small pond” (regular class) but is now “a little fish in a big pond” (gifted class). According to the frame of reference model, academic self-concept will be affected positively with individual achievement but will also be negatively affected by the average achievement of the reference group. Thus the frame of reference model predicts a negative impact of an improvement in peers' achievement on student's own achievement.

Taking this a step further, if the proportion of peers in the top (bottom) 25% increases, then average peer achievement must improve (decline) which will result in a negative (positive) impact on own student achievement. In other words, the frame of reference model predicts that all students are hurt from high-achieving peers and benefit by low-achieving peers (i.e., the

inverse of the monotonicity model). We test models of both the weak and strong frame of reference as follows: for $k = top; middle; bottom$, the weak frame of reference states

$\beta_{k;top} < \beta_{k;bottom}$ and for $k = top; middle; bottom$, the strong frame of reference states

$\beta_{k;top} < \beta_{k;middle}$ and $\beta_{k;middle} < \beta_{k;bottom}$.²²

3.2 Results

3.2.1 Linear-in-Means Results

Column 1 of Panel A and B of Table 3 present our linear-in-means estimations for reading test scores and math test scores, respectively, for the full sample. Specifically, the coefficient estimate on average classroom peer baseline reading achievement is negative and statistically significant (-0.18); a one standard deviation increase in peer achievement is associated with roughly one-ninth of a standard deviation decrease in own endline reading scores. Similarly, the coefficient estimate on average classroom peer baseline math achievement is (-0.24) suggesting that a one standard deviation increase in peer achievement decreases math test scores by around one-ninth of a standard deviation as well.²³

²² Hoxby and Weingarth (2006) outline a number of other potential models including the bad apple model (i.e., one disruptive student has a detrimental effect on the outcomes of all students irrespective of where they are in the achievement distribution); the shining light model (i.e., one excellent student has a positive effect on the outcomes of all students irrespective of where they are in the achievement distribution); the focus model (i.e., homogeneous classrooms are good irrespective of student i 's ability relative to their homogeneous peers); and the rainbow model (i.e., heterogeneous classrooms benefit all students).

²³ For both reading and math test scores, we find very similar results if we exclude teacher characteristics only or if we exclude both student and teacher characteristics. Moreover, we also estimate a specification that adds average block baseline peer achievement to determine whether it is within classroom interactions as opposed to interactions at the school (block) level that matter for student achievement. We find that it is indeed within classroom interactions that matter given adding average block baseline peer achievement does not affect the average classroom baseline peer effects coefficient estimate nor does average block baseline peer achievement have a direct impact on own student's achievement for both math and reading math test scores. Finally, we also tried adding average peer reading and math achievement scores simultaneously to subject-specific achievement equations. Average peer effect coefficients for both subjects on endline scores remain almost intact. They are, however, less precisely estimated. This may not be surprising given the high correlation between these two peer achievement measures. These results are available upon request.

To further explore these findings, we extend our analysis to test for the presence of heterogeneous effects along a number of dimensions. We first focus on student gender. While we continue to find a negative effect of average classroom peer achievement, irrespective of subject or student gender, the magnitude of the peer effect is substantially larger for male students and is imprecisely estimated for female students (see Column 1 of Tables 4 and 5 for female and male students, respectively). These gender differences may stem from the fact that female students tend to be more cooperative and more level even in the presence of ability differences with their peers (see for example, Croson and Gneezy 2009 and Bertrand 2010).

Our next set of results pertains to student free-lunch status, which proxies for family income. We noted earlier that roughly 75 percent of the effective sample is eligible for free-lunch and students from wealthier families are disproportionately distributed at the top of the achievement group. This leaves us with only a limited number of observations at the bottom quartile for students that are not eligible for free-lunch which is particularly important when we allow for non-linearities. As such, we focus our discussion on free-lunch eligible students only. We again find that, irrespective of subject, that average classroom baseline peer achievement adversely influences own endline test scores (see Column 1 of Tables 6).

Our final set of results pertains to student race/ethnicity. As previously noted only 9 percent of our estimation sample comprises white students. Such a small sample prevents us from making a rigorous inference and therefore we only focus on black and Hispanic students. For both black and Hispanic students, we continue to find a negative effect of average classroom peer achievement for reading and math test scores, although the effects are imprecisely estimated (see Column 1 of Tables 7 and 8 for black and Hispanic students, respectively).

Overall, our linear-in-means results sharply contrast with the general consensus in the existing literature of positive (or no) peer effects (see Sacerdote 2011 for a review of the earlier literature and see Imberman, Kugler, Sacerdote 2012 and Burke and Sass 2013 for more recent evidence). One potential explanation for the differences in our results may stem from the fact that we are focusing on students from disadvantaged neighborhoods and peer interactions may differ by socio-economic status and family background (see for example, Ludwig, Duncan and Hirsfield 2001). An alternative potential explanation pertains to the nature of our data. While the existing peer effects literature on grades K through 12 meticulously addresses the identification problems inherent in measuring the causal effect of peer quality on student achievement, usually using grade and school level variation, they may be unable to fully account for all potential confounding effects using survey data. Our use of a well-executed random experiment may then explain the divergence in the results.

3.2.2 Nonlinearities in Peer Effects

In the previous section, we assume the peer effects are linear. There is, however, substantial evidence against the linear-in-means model (see for example, Hoxby and Weingarth 2006; Sacerdote 2011; and Imberman, Kugler, and Sacerdote 2012). Moreover, from a policy point of view, if the peer effects were to be linear, there would be no gain or loss in sorting and tracking students. To examine the potential nonlinearities in peer effects, we first estimate the effect of average classroom subject specific baseline peer achievement on own subject specific endline

test scores separately by the grade and subject-specific baseline achievement quartile k ($k = \text{top 25\%}; \text{middle 25-75\%}; \text{bottom 25\%}$) for student i (see equation 4 in Section 4.1).²⁴

Focusing first on the full sample estimation, we observe a negative and significant impact of average classroom peer baseline reading achievement for students at the bottom quartile of the achievement distribution (-0.45); a one standard deviation increase in peer achievement is associated with roughly one-fourth of a standard deviation decrease in own reading test scores (see Column 3, Panel A of Table 3). The coefficient estimate on peer effects for the middle achievement group, on the other hand, is almost equal to zero in magnitude (-0.06) and is insignificant (see Column 5, Panel A of Table 3), while the effect for the students in the top quartile is negative (-0.26) although imprecisely estimated at conventional levels (see Column 7, Panel A of Table 3). Pair-wise comparisons indicate that the peer effects coefficient for the lowest achievement group is significantly different than the one for middle achievement group (p-value 0.02). Turning to the math test score results (see Column 3, 5, and 7 of Panel B of Table 3), the coefficient estimates are negative and similar in magnitude for all achievement groups, although imprecisely estimated at conventional levels. We fail to reject the null of equality across all pair-wise comparisons of peer effects coefficient estimates.

The patterns for the full sample generally extend to all subgroups under consideration. As such, we only highlight what we believe to be the most interesting findings here (see Columns 3, 5, and 7 of Tables 4-8). For male students, the peer effects for reading is large and more precisely estimated at the top of the achievement distribution (see Column 7, Panel A of Table 5). Similarly, the peer effect for math is larger in all three quartiles of the distribution, although the effect continues to be imprecisely estimated in the top quartile (see Columns 3, 5, and 7,

²⁴ We discuss alternative cut-off points to describe the bottom and top achievement groups in Section 3.2.5. However, we are unable to examine cut-off points based on the top 5 (10) % and bottom 5 (10) % due to data limitations (i.e., our sample size does not allow us to cut the data that finely).

Panel B of Table 5). For free-lunch eligible students, an improvement in class reading achievement hurts the students at the top quartile the most (see Columns 3, 5, and 7, Panel A of Table 6). Moreover, it appears that an increase in peer math quality decreases the achievement level of Hispanic students at the middle quartile (see Column 5, Panel A of Table 8). Finally, we observe a positive coefficient estimate for Hispanic students who are at the top quartile of the math score distribution but this coefficient estimate is imprecisely estimated (see Column 7, Panel A of Table 8).

To further delve into the complexity of the effect of peers on student achievement we replace average classroom baseline achievement with the fraction of bottom 25% and top 25% of peers in classroom c , respectively, based on the grade and subject-specific pre-treatment test score distribution. Due to collinearity, we omit the proportion of middle ability peers in each regression (see equation 5 in Section 4.1). We first replace average peer achievement with these fractions irrespective of student i 's placement in the grade-subject specific baseline achievement distribution (Column 2 of Tables 3-8). We then also allow student i 's placement in the grade-specific baseline achievement distribution to vary as well (Columns 4, 6, and 8 of Tables 3-8). We focus on the results for the full sample and discuss any significant deviations from these results for the subgroups under consideration in the remainder of this section.

