

A Service of

ZBW

Leibniz-Informationszentrum Wirtschaft Leibniz Information Centre for Economics

Bietenbeck, Jan

### Working Paper The Long-Term Impacts of Low-Achieving Childhood Peers: Evidence from Project STAR

Working Paper, No. 2015:35

**Provided in Cooperation with:** Department of Economics, School of Economics and Management, Lund University

*Suggested Citation:* Bietenbeck, Jan (2015) : The Long-Term Impacts of Low-Achieving Childhood Peers: Evidence from Project STAR, Working Paper, No. 2015:35, Lund University, School of Economics and Management, Department of Economics, Lund

This Version is available at: https://hdl.handle.net/10419/260172

#### Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

#### Terms of use:

Documents in EconStor may be saved and copied for your personal and scholarly purposes.

You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.

If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.



### WWW.ECONSTOR.EU

Working Paper 2015:35

Department of Economics School of Economics and Management

## The Long-Term Impacts of Low-Achieving Childhood Peers: Evidence from Project STAR

Jan Bietenbeck

November 2015



### The Long-Term Impacts of Low-Achieving Childhood Peers: Evidence from Project STAR<sup>☆</sup>

#### Jan Bietenbeck

Lund University and IZA

October 2015

#### Abstract

This paper evaluates how sharing a kindergarten classroom with low-achieving repeaters affects the long-term educational performance of regular firsttime kindergarten students. Exploiting random assignment of teachers and students to classes in Project STAR, I document three sets of causal impacts: students who are exposed to repeaters (1) score lower on standardized tests at the end of kindergarten, an effect that fades out in later grades; (2) show persistent improvements in non-cognitive skills such as effort and discipline; and (3) are more likely to graduate from high school and to take a college entrance exam around the age of eighteen. I show that the positive spillovers from repeaters on long-term educational attainment are likely driven by the differential accumulation of non-cognitive skills by repeater-exposed students during childhood. The improvements in these skills are in turn a result of behavioral adjustments by teachers, students, or parents to the presence of low-achieving repeaters in the classroom.

<sup>&</sup>lt;sup>\*</sup>I am greatly indebted to David Dorn for his guidance and support. I thank Joseph Altonji, Manuel Arellano, David Autor, Manuel Bagues, Sascha Becker, Samuel Bentolila, Susan Dynarski, Joshua Goodman, Bryan Graham, Claudio Michelacci, Pedro Mira, Magne Mogstad, Luca Repetto, Jan Stuhler, Simon Wiederhold, and numerous seminar and conference audiences for helpful comments. Diane Schanzenbach generously provided a subset of the data used in this paper. Funding from the Spanish Ministry of Science and Innovation (BES-2011-050947) is gratefully acknowledged. Author contact details: Department of Economics, Lund University, P.O. Box 7082, 220 07 Lund, Sweden; jan.bietenbeck@nek.lu.se

#### 1. Introduction

A large academic literature studies the effects of class composition on student performance in school. Papers in this literature have generally found positive impacts from sharing a classroom with higher-achieving and better-behaved peers (e.g. Hoxby, 2000; Burke and Sass, 2013; Sojourner, 2013) and corresponding negative impacts from sharing a classroom with low-achieving or disruptive peers (e.g. Figlio, 2007; Carrell and Hoekstra, 2010; Lavy, Silva, and Weinhardt, 2012). The vast majority of these papers has focused exclusively on short-term spillover effects such as the impact of kindergarten classmates on test scores at the end of kindergarten. However, in order to judge the overall efficacy of policies that change the student composition of classes and schools, it is important to know how such spillovers play out in the long term.

In this paper, I study how sharing a kindergarten classroom with lowachieving repeaters affects the long-term educational performance of regular first-time kindergarten students. The empirical analysis uses data from the Tennessee Student-Teacher Achievement Ratio experiment (Project STAR), which is uniquely suited for this purpose for three reasons. First, the data allow me to identify kindergarten repeaters as a particularly lowachieving group of peers: by definition, these students have a proven track record of failure, and they are characterized by exceptionally low cognitive and non-cognitive skills.<sup>1</sup> Second, Project STAR randomly assigned teachers and students, including repeaters, to classes within schools. This lets me estimate causal spillover effects from repeaters that are free from selection bias. Finally, the data contain a rich set of medium- and longterm outcomes for students, including measures of non-cognitive skills, high school graduation, and college-test taking.

The main empirical specifications relate regular students' exposure to repeaters in kindergarten, measured as being randomly assigned to a class containing at least one repeater, to their educational performance at dif-

<sup>&</sup>lt;sup>1</sup>Previous studies have mostly identified low-achieving peers by their past academic achievement. In settings different from Project STAR, Lavy, Paserman, and Schlosser (2012), Gottfried (2013), and Hill (2014) document negative short-term spillovers from repeaters on their classmates' test scores.

ferent points in time.<sup>2</sup> Being exposed to repeaters significantly reduces students' test scores at the end of kindergarten, a result which corroborates previous findings of negative short-term spillovers from low-achieving peers. Repeater exposure however substantially increases students' noncognitive skills such as effort and discipline, which are first observed at the beginning of fourth grade. While the negative spillovers from repeaters on test scores fade out rapidly after kindergarten, the gains in non-cognitive skills persist over time. Consistent with this last result, students who shared a kindergarten classroom with repeaters show improved long-term educational attainment as evidenced by higher propensities to graduate from high school and to take a college entrance exam.

In additional analyses, I explore the potential mechanisms behind these results. Motivated by recent findings that non-cognitive skills formed early in life are a key determinant of long-term educational success (e.g. Heckman, Stixrud, and Urzua, 2006), I hypothesize that the positive spillovers on high school graduation and college-test taking are driven by the differential accumulation of such skills by repeater-exposed students. Suggestive evidence from a regression-based test supports this hypothesis. As for how exactly sharing a classroom with repeaters boosts non-cognitive skills, the experimental setup of Project STAR lets me rule out a wide range of potential explanations that involve selection of students or teachers into classes either in kindergarten or in later grades. Ultimately, I narrow the possible mechanisms down to behavioral responses by teachers, students, or parents, although the data do not allow me to pin down the specific response that drives the improvements in non-cognitive skills.

This paper contributes to a large existing literature on peer effects in schools, which was recently reviewed by Sacerdote (2011). This literature includes previous studies based on data from Project STAR, most notably by Whitmore (2005), who examines the effects of classroom gender composition, and by Graham (2008) and Sojourner (2013), who investigate spillovers from peers' academic ability. These papers focus exclusively on short-term spillovers, as does the vast majority of the existing research.

 $<sup>^2{\</sup>rm I}$  use the terms "regular student," "non-repeating student," and "first-time kinder-garten student" interchangeably throughout the paper.

An important exception to this generalization is the study by Gould, Lavy, and Paserman (2009), which investigates how sharing a fifth-grade classroom with immigrants affects the likelihood of natives to graduate from high school in Israel. Furthermore, Bifulco, Fletcher, and Ross (2011) and Black, Devereux, and Salvanes (2013) examine spillovers from high school peers on longer-term educational and labor market outcomes.<sup>3</sup>

The analysis in this paper provides some of the first evidence on longterm spillovers from early childhood peers. Studying these spillovers is important because kindergarten students are at an age where both cognitive and non-cognitive skills are still highly malleable (Kautz et al., 2014). Moreover, the long-term impacts examined here are arguably more relevant than short-term effects for the evaluation of policies as they may translate more directly into changes in labor market outcomes. The importance of studying such long-term impacts is further highlighted by the finding that short- and long-term effects do not necessarily go in the same direction. Finally, my results have important implications for the ongoing policy debate on whether students of different abilities should be separated at an early age. I discuss these implications in detail in the conclusion.

#### 2. The STAR experiment and data

#### 2.1. Background on Project STAR

Project STAR was a randomized experiment designed to study the effects of class size on student achievement. In the beginning of the 1985-86 school year, 6,325 kindergarten students in 79 participating Tennessee schools were randomly assigned to small (target size 13-17 students) or regular-sized (22-25 students) classes within their schools.<sup>4</sup> Students were supposed to stay in their assigned class type (small versus regular-sized)

<sup>&</sup>lt;sup>3</sup>Chetty et al. (2011) show that kindergarten class fixed effects predict earnings at ages 25-27 of participants in Project STAR. As the authors note, these "class effects" combine the impacts of peers, teachers, and any other class-level shocks and therefore cannot be interpreted as peer effects. In a recent paper, Cascio and Schanzenbach (2015) analyze the impacts of students' relative age, measured as the difference between own age and classmates' average age, on short- and long-term outcomes in Project STAR.

<sup>&</sup>lt;sup>4</sup>There was also a third class type: regular-sized classes with a full-time teacher's aide. Like previous analyses of Project STAR, I do not find any differences in treatment effects between regular-sized classes with and without a full-time teacher's aide.

until the end of third grade, after which the experiment ended and they would return to ordinary classes. Students that joined the initial cohort in participating schools after the kindergarten year were also randomly assigned to class types, as were teachers in each grade.

