

The Causal Effect of School Reform: Evidence from California's Quality Education Investment Act*

Paul Burkander[†]

November 1, 2013

Abstract

Beginning in the 2007-08 school year, California's Quality Education Investment Act required schools selected via lottery to institute reforms including class size reduction, increased average teacher experience, and extra professional training. The act provided additional per-pupil funding for schools to meet these requirements. Conditional on known probabilities of selection, which differed across schools, treatment is uncorrelated with potential outcomes, allowing for non-parametric identification of the causal effect by inverse probability weighting. In the first fully-funded year of the program, math scores in 4th grade increased by 0.32 SD in the population of California school-grade averages, and by the second fully-funded year 5th grade math scores improved by 0.36 SD. By the third fully-funded year of the program, math scores in 2nd grade were 0.28 SD higher in the distribution of California school-grade averages, and 0.27 SD higher in 3rd grade. Selected schools did not increase teacher experience, and had 4 to 4.4 fewer students in the first fully-funded year in 4th and 5th grade. In kindergarten through 3rd grade class sizes were reduced later and less dramatically, by 3 to 4 students by the third fully-funded year, due primarily to unselected schools exiting California's previous class size reduction program. The timing of class size reductions and student achievement gains suggests class size was the driving factor.

JEL Classifications: H75, I21, J45.

Keywords: Class Size, Teacher Labor Markets, Education Production Function.

*I am grateful for guidance provided by Gary Solon, Todd Elder, Michael Conlin, and participants at the University of Michigan Causal Inference in Education Research Seminar. I am also grateful to Marsha Porte and Jim Alford at the California Department of Education for their invaluable assistance. The work reported here was supported in part by a Pre-Doctoral Training Grant from the Institute of Education Sciences, U.S. Department of Education (Award # R305B090011) to Michigan State University. The opinions expressed here are mine and do not represent the views of the Institute or the U.S. Department of Education.

[†]Department of Economics, Michigan State University, 110 Marshall Adams Hall, East Lansing, MI 48823. burkande@msu.edu

1 Introduction

Current educational policy in the United States is focused on increasing the proportion of students who meet state-determined proficiency levels on standardized tests. There is disagreement about how to achieve this, with some arguing for additional educational inputs, and others for more efficient use of existing inputs. An extensive literature on educational production functions¹ has attempted to resolve this and related questions. However, despite some random experiments and a plethora of natural experiments, no clear consensus has emerged on the question of whether marginal changes in educational resources have any effect on educational outcomes.

California's Quality Education Investment Act (QEIA) offers a unique opportunity to identify the effect of increased inputs on outcomes. In the 2007-2008² school year the QEIA went into effect, leading to increased funding and obligatory reforms for about 500 selected schools, which were chosen from 1,260 participating schools. Districts were first required to rank all of their participating schools, and then districts were randomly selected to have their highest ranked schools funded. Once selected, funded schools were required to institute several reforms, e.g., they had to reduce average class size, increase average teacher experience, and provide additional professional training to teachers.

Conditional on districts' rankings, selection of schools was random, though schools differed in their probability of selection. The selection process, and therefore the probabilities of selection, are known, and the average treatment effect of QEIA can therefore be non-parametrically identified using inverse probability weighting (IPW). A drawback of QEIA is that the effects of the individual reforms cannot be separately identified. However, bundled reforms are worth studying in their own right: pressure to improve outcomes often leads to concurrent policy changes so QEIA is reflective of how reforms are actually carried out;

¹Summaries of the assumptions and methods employed in the education production function literature can be found in Hanushek (1979), Todd and Wolpin (2003), and Rice and Schwartz (2008).

²Henceforth, school years are referred to by the year in which the Spring semester occurs. For example, the 2007-2008 school year is referred to as 2008.

it may also be the case that interactive effects cause bundled reforms to be more or less effective than the sum of their constituent parts.

Moreover, as I find, QEIA caused a reduction in class size of about 4 students per class by the third fully-funded year of the program, but had no discernible effect on the other main policy lever that I observe, teacher experience. Reportedly, the vast majority of elementary schools eligible to participate in QEIA were already required to meet many of its requirements, with the exception of reduced class size, increased teacher experience, and increased professional training. Also, continued participation in QEIA was contingent on schools meeting achievement growth targets; the evidence therefore suggests that the causal effect of QEIA on standardized test scores occurred through some mix of class size reduction, professional training, and increased pressure to raise test scores.

Indeed, QEIA did cause a statistically significant increase in student achievement, as measured by both California's Average Performance Index, and by grade-level results on California's primary standardized test. The API is a weighted school-level average across all tested subjects, grades, and test types. With respect to the population of all elementary schools, the average treatment effect of QEIA on the API by the third fully-funded year of the program was an increase of 0.33 standard deviations, with larger gains for Hispanic and low-SES students. With respect to the population of grade-level averages across all California schools, by the third fully-funded year of the program standardized math scores increased 0.28 standard deviations in 2nd grade, and by 0.44 standard deviations in 5th grade. QEIA caused more modest gains in English language arts, of 0.19 and 0.22 standard deviations in 2nd and 5th grade, respectively.