If we hold student i 's placement in baseline grade and subject-specific achievement distribution fixed (see Column 2 Table 3), we find that a 1 percentage point increase in the proportion of peers in the top quartile (therefore the proportion of peers in the middle quartile decreases by 1 percentage point) is associated with a 0.048 (0.018) points decrease (increase) in endline reading (math) test scores, although the effects are imprecisely estimated.²⁵ Similarly, if

²⁵ We find that if the proportion of peers in the top quartile increases relative to the middle quartile there is a significant decrease in student i 's endline reading test scores for male students (see Column 2, Panel A of Table 5).

the proportion of peers in the bottom quartile increases by 1 percentage point relative to the middle quartile, then endline reading (math) test scores go up by 0.047 (0.115) points, the effect however is statistically insignificant at conventional levels for reading test scores.

If we now also allow student i 's placement in the baseline grade and subject specific achievement distribution to vary, for students in the bottom quartile we find that if you increase the proportion of peers in the top (bottom) quartile by 1 percentage point relative to the middle quartile, then endline reading tests scores decrease (increase) by 0.06 (0.14) points (see Column 4, Panel A of Table 3), the effect however is statistically insignificant at conventional levels for the top peer quartile relative to the middle quartile.²⁶ For students in the top quartile we find that if you increase the proportion of peers in the top (bottom) quartile by 1 percentage point relative to the middle quartile, then endline reading tests scores decreases (increases) by 0.05 (0.06) points (see Column 8, Panel A of Table 3), although the effects are imprecisely estimated.²⁷ Finally, for students in the middle quartile the effect of changing the proportion of peers in the top (bottom) quartile relative to the middle quartile is much smaller in magnitude and statistically insignificant. The patterns for endline math test scores essentially mirror those found for endline reading test scores with the following exception. Unlike endline reading tests scores, for students in the bottom (middle) quartile we find that if the proportion of peers in the top quartile increases relative to the middle quartile there is an increase, not a decrease, in student i 's endline math test

²⁶ For black students (see Column 4, Panel A of Table 7) however we find that a 1 percentage point increase in the proportion of peers in the top quartile relative to the middle is associated with a 0.15 point increase in student i 's endline reading test scores (the effect is insignificant at conventional levels).

²⁷ For female students (see Column 8, Panel A of Table 4) however we find that if the proportion of peers in the bottom of the distribution increases relative to the middle there is a decrease in student i 's endline reading test scores, although the effects are imprecisely estimated.

scores (see Columns 4 and 6, Panel B of Table 3), although the effects are imprecisely estimated.²⁸

3.2.3 Potential Mechanisms for Peer Effects

Thus far, we have focused on estimating the peer effects on student achievement. An equally important question is the potential mechanisms through which peers in the classroom affect student outcomes. Understanding the underlying mechanism for peer effects is crucial in designing optimal educational policies as different channels may imply different policy prescriptions (Lavy, Paserman, and Schlosser 2012).

To shed light on this, we test the predictions of four models of peer effects outlined in Section 3.1: the monotonicity model (i.e., the effects of peers on student achievement is increasing in peer quality); the invidious comparison model (i.e., higher ability peers adversely influence the outcomes of students who are moved to a lower position in the achievement distribution), ability grouping (boutique) model (i.e., student performance is highest when their peers are similar to themselves), and frame of reference model (i.e., higher ability peers adversely influence the outcomes of students due to a lower academic self-concept).

Specifically, based on the non-linear results for the full sample where both own student and peer effects are allowed to vary by placement in the subject specific achievement distribution (i.e., Columns 4, 6, and 8 of Table 3), we count the number of tests that are significant at the 10% level in the direction predicted by the model under question, the number of tests that are

²⁸ For black students (see Column 8, Panel A of Table 7) and Hispanic students (see Column 8, Panel A of Table 8) in the top quartile we find that a 1 percentage point increase in the proportion of peers in the top quartile relative to the middle is associated with a 0.16 and 0.05 point increase, respectively, in student i 's endline math test scores and for Hispanic students (see Column 8, Panel A of Table 8) in the top quartile we find that a 1 percentage point increase in the proportion of peers in the bottom quartile relative to the middle is associated with a 0.08 point decrease in student i 's endline math test scores. The effects however are insignificant at conventional levels.

significant at the 10% level in the opposite direction predicted by the model under question, and the number of tests that are statistically insignificant. This is very similar to the inference procedure in Imberman, Kugler, and Sacerdote (2012).²⁹ We take any tests indicating significant in the correct (opposite) direction to be consistent (inconsistent) with the model. If there are no significant tests, irrespective of the direction predicted by the model, we find no support for the model. The results are presented in Table 10.

We find evidence against both the weak and strong monotonicity models for reading test scores. For math test scores while we also find evidence against the strong monotonicity model we do not find evidence for or against the weak monotonicity model (i.e., no tests are significant in the direction of or in the opposite direction the model suggest). For instance, consider the tests for the strong monotonicity model, 1 (2) out of 6 tests are significant in the opposite direction the model suggests for reading (math) test scores, while there are no significant tests in the direction predicted by the model. The evidence with respect to the ability grouping model is mixed. Specifically, we find evidence against the ability grouping model based on math test scores (i.e., 1 out of 6 tests are significant in the opposite direction the model suggests) and evidence for the ability grouping model based on reading test scores (i.e., 2 out of 6 tests are significant in the direction the model suggests). We do however find stronger support for both the invidious comparison model and the frame of reference model. Specifically, we never find significant tests in the opposite direction predicted by either invidious comparison model or the frame of reference models. For the invidious comparison model, we find 0 (2) out of 4 tests are significant in the direction the model suggests for reading (math) test scores. For the weak (strong) frame of reference model, 1 (1) out of 3 (6) tests predict the model in the correct

²⁹ We similarly test the peer effect models for each sub-group (see Appendix Table A2). The patterns for the subgroups mostly coincide with those of the full sample.

direction for reading test scores and 0 (2) out of 3 (6) tests predict the model in the correct direction for math test scores.

We also offer some suggestive evidence to further support the patterns above. Specifically, if we break up the insignificant results into those that predict the model in the correct and the incorrect direction, we find that all the tests go in the opposite direction predicted by the weak monotonicity models for both reading and math test scores. While the same is true for the strong monotonicity model for reading, for math all but 2 of the 6 tests go in the opposite direction predicted by the model. For the invidious comparison model we find that 4 (2) out of the 4 tests go in the direction predicted by the model for reading (math) test scores. For the ability grouping model we find that 3 (4) out of 6 tests go in the opposite direction predicted by the model for reading (math) test scores. It is also important to note that the evidence for the ability grouping model only comes from the bottom quartiles of the reading achievement distribution both in terms of sign and statistical significance. Finally, for the weak frame of reference model all the tests go in the direction predicted by the model irrespective of test subject while for the strong frame of reference model all (4) of the 6 tests go in the direction predicted by the model for reading (math) test scores.³⁰

Taken together, we appear to find no support for the weak and strong monotonicity models and little evidence in favor of the ability grouping model. However, we appear to find stronger support for the invidious comparison and frame of reference models. Our test results appear to be at odds with Imberman, Kugler, and Sacerdote (2012) who use the inflow of student

³⁰ While we get similar results based on reading test scores if we allow for interdependence between tests (i.e., do the Bonferroni correction), for math test scores the results with the Bonferroni correction do not allow us to either reject or support any of the peer effect models (see Appendix Table A3). However, we argue that this is largely an artifact of small sample sizes combined with the marginal level of significance found for the math test scores. Moreover, it is a well known fact that the Bonferroni tests are too conservative and may lack power for correlated tests.

evacuees from Hurricane Katrina into the non-evacuated school districts in Louisiana as their measure of peer quality and examine the effects of evacuees with a wide range of achievement on incumbent achievement. While the authors find evidence for both weak and strong monotonicity models, they reject the ability grouping and invidious comparison models. One potential explanation for the conflicting findings is that we focus on primary school students whereas Imberman, Kugler, and Sacerdote (2012) focus on middle and high school students. The peer dynamics may differ at different levels of education. An alternative and more plausible explanation is that the sample in Imberman, Kugler, and Sacerdote (2012) is less disadvantaged and more closely mirrors the nation as a whole compared to our sample. For instance, roughly 54 (75) % of the Imberman, Kugler, and Sacerdote (2012) sample (our sample) is eligible for free-lunch compared to 40% nationwide and 54% of their sample is white compared to 9% of our sample.

Finally our findings may also shed some light on the channels through which peer effects in primary school operate. It may be the case that high quality peers in the classroom depresses the academic performance of all students presumably through the frame of reference model or the invidious comparison model. In this case, negative peer effects result from the interactions across students. Alternatively, as noted in Lavy, Paserman, and Schlosser (2012), an increase in the overall achievement of the classroom may force teachers to raise the level more towards higher ability students and this may hurt some students. The negative peer effects here result from a change in teaching practices and methods.³¹ With our data, it is not possible to directly see whether the teacher raises the level of teaching more towards higher ability students or not. That said, however, if the peer effects were to stem from changes in teachers' pedagogical

³¹ Burke and Saas (2013) argue that the effect of peers in their analysis based on public school students in Florida (1999-2005) largely stems from changes in the level of teaching towards the higher ability students as opposed to interactions across students.

practices, one would not observe any negative peer effects at the top of the achievement distribution as we do. Stated somewhat differently, a teacher raising the bar so high that even students in the top quartile are hurt does not appear to be a valid explanation because for this to be the case we would also expect to observe negative peer effects for middle achievement students and we do not.