This study exploits the fact that kindergarten students, including repeaters, and teachers were randomly assigned not only to class type, but also to a particular class within each type (50 schools in the experiment had multiple classes per type). Early analyses of Project STAR were reluctant to conclude that this was indeed the case, mainly because the STAR Technical Report (Word et al., 1990) does not describe the exact procedure by which students were allocated to specific classes. However, several more recent studies (Chetty et al., 2011; Sojourner, 2013; Cascio and Schanzenbach, 2015) also rely on random assignment of students and teachers to classes in Project STAR and provide new evidence in support of this assumption. Section 3 revisits some of this evidence and provides additional statistical support for the claim that repeaters were randomly assigned to kindergarten classes within schools.

The eventual implementation of Project STAR differed somewhat from the original experimental design. Three of these differences are particularly important in the context of this paper. First, as the initial cohort of students advanced from kindergarten to third grade, there was substantial attrition due to students moving to other schools or being retained in grade. Thus, by the time the cohort reached third grade, 49% of students who had participated in the experiment in kindergarten had left the sample. Second, because of complaints by some parents about their children's initial assignment, students in regular-sized classes were re-randomized at the beginning of first grade. Third, while compliance with treatment assignment was nearly perfect in kindergarten, approximately 10% of students managed to switch between small and regular-sized classes in each of the subsequent grades (Krueger, 1999).

Due to the focus on spillovers from repeaters in kindergarten, noncompliance with class assignment in the later grades does not affect the (reduced-form) causal interpretation of results in this paper. In contrast, sample attrition could potentially confound some of the estimates on longterm outcomes. I test for selective attrition in Section 6 below, but do not find it to be a problem for the large majority of outcomes studied here. Finally, the three aspects of the implementation mentioned in the previous paragraph change the total amount of time that students spent in class with a kindergarten repeater. This affects the interpretation of the repeater-exposure treatment, a point that I discuss in more detail in the following subsection.<sup>5</sup>

#### 2.2. Variable definitions

Data for students participating in Project STAR were collected by various research teams and organizations both during the experiment and in several rounds after the experiment ended. The Project STAR public use file, on which the empirical analysis below is based, combines these data such that students can be followed throughout their scholastic careers until the end of high school.<sup>6</sup> This subsection gives a brief overview of the dependent and independent variables used in the empirical analysis. Online appendix A provides additional details on data collection procedures and on the construction of outcome variables.

Demographic characteristics. The data contain information on students' gender, race, eligibility for free or reduced-price lunch, and exact date of birth. Children in Tennessee are supposed to enter kindergarten if they are five years or older on September 30 of a given year, and I use this rule to construct an old-for-grade indicator which takes value 1 if the student was six years or older on September 30, 1985, and 0 otherwise. Students in Project STAR may be old for grade either because they entered school late (the so-called "red-shirting") or because they were repeating kindergarten.<sup>7</sup>

*Kindergarten repeaters.* The data include an indicator for whether each student was repeating kindergarten in the 1985-86 school year. There are 253 repeaters in the sample, 193 of whom are old for grade. Note that

<sup>&</sup>lt;sup>5</sup>Additional details regarding the design and implementation of Project STAR can be found in Word et al. (1990), Krueger (1999), and Finn et al. (2007).

<sup>&</sup>lt;sup>6</sup>Data on some of the outcomes studied in this paper were generously provided to me by Diane Schanzenbach; see Online appendix A for details.

<sup>&</sup>lt;sup>7</sup>See Deming and Dynarski (2008) for an analysis of the red-shirting phenomenon in the United States.

all repeaters would be expected to be old for grade if they had entered kindergarten in accordance with Tennessee's school entry rules during one of the previous school years. Therefore, the 60 repeaters who were not old for grade must have entered school early. The empirical analysis below focuses on spillover effects from the 193 old-for-grade repeaters, who first entered kindergarten at the regular entry age. While the data do not contain information on the exact reason for their retention, these students had likely been identified by principals or teachers as having cognitive or behavioral deficiencies that would have put them at a disadvantage had they been promoted to first grade. The same is not necessarily true for the 60 other repeaters, who may have stayed in kindergarten only because they were too young to enter first grade.<sup>8</sup>

Repeater exposure. Figure 1 shows the distribution of repeaters across classes in schools with at least one repeater.<sup>9</sup> 126 of the 254 classes in this subsample contain no repeater, 81 contain one repeater, and only 47 contain two or more repeaters. In view of this heavily skewed distribution, the main specifications of the empirical analysis will distinguish just between classes with and without repeaters. As a robustness check, I also measure repeater exposure as the actual number of repeaters in class, or as the share of repeaters in class. Results from these alternative specifications suggest that outcomes are similar for students who are exposed to one or to several repeaters, which implies that the main specifications using a

<sup>&</sup>lt;sup>8</sup>Children are required to be six years old on September 30 of the year they start first grade. It seems reasonable to assume that this rule was enforced more strictly than the kindergarten entry rule since kindergarten attendance was not mandatory in Tennessee at the time of the experiment. Empirically, the 60 "young" repeaters come from more favorable demographic backgrounds and exhibit better cognitive and non-cognitive outcomes than the 193 old-for-grade repeaters. If all 253 repeaters are used as treatment in the empirical analysis, the estimated spillover effects are somewhat attenuated compared to the ones reported in the paper.

<sup>&</sup>lt;sup>9</sup>Out of the 79 participating schools, 60 contain at least one repeater. The 19 schools without repeaters do not contribute to the identification of spillover effects in this paper, which is based on between-class variation in the number of repeaters within schools. Compared to schools without repeaters, schools with positive numbers of repeaters are slightly smaller (average enrollment of 73 students versus 83 students), are less likely to be located in the inner city (12% versus 47% of schools), and contain lower fractions of black students (20% versus 61%) and low-income students (41% versus 67%) on average.

dummy variable for the presence of at least one repeater in class do not unduly miss heterogeneous treatment effects.

An important question for the interpretation of results is whether the spillovers on long-term outcomes documented in this paper arise from exposure to repeaters during kindergarten or from exposure over a longer time horizon. If all children had stayed in their assigned kindergarten classes until the end of the experiment, regular students would have been exposed to repeaters either for four years or not at all until third grade. In practice, however, due to the various deviations from the original experimental design described above, students who were exposed to repeaters in kindergarten and who had not left the experiment by third grade ended up being in class with at least one of these repeaters for 2 years on average, whereas students not exposed to repeaters in kindergarten ended up being in class with repeaters for an average of 0.6 years.<sup>10</sup> The treatment studied in this paper thus consists of exposure to repeaters during kindergarten and an additional six months of differential exposure during grades 1-3.

*Outcomes.* At the end of each grade level from kindergarten through third grade, students were administered the grade-appropriate version of the Stanford Achievement Test. Moreover, in the spring of grades 5-8, all participants still attending public school in Tennessee took the Comprehensive Test of Basic Skills as part of a statewide student assessment program. Both tests are standardized multiple-choice assessments with components in mathematics and reading. The empirical analysis below studies the effects of repeater exposure in kindergarten on student performance on these tests at each grade level.

In November 1989, when participants were in fourth grade, teachers in

<sup>&</sup>lt;sup>10</sup>These figures come from a regression of cumulative years of exposure at the end of third grade on a constant, an indicator for repeater exposure in kindergarten, and school fixed effects. Further analysis showed that cumulative years of exposure are very similar for students assigned to small and to regular-sized kindergarten classes. Note that these figures measure exposure to at least one of the 193 *original* repeaters for students who did not attrit from the experiment. A complete history of exposure to *any* repeaters cannot be determined for participants in Project STAR because class composition is no longer observed for students who leave the experiment and because repeater status was not recorded for students who entered the experiment after kindergarten.

the STAR schools were asked to evaluate a random subset of their students on a set of behavioral measures. Teacher ratings were recorded on a scale from 1-5 and were consolidated into four indices. The effort index is based on such items as whether a student completes her homework and whether she is persistent when confronted with difficult problems. The initiative index captures such characteristics as whether a student actively participates in classroom discussions. The value index measures how much a student appreciates the school learning environment. Finally, the discipline index is based on such items as whether a student often acts restless and whether she interferes with her peers' work. In eighth grade, math and English teachers were asked to rate a different random subset of STAR participants on similar questions, the answers to which were consolidated into the same four indices. The total of eight fourth- and eighth-grade indices derived from teacher ratings serve as measures of non-cognitive skills in the empirical analysis below.

Most STAR participants graduated from high school in 1998, and transcripts including information on high school grade point average (GPA) and graduation status were collected from selected high schools in 1999 and 2000. Colleges and universities in the United States typically require applying students to report results from either the ACT or the SAT test. In 1998, Krueger and Whitmore (2001) matched all STAR students to the administrative records of the two companies responsible for these tests. The outcome of this process is an indicator that takes value 1 if a student took either of these college entrance exams in 1998 and 0 otherwise. Together, high school GPA, high school graduation, and college-test taking are the measures of long-term educational attainment studied in this paper.