In what follows, section 2 reviews the relevant literature; section 3 describes QEIA in greater detail; section 4 reviews the data used in this analysis; section 5 outlines the identification strategy; section 6 presents the results. Section 7 concludes.

2 Literature Review

This analysis contributes causal evidence to the aforementioned extensive literature on education production functions, which generally has found mixed results. Meta-analyses that find no clear evidence of an effect of increased school inputs on student outcomes include Hanushek (1986), and Hanushek (1997), though the methods employed in those analyses are criticized by Krueger (2002). In contrast, Greenwald et al. (1996) provide a meta-analysis that finds many school inputs do have positive effects, though their methods are criticized by Hanushek (1996).

Within the education production function literature, this paper contributes most to those strands concerned with the effect of reducing class size, providing professional training to teachers, and increasing accountability. The study of class size effects on student achievement has a rich history, dating back at least a century. As noted by Rockoff (2009), early waves of the literature, which include field experiments as early as the 1920s, tended to find no effect from a reduction of class size.

Recent studies of variation in class size tend to be quasi-experimental, with the notable exception of Tennessee's Project STAR (Student/Teacher Achievement Ratio). Project STAR was a randomized control trial that assigned students to either small classes (13-17 students per class), regular classes (22-26 students per class), or regular classes with a teaching aide. The reduced class size treatment of Project STAR has generally been found to have had positive effects in the short run (Nye et al., 1999, Krueger, 1999), and in longer run outcomes (Krueger and Whitmore, 2001, Chetty et al., 2011a) though non-random attrition, lack of baseline measure of student performance, and little information about teachers and how they were randomized should give us pause in interpreting results (Hanushek, 1999).

Notable natural experiments include Angrist and Lavy (1999), which uses a regression discontinuity design based on a class size rule in Israel. They find significant returns to achievement from class size reduction for math and reading scores for 5th graders. Hoxby (2000) also uses variation in class size generated by class size caps, and exploits exogenous

variation in population, to analyze the effect of class size in Connecticut. She finds no returns to class size reduction, and in fact rules out even modest returns to class size reduction.

The results in Hoxby (2000) are questioned by Jepsen and Rivkin (2009), who note that, in using test scores from the following year, estimates may be attenuated. Jepsen and Rivkin (2009) analyze a previous class size reduction program in California, which was first implemented in 1996. Using a fixed-effects analysis, and with a school-level measure of achievement as the outcome variable, the authors find that a ten-student reduction in class size led to a 0.06 to 0.1 standard deviation improvement in Math, and a 0.04 to 0.6 standard deviation improvement in Reading. An important contribution by Jepsen and Rivkin (2009) is that, unlike previous class size analyses, they explore changes in teacher quality that result from rapidly reducing class size, finding that in the early years of the program class size returns were offset by losses from reduced teacher quality. This issue is explored further by Dieterle (2013), who finds that the reduction in teacher quality was not large enough to account for only modest returns to reduced class size reduction in an anonymous state. Chingos (2012) examines another class size reduction policy implemented in Florida, and finds no effect using a comparative interrupted time series design. Chingos (2012) exploits the fact that many districts already met the class size requirements of Florida's law when it was implemented. Districts that already met the requirement received the same increase in funding, so the counterfactual to increased funding and class size reduction is an unconstrained increase in funding.

The literature on effects of teacher professional development is less-well developed. Yoon et al. (2007) reviews over 1,300 studies conducted between 1986 and 2006, and found only nine to rigorously examine the effect of teacher professional development on student achievement. Among these, the range of effects was quite large, from -0.53 to 2.39 standard deviations,³ with the smallest and largest effect coming from the same study. The student assessment tools were generally closely aligned with teacher training, and within study there was wide

³Presumably the standard deviations are with respect to the population of students studied, though neither Yoon et al. (2007) nor the source material clarifies the point.

variation of effects. In a more recent study, Barrett et al. (2012) uses a propensity score model to test whether less effective teachers are more likely to select into professional development, and whether accounting for this selection affects estimates of effectiveness of such programs. They find that pre-treatment value added scores are an important predictor of participation, and controlling for selection professional development increases student test scores by 0.08 standard deviations in elementary school.