3.2.4 Robustness Checks

We undertake several sensitivity checks to examine the robustness of our results. First, following Foster (2006) we replace the average baseline peer achievement with the median baseline achievement level of the class and re-run all the specifications. The results from this exercise are qualitatively similar to those presented in the paper.³² Second, to examine the potential differential effects of peer quality at different grades, we divide the sample into lower grades (kindergarten and first grade) and upper grades (second through fifth). The peer effect coefficient estimates from the lower versus the upper grades do not indicate any discernible pattern. Third, rather than splitting the sample within each achievement group based on selected student characteristics, we run fully interacted models (e.g., female dummy interacted with all covariates) within each achievement group. We also repeated a similar exercise within each subgroup and run fully interacted models in ability (i.e., baseline achievement indicators-top, middle, and bottom-interacted with all covariates). The precision of our results from these robustness checks are very similar to those presented in the paper. Finally, we choose different cut-off points to describe the bottom and top achievement groups (i.e., one-third). Doing so does not alter our conclusions. All these results are available upon request.

³² We also tried adding the standard deviation of baseline peer achievement along with average baseline peer achievement. The coefficient estimates on this additional measure are not different than zero in any of the specifications.

3.2.5 Peer Effects and Teacher Performance

We take our analysis one step further and focus on the potential relationship between peer achievement and teacher performance. Aside from the impact peers may have on the outcomes of their classmates, they may also influence teacher performance. Understanding this relationship is particularly important in light of the Obama Administration's commitment to improve student achievement by focusing on enhancing teacher quality and accountability, particularly for students in disadvantaged neighborhoods with lower achieving students.

To the best of our knowledge, this is the first study that examines the relationship between peer achievement and teacher performance at the classroom level using U.S. data.³³ Specifically, the TTTDR data set includes a measure of teacher performance that is based on the Vermont Classroom Observational Tool (VCOT).³⁴ Specifically, the survey administrators directly observed teachers teaching two reading and two math lessons and they rated their performance using VCOT. We use two domains of VCOT: (i) lesson implementation (i.e., designed to measure the use of best practices and teacher pacing), and (ii) classroom culture (i.e., designed to measure clarity and consistency of classroom routines). The first domain consists of five questionnaires while the second domain consists of seven questionnaires. Each questionnaire is measured on a five-point scale ranging from no evidence to extensive evidence. A larger value implies a better rating of the teacher.

To examine the association between teachers' ratings and peer achievement, we estimate the following two equations:

³³ In this context, the closest study to ours is Lavy, Paserman, and Schlosser (2012). Using Israeli data and the proportion of grade repeaters as their peer measure, the authors examine the association between peer effects and high school students' opinions about their teacher performance. The study finds negative effects of the proportion of grade repeaters on teachers' pedagogical practices, as well as on classroom environment.

³⁴ The TTTDR data also includes a measure of teacher performance based on the ratings of each school principal. We argue that this measure is more likely to be subjective and biased. Moreover, there is no information on when this rating is performed (i.e., prior to or after the endline test). As such, we do not use this measure.

$$Rating_{cb} = \delta_0 + \delta_1 \overline{TS}_{-i,cb}^{base} + TC'_{cb} \lambda + \eta_b + \varepsilon_{cb} \quad (6)$$

and

$$Rating_{cb} = \delta_0 + \delta_1 P_{-i,cb}^{bottom} + \delta_2 P_{-i,cb}^{top} + TC'_{cb} \lambda + \eta_b + \varepsilon_{cb} \quad (7)$$

where $Rating_{cb}$ is the subject-specific VCOT rating of the teacher (in class c and block b) (standardized at the grade level to have a mean of zero and standard deviation of one) and all other variables are defined as previously.

Columns 1 and 3, Panel A (Panel B) of Table 11 presents the estimated effect of peers on teacher performance based on reading (math) VCOT ratings using average peer baseline reading (math) test scores as our measure of peer achievement [see equation 6] for lesson implementation and classroom culture, respectively. Irrespective of subject-specific peer measure, we find a positive effect of average peer quality on teaching performance, although the effect is imprecisely estimated for the reading VCOT rating.

Columns 2 and 4, Panel A (Panel B) of Table 11 present the estimated effect of peers on teacher performance based on reading (math) VCOT ratings replacing average peer baseline reading (math) test scores as our measure of peer achievement with the proportion of peers in the top 25% and the bottom 25% of the grade baseline reading achievement distribution [see equation 7] for lesson implementation and classroom culture, respectively. Our results show no apparent difference between any of the quartiles if we base our peer achievement measure on baseline reading and our VCOT rating on reading instruction. If, however, we base our peer achievement measure on baseline math and our VCOT rating on math instruction, we find some evidence suggesting that if the proportion of peers in the top 25% is increased relative to the middle peer quartile there is a positive effect on teacher performance irrespective of VCOT rating measure but there is no apparent difference between the bottom and middle quartile.

Overall, the evidence in this section suggests that a higher achieving classroom is likely to improve the way the teacher delivers instruction in disadvantaged neighborhoods. In addition, the positive impact of peer achievement on classroom management/culture may suggest less crowding out of teacher resources in terms of time and energy which may then lead to a more conducive environment of learning.

3.2.6 Policy Experiment

What might the consequences of manipulating peer effects on student achievement be? To answer this question we conduct a simple policy exercise similar to that of Lavy, Silva, and Weinhardt (2012). In particular, we take an average class in our sample and rank students in the class from the best to the worst based on their baseline reading (math) test scores. We then group all of the students with below-median reading (math) test scores in a “low achievement” class and the remaining students in a “high achievement” class. The *direct effect* of our tracking policy is the proportion of students in the bottom quartile of the baseline reading (math) achievement doubles from 24 (25) percent to 48 (50) percent in the low achievement reading (math) class and decreases from 24 (25) percent to 0 in the high achievement reading (math) class and the proportion of students in the top quartile of the baseline reading (math) achievement decreases from 26 (25) percent to 0 in the low achievement reading (math) class and doubles from 26 (25) percent to 52 (50) percent in the high achievement reading (math) class. We also consider the *indirect effect* of the tracking policy through the effect of peers on teacher performance which itself is very likely to be correlated with student achievement.

In doing this policy exercise we are making the following three somewhat strong assumptions: (i) the peer effect dynamics in the new classrooms do not change; (ii) peer effects

are the same for each student within an achievement group, but differ across each achievement group; and (iii) the effect of teacher performance on student test scores is homogenous regardless of any given student's baseline achievement. It is important to note however that the predictions of a large policy intervention like the one we are proposing should be viewed with caution given the potential unintended consequences (i.e, endogeneous sorting) of the intervention (see Carell, Sacerdote, and West 2013). We discuss the impacts of our tracking policy on low achievement and high achievement classes in turn.

3.2.6.1 Low achievement class

Let us begin with calculating the *indirect effect* of our tracking policy in the low achievement reading and math classes, respectively. For reading test scores, we can safely conclude that the indirect effect of the tracking policy is zero given the coefficients in Column 2, Panel A of Table 11 (reading lesson implementation) almost cancel each other out. For the math test scores, based on the results presented in Column 2, Panel B of Table 11 (math lesson implementation), the doubling of students in the bottom quartile from 25 percent to 50 percent hurts teacher performance in math by -0.075 (25×-0.003) of a standard deviation. Moreover, the reduction of students in the top quartile from 25 percent to zero also decreases the teacher performance by -0.550 (-25×0.022) of a standard deviation. Thus, the total indirect effect of the tracking policy on teacher performance in math is a decline of -0.625 of a standard deviation.

Turning to the *direct effect* of our tracking policy, we find that for the students in the bottom quartile, the 24 (25) percentage point increase in their proportion in the new class *increases* each students' reading test scores by 3.36 (0.14×24) points and math test scores by 2.925 (0.117×25) points (see Column 4, Panels A and B of Table 3, respectively). Moreover,

since the proportion of students in the top quartile of the baseline reading (math) achievement changes from 26 (25) percent to 0 in the low achievement class, the reading test scores of each student in the bottom quartile increases by 1.56 points (-0.06×-26) and the math test scores of each student in the bottom quartile decreases by 1.45 (0.058×-25). Therefore, the total *direct effect* of tracking on each student in the bottom quartile is 4.92 (1.475) points, which is a gain of 0.25 (0.07) of a standard deviation in reading (math) test scores. The effect of tracking for students in the middle achievement group is a gain of 1.282 (0.925) points in reading (math) achievement (see Column 6, Panel A (B) of Table 3), or roughly 0.06 (0.046) of a standard deviation.