#### 2.3. Sample selection and descriptive statistics

The full sample includes 6,325 kindergarten students in 127 small and 198 regular-sized classes in 79 schools. I exclude 28 students for whom repeater status is not observed and five students with missing demographic characteristics from this sample. I further drop the 60 repeaters who are not old for grade as they had likely been in class with one of the oldfor-grade repeaters during the previous (1984-85) school year and are thus subject to a fundamentally different treatment. Finally, while schools without repeaters do not contribute to the identification of spillover effects in this paper, they are kept in the sample in order to increase the precision of the estimated impacts of other covariates included in the regressions. The final estimation sample thus consists of 6,232 students, 193 of whom are repeaters. Results in this paper are robust to relaxing the sample restrictions discussed in this paragraph.

Table 1 shows descriptive statistics for demographic characteristics, repeater exposure, and key outcome variables separately for non-repeating and repeating kindergarten students in the estimation sample. Students in general exhibit lower socioeconomic characteristics than the student populations in Tennessee and the United States as a whole because Project STAR oversampled schools in low-income neighborhoods (Krueger and Whitmore, 2001). Repeaters are predominantly male and are more likely to be eligible for free or reduced-price lunch than non-repeating students. Repeaters are also older than non-repeating students by definition. Since low-income schools with primarily black student populations have lower repeater shares on average, repeating students in the sample are less likely to be black. Finally, only three percent of non-repeating students are old for grade, which shows that red-shirting was not common in the schools participating in Project STAR at the time of the experiment.

In order to facilitate easy comparison between the outcomes of regular students and repeaters, I standardize all test scores and non-cognitive skill measures to have mean 0 and standard deviation 1 across non-repeating students in the estimation sample. Consistent with the idea that repeaters were retained in grade because of some cognitive or behavioral deficiencies, Table 1 shows that they tend to perform substantially worse than regular students in school. For instance, repeaters score half a standard deviation below non-repeating students on the end-of-kindergarten reading test, a gap that widens to almost a full standard deviation by eighth grade. Repeaters are also rated considerably worse on measures of effort, initiative, value, and discipline by their teachers. Therefore, by focusing on repeater exposure as my treatment, I capture the impacts of sharing a classroom with students with exceptionally low cognitive and non-cognitive skills.<sup>11</sup>

<sup>&</sup>lt;sup>11</sup>Repeaters' measured cognitive and non-cognitive skills are also significantly below

#### 3. Identification strategy and validity of the experimental design

#### 3.1. Identification based on between-class variation in repeater exposure

Identification of spillovers from repeaters in this paper is based on between-class variation in repeater exposure within schools. The regression framework, which is described in detail below, thus compares the outcomes of regular students who attend kindergarten in the same school but who are randomly assigned to classes with and without repeating schoolmates. This identification strategy requires that these classes do not systematically differ from each other in any other dimension. In non-experimental data, this requirement will not be met if, for example, school principals assign low-achieving repeaters to classes with high-achieving other students or teachers. In contrast, random assignment in Project STAR ensures that classes with and without repeaters are balanced on characteristics of regular students and teachers.

One challenge to identification arises because repeater exposure is positively correlated with class size. In particular, repeaters are more likely to be observed in regular-sized classes because (i) larger classes are more likely to contain at least one repeater when students are randomly assigned to classes, and (ii) the sample contains more regular-sized classes than small classes.<sup>12</sup> Previous analyses of Project STAR have documented large negative effects of class size on student outcomes (see Schanzenbach (2006) for an overview of these findings). Therefore, a regression of student performance on repeater exposure that does not control for class size will yield an estimate that is negatively biased. I avoid such bias by controlling for class size in all of my regressions. In Section 6 below, I also present results from specifications that allow the effects of repeater exposure to vary with class size.

those of male students, black students, and students eligible for free or reduced-price lunch (results are available upon request). This suggests that by focusing on repeaters, I may be more successful in identifying truly low-achieving peers than by simply categorizing students as low achievers based on their demographic background.

 $<sup>^{12}</sup>$ Consider, for example, a school with the typical configuration of one small class of 15 students and two regular-sized classes of 23 students. If this school contains one repeater (the mode among schools with positive numbers of repeaters), this repeater has a 46/61 probability of being assigned to a regular-sized class and a 15/61 probability of being assigned to the small class.

Section 4 reports estimates of the following empirical model:

$$y_{ics} = \alpha_s + \beta_1 \text{EXPOSURE}_{cs} + \beta_2 \text{SMALL}_{cs} + X_{ics} \gamma + \varepsilon_{ics}, \qquad (1)$$

where  $y_{ics}$  is a kindergarten or long-term outcome for non-repeating student i randomly assigned to kindergarten class c in school s, EXPOSURE<sub>cs</sub> is an indicator for whether student i's class contains at least one repeater, SMALL<sub>cs</sub> is an indicator for small class in kindergarten, and  $X_{ics}$  is a vector containing the five student demographic characteristics shown in Table 1. Because random assignment to classes took place within schools, the model also controls for a vector of school fixed effects ( $\alpha_s$ ).

#### 3.2. Evidence on random assignment of repeaters

The key identification assumption underlying the specification in equation 1 is that conditional on class size and school fixed effects, classes with and without repeaters do not differ systematically in any other dimension. Intuitively, this assumption holds here because of the random assignment of students and teachers to classes in Project STAR. This intuition is supported by evidence from previous studies of the experiment (e.g. Chetty et al., 2011; Cascio and Schanzenbach, 2015), which show that classes are balanced on a wide range of student demographics and teacher characteristics. Here, I complement this evidence by evaluating whether repeaters were indeed randomly assigned to classes within schools.

As a first test for random assignment, I checked whether the withinschool variation in repeater exposure observed in the data is consistent with a random allocation process. To that end, I performed a Monte Carlo simulation in which students were randomly assigned to classes within schools and in which the number and size of classes and the number of repeaters in each school were based on the actual data. I then computed the withinschool standard deviation in repeater exposure, which is a summary measure of the identifying variation used in this paper, in the re-randomized data. Across 1,000 replications, the median standard deviation was 0.381 with a narrow 90% empirical confidence interval of [0.369, 0.391]. This confidence interval comfortably contains the within-school standard deviation of 0.383 observed in the actual data. As a second test for random assignment, I regressed an indicator taking value 1 if the student is a repeater and 0 otherwise on school and class fixed effects (omitting one class per school to avoid collinearity). Following the intuition described in Chetty et al. (2011), if assignment to classes was indeed random, then class indicators should not predict predetermined repeater status in this regression. Consistent with this idea, the *p*-value from an *F*-test for the joint significance of the class fixed effects was 0.65, suggesting that repeater status is indeed balanced across classes.

Finally, I tested whether being exposed to a repeater predicts nonrepeating students' demographic characteristics. Online appendix Table B.1 reports results from regressions of the five demographic characteristics available in the data on the repeater-exposure dummy (panel A) and on the number of repeaters in class (panel B). All specifications in this table also control for school fixed effects. Across the ten regressions, the estimated coefficients on the measures of repeater exposure are small and, with one exception, not statistically significant at conventional levels. Overall, the evidence presented here strongly suggests that repeaters were indeed randomly assigned to classes within schools in Project STAR.

#### 4. Main results

#### 4.1. Effects on end-of-kindergarten test scores

I begin the empirical analysis by estimating the impact of repeater exposure on regular students' math and reading performance at the end of kindergarten. These short-term estimates serve as a benchmark for comparison with findings from the previous literature and with the estimates for long-term outcomes reported later on. Column 1 of Table 2 shows that being exposed to repeaters reduces regular students' math scores by 9.0% of a standard deviation on average. Column 2 adds controls for students' demographic background to this regression. Due to the random assignment of students to classes, these controls do not change the coefficient estimate for the repeater-exposure treatment, but they slightly improve its precision. Columns 3 and 4 show the corresponding results for reading scores. The estimated impact of repeater exposure in these specifications is also negative, but it is substantially smaller than that in the math regressions and not statistically significant at conventional levels.

The estimates in Table 2 reveal that sharing a kindergarten classroom with repeaters has detrimental effects on regular students' test performance in the short term. This finding is in line with the results reported in the previous literature, which documents negative spillover effects from low-achieving classmates (e.g. Figlio, 2007; Carrell and Hoekstra, 2010; Lavy, Silva, and Weinhardt, 2012). Taken at face value, it suggests that policies which separate low-achieving repeaters from first-time kindergarten students would greatly benefit the latter, who make up the vast majority of the student population in schools.

#### 4.2. Effects on post-kindergarten test scores

Previous analyses of peer effects in schools have focused almost exclusively on short-term spillovers like the ones reported in Table 2. However, in order to judge the overall efficacy of policies that change the student composition of classes and schools, it is important to know how these spillovers play out in the long term. The STAR data provide me with the unique opportunity to analyze such long-term impacts on students' test scores, their non-cognitive skills, and their long-term educational attainment. In this subsection, I present the results for test scores.

Figure 2 plots the estimated impacts of repeater exposure on regular student' math and reading scores for each grade level from kindergarten to eighth grade. Panel A reveals a rapid fade-out of the negative spillover effect from repeaters on math scores: already one year after kindergarten, the estimated impact turns slightly positive, and it never falls below zero again afterwards. Indeed, the magnitude of the repeater-exposure effect seems to rise over time, culminating in an estimate of a 6.0% of a standard deviation increase in eighth-grade math scores which is marginally statistically significant. Panel B shows point estimates for reading scores that are qualitatively similar, though generally smaller in size. Overall, the results in Figure 2 therefore point to an interesting pattern of substantial negative impacts (at least in math) of repeater exposure on test scores in the short term and a rapid fade-out of these effects later on.