Linking incentives to test scores has been shown to improve medium term math outcomes (Chiang, 2009), and long-term outcomes of low-performing students (Cohodes et al., 2013). A growing body of literature considers potential erosion of signal quality of student assessments when those assessments are linked to incentives. There is evidence that schools manipulate the population of test takers (Figlio and Getzler, 2006, Jacob, 2005), shift resources towards marginal students (Neal and Schanzenbach, 2010), and that teachers manipulate test results (Jacob and Levitt, 2003). The QEIA link between incentives and test scores differs from those studied above in at least two ways: using the API, an average of scores across all students, instead of percent proficient removes the incentive to teach to the marginal student, and scores on tests for cognitively impaired students count toward a school's API. I nonetheless test below whether the population of test takers changes.

With conflicting results in meta-analyses and natural experiments, and few randomized control trials, it seems clear that after a century of research into education production functions, more research is needed. States such as California in 1996, Florida in 2003, and Ohio in 2009 have passed class size reduction laws, devoting resources toward increasing inputs that may or may not improve outcomes. If increased inputs can lead to improved output, it must be determined how to move closer to an optimal mix of inputs. To address these questions, more natural experiments with credible exogenous variation are needed. QEIA provides such credible evidence.

3 Policy Description

The QEIA was preceded in California by a larger and more ambitious class size reduction policy. That policy was enacted in 1996, and continues nominally to this day. Participation was voluntary, but incentivized: districts received \$650 in the first year⁴ per student in a K-3 class of 20 or fewer students. However, this incentive has diminished twice over time. In 2004 the maximum qualifying average class size was increased to just under 22, and as of February 2009 classes of 25 or more students are still eligible for 70% of the per-pupil funds, though funding is given for no more than 20 students.⁵

The QEIA came about as the consequence of litigation against then California Governor Arnold Schwarzenegger. The plaintiffs in the case argued successfully that the state paid less than the legislated minimum amount to kindergarten through 12th grade public schools in the 2005⁶ and 2006 school years. As a result, the state was required to pay back approximately \$2.7 billion to K-12 schools.

Rather than distribute the money equally across all schools, legislators decided to focus on a subset of low-performing schools. The subset was chosen on a semi-random basis, and the number of schools was chosen such that per-student funding would increase by \$500 in grades K-3,⁷ \$900 in grades 4-8, and \$1,000 in high school from 2009-2014, and by half as much in 2008.⁸

Schools were deemed eligible to participate in QEIA if they were in the bottom quintile of the state's 2005 academic performance distribution, as determined by the API.⁹ Eligible schools had to commit to meeting the requirements of QEIA before they could participate

⁴This number was adjusted for inflation in subsequent years.

⁵For instance, a class of 25 or more students would receive $0.70 \times 20 \times \text{Full Per-Pupil Amount}$

⁶Governor Schwarzenegger had reached an agreement with a coalition including the California Teachers Association to underfund education by \$2 billion below the amount guaranteed by Proposition 98, which pegs education funding to growth in general funds. However, state revenue exceeded expectations, and education funding was not updated to reflect this. For more information, see Bluth (2005).

⁷For comparison, in the first full year of QEIA funding the per-pupil funding for participation in California's existing class size reduction program was \$1,071.

⁸The reduced amount in 2008 was intended to give schools a chance to prepare for full implementation of reforms by 2009.

⁹Very small schools, whose API scores were considered unreliable, were excluded.

in the selection process. Schools could choose to participate in the regular QEIA program, or an alternative program. Schools in the alternative program, which are excluded from this analysis, were able to design their own reform plans, which had to be approved as part of the application to participate in QEIA. Of the 1,455 schools eligible to participate in QEIA, 1,260 chose to do so, and 88 of these chose to participate in the alternative program.

Each district with more than one participating school was required to rank its schools. It was permissible to give multiple schools the same rank, and indeed several districts did so. Districts received as many random numbers as they had participating schools, and these random numbers were assigned to each district's schools based on the district's rankings. For example, if a district with two schools received random numbers 213 and 314, the highest ranked school was assigned 213 and the second was assigned 314. If a district assigned the same ranking to multiple schools, the order was determined randomly within that ranking, and was done so by the California Department of Education. The selection then proceeded in four stages.

First, schools for the alternative program were selected. High schools were given priority for this program, and the number of schools was chosen such that no more than 15% of the anticipated number of students in funded schools would be in the alternative program. The high schools with the lowest random numbers in the alternative program were funded until this target was reached.¹⁰

Second, to ensure geographic diversity, in each county without a funded school from the first stage, the school with the lowest random number was selected. Districts were told that after schools were selected for the alternative program and geographic diversity, all schools with the lowest random numbers would be funded until funds were exhausted.¹¹

In fact, high schools were selected separately in the third stage: to ensure the legislatively

¹⁰Several middle and elementary schools applied for the alternative program, but given the number of high schools that applied they effectively had zero probability of being chosen.