The *overall effect* of the tracking policy needs to incorporate both the *indirect* and *direct* effects. For reading achievement the *overall effect* of the tracking policy is roughly equal to the *direct effect* and results in an improvement in the reading achievement of all students in the low achievement class since *the indirect effect* of the tracking policy through its impact on teacher performance is essentially zero. For math achievement, we need to take into account that the *indirect effect* of the tracking policy hurts teacher performance in the low achievement class by 0.625 of a standard deviation. Specifically, to eliminate the positive *direct effect* of tracking on each student in the bottom and middle quartiles in the low achievement class, *the indirect effect* of teacher performance on student math achievement, P , must be equal to 0.112 of a standard deviation ($-0.07 = -0.625 \times P$) and 0.07 of a standard deviation ($-0.046 = -0.625 \times P$), respectively. Therefore, for this tracking policy to be a Pareto optimum in the low achievement math class, *the indirect effect* of teacher performance on student math test scores should be less than 0.07 of a standard deviation. If the *indirect effect* is between 0.07 and 0.112 of a standard deviation, only

the students in the bottom quartile benefit from the policy. If it is larger than 0.112 of a standard deviation, then tracking hurts all students in the low achievement math class.

3.2.6.2 High Achievement Class

Since we assume a homogeneous impact of teacher performance on students irrespective of their achievement group, the *indirect effect* in the high achievement class is zero and 0.625 of a standard deviation improvement in reading and math, respectively. The total *direct effect* of our tracking policy on reading (math) test scores is a *decrease* of -2.764 (-4.8) points, or about -0.14 (-0.24) of a standard deviation for each student in the top quartile of the baseline reading (math) achievement distribution. The total *direct effect* of our tracking policy is a *decrease* of -1.282 (-0.7) points, or about -0.06 (-0.046) of a standard deviation for each student in the middle quartile of the baseline reading (math) achievement distribution (see Column 8, Panels A and B of Table 3).

As with the low achievement class, the *overall effect* is roughly equal to the *direct effect* and our tracking policy *harms* all students in the high achievement reading class. For high achievement math class, for our tracking policy to be a Pareto optimum, the effect of teacher performance on student math test scores should be greater than 0.38 of a standard deviation. If the indirect effect is in the range of 0.07 and 0.38, only the students in the middle of the math achievement distribution benefit from the policy. If the effect is smaller than 0.07, tracking hurts all students in the high achievement class.

What should we expect the effect of teacher performance on student math achievement, P , to be in our data? To answer this, we run a regression of teachers' math implementation scores on students' endline math test scores, controlling for our usual set of variables except block fixed

effects. We exclude block fixed effects to obtain an upper bound estimate of the math implementation coefficient. We find that a one standard deviation increase in teacher performance (math implementation score) increases math test scores by 0.06 of a standard deviation. If we take this 0.06 as our benchmark estimate, the low ability math class appears to benefit from the tracking policy while the high ability math class appears to be hurt.

4. Conclusion

For decades, there has been a flurry of research by social scientists trying to pinpoint the underlying determinants of student achievement, particularly since the Coleman Report was released in 1966. Despite this, we still know very little about the impact of specific education policies on student achievement outcomes. This paper further analyzes how to improve student achievement in disadvantaged neighborhoods with a particular interest on the effect of peers on both student achievement and teacher performance. This focus should not only help pinpoint which types of classrooms enhance teacher performance in disadvantaged neighborhoods—which is a key goal of the Obama Administration's education policy, *Race to the Top* (RTTT)—but also whether policy makers should consider complementary policies such as tracking and/or any demographic (e.g., gender, racial/ethnic, and economic) desegregation.³⁵

We use data from a well-executed randomized experiment which allows us to measure the causal effect of peer quality, as well as affords us a large sample of primary schools, students, teachers, and states. Furthermore, our data comes from a disadvantaged part of the student population which allows us to take a closer look at peer effects in a setting where the influence of family background may be particularly less pronounced and the problems with the education

³⁵ In a companion paper, we examine in details the contextual peer effects (see Antecol, Eren, and Ozbeklik 2013).

system in the United States are most evident and arguably more important from a policy perspective.

Unlike the existing literature which generally finds positive and significant effects or small positive to no effects, we find that the average classroom baseline peer achievement *adversely* influences student's own endline achievement in the linear-in-means regressions. The linear-in-means model, however, masks a great deal of information. We therefore extend our analysis to take into account non-linearities in peer effects which reveals substantial heterogeneity across the achievement distribution. Specifically, we consistently find negative peer effects at the bottom of the reading achievement distribution for the full sample, as well as for all subgroups, although some effects are imprecisely estimated. We also find negative peer effects at the top of the reading achievement distribution for the full sample and all subgroups however for certain subgroups the effects are more precisely estimated (particularly at the top of the distribution). Peer effects estimates on reading achievement for middle ability students, on the other hand, are essentially zero for the full sample and across all subgroups of interest. Turning to math test scores, we find that average peer quality adversely affects student achievement for the full sample over the entire achievement distribution, although the effects are imprecisely estimated. The full sample pattern for math achievement is essentially mirrored for two subgroups only. Specifically, we observe negative (and significant) peer effects for male students at almost all achievement levels and for Hispanic students in the middle achievement group. Taken altogether, direct peer effects as opposed to teacher responses to student compositional changes appear to be driving our results.

In an attempt to further understand the underlying mechanisms behind our peer effect results, we test several peer effect models. While we find no evidence to support the

monotonicity model and little evidence in favor of the ability grouping model, we find stronger evidence to support the frame of reference and the invidious comparison models. We also find that higher achieving classes improve teaching performance in math in disadvantaged neighborhoods. Finally, we perform a simple policy experiment where we rank students in an average class and group them into new low achievement and high achievement classes based on their baseline achievement. We find that this tracking policy appears to benefit students in the low achievement class and to hurt the students in the high achievement class.

References

- Ammermueller, Andreas, and Jorn-Steffen Pischke. 2009. "Peer Effects in European Primary Schools: Evidence from PIRLS," *Journal of Labor Economics*, 27(3): 315-348.
- Angrist, Joshua D., and Kevin Lang. 2004. "Does School Integration Generate Peer Effects? Evidence from Boston's Metco Program," *American Economic Review*, 94(5): 1613-1634.
- Antecol, Heather, Ozkan Eren, and Serkan Ozbeklik. 2013. "Gender and Racial Composition of Peers and Student Achievement in Disadvantaged Primary Schools," Unpublished Manuscript.
- Arcidiacono, Peter, and Sean Nicholson. 2005. "Peer Effects in Medical School," *Journal of Public Economics*, 89(2-3): 327-350.
- Bandiera, Oriana, Iwan Barankay, and Imran Rasul. 2005. "Social Preferences and the Response to Incentives: Evidence from Personnel Data," *Quarterly Journal of Economics*, 120(3): 917-62.
- Bertrand, Marianne. 2010. "New Perspectives in Gender," in *Handbook of Labor Economics* Volume 4B, eds. Orley Ashenfelter and David Card, 1545-1592, Elsevier.
- Betts, Julian R., and Andrew Zau. 2004. "Peer Groups and Academic Achievement: Panel Evidence from Administrative Data," Unpublished Manuscript.
- Boozer, Michael A., and Stephen E. Cacciola. 2001. "Inside the Black Box of Project Star: Estimation of Peer Effects," Economic Growth Center Discussion Paper No. 832.
- Brown, Jennifer. 2011. "Quitters Never Win: The (Adverse) Incentive Effects of Competing with Superstars," *Journal of Political Economy*, 119(5): 982-1013.
- Bui, Sa A., Steven G. Craig and Scott A. Imberman. (2012). "Is Gifted Education a Bright Idea? Assessing the Impact of Gifted and Talented Programs on Students," Unpublished Manuscript.
- Burke, Mary A., and Tim R. Saas. 2013. "Classroom Peer Effects and Student Achievement," *Journal of Labor Economics*, 31(1): 51-82.
- Carrell, Scott E., Richard L. Fullerton, and James E. West. 2009. "Does Your Cohort Matter? Measuring Peer Effects in College Achievement," *Journal of Labor Economics*, 27(3): 439-464.
- Carell, Scott E., Bruce I. Sacerdote, and James E. West. 2013. "From Natural Variation to Optimal Policy? The Importance of Endogenous Peer Group Formation," *Econometrica*, 81(3): 855-882.
- Carman, Katherine G., and Lei Zhang. 2012. "Classroom Peer Effects and Academic Achievement: Evidence from a Chinese Middle School," *Chinese Economic Review*, 23, 223-237.

Case, Anne, and Larry F. Katz. 1991. "The Company You Keep: The Effects of Family and Neighborhood on Disadvantaged Youths," NBER Working Paper No. 3705.

Clotfelter, Charles T., Helen F. Ladd, and Jacob L. Vigdor. 2006. "Teacher-Student Matching and the Assessment of Teacher Effectiveness," *Journal of Human Resources*, 41(4): 778-820.

Coleman, James.S. 1966. "Equality of Educational Opportunity," Washington: U.S. Government Printing Office, 1966 [summary report].

Constantine, Jill, Daniel Player, Tim Silva, Kristin Hallgren, Mary Grider, John Deke, and Elizabeth Warner. 2009. "An Evaluation of Teachers Trained Through Different Routes to Certification: Final Report," NCEE 2009-4043. Institute of Education Sciences. Department of Education.