#### 4.3. Effects on non-cognitive skills

A growing literature in economics documents the importance of noncognitive skills for success in life and argues that such skills are partly formed in school (e.g. Heckman, Stixrud, and Urzua, 2006; Chetty et al., 2011; Heckman, Pinto, and Savelyev, 2012). I analyze the impacts of repeater exposure in kindergarten on non-cognitive skills in Table 3. In stark contrast to the negative short-term effects on test scores discussed above, panel A shows large positive spillovers from repeaters on regular students' effort, initiative, value, and discipline in fourth grade, when these skills are first measured. Panel B reveals that these effects persist into eighth grade, the second and last point of measurement of these outcomes.

Panel C shows the estimated effect of repeater exposure on a summary index of non-cognitive skills. Following Kling, Liebman, and Katz (2007), this index is constructed by averaging the eight standardized fourth- and eighth-grade indices for each student and normalizing the resulting composite to have mean 0 and standard deviation 1.<sup>13</sup> Repeater exposure raises non-cognitive skills, as measured by the summary index, by a highly significant 11.7% of a standard deviation. In comparison, being assigned to a small rather than to a regular-sized kindergarten class is estimated to increase the index by only 4.3% of a standard deviation, an effect that is not statistically significant at conventional levels (not shown in table). Taken together, the results in Table 3 point to lasting positive impacts of repeater exposure in kindergarten on regular students' non-cognitive skills.

#### 4.4. Effects on long-term educational attainment

The scholastic outcomes of participants in Project STAR were last tracked at the end of high school through collection of data on high school GPA, high school graduation, and college-test taking. Table 4 reports estimates from regressions that relate these measures of long-term educational attainment to students' exposure to repeaters in kindergarten. Sharing a classroom with repeaters raises regular students' high school GPA by 0.6 points on a scale of 100 points (column 1) and increases their likelihood to graduate from high school by 2.1 percentage points (column 2). Strikingly, repeater exposure also increases their likelihood of taking a college entrance exam by 3.3 percentage points (column 3), which corresponds to a sizable 8% increase over the base rate of 41%. Finally, column 4 shows a

 $<sup>^{13}{\</sup>rm If}$  only fourth-grade or only eighth-grade non-cognitive skills are observed for a student, the average of the available skill variables is used.

highly significant positive impact of repeater exposure on a summary index of these three long-term outcomes.<sup>14</sup> Overall, the findings in Table 4 thus suggest that mixing repeaters and regular students in kindergarten benefits the latter in the long term, a conclusion that is very different from the one drawn based on short-term spillovers on test scores only.

#### 5. Discussion and mechanisms

#### 5.1. Summary of main results

The results in Section 4 reveal important spillovers from repeaters on the outcomes of their non-repeating kindergarten classmates. Students who are exposed to repeaters score worse on standardized tests in math and reading in the short term, but they catch up rapidly with their non-exposed peers after kindergarten. In contrast, there are positive effects of repeater exposure on non-cognitive skills both when these are first measured about three years after kindergarten and later on in eighth grade, pointing to persistent impacts.<sup>15</sup> Perhaps most strikingly, students who share a kindergarten classroom with repeaters show improved long-term educational attainment as measured by high school grades, high school graduation, and college-test taking. I now discuss possible mechanisms behind these results.

#### 5.2. Non-cognitive skills as a channel for long-term impacts

There is mounting evidence that non-cognitive skills formed early in life are a key determinant of long-term educational success (e.g. Heckman, Stixrud, and Urzua, 2006; Chetty et al., 2011; Heckman, Pinto, and Savelyev, 2012). Therefore, it is natural to hypothesize that the greater accumulation of such skills by repeater-exposed students is responsible for their improved long-term outcomes. An informal test of this hypothesis

<sup>&</sup>lt;sup>14</sup>To construct the summary index, each of the long-term outcomes is first standardized by subtracting its mean and dividing by its standard deviation. In a second step, the average of these standardized outcomes is then normalized to have mean 0 and standard deviation 1 across non-repeating students in the estimation sample. All available long-term outcomes are used for each student.

<sup>&</sup>lt;sup>15</sup>Similar patterns of only temporary effects on cognitive skills but permanent impacts on non-cognitive skills have previously been documented for a variety of other childhood interventions (e.g. Chetty et al., 2011; Heckman, Pinto, and Savelyev, 2012).

is to control for intermediate non-cognitive skills in a regression of longterm educational attainment on repeater exposure in kindergarten. If noncognitive skills are indeed an important channel for the long-term spillovers documented above, then the coefficient on repeater exposure should be substantially attenuated in this augmented regression. I present estimates from such an informal test in Table 5.

Column 1 shows that non-cognitive skills measured in fourth and eighth grade are highly predictive of long-term educational attainment, corroborating findings from previous studies.<sup>16</sup> Column 2 reports the estimated effect of repeater exposure on long-term educational attainment for the subsample of students observed with non-cognitive skills. The coefficient in this regression is similar to the one reported for the unrestricted sample in Table 4, but due to the small sample size it is imprecisely estimated. Column 3 shows that this coefficient is reduced by 80% when non-cognitive skills are added to the regression as a control. Finally, column 4 reports that a Wald test rejects the null of equal coefficients in columns 2 and 3 with p<0.01. The evidence in Table 5 thus supports the hypothesis that the impact of repeater exposure on long-term educational attainment works mainly through the differential accumulation of non-cognitive skills by exposed students during childhood.<sup>17</sup>

#### 5.3. Mechanisms for impacts on non-cognitive skills

How exactly does exposure to repeaters in kindergarten raise regular students' non-cognitive skills? I consider two broad classes of explanations for these impacts. First, a systematic pairing of repeaters with particular students or teachers, either in kindergarten or in later grades, might mean that repeater-exposed students mechanically exhibit higher non-cognitive skills. Indeed, perhaps the most obvious interpretation of the results found above is that principals assign low-achieving and undisciplined repeaters to

<sup>&</sup>lt;sup>16</sup>For the sake of brevity, Table 5 and subsequent tables present estimates from regressions in which non-cognitive skills and long-term educational attainment are measured by the respective summary indices. Results are qualitatively similar if individual non-cognitive skill measures or long-term outcomes are used instead.

<sup>&</sup>lt;sup>17</sup>Note that the results do not let me rule out that repeater exposure improves longterm outcomes via raising different unobserved skills that are correlated with noncognitive skills. The evidence in Table 5 should therefore be interpreted as suggestive.

kindergarten teachers who are relatively better at teaching non-cognitive skills. Such a systematic allocation could notably explain both the temporary drop in test scores and the persistent improvements in skills such as effort and discipline. However, the random assignment of students and teachers to kindergarten classes in Project STAR means that this mechanism cannot be behind the results in this paper.

A closely related explanation involves a systematic selection of students into classes after kindergarten. For example, if repeater-exposed students consistently attend classes with better peers or teachers during grades 1-3, this could explain their differentially better non-cognitive skills at the beginning of fourth grade. Because students and teachers in Project STAR were randomly assigned to classes until the end of third grade, the main way for such systematic sorting to happen was via switching to another school (recall that 49% of students left the STAR schools before the end of the experiment). To test for differential switching by repeater-exposed students, I estimated the impact of repeater exposure on an indicator that takes value 1 if a student attrited from the experiment and 0 otherwise. The resulting coefficient of 0.004 indicates that repeater-exposed students were not more likely to switch schools during grades 1-3. Overall, there is thus little support for the idea that the reported gains in non-cognitive skills come about because of a systematic pairing of repeaters.

A second broad class of explanations relates to behavioral responses by teachers, students, or parents that affect the path of skill accumulation of repeater-exposed students. For example, teachers who are randomly assigned to kindergarten classes with undisciplined repeaters may react to this by teaching students primarily non-cognitive skills. Alternatively, students themselves may develop skills such as discipline and resilience as a response to classroom disruption by their repeating classmates. Finally, parents who are concerned about their children's learning progress in school may compensate for a worse classroom environment by helping their children at home. Unfortunately, the data do not let me distinguish between these individual mechanisms. However, given the evidence against a systematic pairing of repeaters, the general class of behavioral explanations seems to best fit the observed pattern of results.

#### 6. Additional results and robustness

#### 6.1. Heterogeneous effects by class size

One interesting question is whether the spillovers from repeaters documented in Section 4 vary with class size. For example, to the extent that smaller classes allow teachers to better respond to the individual needs of each student, the negative short-term impact of repeater exposure on test scores might be attenuated in these classes.<sup>18</sup> I explore such heterogeneity in Online appendix Table B.2, which reports estimates from regressions of four main outcomes in which the repeater-exposure treatment is interacted with the small-class indicator. Across all specifications, the estimated coefficients on the interaction term are imprecisely estimated and, with one exception, small relative to the main repeater-exposure effect. Thus, there is little evidence that spillovers from repeaters differ between small and regular-sized classes.<sup>19</sup>

#### 6.2. Alternative measures of repeater exposure

The main analysis of this paper distinguishes between classes with and without repeaters, but does not further differentiate classes according to the actual number of repeaters. In Online appendix Table B.3, I explore whether the results are sensitive to this particular definition of treatment. Panel A shows estimates from regressions of four main outcomes on separate indicators for being in class with one, two, and three to five repeaters. Across all specifications, the estimated impacts of exposure to one and exposure to two repeaters are qualitatively and quantitatively similar to the main effects reported in Section 4. While the coefficients on exposure to three to five repeaters are smaller in absolute value, they are very imprecisely estimated and not statistically different from these effects either.