¹¹The actual selection differed somewhat, as described below. That districts were told this simplified version is evidenced in contemporaneous school board minutes (Santa Rosa City Schools, 2007), and CDE presentations (Balcom, 2007). This is also the depiction in the report to CA legislature (CDE, 2010), written 3 years after the selection process.

mandated fair representation of grade spans, a target number of high school students was selected so that the proportion of high school students in funded schools would be roughly equivalent to the proportion of high school students in all participating schools. The high schools with the lowest random numbers were selected until this target was reached. Any school with at least one high school student in 2007 was considered a high school for this purpose. Finally, the elementary and middle schools with the lowest random numbers were selected until QEIA funds were exhausted.¹²

At the conclusion of the selection process, 25 schools had been selected for the alternative program and 463 for the regular program. One school that was selected immediately withdrew from the program, and in subsequent years 13 schools were added. For the purpose of this analysis, I consider all schools initially selected to be treated, and all participating schools not selected to be the control group. Additionally, I restrict the sample to elementary schools,¹³ which account for over 70% of schools participating in the regular QEIA program.

Funded schools were required to implement the following: reduce class size; align average teacher experience with their district average; ensure that all teachers in the school be considered Highly Qualified Teachers (HQT) under California's Elementary and Secondary Education Act; satisfy the requirements of the *Williams* settlement, which required schools to provide qualified teachers and safe, well-maintained facilities; provide professional training to teachers and paraprofessionals; and, for high schools, increase the counselor-student ratio.

According to CDE (2010), the vast majority of schools eligible to participate in QEIA were already required to meet the HQT standard and the requirements of the *Williams* settlement, regardless of whether they were selected to be funded. This claim is substantiated by Table

¹²As a result of this process, and unbeknownst to districts prior to selection, the funding results did not always follow district rankings. For instance, a highly ranked high school could go unfunded, and a lower ranked elementary school could be funded. Ranking a high school ahead of an elementary school could also lead to both not being funded, while if the elementary were ranked higher it would be funded, if the difference in random numbers is sufficiently large.

¹³This restriction has two motivations: there are additional QEIA requirements for high schools, complicating the interpretation of the treatment, and beginning in 6th grade, students are sorted into various math examinations, thus compositional changes may be confounded with changes in achievement. Results that include middle and high schools are qualitatively quite similar, and are available in the online appendix.

1, which shows that the typical participating elementary school had 94% of its teachers classified as Highly Qualified, and 95% of participating schools were required to satisfy the terms of the *Williams* settlement.

The class size reduction requirement stipulated that funded schools reduce class size to 20 students per class in grades K-3.¹⁴ In grades 4-12, class sizes had to be reduced from their baseline level¹⁵ by 5 students, or to 25 students per class, whichever was lower. In each of the first three years of QEIA, schools were required to reduce the difference between the pre-QEIA average class size and QEIA target class size by 1/3. For some schools, the average in 2007 was quite low, which was particularly strenuous for small schools with a single classroom per grade. As such, many schools applied for and were granted waivers from this requirement, and instead met a higher minimum class size requirement.

Under QEIA, teacher experience is measured by the Teacher Experience Index (TEI). Teachers with more than 10 years of experience are assigned 10 years in calculating the average. Part-time teachers are given full weight in the calculation, and teachers teaching at multiple schools count towards each school's average. Funded schools are required to exceed the district average TEI.

Districts selected for QEIA are required to provide professional development opportunities for teachers, administrators, and paraprofessionals, e.g., teaching assistants. Funded schools are required to build and maintain a system for tracking participation in professional development programs, and districts are required to ensure that funded schools are in fact meeting the requirements. Participation requirements for teachers are clearly spelled out by QEIA, e.g., each year at least one third of teachers in a QEIA funded school must participate in training, but the specifics of the training program are largely left to the schools and districts.

In addition to these reforms, participation in QEIA was contingent on meeting accelerated

¹⁴This is precisely the original requirement of California's 1996 class size reduction policy, the maximum cap of which increased over time.

¹⁵The baseline was the grade-level average class size in 2006, unless that average was greater than 25, in which case 2007 was used

student achievement growth targets, as measured by California's API. The target API for all schools in California is 800; all California schools below that target have a growth target of 5% of the difference between their API and 800, or 1 point, whichever is greater. By the third year of QEIA, funded schools are required to have exceeded growth targets on average over those first three years. A school is permitted to fall short of its growth target in the first two years of full funding without repercussions, then after the third fully-funded year schools whose average growth did not exceed average growth targets lost QEIA funding.

4 Data

This analysis relies on several publicly available data sets produced by the California Department of Education. These include school-level data on API scores, a data set that includes a rich set of school demographics; teacher-level data; assignment-level data, including, e.g., the number of classes assigned to a teacher and the number of students in each of those classes; and subject-grade-level data on California standardized tests. Though these data sets are available for earlier years as well, I rely primarily on data from 2005-2011 with one exception: the assignment-level data were not collected in 2010 due to budget constraints. As a result, I am unable to calculate average class size, proportion of teachers classified as HQT, or TEI for 2010. In addition to these publicly available data, I have obtained from CDE the rankings of participating schools submitted by each district, which include a variable for whether the school applied for the regular or alternative program.