Crosen, Rachel, and Uri Gneezy. 2009. "Gender Differences in Preferences," *Journal of Economic Literature*, 47(2): 1-27.

Darling-Hammond, Linda, Deborah J. Holtzman, Su J. Gatlin, and Julian V. Heilig. 2005. "Does Teacher Preparation Matter? Evidence about Teacher Certification, Teach for America, and Teacher Effectiveness," *Education Policy Analysis Archives*, 13(42): 1-51.

Decker, Paul.T., Daniel P. Mayer, and Steven Glazerman. 2004. "The Effects of Teach for America on Students: Findings from a National Evaluation," Mathematica Policy Research. Report 8792-8750, New York.

Duflo, Esther, Pascaline Dupas and Michael Kremer. 2011. "Peer Effects, Teacher Incentives, and the Impact of Tracking: Evidence from a Randomized Evaluation in Kenya," *American Economic Review*, 101(5): 1739-1774.

Duncan, Greg J., Johanne Boisjoly, Michael Kremer, Dan M. Levy, and Jacque Eccles. 2005. "Peer Effects in Drug Use and Sex among College Students," *Journal of Abnormal Child Psychology*, 33(3): 375-385.

Epple, Dennis, and Richard E. Romano. 2011. "Peer Effects in Education: A Survey of the Theory and Evidence," in *Handbook of Social Economics* Vol. 1, eds. Jess Benhabib, Alberto Bisin and Matthew O. Jackson, 1053-1163. Amsterdam, North-Holland.

Falk, Armin, and Andrea Ichino. 2006. "Clear Evidence on Peer Effects," *Journal of Labor Economics* 24(1): 39-57.

Fleischman, Howard L., Paul J. Hopstock, Pelczar, Marisa P. Pelczar, and Brooke E. Shelley. 2010. "Highlights from PISA 2009: Performance of U.S. 15-Year-Old Students in Reading," Mathematics, and Science Literacy in an International Context. NCES 2011-004.

Frölich, Markus, and Blaise Melly. 2013. "Unconditional Quartile Treatment Effects under Endogeneity," *Journal of Business and Economic Statistics*, 31(3): 346-357.

- Foster, Gigi. 2006. "It's Not Your Peers, and It's Not Your friends: Some Progress toward Understanding the Educational Peer Effect Mechanism," *Journal of Public Economics*, 90(10-11): 1455-1475.
- Gaviria, Alejandro and Steven Raphael. 2001. "School-Based Peer Effects and Juvenile Behavior," *Review of Economics and Statistics*, 83(2): 257–268.
- Guryan, Jonathan, Kory Kroft, and Matthew J. Notowidigdo. 2009. "Peer Effects in the Workplace: Evidence from Random Groupings in Professional Golf Tournaments," *American Economic Journal: Applied Economics*, 1(4): 34-68.
- Hanushek, Eric A., John F. Kain, Jacob M. Markman, and Steven G. Rivkin. 2003. "Does Peer Ability Affect Student Achievement," *Journal of Applied Econometrics*, 18(5): 527-544.
- Hanushek, Eric. 2006. "School Resources," in *Handbook of the Economics of Education* Vol. 2, eds. Eric A Hanushek and Finis Welch, 865-908, Elsevier.
- Hanushek, Eric A., and Steven G. Rivkin. 2006. "Teacher Quality," in *Handbook of the Economics of Education* Vol. 2, eds. Eric A. Hanushek and Finis Welch, 1051-1162, Elsevier.
- Hanushek, Eric A. and , Steve G. Rivkin. 2009. "Harming the Best: How Schools Affect the Black-White Achievement Gap," *Journal of Policy Analysis and. Management*, 28 (Summer): 366–393.
- Hoxby, Caroline M. 2000. "Peer Effects in the Classroom: Learning from Gender and Race Variation," NBER Working Paper No. 7867.
- Hoxby, Caroline M. and Gretchen Weingarth. 2006. "Taking Race Out of the Equation: School Reassignment and the Structure of Peer Effects." unpublished manuscript.
- Imberman, Scott A., Adriana D. Kugler, and Bruce I. Sacerdote. 2012. "Katrina's Children: Evidence on the Structure of Peer Effects from Hurricane Evacuees," *American Economic Review*, 102(5): 2048-2082.
- Jonsson, Jan O. and Carina Mood. 2008. "Choice by Contrast in Swedish Schools: How Peers' Achievement Affect Educational Choice," *Social Forces*, 87(2): 741-765.
- Kane, Thomas J., Eric S. Taylor, John H. Tyler, and Amy L. Wooten. 2011. "Identifying Effective Classroom Practices Using Student Achievement Data," *Journal of Human Resources*, 46(3): 587-613.
- Kling, Jeffrey R., Jens Ludwig, and Larry F. Katz. 2005. "Neighborhood Effects on Crime for Female and Male Youth: Evidence from a Randomized Housing Voucher Experiment," *Quarterly Journal of Economics*, 120(1): 87–130.

- Krueger, Alan B., and Diane M. Whitmore. 2001. "The Effect of Attending a Small Class in the Early Grades on College-Test Taking and Middle School Results: Evidence from Project STAR," *Economic Journal*, 111(468): 1-28.
- Lavy, Victor, Daniele Paserman, and Analia Schlosser. 2012. "Inside the Black Box of Ability Peer Effects: Evidence from Variation in the Proportion of Low Achievers in the Classroom," *Economic Journal*, 122(559): 208-237.
- Lavy, Victor, Olma Silva, and Felix Weinhardt. 2012. "The Good, the Bad, and the Average: Evidence on Ability Peer Effects in Schools," *Journal of Labor Economics*, 30(2): 367-414.
- Lefgren, Lars. 2004. "Educational Peer Effects and the Chicago Public Schools," *Journal of Urban Economics*, 56(2): 169-191.
- Ludwig, Jens, Greg J. Duncan, and Paul Hirschfield. 2001. "Urban Poverty and Juvenile Crime: Evidence from a Randomized Housing-Mobility," *Quarterly Journal of Economics*, 116(2): 655-679.
- Lyle, David S. 2007. "Estimating and Interpreting Peer and Role Model Effects from Randomly Assigned Social Groups at West Point," *Review of Economics and Statistics*, 89(2): 289-299.
- Manski, Charles F. 1993. "Identification of Endogenous Social Effects: The Reflection Problem," *Review of Economic Studies*, 60(3): 531-542.
- Marsh, Herbert W., and Parker, John W. 1984. "Determinants of Student Self-Concept: Is it Better to be a Relatively Large Fish in a Small Pond Even If You Don't Learn to Swim as Well?" *Journal of Personality and Social Psychology*, 47(1): 213-231.
- Marsh, Herbert W., 1987. "The Big Fish Little Pond Effect on Academic Self-Concept," *Journal of Educational Psychology*, 79(3): 280-295.
- Marsh, Herbert W., and Hau, Kit-Tai. 2003. "Big Fish Little Pond Effect on Academic Self-Concept: A Cross Cultural (26-Country) Test of the Negative Effects of Academic Selective Schools," *American Psychologist*, 58(5): 364-376.
- Mass, Alexandre and Enrico Moretti. 2009. "Peers at Work," *American Economic Review*, 99(1): 112-145.
- Moffitt, Robert A. 2001. "Policy Interventions, Low-Level Equilibria, and Social Interactions," in *Social Dynamics*, eds. Steven N. Durlauf and H. Peyton Young, 45-82. Cambridge, MA: MIT Press.
- Pop-Eleches, Christian and Miguel Urquiola. 2013. "Going to a Better School: Effects and Behavioral Responses," *American Economic Review*, 103(4): 1289-1324.

Sacerdote, Bruce I. 2001. "Peer Effects with Random Assignment: Results for Dartmouth Roommates," *Quarterly Journal of Economics*, 116(2): 681-704.

Sacerdote, Bruce I. 2011. "Peer Effects in Education: How Might They Work, How Big Are They and How Much Do We Know Thus Far?" in *Handbook of the Economics of Education* Vol. 3, eds. Erik Hanushek, Stephen Machin and Ludger Woessmann, 249-277. Elsevier.

Stinebrickner, Ralph, and Todd R. Stinebrickner. 2006. "What Can Be Learned about Peer Effects Using College Roommates? Evidence from New Survey Data and Students from Disadvantaged Backgrounds," *Journal of Public Economics*, 90(8-9): 1435-1454.

US Dept of Education, 2009. Race to the Top Executive Summary. Retrieved from <http://www2.ed.gov/programs/racetothetop/>

Vigdor, Jacob L. 2006. "Peer Effects in Neighborhoods and Housing," in *Deviant Peer Influences in Programs for Youth: Problems and Solutions*, eds. Kenneth A. Dodge., Thomas J. Dishion, Jennifer E. Lansford, Guilford Press.

Vigdor, Jacob L., and Thomas Nechyba. 2007. "Peer Effects in North Carolina Public Schools," in *Schools and the Equal Opportunity Problem*, eds. Ludger Woessmann and Paul E. Peterson, 73-102. Cambridge, MA: MIT Press.