<sup>&</sup>lt;sup>18</sup>Such an attenuated spillover effect in small classes would notably be consistent with the predictions of the widely cited theoretical paper by Lazear (2001).

<sup>&</sup>lt;sup>19</sup>I also confirmed that estimates are qualitatively similar, though less precise, when the entire empirical analysis is conducted separately for small and for regular-sized classes. Moreover, I examined whether spillover effects differ by regular students' demographic background. In regressions that interacted repeater exposure with demographic characteristics, there was suggestive evidence that black students, male students, and students on free lunch suffered larger declines in test scores in the short term if exposed to a repeater, and that they benefited less from this treatment in the long term. However, all of these interactions were imprecisely estimated.

Panel B shows that using the class share of repeaters as treatment again yields results that are qualitatively similar to the main results. Overall, the estimates in this paper are therefore not very sensitive to the particular definition of the treatment variable.

#### 6.3. Robustness to selective attrition based on repeater exposure

One potential concern with the results in Section 4 is that some of the long-term outcomes are observed only for a subset of the students who attended kindergarten in a STAR school. If the follow-up rates for these outcomes differ between students who were exposed to repeaters and those who were not, this could lead to a mechanical bias in the corresponding estimates. Particularly worrisome is the possibility of a "healthy survivor effect:" if students who were negatively affected by repeaters in the short term are less likely to be observed with long-term outcomes, this could explain some of the positive long-term spillover effects documented above.<sup>20</sup> In Online appendix Table B.4, I provide evidence that this mechanism is unlikely to drive my results.

As a first test for selective attrition, I estimated the effect of repeater exposure on indicators for being observed with six key outcomes. As panel A shows, the resulting follow-up differentials are small and not significantly different from zero in all specifications, suggesting that repeater-exposed students are not more likely to attrit from the sample. The results in Section 4 might still be biased, however, if the composition of the observed exposed or non-exposed students changed over time. I tested for such compositional changes by adding interactions between the repeater-exposure treatment and the five demographic characteristics to the specifications in panel A. The results, which are reported in panel B, show that the estimated coefficients on both the main exposure effect and the interaction terms tend to be small and are always jointly insignificant. This suggests that there are no systematic differences between exposed and non-exposed students observed with long-term outcomes.

Finally, I checked whether the negative impact of repeater exposure on kindergarten math scores can be replicated in the subsamples of students

 $<sup>^{20}</sup>$ Notably, such selective attrition cannot explain the positive impact of repeater exposure on college-test taking, which does not have missing values by construction.

observed with each of the outcomes studied in this table. Panel C presents the corresponding results. Across the six specifications, the coefficients on the repeater-exposure treatment are qualitatively and quantitatively similar to that in column 2 of Table 2, even though some of them are less precisely estimated due to reduced sample sizes. Overall, the results in this table do not support the notion that selective attrition based on repeater exposure in kindergarten biases the results in this paper.

#### 6.4. Robustness to relative measurement of non-cognitive skills

Table 3 reports positive impacts from repeater exposure in kindergarten on regular students' non-cognitive skills. A potential concern with these findings is that these improvements might simply reflect higher teacher ratings of students' behavior *relative* to the behavior of repeaters in the same class. I address this concern in Online appendix Table B.5. In panel A, I re-estimate the impacts of repeater exposure on fourth-grade non-cognitive skills for the subsample of students whose fourth-grade classes did not contain any of the 193 original kindergarten repeaters. The effects of repeater exposure in these regressions are somewhat attenuated compared to those reported in Table 3 but qualitatively similar.

The data do not allow me to observe classroom composition during eighth grade. However, I can restrict the sample to students who at that time attended a *school* that did not contain any of the original repeaters (most students had switched to a different middle school by eighth grade). Panel B shows that the estimated impacts of repeater exposure on noncognitive skills in this restricted sample are very similar to the ones reported in Table 3. Therefore, the evidence does not support the idea that the positive impacts of repeater exposure on non-cognitive skills capture purely mechanical effects due to relative teacher ratings.

#### 6.5. Testing for mechanical spillover effects

In a recent paper, Angrist (2014) documents a mechanical bias in peereffects regressions that arises if students both provide treatment for other students and are subject to treatment from these other students themselves. Intuitively, this bias is avoided in this paper due to the clear separation of initiators and recipients of spillover effects. I confirmed this intuition in a simulation-based falsification test similar to the one developed by Feld and Zoelitz (2014). In particular, I exchanged each student's classmates with a new set of peers randomly drawn from other classes in the same school. In this way, all students were assigned to a group of placebo classmates with whom they did not interact in their real-world classroom. I then re-estimated the effect of repeater exposure, measured using the placebo classmates, on kindergarten math scores. Any effect of repeater exposure in this regression reflects purely mechanical forces. In 1,000 replications of this exercise, the median coefficient on repeater exposure was 0.017 with a 90% empirical confidence interval of [-0.028, 0.063], which excludes the coefficient of -0.090 found in the actual data. Thus, the mechanical forces described by Angrist (2014) do not bias the results in this paper.

#### 7. Conclusion

Many education policies change the grouping of students into classes and schools, but little is known about the long-term impacts of school peers. This paper provides some of the first evidence on such impacts by evaluating how sharing a kindergarten classroom with low-achieving repeaters affects first-time students' test scores, their non-cognitive skills, and their long-term educational attainment.

The empirical analysis exploits the random assignment of teachers and students to classes in Project STAR in order to estimate causal spillover effects. Regular students who are exposed to repeaters in their kindergarten class perform worse on standardized tests at the end of kindergarten. However, these students display substantially improved non-cognitive skills, such as effort and discipline, when these are first measured at the beginning of fourth grade. While the negative spillovers from repeaters on test scores fade out rapidly after kindergarten, the positive spillovers on non-cognitive skills persist over time. The favorable development of repeater-exposed students culminates in significantly raised propensities to graduate from high school and to take a college entrance exam around the age of eighteen. My analysis suggests that these positive long-term impacts are likely due to the differential accumulation of non-cognitive skills by exposed students, which in turn appears to be a result of behavioral adjustments by teachers, students, or their parents. The striking divergence of the impacts of repeater exposure on shortterm test scores and long-term educational attainment highlights the importance of studying the long-term effects of educational interventions. By themselves, the negative short-term spillovers on test scores would have suggested that policies which separate low-achieving repeaters from regular first-time students would greatly benefit the latter, who make up the vast majority of the student population in schools. However, this conclusion has to be reversed once long-term impacts are taken into account. Indeed, the overall results show that mixing students of very different abilities at an early age can be beneficial for most students in the long term.

#### References

- Angrist, J.D. 2014. "The Perils of Peer Effects." Labour Economics 30:98– 108.
- Bifulco, R., J. Fletcher, and S. Ross. 2011. "The Effect of Classmate Characteristics on Post-Secondary Outcomes: Evidence from the Add Health." *American Economic Journal: Economic Policy* 3:25–53.
- Black, S.E., P.J. Devereux, and K.G. Salvanes. 2013. "Under Pressure? The Effect of Peers on Outcomes of Young Adults." *Journal of Labor Economics* 31:119–153.
- Burke, M.A., and T.R. Sass. 2013. "Classroom Peer Effects and Student Achievement." *Journal of Labor Economics* 31:51–82.
- Carrell, S.E., and M.L. Hoekstra. 2010. "Externalities in the Classroom: How Children Exposed to Domestic Violence Affect Everyone's Kids." American Economic Journal: Applied Economics 2:211–228.
- Cascio, E.U., and D.W. Schanzenbach. 2015. "First in the class? Age and the Education Production Function." *Education Finance and Policy* forthcoming.
- Chetty, R., J.N. Friedman, N. Hilger, E. Saez, D.W. Schanzenbach, and D. Yagan. 2011. "How Does Your Kindergarten Classroom Affect Your Earnings? Evidence from Project STAR." *The Quarterly Journal of Economics* 126:1593–1660.
- Deming, D., and S. Dynarski. 2008. "The Lengthening of Childhood." Journal of Economic Perspectives 22(3):71–92.