The teacher-level data are not linked from year to year. Instead, in each year teachers are assigned a new ID. The purpose of the ID is to facilitate linking teacher-level data to assignment-level data. The teacher-level data do however contain a number of teacher characteristics, such as years of teaching experience, years teaching in the same district, ethnicity, gender and education. I use these characteristics to link teachers across years within a school. If multiple teachers at a school are observationally similar, I randomly link

them across years.

As a result, I do not reliably observe duration of employment spells at a school. Similarly, if a teacher leaves, and in the following year a new observationally identical teacher enters the school, I do not observe a change in the composition of teachers. The data can however be used to reliably determine net changes in the characteristics of the teacher workforce at a school. I use these data to measure average teacher experience, the proportion of teachers new to a school, and the proportion of teachers new to a school who are either new to teaching, or experienced but new to the district.

For my measure of class size, I restrict the set of classes to math, English, science, and self-contained classes. Self-contained classes are those in which subjects such as math and English are taught by the same teacher, and are the most common class type in elementary schools. This analysis excludes special education courses, vocational courses, and other electives.¹⁶

It has become common in the education production function literature to use student performance on standardized tests as a measure of the output of this production process. Standardized tests surely fail to capture a number of cognitive and non-cognitive skills that an educational system is expected to impart on students. However, there is evidence that variation in school inputs that increase test scores also have a positive impact on a number of later-life outcomes, such as probability of attending college, selectivity of college, and income (Chetty et al., 2011a, Chetty et al., 2011b). Often a student's performance on standardized tests is the outcome in a regression including measures of scholastic inputs and the student's performance in previous years as controls. The use of California's API in this analysis is similar, but it differs from student level assessment scores in important ways.

Notably, the API is an average of performance not just across students, but across subjects and even test types. For instance, an elementary school in 2010 would have administered

¹⁶Teachers for these excluded classes are included in the teacher experience category, in part so my of experience is not dependent on data missing in 2010. The TEI is based on a subset of classes similar to that which I use to calculate average class size.

an English and language arts test in grades 2-5, a math test in grades 2-5, and a science exam in grade 5. Additionally, two alternative exams, the California Modified Assessment and the California Alternative Performance Assessment would have been administered to students with varying degrees of cognitive impairment. The API for that school is a weighted average across all these tests, subjects, and grades.¹⁷ Nonetheless, the API is California's primary tool for assessing academic performance, and the goal of QEIA was to improve API scores, so I include it in my analysis. I standardize API scores within years with respect to the distribution of all elementary school APIs.

I supplement this measure of student achievement with grade-subject-level data on California's primary standardized test, the eponymous CST. California makes publicly available the mean scaled score,¹⁸ and percent of students whose scores fall into particular bins, referred to as proficiency levels. I use these data for 2nd-5th grade math and English language arts tests.

Table 1 lists descriptive statistics for all funded and unfunded elementary schools in 2007, and those for which p_i , the probability a school is selected, lies between 0.10 and 0.90, as well as descriptive statistics for the restricted sample in 2011. The restricted sample is similar to the full sample with two notable exceptions: a much smaller proportion of schools in the restricted sample are in Los Angeles, and funded schools in the restricted sample have a higher TEI.

Both funded and unfunded schools in QEIA typically had high proportions of students who were Hispanic, English language learners, eligible for free and reduced price lunch, and whose parents did not have a college degree. In 2007, a typical school in my sample had at least 1/3 of its teacher work force that was not in the school the prior year.¹⁹ New teachers

¹⁷The average is weighted by the proportion of students for whom there is a valid score, and each subject and test receives an additional weight.

¹⁸The scaling of scores takes into account changes in difficulty of tests across years, and therefore makes yearly comparisons more meaningful.

¹⁹Recall that my teacher-level data set can only distinguish net changes in teacher characteristics. New teachers who are observationally identical to departing teachers from the previous year are not recorded as new, and thus the one third estimate is a lower bound.

did however tend to have nearly three years of experience.

From Table 1 it is apparent that the reduction in class size in grades kindergarten through third is driven by increased class sizes in unfunded schools, while the reduction in class size in fourth and fifth grade is driven by smaller class sizes in funded schools. The QEIA requirement for class sizes in kindergarten through third grade replicated that of California's prior class size reduction policy, the incentives for which were drastically weakened in the first year of QEIA. This weakened incentive led many unfunded schools to gradually increase class sizes in lower grades.