Zimmerman, David J. 2003. "Peer Effects in Academic Outcomes: Evidence from a Natural Experiment," *Review of Economics and Statistics*, 85(1): 9-23.

Table 1: Student Summary Statistics and Basic Randomization Regressions

	TTTDR Students	AC Students	TC (Control) Students	Coefficient (Standard Error)
	Mean (Standard Error)	Mean (Standard Error)	Mean (Standard Error)	AC
Endline Reading Test Score (NCE)	38.59 (20.21)	39.07 (20.38)	38.13 (20.03)
Endline Math Test Score (NCE)	42.43 (22.66)	42.59 (22.83)	42.28 (22.51)
Baseline Reading Test Score (NCE)	38.93 (20.87)	39.86 (20.92)	38.05 (20.80)	0.00 (0.00)
Baseline Math Test Score (NCE)	42.55 (21.28)	42.93 (21.06)	42.18 (21.49)	-0.00 (0.00)
Female (1=Yes)	0.45 (0.49)	0.46 (0.49)	0.44 (0.49)	0.01 (0.01)
Race				
White	0.09 (0.28)	0.09 (0.29)	0.08 (0.28)	-0.02 (0.03)
Black	0.35 (0.47)	0.34 (0.47)	0.35 (0.47)	-0.03 (0.02)
Hispanic	0.47 (0.49)	0.47 (0.49)	0.47 (0.49)	0.01 (0.02)
Free/Reduced Lunch (%)	0.75 (0.42)	0.74 (0.43)	0.77 (0.41)	-0.04 (0.03)
Sample Size	2,610	1,280	1,340	

NOTES: All test scores are expressed in NCEs. NCE scale has a mean 50 and standard deviation 21.06 nationally. Randomization regression tests control for block fixed effects. The standard errors clustered at the block level are reported. AC indicator takes the value of one if the student is taught by a AC teacher and takes the value of zero if the student is taught by a TC teacher. Due to confidential nature of the data, the sample sizes are rounded to the nearest tenth.

Table 2: Randomization Tests

Dependent Variable: Baseline Test Scores	Coefficients (Standard Error)			
	Reading	Math	Reading	Math
Average Peer Baseline Reading Achievement	-0.864*** (0.169)	0.018 (0.034)
Average Peer Baseline Math Achievement	-0.722*** (0.168)	-0.063 (0.053)
Average Block Baseline Reading Achievement	-25.478*** (0.716)
Average Block Baseline Math Achievement	-25.783*** (0.754)

NOTES: All test scores are expressed in NCEs. Standard errors clustered at the block level are reported. Randomization regressions control for block fixed effects. Average peer subject-specific baseline achievement measured at the classroom level.

** significant at 5%, *** significant at 1%.

Table 3: Estimates of Peer Effects by Own Student Baseline Achievement at the Grade Level

Dependent Variable: Endline Test Scores	Coefficients (Standard Error)							
	Own Student Baseline Achievement							
Panel A: Reading Test Scores	All		Bottom 25%		Middle 25%-75		Top 25%	
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Average Peer Baseline Achievement	-0.182*** (0.069)		-0.453*** (0.164)		-0.067 (0.089)		-0.269 (0.196)	
Proportion of Peers in Top 25%		-0.048 (0.032)		-0.060 (0.086)		-0.029 (0.040)		-0.050 (0.077)
Proportion of Peers in Bottom 25%		0.047 (0.030)		0.140*** (0.047)		0.022 (0.046)		0.061 (0.080)
Baseline Test Score	0.632*** (0.020)	0.634*** (0.020)	0.544*** (0.080)	0.562*** (0.080)	0.726*** (0.057)	0.728*** (0.057)	0.552*** (0.052)	0.564*** (0.052)
Bottom vs. Middle (p-value)	0.02							
Bottom vs. Top (p-value)	0.44							
Middle vs. Top (p-value)	0.33							
Top Proportion vs. Bottom Proportion (p-value)		0.03		0.04		0.40		0.31
Sample Size	2,610		640		1,290		680	
Panel B: Math Test Scores	All		Bottom 25%		Middle 25%-75%		Top 25%	
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Average Peer Baseline Achievement	-0.236** (0.103)		-0.281 (0.181)		-0.194 (0.119)		-0.198 (0.192)	
Proportion of Peers in Top 25%		0.018 (0.036)		0.058 (0.086)		0.052 (0.054)		-0.032 (0.089)
Proportion of Peers in Bottom 25%		0.115** (0.045)		0.117 (0.083)		0.089* (0.051)		0.160* (0.093)
Baseline Test Score	0.627*** (0.020)	0.630*** (0.019)	0.555*** (0.068)	0.565*** (0.066)	0.654*** (0.072)	0.661*** (0.072)	0.520*** (0.080)	0.530*** (0.080)
Bottom vs. Middle (p-value)	0.67							
Bottom vs. Top (p-value)	0.73							
Middle vs. Top (p-value)	0.99							
Top Proportion vs. Bottom Proportion (p-value)		0.09		0.62		0.61		0.13
Sample Size	2,580		670		1,250		660	

NOTES: All test scores are expressed in NCEs. Standard errors are clustered at the block level. All specifications control for block fixed effects. Student controls include gender, race/ethnicity and eligibility for free lunch. Teacher controls include teacher's type: AC or TC, gender, race/ethnicity and teaching experience. Average peer subject-specific achievement measured at the classroom level. The proportion of top 25% and bottom 25% of peers in a classroom are based on the grade and subject-specific baseline test score distribution. Due to confidential nature of the data, the sample sizes are rounded to the nearest tenth.

* significant at 10%, ** significant at 5%, *** significant at 1%.

Table 4: Estimates of Peer Effects by Own Student Baseline Achievement at the Grade Level for Female Students

Dependent Variable: Endline Test Scores	Coefficients (Standard Error)							
	Own Student Baseline Achievement							
Panel A: Reading Test Scores	All	Bottom 25%		Middle 25%-75		Top 25%		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Average Peer Baseline Achievement	-0.118 (0.094)		-0.428 (0.265)		-0.025 (0.153)		-0.207 (0.247)	
Proportion of Peers in Top 25%		-0.012 (0.037)		-0.071 (0.124)		0.006 (0.055)		-0.085 (0.099)
Proportion of Peers in Bottom 25%		0.048 (0.041)		0.133* (0.079)		0.037 (0.086)		-0.036 (0.109)
Baseline Test Score	0.608*** (0.023)	0.609*** (0.023)	0.562*** (0.112)	0.573*** (0.113)	0.623*** (0.080)	0.621*** (0.081)	0.452*** (0.076)	0.460*** (0.075)
Bottom vs. Middle (p-value)	0.18							
Bottom vs. Top (p-value)	0.53							
Middle vs. Top (p-value)	0.52							
Top Proportion vs. Bottom Proportion (p-value)		0.27		0.16		0.76		0.73
Sample Size	1,190		290		580		330	
Panel B: Math Test Scores	All	Bottom 25%		Middle 25%-75%		Top 25%		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Average Peer Baseline Achievement	-0.037 (0.124)		0.144 (0.221)		0.009 (0.182)		-0.343 (0.332)	
Proportion of Peers in Top 25%		0.030 (0.057)		0.129 (0.161)		0.067 (0.089)		-0.158 (0.158)
Proportion of Peers in Bottom 25%		0.074 (0.062)		0.009 (0.138)		0.051 (0.087)		0.146 (0.168)
Baseline Test Score	0.641*** (0.027)	0.638*** (0.027)	0.625*** (0.099)	0.613*** (0.009)	0.641*** (0.124)	0.636*** (0.123)	0.552*** (0.128)	0.563*** (0.127)
Bottom vs. Middle (p-value)	0.62							
Bottom vs. Top (p-value)	0.22							
Middle vs. Top (p-value)	0.36							
Top Proportion vs. Bottom Proportion (p-value)		0.60		0.57		0.89		0.18
Sample Size	1,170		310		580		280	

NOTES: All test scores are expressed in NCEs. Standard errors are clustered at the block level. All specifications control for block fixed effects. Student controls include race/ethnicity and eligibility for free lunch. Teacher controls include teacher's type: AC or TC, gender, race/ethnicity and teaching experience. Average peer subject-specific achievement measured at the classroom level. The proportion of top 25% and bottom 25% of peers in a classroom are based on the grade and subject-specific baseline test score distribution. Due to confidential nature of the data, the sample sizes are rounded to the nearest tenth.

* significant at 10%, ** significant at 5%, *** significant at 1%.