- Feld, J., and U. Zoelitz. 2014. "Understanding Peer Effects: On the Nature, Estimation and Channels of Peer Effects." mimeo.
- Figlio, D.N. 2007. "Boys Named Sue: Disruptive Children and Their Peers." *Education Finance and Policy* 2:376–394.
- Finn, J.D., J. Boyd-Zaharias, R.M. Fish, and S.B. Gerber. 2007. "Project STAR and Beyond: Database User's Guide." Report, HEROS Incorporated.
- Gottfried, M.A. 2013. "The Spillover Effects of Grade-Retained Classmates: Evidence from Urban Elementary Schools." American Journal of Education 119:1–64.
- Gould, E.D., V. Lavy, and M.D. Paserman. 2009. "Does Immigration Affect the Long-Term Educational Outcomes of Natives? Quasi-Experimental Evidence." *The Economic Journal* 119:1243–1269.
- Graham, B. 2008. "Identifying Social Interactions through Conditional Variance Restrictions." *Econometrica* 76:643–660.
- Heckman, J., R. Pinto, and P. Savelyev. 2012. "Understanding the mechanisms through which an influential early childhood program boosted adult outcomes." *American Economic Review* 103:2052–2086.
- Heckman, J., J. Stixrud, and S. Urzua. 2006. "The Effects of Cognitive and Noncognitive Abilities on Labor Market Outcomes and Social Behavior." *Journal of Labor Economics* 24:411–482.
- Hill, A.J. 2014. "The Costs of Failure: Negative Externalities in High School Course Repetition." *Economics of Education Review* 43:91–105.
- Hoxby, C. 2000. "Peer Effects in the Classroom: Learning from Gender and Race Variation." NBER Working Paper No. 7867.
- Kautz, T., J.J. Heckman, R. Diris, B.t. Weel, and L. Borghans. 2014. "Fostering and Measuring Skills: Improving Cognitive and Non-Cognitive Skills to Promote Lifetime Success." OECD Education Working Paper No. 110.
- Kling, J.R., J.B. Liebman, and L.F. Katz. 2007. "Experimental Analysis of Neighborhood Effects." *Econometrica* 75:83–119.
- Krueger, A.B. 1999. "Experimental Estimates of Education Production Functions." The Quarterly Journal of Economics 114:497–532.

- Krueger, A.B., and D.M. Whitmore. 2001. "The Effect of Attending a Small Class in the Early Grades on College-Test Taking and Middle School Test Results: Evidence from Project STAR." *The Economic Journal* 111:1– 28.
- Lavy, V., M.D. Paserman, and A. Schlosser. 2012. "Inside the Black Box of Ability Peer Effects: Evidence from Variation in the Proportion of Low Achievers in the Classroom." *The Economic Journal* 122:208–237.
- Lavy, V., O. Silva, and F. Weinhardt. 2012. "The Good, the Bad, and the Average: Evidence on Ability Peer Effects in Schools." *Journal of Labor Economics* 30:367–414.
- Lazear, E.P. 2001. "Educational Production." The Quarterly Journal of Economics 116:777–803.
- Sacerdote, B. 2011. "Peer Effects in Education: How Might They Work, How Big Are They and How Much Do We Know Thus Far?" Elsevier, vol. 3 of *Handbook of the Economics of Education*, pp. 249 – 277.
- Schanzenbach, D.W. 2006. "What have researchers learned from Project STAR?" Brookings Papers on Education Policy 9:205–228.
- Sojourner, A. 2013. "Identification of Peer Effects with Missing Peer Data: Evidence from Project STAR." *The Economic Journal* 123:574–605.
- Whitmore, D. 2005. "Resource and peer impacts on girls' academic achievement: Evidence from a randomized experiment." American Economic Review, pp. 199–204.
- Word, E., J. Johnston, H.P. Bain, D.B. Fulton, C.M. Achilles, M.N. Lintz, J. Folger, and C. Breda. 1990. "The State of Tennessee's Student/Teacher Achievement Ratio (STAR) Project: Technical Report 1985-1990." Report, Tennessee State Department of Education.

Figure 1 Distribution of repeaters across classes



*Notes:* The figure displays a histogram of the number of repeaters in class. The sample includes only the 60 (out of 79) schools with at least one repeater in kindergarten. There are 254 classes in this sample, with a mean (median) number of repeaters of 0.76 (1).

Figure 2 Repeater exposure in kindergarten and post-kindergarten test scores



*Notes:* The figure plots regression coefficients and 90% confidence intervals obtained from 16 variations of the specifications in columns 2 and 4 of Table 2. The dependent variable in each regression is the math score (panel A) or reading score (panel B) in the grade indicated on the horizontal axis. No results are reported for fourth grade because test scores are available for only a small fraction of students in that grade; see Online appendix A for details. See the notes to Table 2 for information about additional controls included in each of the regressions.

	No	on-repeat	ers		Repeaters	3
	N	Mean	SD	N	Mean	SD
Demographic characteristics						
Male	$6,\!039$	0.51	0.50	193	0.70	0.46
Black	$6,\!039$	0.33	0.47	193	0.17	0.38
Free lunch	$6,\!039$	0.48	0.50	193	0.65	0.48
Age in years	$6,\!039$	5.48	0.31	193	6.39	0.31
Old for grade	$6,\!039$	0.03	0.17	193	1.00	0.00
Repeater exposure						
At least 1 repeater in class	$6,\!039$	0.39	0.49	_	—	_
Standardized test scores						
Kindergarten math score	$5,\!614$	0.00	1.00	175	-0.36	0.80
Kindergarten reading score	$5,\!535$	0.00	1.00	173	-0.47	0.69
8th-grade math score	$4,\!353$	0.00	1.00	102	-0.88	1.09
8th-grade reading score	$4,\!364$	0.00	1.00	108	-0.93	1.15
Non-cognitive skills						
4th-grade effort	$1,\!628$	0.00	1.00	32	-1.13	1.24
4th-grade initiative	$1,\!628$	0.00	1.00	32	-1.01	1.01
4th-grade value	$1,\!628$	0.00	1.00	32	-0.83	1.25
4th-grade discipline	$1,\!628$	0.00	1.00	32	-0.32	1.20
8th-grade effort	1,731	0.00	1.00	37	-0.50	1.09
8th-grade initiative	1,731	0.00	1.00	37	-0.43	0.91
8th-grade value	1,731	0.00	1.00	37	-0.36	1.17
8th-grade discipline	1,731	0.00	1.00	37	-0.29	1.06
Long-term outcomes						
High school GPA	$^{2,438}$	84.20	7.42	40	81.82	7.35
High school graduation	$2,\!955$	0.87	0.34	60	0.67	0.48
${\rm Took}{\rm ACT}/{\rm SAT}$	$6,\!039$	0.41	0.49	193	0.12	0.32

Table	e 1
Descriptive	statistics

*Notes:* The table reports descriptive statistics of key variables separately for the 6,039 non-repeating students and the 193 repeaters in the estimation sample. A student is considered old for grade if based on her age and Tennessee's kindergarten entry cutoff date of September 30 she would be expected to attend at least first grade in the 1985-86 school year. Repeater exposure is measured as an indicator taking value 1 if the student's kindergarten class contains at least one repeater and 0 otherwise. Repeater exposure is not defined for repeaters because this paper studies spillovers from repeaters on non-repeating students. The non-cognitive skill measures are indices summarizing teacher ratings of student behavior in four areas: effort, initiative, value, and discipline. All test scores and measures of non-cognitive skills are standardized to have mean 0 and standard deviation 1 across non-repeating students. High school GPA is measured on a scale from 0-100. Took ACT/SAT is an indicator for whether the student took either of these tests in 1998, when most students were in their final year of high school.

	Math	Math	Reading	Reading
	(1)	(2)	(3)	(4)
Repeater exposure	-0.090**	-0.090**	-0.014	-0.014
	(0.043)	(0.041)	(0.046)	(0.044)
Male		$-0.144^{***}$		$-0.175^{***}$
		(0.024)		(0.025)
Black		$-0.355^{***}$		$-0.249^{***}$
		(0.051)		(0.053)
Free lunch		$-0.411^{***}$		$-0.450^{***}$
		(0.029)		(0.029)
Age in years		$0.550^{***}$		$0.408^{***}$
		(0.044)		(0.048)
Old for grade		$-0.411^{***}$		$-0.346^{***}$
		(0.081)		(0.074)
Small class	$0.169^{***}$	$0.158^{***}$	$0.194^{***}$	$0.185^{***}$
	(0.045)	(0.043)	(0.043)	(0.042)
Observations	$5,\!614$	$5,\!614$	$5,\!535$	$5,\!535$

Table 2Repeater exposure in kindergarten and end-of-kindergarten test scores

Notes: The table reports estimates from regressions of end-of-kindergarten math and reading scores on the variables listed in rows and school fixed effects. Test scores are standardized to have mean 0 and standard deviation 1 across non-repeating students in the estimation sample. Repeater exposure is measured as an indicator taking value 1 if the student's class contains at least one repeater and 0 otherwise. Standard errors in parentheses allow for clustering at the class level. \* p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01.