5 Identification

Estimating the causal effect of QEIA is complicated by two facts: districts ranked schools according to unobserved objectives, and districts with more participating schools were more likely to be chosen at least once. A simple comparison of funded and unfunded schools within a district would surely be biased, though the direction of bias would depend on the district's objective functions. A comparison across even just the highest ranked schools in each district would also likely be biased, since larger districts, e.g., Los Angeles Unified, were almost certain to have their highest ranked schools funded, and the size of a district could be correlated with potential outcomes. Even within school over time, potential outcomes might be correlated with treatment if districts gave higher rankings to schools poised to improve.²⁰

Instead, my estimation strategy relies on the following intuition: if we were to compare only schools that had an equal probability of being funded, e.g., 50%, then within that group treatment is random, and an OLS estimate would be consistent and unbiased. For each probability we could repeat this exercise, yielding treatment effects conditional on each probability. By the Law of Iterated Expectations, the unconditional average treatment effect could then be recovered. As Wooldridge (2004) shows, the result of an exercise like this is

²⁰There is anecdotal evidence that this did in fact happen: in personal communication with a CDE employee who was on Sacramento's school board when schools were ranked, I was told that one school ranked highly because, with a new and talented principal soon starting there, it was poised to improve.

equivalent to the following population specification for τ_{ATE} , the average treatment effect:²¹

$$\tau_{\text{ATE}} = \text{E} \left(\frac{T_i y_i}{p_i} - \frac{(1 - T_i) y_i}{(1 - p_i)} \right) = \text{E} \left(\frac{(T_i - p_i) y_i}{(1 - p_i) p_i} \right) = \text{E}(y_{i1} - y_{i0}) \quad (1)$$

where T_i is an indicator for treatment, in this case being funded by QEIA; y_i is an outcome measure; y_{i0} is the outcome for school i if it is not selected, and y_{i1} is the outcome for school i if it is selected; and p_i is the probability of selection, i.e., the propensity score. $\text{E}(y_{i1} - y_{i0})$ is the Average Treatment Effect: it captures the average change in outcome caused by QEIA. The parameter in equation (1) can be estimated using its sample analog.

Given the selection mechanisms, determining the functional form of p_i is complex.²² However, given the rules of selection and districts' rankings, I am able to determine the true propensity score (up to an arbitrarily small error) by simulation; I do so by randomly assigning the numbers 1-1,260 to districts and replicating the selection process 10,000 times.²³

This method allows for the causal effect of QEIA to be non-parametrically identified if two assumptions are satisfied. First, treatment must be mean independent of the potential outcomes conditional on the propensity score (i.e., $\text{E}(y_j | T, p) = \text{E}(y_j | p)$, $j \in \{0, 1\}$). This requirement is satisfied by the nature of the selection process. The second requirement is that there can be no schools for which $p_i = 1$ or $p_i = 0$. The intuition for this requirement is straightforward. Among the schools for which $p_i = 0$ or $p_i = 1$, there is no variation in treatment, and so these schools contribute nothing to identification. Among schools participating in QEIA, some were in counties with only one participating school, and that

²¹Wooldridge (2004) shows that the coefficient in a “conditional linear projection” of outcome on treatment, where the conditioning is on probability of selection, can be averaged across probabilities to yield this form of the average treatment effect. He also notes several alternative and asymptotically equivalent forms of the estimator. The estimator is similar to that used in Horvitz and Thompson (1952). See also Imbens and Wooldridge (2009) and Wooldridge (2010).

²²Were the total number of high schools and elementary or middle schools predetermined, the problem would be considerably simpler, and p_i would be based on a summation of hypergeometric functions, weighted by the probability that the district has a school selected for geographic diversity. Since the number of schools selected depended ultimately on the number of students in each grade level in each school, the problem is considerably more complicated.

²³The actual random numbers assigned to districts are also made publicly available by CDE. Using these, and the district rankings, my simulation of the selection process perfectly predicts funded schools.

school was therefore selected with probability one. Conversely, the middle and elementary schools that applied for the alternative program had zero probability of being selected. There were also many schools, e.g., Los Angeles Unified’s highest ranked schools, whose probability of selection was near one, and many, e.g., Los Angeles Unified’s lowest ranked schools, whose probability of selection was very near zero.

In practice, researchers drop observations with probability of treatment “close” to zero or one. Crump et al. (2009) suggest discarding observations less than α away from zero or one, where α satisfies the following:

$$\frac{1}{\alpha(1-\alpha)} = 2 * E \left[\frac{1}{p_i(1-p_i)} \mid \frac{1}{p_i(1-p_i)} < \frac{1}{\alpha(1-\alpha)} \right]. \quad (2)$$

As a general rule of thumb, Crump et al. (2009) suggest using $\alpha = .10$. After dropping schools for which p_i is identically zero or one, I am able to calculate (2), and $\alpha = .10$ nearly satisfies this requirement exactly. I therefore restrict the sample to schools for which $p_i \in [0.10, 0.90]$.