Table 5: Estimates of Peer Effects by Own Student Baseline Achievement at the Grade Level for Male Students

Dependent Variable: Endline Test Scores	Coefficients (Standard Error)							
	Own Student Baseline Achievement							
Panel A: Reading Test Scores	All	Bottom 25%		Middle 25%-75		Top 25%		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Average Peer Baseline Achievement	-0.224** (0.094)		-0.446** (0.198)		-0.171 (0.133)		-0.596* (0.344)	
Proportion of Peers in Top 25%		-0.080** (0.040)		-0.092 (0.121)		-0.063 (0.061)		-0.057 (0.118)
Proportion of Peers in Bottom 25%			0.055 (0.040)	0.135* (0.079)		0.043 (0.060)		0.142 (0.133)
Baseline Test Score	0.642*** (0.028)	0.643*** (0.027)	0.498*** (0.127)	0.515*** (0.123)	0.759*** (0.082)	0.765*** (0.079)	0.535*** (0.077)	0.572*** (0.073)
Bottom vs. Middle (p-value)	0.25							
Bottom vs. Top (p-value)	0.70							
Middle vs. Top (p-value)	0.24							
Top Proportion vs. Bottom Proportion (p-value)		0.01		0.11		0.21		0.26
Sample Size	1,420		360		710		350	
Panel B: Math Test Scores	All	Bottom 25%		Middle 25%-75%		Top 25%		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Average Peer Baseline Achievement	-0.441*** (0.130)		-0.574** (0.241)		-0.360** (0.180)		-0.454 (0.299)	
Proportion of Peers in Top 25%		-0.011 (0.049)		0.010 (0.120)		0.062 (0.073)		-0.041 (0.138)
Proportion of Peers in Bottom 25%			0.151*** (0.056)	0.186* (0.108)		0.133* (0.071)		0.260** (0.130)
Baseline Test Score	0.612*** (0.027)	0.620*** (0.027)	0.489*** (0.118)	0.512*** (0.118)	0.654*** (0.086)	0.666*** (0.085)	0.436*** (0.106)	0.460*** (0.110)
Bottom vs. Middle (p-value)	0.48							
Bottom vs. Top (p-value)	0.74							
Middle vs. Top (p-value)	0.79							
Top Proportion vs. Bottom Proportion (p-value)		0.02		0.27		0.48		0.11
Sample Size	1,410		360		670		380	

NOTES: All test scores are expressed in NCEs. Standard errors are clustered at the block level. All specifications control for block fixed effects. Student controls include race/ethnicity and eligibility for free lunch. Teacher controls include teacher's type: AC or TC, gender, race/ethnicity and teaching experience. Average peer subject-specific achievement measured at the classroom level. The proportion of top 25% and bottom 25% of peers in a classroom are based on the grade and subject-specific baseline test score distribution. Due to confidential nature of the data, the sample sizes are rounded to the nearest tenth. * significant at 10%, ** significant at 5%, *** significant at 1%.

Table 6: Estimates of Peer Effects by Own Student Baseline Achievement at the Grade Level for Free-Lunch Eligible Students

Dependent Variable: Endline Test Scores	Coefficients (Standard Error)							
	Own Student Baseline Achievement							
Panel A: Reading Test Scores	All		Bottom 25%	Middle 25%-75		Top 25%		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Average Peer Baseline Achievement	-0.157*		-0.459**		0.005		-0.546**	
	(0.096)		(0.191)		(0.123)		(0.256)	
Proportion of Peers in Top 25%		-0.019		-0.045		0.003		-0.036
		(0.038)		(0.094)		(0.049)		(0.113)
Proportion of Peers in Bottom 25%		0.061*		0.146***		0.028		0.163
		(0.033)		(0.053)		(0.050)		(0.110)
Baseline Test Score	0.633***	0.634***	0.544***	0.563***	0.797***	0.797***	0.510***	0.544***
	(0.025)	(0.024)	(0.081)	(0.080)	(0.055)	(0.055)	(0.105)	(0.105)
Bottom vs. Middle (p-value)	0.05							
Bottom vs. Top (p-value)	0.77							
Middle vs. Top (p-value)	0.05							
Top Proportion vs. Bottom Proportion (p-value)	0.11		0.07		0.72		0.20	
Sample Size	1,980		570		1,040		360	
<hr/>								
Panel B: Math Test Scores	All		Bottom 25%	Middle 25%-75%		Top 25%		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Average Peer Baseline Achievement	-0.218*		-0.243		-0.060		-0.236	
	(0.128)		(0.212)		(0.123)		(0.281)	
Proportion of Peers in Top 25%		0.016		0.038		0.079		-0.016
		(0.048)		(0.093)		(0.067)		(0.185)
Proportion of Peers in Bottom 25%		0.112**		0.086		0.078		0.185
		(0.049)		(0.091)		(0.051)		(0.122)
Baseline Test Score	0.639***	0.641***	0.538***	0.546***	0.619***	0.622**8	0.527***	0.539***
	(0.024)	(0.023)	(0.075)	(0.072)	(0.083)	(0.082)	(0.119)	(0.120)
Bottom vs. Middle (p-value)	0.45							
Bottom vs. Top (p-value)	0.97							
Middle vs. Top (p-value)	0.57							
Top Proportion vs. Bottom Proportion (p-value)	0.16		0.71		0.99		0.36	
Sample Size	1,950		570		990		400	

NOTES: All test scores are expressed in NCEs. Standard errors are clustered at the block level. All specifications control for block fixed effects. Student controls include gender and race/ethnicity. Teacher controls include teacher's type: AC or TC, gender, race/ethnicity and teaching experience. Average peer subject-specific achievement measured at the classroom level. The proportion of top 25% and bottom 25% of peers in a classroom are based on the grade and subject-specific baseline test score distribution. Due to confidential nature of the data, the sample sizes are rounded to the nearest tenth.

* significant at 10%, ** significant at 5%, *** significant at 1%.

Table 7: Estimates of Peer Effects by Own Student Baseline Achievement at the Grade Level for Black Students

Dependent Variable: Endline Test Scores	Coefficients (Standard Error)							
	Own Student Baseline Achievement							
Panel A: Reading Test Scores	All	Bottom 25%		Middle 25%-75		Top 25%		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Average Peer Baseline Achievement	-0.082 (0.108)		-0.319 (0.313)		0.051 (0.141)		-0.351 (0.291)	
Proportion of Peers in Top 25%		0.032 (0.056)		0.150 (0.137)		0.034 (0.069)		-0.010 (0.118)
Proportion of Peers in Bottom 25%		0.072** (0.032)		0.198** (0.082)		0.047 (0.075)		0.110 (0.153)
Baseline Test Score	0.652*** (0.028)	0.653*** (0.028)	0.747*** (0.124)	0.760*** (0.122)	0.678*** (0.100)	0.675*** (0.101)	0.460*** (0.102)	0.484*** (0.099)
Bottom vs. Middle (p-value)	0.28							
Bottom vs. Top (p-value)	0.94							
Middle vs. Top (p-value)	0.21							
Top Proportion vs. Bottom Proportion (p-value)		0.55		0.76		0.89		0.53
Sample Size	900		180		500		220	
Panel B: Math Test Scores	All	Bottom 25%		Middle 25%-75%		Top 25%		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Average Peer Baseline Achievement	-0.240 (0.186)		-0.367 (0.310)		-0.050 (0.208)		0.062 (0.489)	
Proportion of Peers in Top 25%		0.029 (0.059)		0.070 (0.131)		0.040 (0.098)		0.162 (0.170)
Proportion of Peers in Bottom 25%		0.110 (0.072)		0.143 (0.140)		0.029 (0.070)		0.126 (0.161)
Baseline Test Score	0.668*** (0.035)	0.672*** (0.035)	0.514*** (0.170)	0.520*** (0.170)	0.625*** (0.132)	0.629*** (0.132)	0.633*** (0.165)	0.614*** (0.162)
Bottom vs. Middle (p-value)	0.39							
Bottom vs. Top (p-value)	0.45							
Middle vs. Top (p-value)	0.83							
Top Proportion vs. Bottom Proportion (p-value)		0.38		0.70		0.92		0.87
Sample Size	900		230		470		210	

NOTES: All test scores are expressed in NCEs. Standard errors are clustered at the block level. All specifications control for block fixed effects. Student controls include gender and eligibility for free lunch. Teacher controls include teacher's type: AC or TC, gender, race/ethnicity and teaching experience. Average peer subject-specific achievement measured at the classroom level. The proportion of top 25% and bottom 25% of peers in a classroom are based on the grade and subject-specific baseline test score distribution. Due to confidential nature of the data, the sample sizes are rounded to the nearest tenth. * significant at 10%, ** significant at 5%, *** significant at 1%.