	Effort (1)	Initiative (2)	Value (3)	Discipline (4)
		Panel A:	4th grade	
– Repeater exposure	$0.104^{*}$	0.025	$0.124^{**}$	0.142***
	(0.054)	(0.056)	(0.053)	(0.054)
Observations	$1,\!628$	$1,\!628$	$1,\!628$	$1,\!628$
		Panel B:	8th grade	
– Repeater exposure	0.169***	$0.105^{*}$	0.160***	$0.194^{***}$
	(0.054)	(0.056)	(0.051)	(0.052)
Observations	1,731	1,731	1,731	1,731
		Panel C: sur	nmary index	
Repeater exposure		0.11	7***	
		(0.0)	41)	
Observations		$2,\!5$	89	

## Table 3Repeater exposure in kindergarten and non-cognitive skills

Notes: The table reports estimates from regressions that relate students' non-cognitive skills in fourth and eighth grade to their exposure to repeaters in kindergarten. The outcome variables in panels A and B are indices summarizing teacher ratings of student behavior in four areas: effort, initiative, value, and discipline. The indices are standardized to have mean 0 and standard deviation 1 across non-repeating students in the estimation sample. The outcome variable in panel C is a summary index of non-cognitive skills that combines the available information from fourth and eighth grade for each student; see text for details on how this index is constructed. Repeater exposure is measured as an indicator taking value 1 if the student's kindergarten class contains at least one repeater and 0 otherwise. All specifications control for students' demographic background, an indicator for small class in kindergarten, and kindergarten school fixed effects. Standard errors in parentheses allow for clustering at the kindergarten class level. \* p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01.

	Table 4		
Repeater exposure in	kindergarten and long-term	educational	attainment

	High school GPA (1)	$\begin{array}{c} {\rm High\ school}\\ {\rm graduation}\\ (2) \end{array}$	${ m Took} \ { m ACT/SAT} \ { m (3)}$	Summary index (4)
Repeater exposure	$0.552^{*}$ (0.308)	$0.021^{*}$ (0.013)	$0.033^{**}$ (0.015)	$0.074^{***}$ (0.028)
Observations	2,438	2,955	6,039	6,039

Notes: The table reports estimates from regressions that relate students' educational attainment, measured at the end of high school, to their exposure to repeaters in kindergarten. See the notes to Table 1 for descriptions of the outcome variables in columns 1-3. See text for details on the construction of the summary index of long-term educational attainment used as outcome in column 4. Repeater exposure is measured as an indicator taking value 1 if the student's kindergarten class contains at least one repeater and 0 otherwise. All specifications control for students' demographic background, an indicator for small class in kindergarten, and kindergarten school fixed effects. Standard errors in parentheses allow for clustering at the kindergarten class level. \* p<0.10, \*\* p<0.05, \*\*\* p<0.01.

		Table 5			
Repeater expo	osure, non-cogniti	ve skills, and	long-term	educational	attainment

	Summary in	dex of long-term	attainment	Difference
_	(1)	(2)	(3)	[(2)-(3)]
Repeater exposure		0.060	0.012	$0.048^{***}$
		(0.042)	(0.038)	$[p{=}0.004]$
Non- $\cos$ . skills (index)	$0.408^{***}$		$0.408^{***}$	
	(0.019)		(0.019)	
Observations	$2,\!589$	$2,\!589$	$2,\!589$	

Notes: The table reports estimates from regressions that relate students' educational attainment to their exposure to repeaters in kindergarten and to their non-cognitive skills measured in fourth and eighth grade (columns 1-3). Repeater exposure is measured as an indicator taking value 1 if the student's kindergarten class contains at least one repeater and 0 otherwise. See text for descriptions of the summary indices of long-term educational attainment and non-cognitive skills. All specifications control for students' demographic background, an indicator for small class in kindergarten, and kindergarten school fixed effects. Regressions include all non-repeating students for whom non-cognitive skills are observed. Standard errors in parentheses allow for clustering at the kindergarten class level. The rightmost column reports results from a test of the null hypothesis that the coefficients on repeater exposure in columns 2 and 3 are equal. The *p*-value in brackets is based on a Wald test conducted after re-estimating the specifications in columns 2 and 3 using seemingly unrelated regression. \* p<0.10, \*\* p<0.05, \*\*\* p<0.01.

### ONLINE APPENDIX TO:

### The Long-Term Impacts of Low-Achieving Childhood Peers: Evidence from Project STAR

Jan Bietenbeck Lund University and IZA

October 2015

#### A. Data appendix

The Tennessee State Department of Education entrusted a consortium of researchers from four Tennessee universities and various state institutions with the planning and implementation of Project STAR. After the experiment ended, some researchers continued to collect data on outcomes of participating students. Finn et al. (2007) provide a detailed account of these data collection efforts. The Project STAR public use file, on which the empirical analysis in this paper is based, combines these data such that students can be followed throughout their scholastic careers until the end of high school. Additional data on test scores in grades 5-8 were generously provided to me by Diane Schanzenbach. In what follows, I discuss in detail how I constructed the outcome variables used in the empirical analysis.

Test scores. At the end of each school year from kindergarten through third grade, students in Project STAR were administered the grade-specific version of the Stanford Achievement Test. From fifth grade through eighth grade, students who were still residing in Tennessee took the Comprehensive Test of Basic Skills (CTBS) as part of a statewide testing program.<sup>1</sup> Both tests are standardized multiple-choice assessments with components in mathematics and reading and are graded centrally.

The public use file contains Stanford Achievement Test scores for all students who took these tests. However, it contains CTBS scores only for students who were on grade level, i.e. students who attended grade 5/6/7/8 in 1991/1992/1993/1994, respectively. This implies that test scores are not observed for a number of students who had been retained in grade by those years.<sup>2</sup> In contrast, the data supplied by Diane Schanzenbach contain CTBS scores for students who attended grades 5-8 in Tennessee in any year between 1990 and 1997. Test scores are provided as scale scores, which are

<sup>&</sup>lt;sup>1</sup>An unrepresentative subsample of students took the CTBS also in fourth grade, see Finn et al. (2007). Due to the selective nature of this subsample, I chose not to use fourth-grade test scores in the empirical analysis.

<sup>&</sup>lt;sup>2</sup>Note that students who were retained in grade at any point between kindergarten and third grade dropped out of the STAR cohort and therefore did not write the subsequent Stanford Achievement Tests. However, these students did write the CTBS in later grades as long as they stayed in Tennessee.

comparable across grade levels (Finn et al., 2007). In order to increase sample size, I define test scores for a given grade level as scores obtained in the school year in which participating students were supposed to be in that grade (e.g., eighth-grade scores are defined as scores obtained in 1994, even though some students were attending seventh grade in that year). Results are however robust to using only the test scores available in the public use file. I standardize all test scores to have mean 0 and standard deviation 1 across non-repeating kindergarten students in the estimation sample.

*Non-cognitive skills.* I obtain fourth-grade non-cognitive skill measures from a questionnaire administered to teachers of a random sample of participating students in November 1989. The questionnaire asked teachers to rate how often each student had engaged in 31 different behaviors over the last two to three months. Ratings were recorded on a scale from 1 ("never") to 5 ("always"), and ratings of 28 of these behaviors were consolidated into four indices. The effort index includes items such as whether a student is persistent when confronted with difficult problems, whether she completes her homework, and whether she gets discouraged easily when encountering an obstacle in schoolwork. The initiative index is based on such items as whether a student participates actively in classroom discussions, whether she does more than just the assigned work, and whether she often asks questions. The value index measures how much a student appreciates the school learning environment. Finally, the discipline index captures such characteristics as whether a student often acts restless, whether she needs reprimanding, and whether she interferes with peers' work.<sup>3</sup>

During the 1993-94 school year, eighth-grade math and English teachers of a different random subset of participants were asked about student behaviors on a similar though shorter questionnaire. Thirteen of these behaviors were again consolidated into four indices measuring each student's effort, initiative, value, and discipline. I first average these indices across math and English for each student, and then normalize each of the eight

<sup>&</sup>lt;sup>3</sup>Note that what the paper refers to as the "discipline index" is the inverse of the "index of non-participatory behavior" in the original data. See Finn et al. (2007) for a complete listing of the behaviors included in each of the indices.

fourth- and eighth-grade indices by subtracting its mean and dividing by its standard deviation (computed across non-repeating students in the estimation sample). Finally, I construct the summary index of non-cognitive skills by averaging the available normalized indices for each student and normalizing the resulting composite.

High school grade point average and graduation. Most students in Project STAR graduated from high school in 1998, and transcripts were gathered from selected high schools in 1999 and 2000. High schools were chosen for data collection based on the likelihood that STAR participants would attend them given the locations of students' last known middle schools. Course grades from transcripts were transferred to a scale from 0-100 if necessary, and separate GPAs for math, science, and foreign languages were computed and are available in the data. The empirical analysis in this paper uses overall GPA, defined as the average of the these three subject-specific GPAs, as an outcome variable.

Information on high school graduation was also derived from transcripts and cross-checked with data from the Tennessee State Department of Education in ambiguous cases. Nevertheless, graduation status could not be determined with certainty for all students. In these cases, which comprise 7% of the non-repeating students in the estimation sample, the data collectors made a best guess whether a student "probably graduated" or "probably dropped out" based on the available course grades, information on attendance, and additional information from the Tennessee State Department of Education. The variable used in the empirical analysis codes 2,378 students who graduated, 98 students who probably graduated, and 82 students who received a General Educational Development certificate as graduates, and 296 students who dropped out and 101 students who probably dropped out as dropouts.