To examine whether funded and unfunded schools share a common support across p_i , Figure 1 graphs the number of elementary schools that are funded and unfunded by bins of p_i . Since schools are not uniformly distributed within each bin, we should not necessarily expect the proportion of schools funded to be the midpoint in each bin even in the population.

Though consistent, the sample analog to (1) is not efficient: as Hahn (1998) shows, it fails to achieve the semiparametric efficiency bound. Hirano, Imbens and Ridder (2003) show that a two-step estimator, in which the first step estimates the probability of treatment using a logit series estimator, does achieve the semiparametric efficiency bound, even when the true probability is known. This puzzle is well known in the econometric literature (Henmi and Eguchi, 2004, Hitomi et al., 2008, Prokhorov and Schmidt, 2009, Han and Kim, 2011), though as far as I know the result has never been applied empirically, presumably because probability of treatment is rarely known, as it is in this case.

Though seemingly counter-intuitive, this result rests on a well-known fact: even under

exogenous treatment, if variation in the outcome can be explained by variation in other observables, partialling out this variation results in more efficient estimation. This same principle leads to the inclusion of covariates in an OLS estimate with random and dichotomous treatment. An OLS estimate of the causal effect is consistent and unbiased without covariates, but is more precise when covariates that explain variation in the outcome are included.

Wooldridge (2010) makes explicit the application of this intuition. Consider $k_i = [(T_i - p_i)y_i]/[p_i(1 - p_i)]$, where $E(k_i) = \tau$, my population parameter of interest. We would of course estimate τ using the sample average of k_i , but doing so treats variation in k_i that is explained by variation in covariates as noise, leading to inefficient estimation.

If instead we were to estimate p_i in a first stage using a logit model, as Hirano et al. (2003) suggest, this would be equivalent to regressing $\hat{k}_i = [(T_i - \hat{p}_i)y_i]/[\hat{p}_i(1 - \hat{p}_i)]$, where \hat{p}_i is the predicted probability from the first stage, on a constant and $\hat{d}_i = \mathbf{X}_i(T_i - \hat{p}_i)$: the constant would be an estimate of τ , and the residuals can be used to estimate the variance of $\hat{\tau}$. To the degree that \hat{d}_i explains variation in \hat{k}_i , $\hat{\tau}$ will gain efficiency. Another way to reach the same conclusion is to note that $E(k_i - \tau) = 0$ and $E(d_i) = 0$ are moment conditions. Estimating τ treating p_i as known disregards the second moment condition, which, so long as it is correlated with the first moment condition, contains useful information that we incorporate in estimation by treating p_i as unknown.

With known p_i , the gains in efficiency can be achieved by regressing k_i on $d_i = \mathbf{X}_i(T_i - p_i)$. I provide results using the sample analog of (1), which I refer to as those with one moment condition, and results that regress k_i on d_i , where \mathbf{X}_i includes an indicator for having met the growth target, proportion of students eligible for Free and Reduced Price Lunches, enrollment, Standardized API, percent of students who are Hispanic, English language learners, and migrant. I use the value of these variables in 2007. As illustrated in Table 2, the sample moment conditions implied by $E(d_i) = 0$ are all quite close to zero.

6 Results

6.1 Regression Results

For comparison, I first present results based on various regression specifications with the main QEIA requirements as outcome variables. It's important to note that, unlike the IPW estimator, consistent estimation of average treatment effects by the regression models depends on untestable assumptions. Each regression has a full set of year dummies and interactions between a treatment indicator and the year dummies. For expositional purposes the table includes only the coefficient on the interaction between the treatment indicator and the dummy for 2011.²⁴ I present results from regressions on the full sample as well as regressions on the restricted sample such that $p_i \in [.10, .90]$.

For each main QEIA requirement, I present five regression models. Model 1 includes only the year dummies and interactions between a treatment indicator and the year dummies. This model consistently estimates the effect of QEIA on outcomes only if treatment is uncorrelated with potential outcomes. Since the probability of treatment is dependent on district rankings, as well as on the size of the district, we wouldn't expect this assumption to be satisfied.

The second regression model adds to Model 1 an interaction between the year dummies and the probability of treatment, p_i . If there is no heterogeneity in the treatment effect, Model 2 will consistently estimate it. If there is heterogeneity, then consistent estimation of the average treatment effect requires $\text{Var}(T|p)$ to be uncorrelated with potential outcomes (Wooldridge, 2004). There is of course no way to know whether this condition is satisfied.

Model 3 also nests Model 1, and includes as covariates an indicator for whether the school met its growth target in 2007, whether the school is in Los Angeles Unified, the percent of students eligible for free and reduced price lunch in 2007, and the enrollment in 2007. If conditional on these covariates treatment is uncorrelated with potential outcomes,

²⁴The full set of results is available in an online appendix

the average treatment effect will be consistently estimated.