Table 8: Estimates of Peer Effects by Own Student Baseline Achievement at the Grade Level for Hispanic Students

Dependent Variable: Endline Test Scores	Coefficients (Standard Error)							
	Own Student Baseline Achievement							
Panel A: Reading Test Scores	All	Bottom 25%	Middle 25%-75	Top 25%				
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Average Peer Baseline Achievement	-0.080 (0.160)		-0.361 (0.341)		0.044 (0.223)		-0.204 (0.438)	
Proportion of Peers in Top 25%		-0.048 (0.051)		-0.096 (0.122)		-0.034 (0.069)		-0.054 (0.121)
Proportion of Peers in Bottom 25%		0.048 (0.056)		0.135* (0.074)		0.017 (0.079)		0.055 (0.191)
Baseline Test Score	0.621*** (0.032)	0.618*** (0.031)	0.480*** (0.101)	0.497*** (0.099)	0.812*** (0.081)	0.806*** (0.080)	0.402*** (0.112)	0.416*** (0.113)
Bottom vs. Middle (p-value)	0.32							
Bottom vs. Top (p-value)	0.77							
Middle vs. Top (p-value)	0.61							
Top Proportion vs. Bottom Proportion (p-value)		0.20		0.10		0.62		0.62
Sample Size	1,230		400		610		220	
Panel B: Math Test Scores	All	Bottom 25%	Middle 25%-75%	Top 25%				
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Average Peer Baseline Achievement	-0.205 (0.128)		-0.219 (0.300)		-0.370** (0.185)		0.365 (0.450)	
Proportion of Peers in Top 25%		-0.012 (0.057)		0.063 (0.156)		-0.041 (0.095)		0.049 (0.256)
Proportion of Peers in Bottom 25%		0.089 (0.065)		0.094 (0.128)		0.152* (0.085)		-0.079 (0.245)
Baseline Test Score	0.609*** (0.026)	0.611*** (0.025)	0.551*** (0.076)	0.562*** (0.075)	0.768*** (0.097)	0.784*** (0.098)	0.690*** (0.181)	0.676*** (0.180)
Bottom vs. Middle (p-value)	0.66							
Bottom vs. Top (p-value)	0.28							
Middle vs. Top (p-value)	0.13							
Top Proportion vs. Bottom Proportion (p-value)		0.24		0.87		0.13		0.71
Sample Size	1,200		370		580		250	

NOTES: All test scores are expressed in NCEs. Standard errors are clustered at the block level. All specifications control for block fixed effects. Student controls include gender and eligibility for free lunch. Teacher controls include teacher's type: AC or TC, gender, race/ethnicity and teaching experience. Average peer subject-specific achievement measured at the classroom level. The proportion of top 25% and bottom 25% of peers in a classroom are based on the grade and subject-specific baseline test score distribution. Due to confidential nature of the data, the sample sizes are rounded to the nearest tenth. * significant at 10%, ** significant at 5%, *** significant at 1%.

Table 9: Tests of Peer Effect Models-Full Sample

Panel A: Endline Reading Test Scores		Tests		
Models (Number of Tests)	Number of Estimates Significant in Direction the Model Suggest	Number of Estimates Insignificant in Direction the Model Suggest	Number of Estimates Significant in Opposite Direction	Number of Estimates Insignificant in Opposite Direction
Weak Monotonicity (3)	0	0	1	2
Strong Monotonicity (6)	0	0	1	5
Invidious Comparison (4)	0	4	0	0
Ability Grouping (6)	2	1	0	3
Weak Frame of Reference (3)	1	2	0	0
Strong Frame of Reference (6)	1	5	0	0

Panel B: Endline Math Test Scores		Tests		
Models (Number of Tests)	Number of Estimates Significant in Direction the Model Suggest	Number of Estimates Insignificant in Direction the Model Suggest	Number of Estimates Significant in Opposite Direction	Number of Estimates Insignificant in Opposite Direction
Weak Monotonicity (3)	0	0	0	3
Strong Monotonicity (6)	0	2	2	2
Invidious Comparison (4)	2	0	0	2
Ability Grouping (6)	0	2	1	3
Weak Frame of Reference (3)	0	3	0	0
Strong Frame of Reference (6)	2	2	0	2

NOTES: See text for further details on the models and the tests conducted. The tests are based on the coefficients on the proportion of peers in the top 25% and bottom 25% from the full sample when own student varies by grade and subject specific placement in the baseline distribution.

Table 10: Peer Effects and Teachers' Performance Based VCOT Ratings

Dependent Variable: VCOT Ratings	Coefficients (Standard Error)			
	Reading Lesson Implementation		Reading Classroom Culture	
Panel A: Reading VCOT Ratings	(1)	(2)	(3)	(4)
Average Peer Baseline Reading Achievement	0.000 (0.016)		0.017 (0.018)	
Proportion of Peers in Top 25%		-0.002 (0.005)		-0.001 (0.006)
Proportion of Peers in Bottom 25%		-0.002 (0.006)		-0.007 (0.008)
Top Proportion vs. Bottom Proportion (p-value)		0.99		0.54
Panel B: Math VCOT Ratings	Math Lesson Implementation		Math Classroom Culture	
	(1)	(2)	(3)	(4)
Average Peer Baseline Math Achievement	0.043*** (0.014)		0.036** (0.015)	
Proportion of Peers in Top 25%		0.022*** (0.006)		0.021*** (0.006)
Proportion of Peers in Bottom 25%		-0.003 (0.005)		0.001 (0.006)
Top Proportion vs. Bottom Proportion (p-value)		0.00		0.01

NOTES: VCOT ratings are standardized at the grade level to have a mean of zero and standard deviation one. All test scores are expressed in NCEs. Standard errors are clustered at the block level. All specifications control for block fixed effects and teacher characteristics (teacher's type: AC or TC, gender, race/ethnicity, and teaching experience). Average peer subject-specific baseline achievement measured at the classroom level. The proportion of top 25% and bottom 25% of peers in a classroom are on the grade and subject-specific baseline test score distribution.

*significant at 10%, ** significant at 5%, *** significant at 1%.

Appendix:

Table A1: Teacher Summary Statistics

	All Teachers	AC Teachers	TC Teachers
	Mean (Standard Error)	Mean (Standard Error)	Mean (Standard Error)
Female	0.90 (0.29)	0.87 (0.33)	0.93 (0.25)
Race			
White	0.59 (0.49)	0.45 (0.50)	0.72 (0.44)
Black	0.24 (0.43)	0.35 (0.48)	0.12 (0.33)
Hispanic	0.17 (0.38)	0.20 (0.40)	0.15 (0.36)
Experience	3.29 (1.59)	3.10 (1.59)	3.36 (1.59)
Hours of Instruction for Certification			
Total	462.06 (253.57)	296.91 (150.76)	623.69 (228.64)
Reading	89.28 (57.31)	61.32 (44.95)	116.65 (55.02)
Math	33.15 (23.18)	25.44 (23.14)	40.85 (20.61)
SAT (or SAT Equivalent) Composite Score	972.40 (161.98)	960.76 (179.35)	982.91 (145.01)
# of Teachers	180	90	90

NOTES: Due to confidential nature of the data, the sample sizes are rounded to the nearest tenth.

Table A2: Tests of Peer Effect Models-Subgroups

Panel A: Endline Reading Test Scores		Tests		
Models (Number of Tests)	Number of Estimates Significant in Direction the Model Suggest	Number of Estimates Insignificant in Direction the Model Suggest	Number of Estimates Significant in Opposite Direction	Number of Estimates Insignificant in Opposite Direction
Weak Monotonicity (15)	0	0	2	13
Strong Monotonicity (30)	0	5	5	20
Invidious Comparison (20)	0	16	0	4
Ability Grouping (30)	7	5	0	18
Weak Frame of Reference (15)	2	13	0	0
Strong Frame of Reference (30)	5	20	0	5

Panel B: Endline Math Test Scores		Tests		
Models (Number of Tests)	Number of Estimates Significant in Direction the Model Suggest	Number of Estimates Insignificant in Direction the Model Suggest	Number of Estimates Significant in Opposite Direction	Number of Estimates Insignificant in Opposite Direction
Weak Monotonicity (15)	0	6	0	9
Strong Monotonicity (30)	0	12	4	14
Invidious Comparison (20)	3	7	0	10
Ability Grouping (30)	1	13	2	14
Weak Frame of Reference (15)	0	9	0	6
Strong Frame of Reference (30)	4	14	0	12

NOTES: See text for further details on the models and the tests conducted. The tests are based on the coefficients on the proportion of peers in the top 25% and bottom 25% from the subgroups (female students, male students, free-lunch eligible students, black students, and Hispanic students) when own student varies by grade and subject specific placement in pre-treatment distribution.

Table A3: Tests of Peer Effect Models-Full Sample with Bonferroni Corrections

Panel A: Endline Reading Test Scores		Tests		
Models (Number of Tests)	Number of Estimates Significant in Direction the Model Suggest	Number of Estimates Significant in Opposite Direction	Number of Insignificant Estimates	
Weak Monotonicity (3)	0	1	2	
Strong Monotonicity (6)	0	1	5	
Invidious Comparison (4)	0	0	4	
Ability Grouping (6)	2	0	4	
Weak Frame of Reference (3)	1	0	2	
Strong Frame of Reference (6)	1	0	5	

Panel B: Endline Math Test Scores		Tests		
Models (Number of Tests)	Number of Estimates Significant in Direction the Model Suggest	Number of Estimates Significant in Opposite Direction	Number of Insignificant Estimates	
Weak Monotonicity (3)	0	0	3	
Strong Monotonicity (6)	0	0	6	
Invidious Comparison (4)	0	0	4	
Ability Grouping (6)	0	0	6	
Weak Frame of Reference (3)	0	0	3	
Strong Frame of Reference (6)	0	0	6	

NOTES: See text for further details on the models and the tests conducted. The tests are based on the coefficients on the proportion of peers in the top 25% and bottom 25% from the full sample when own student varies by grade and subject specific placement in the baseline distribution.