College-test taking and summary index of long-term educational attainment. ACT/SAT test taking was recorded by Krueger and Whitmore (2001), who matched all students in STAR to the administrative records of the two companies responsible for these tests in 1998. The outcome variable used in the empirical analysis is an indicator that takes value 1 if a student took either of these college entrance exams in 1998 and 0 otherwise. The summary index of long-term educational attainment combines information on high school GPA and graduation and college-test taking by first standardizing each of these variables to have mean 0 and standard deviation 1 across non-repeating students in the estimation sample. The average of these standardized variables is then normalized by subtracting its mean and dividing by its standard deviation.

#### References

- Finn, J.D., J. Boyd-Zaharias, R.M. Fish, and S.B. Gerber. 2007. "Project STAR and Beyond: Database User's Guide." Report, HEROS Incorporated.
- Krueger, A.B., and D.M. Whitmore. 2001. "The Effect of Attending a Small Class in the Early Grades on College-Test Taking and Middle School Test Results: Evidence from Project STAR." *The Economic Journal* 111:1– 28.

#### **B.** Appendix tables

	Male	Black	Free	Age in	Old for
	(1)	$\langle 0 \rangle$	lunch	years	grade
	(1)	(2)	(3)	(4)	(5)
		Panel A: re	peater-exposi	ıre dummy	
Repeater exposure	-0.005	-0.001	0.004	0.001	-0.003
	(0.015)	(0.007)	(0.015)	(0.009)	(0.005)
Observations	$6,\!039$	$6,\!039$	$6,\!039$	$6,\!039$	6,039
		Panel B: nui	mber of repea	ters in class	
No. of repeaters / $100$	-0.910	$-0.791^{**}$	0.303	-0.008	-0.417
	(0.829)	(0.333)	(0.841)	(0.509)	(0.273)
Observations	$6,\!039$	$6,\!039$	$6,\!039$	$6,\!039$	6,039

## Table B.1Randomization tests

Notes: The table reports estimates from regressions that relate non-repeating students' demographic characteristics to their exposure to repeaters in kindergarten. Repeater exposure in panel A is measured as an indicator taking value 1 if the student's kindergarten class contains at least one repeater and 0 otherwise. Specifications in panel B include the number of repeaters in a student's class as treatment instead. All regressions also control for kindergarten school fixed effects. Standard errors in parentheses allow for clustering at the kindergarten class level. \* p<0.10, \*\* p<0.05, \*\*\* p<0.01.

	Kindergarten math score (1)	8th-grade math score (2)	Non-cog. index (3)	Long-term index (4)
Repeater exposure	-0.078	0.089**	0.119**	0.065**
	(0.048)	(0.039)	(0.050)	(0.032)
$\times$ small class	-0.043	-0.099	-0.007	0.031
	(0.095)	(0.064)	(0.084)	(0.056)
Observations	$5,\!614$	4,353	$2,\!589$	6,039

## Table B.2Heterogeneous effects by class size

Notes: The table reports estimates from regressions that probe for heterogeneous effects of repeater exposure by kindergarten class size. Repeater exposure is measured as an indicator taking value 1 if the student's kindergarten class contains at least one repeater and 0 otherwise; see text for details on the construction of the outcome variables. All specifications control for students' demographic background, an indicator for small class in kindergarten, and kindergarten school fixed effects. Standard errors in parentheses allow for clustering at the kindergarten class level. \* p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01.

	Kindergarten math score (1)	8th-grade math score (2)	$egin{array}{c} { m Non-cog.} \ { m index} \ (3) \end{array}$	Long-term index (4)
	Panel A: indic	ators for differen	t numbers of rep	eaters in class
1 repeater in class	-0.096**	0.070*	0.120***	0.072**
	(0.046)	(0.036)	(0.046)	(0.031)
2 repeaters in class	-0.092	0.039	$0.134^{**}$	$0.093^{**}$
	(0.070)	(0.053)	(0.065)	(0.041)
3-5 repeaters in class	-0.021	0.019	0.025	0.030
	(0.103)	(0.090)	(0.090)	(0.063)
Observations	$5,\!614$	$4,\!353$	$2,\!589$	$6,\!039$
	Pan	el B: linear share	of repeaters in o	class
Share of repeaters	-0.601	$0.3\overline{70}$	$1.045^{**}$	0.781**
	(0.483)	(0.406)	(0.445)	(0.310)
Observations	$5,\!614$	$4,\!353$	$2,\!589$	$6,\!039$

Table B.3Alternative measures of repeater exposure

Notes: The table reports estimates from regressions that probe the robustness of results to using alternative measures of repeater exposure. In panel A, the repeater-exposure dummy is replaced by dummies for 1, 2, and 3-5 repeaters in class. Specifications in panel B include the class share of repeaters as treatment instead. See text for details on the construction of the outcome variables. All regressions control for students' demographic background, an indicator for small class in kindergarten, and kindergarten school fixed effects. Standard errors in parentheses allow for clustering at the kindergarten class level. \* p<0.10, \*\* p<0.05, \*\*\* p<0.01.

	Robustness	s to selective attı	rition based on r	epeater exposure		
	Kindergarten math score (1)	8th-grade math score (2)	4th-grade non-cog. (3)	8th-grade non-cog. (4)	High school GPA (5)	High school graduation (6)
	Panel 4	A: outcome is an in	dicator for being o	bserved with the va	ariable in the colum	ı head
Repeater exposure	-0.011 (0.008)	0.009 (0.012)	-0.012 (0.015)	-0.020 (0.013)	0.011 (0.013)	0.007 (0.014)
	Panel 1	B: outcome is an in	dicator for being o	oserved with the va	ariable in the column	ı head
Repeater exposure	-0.016	0.033	0.016	-0.016	$0.040^{*}$	0.023
	(0.013)	(0.020)	(0.023)	(0.024)	(0.023)	(0.024)
$\times$ male	0.005		-0.013	-0.002	-0.017	0.014
-	(0.014)	(0.025)	(0.021)	(0.025)	(0.025)	(0.027)
× black	-0.021 (0.018)	-0.004 (0.032)	-0.017 (80 0)	0.034 (0.098)	0.026 (0.033)	-0.036 (0.034)
$\times$ free lunch	0.016	$-0.050^{*}$	-0.036	-0.024	$-0.068^{**}$	-0.040
	(0.016)	(0.028)	(0.027)	(0.030)	(0.029)	(0.030)
× age in years	0.001	-0.025	0.014	0.007	-0.042	-0.035
	(0.026)	(0.042)	(0.039)	(0.042)	(0.041)	(0.043)
$\times$ old for grade	0.049	0.012	-0.068	-0.048	0.023	0.080
	(0.046)	(0.079)	(0.080)	(0.077)	(0.074)	(0.082)
p-value (joint significance)	0.47	0.52	0.40	0.60	0.26	0.48
	Panel C: outc	ome is the kinderg	arten math score, c	uly students observ	ved with variable in	column head
Repeater exposure	-0.090**	-0.088**	-0.075	-0.046	-0.067	$-0.131^{***}$
4	(0.041)	(0.044)	(0.061)	(0.057)	(0.054)	(0.050)
<i>Notes:</i> The table reports results the dependent variable is an indic panels include the 6,039 non-repe- repeater-exposure dummy and the include all non-repeating students Regressions in all three panels als effects. Standard errors in parenth	from a series of regre- ator taking value 1 if ating students in the b five interaction term in the estimation sa o control for students eses allow for clusteri	ssion-based tests for the outcome in the estimation sample. is. In panel C, the d mple for whom the c s' demographic back and at the kindergart	: selective attrition $l$ column head is obse The last row in pan lependent variable is outcome in the colur ground, an indicator en class level. * $p<0$	y exposure to repeative for a given stud- el B reports <i>p</i> -value the end-of-kinderga in head and the end for small class in ki 10, ** p<0.05, *** p	there in kindergarten. lent and 0 otherwise. s from $F$ -tests for joi tren math score. Reg l-of-kindergarten math indergarten, and kind <0.01.	In panels A and B, Regressions in these nt significance of the ressions in this panel a score are observed. ergarten school fixed

Table B.4

	Effort (1)	Initiative (2)	$\begin{array}{c} \text{Value} \\ (3) \end{array}$	Discipline (4)
	Panel A: 4th grade, no repeaters in class			
Repeater exposure	$0.072 \\ (0.070)$	-0.006 (0.074)	$0.031 \\ (0.071)$	0.081 (0.068)
Observations	1,037	1,037	1,037	1,037
	Panel B: 8th grade, no repeaters in school			
Repeater exposure	$0.161^{*}$ (0.082)	$0.074 \\ (0.086)$	$0.170^{**}$ (0.071)	$0.227^{***}$ (0.084)
Observations	866	866	866	866

# Table B.5Robustness to relative measurement of non-cognitive skills

Notes: The table reports estimates from regressions that probe for measurement of non-cognitive skills relative to repeaters. In panel A (panel B), the sample is restricted to students whose fourth-grade class (eighth-grade school) did not contain any of the 193 original kindergarten repeaters. Repeater exposure is measured as an indicator taking value 1 if the student's kindergarten class contains at least one repeater and 0 otherwise. See the notes to Table 3 for descriptions of the outcome variables. All regressions control for students' demographic background, an indicator for small class in kindergarten, and kindergarten school fixed effects. Standard errors in parentheses allow for clustering at the kindergarten class level. \* p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01.