Models 4 and 5 build on Model 1 by including fixed effects in the former, and the above mentioned covariates with fixed effects in the latter. Fixed effects estimation requires treatment to be uncorrelated with trends in potential outcomes. This assumption would be violated if districts ranked highly those schools that were primed for improvement.

As Table 3 illustrates, the estimated effect of QEIA on class size is robust to a broad range of specifications and to the sample restriction. Average class size is estimated to have decreased in selected schools by about 4.5 students per class. In the full sample, estimates of the effect of QEIA on teacher experience are similarly robust to a broad range of specifications. Average experience appears to have decreased by 0.73 to 0.98 years, suggesting that funded schools were not able to reduce class size by hiring more experienced teachers. In the restricted sample, the standard errors are generally larger and the effects are smaller in each model, suggesting at most a 0.74 reduction in average teacher experience, significant at the 10% level.

The regression estimates of the effect of QEIA on schools' API vary across models. Controlling for the probability of treatment, the API in funded schools is estimated to have increased by 0.41 standard deviations ($p < 0.001$) in the distribution of APIs across all elementary schools in California. At the other extreme, controlling for covariates suggests QEIA had no effect on schools' API. Estimates from the restricted sample are precise and less widely dispersed, with a maximum of 0.4 standard deviations and a minimum of 0.26 standard deviations.

Results for 5th grade assessments in math and English language arts vary across specifications, with effects for math being larger across specifications in the full and restricted sample. In the full sample, the effect of QEIA on math scores varies from 0.25 ($p < 0.01$) to 0.46 ($p < 0.001$) standard deviations in the population of all school-level averages in California, and ELA scores range from an insignificant 0.06 to 0.30 ($p < 0.001$). The estimates are more uniform in the restricted sample, varying from 0.41 ($p < 0.001$) to 0.50 standard

deviations ($p < 0.001$) in 9math, and from 0.31 ($p < 0.001$) to 0.19 ($p < 0.01$) in ELA.

There is some evidence from the regression results of an increase in persistence in funded schools, measured by the percent of students who were in the school the previous year. A causal effect of QEIA on the composition of students in a school could suggest that it did not benefit particular students, but rather better students were attracted to the QEIA schools. However, the estimated effects are small relative to the baseline of about 90% (see Table 1), precise only in some specifications, and then only precise at the 5% level. More importantly, as indicated below, results from IPW estimates suggest no change in student characteristics.

6.2 Main Results

Unlike the regression results above, IPW estimates depend only on the assumption that the randomization was carried out correctly, and by all accounts it was. The remainder of the paper therefore focuses on these estimates, presenting those that depend on one and two sets of moment conditions.²⁵

Table 4 shows the causal effect of QEIA on average class size, the percent of teachers classified as highly qualified, average teacher experience, and the TEI. The point estimate on class size in 2009, the first full year of QEIA funding, suggests that QEIA reduced class size, but the effect is imprecisely measured. The standard errors are much smaller using two sets of moment conditions, though they are still larger than those from the regression. Using estimates based on two moment conditions, in the final year for which class size data are available, QEIA reduced class size by 3.95 students per class, an estimate that is significant at the 0.001 level.

Consistent with the claim that both funded and unfunded schools were required to have high proportions of HQT teachers, being funded had no causal impact on the proportion of HQT teachers. The estimates that rely on two sets of moment conditions are all practically

²⁵Average treatment effects on the treated are available in the online appendix. The point estimates are quite similar to the average treatment effects, and are too noisy to distinguish from the average treatment effects.

small, precisely measured, and not statistically discernible from zero.

Similarly, teacher experience does not appear to have been affected by QEIA. From 2009 through 2011 the point estimates using both one and two sets of moment conditions are neither positive nor statistically discernible from zero, as evidenced by Table 4. The point estimates based on two moment conditions are smaller in absolute value than the regression estimates, both of which are small relative to the 2007 baseline of 11.81. The marginally significant effects on TEI in 2007 are presumably spurious, and cast doubt on the significant differences in 2008 and 2009, using one moment condition, and 2008, using two moment conditions. Even taking the point estimates at face value, QEIA appears to have reduced teacher experience, measured by years or by the TEI.

The estimated effect of QEIA on student achievement, as measured by California's API, is quite similar to the regression estimates based on the restricted sample, as shown in Table 5. Using two sets of moment conditions,²⁶ QEIA increased API scores in funded schools by 0.33 standard deviations ($p < 0.001$) by 2011, with respect to the population of all elementary schools. The effect for Hispanic students is significantly larger than for all students by 2011 ($p = 0.046$), as is the effect for low-SES students ($p = 0.051$). From 2008 onward there is a clear pattern of funded schools improving over unfunded schools.

These estimates capture the causal effect of QEIA at the school level. However, it is possible that these results are driven partly by changes in the composition of students in response to QEIA. For instance, it may be that especially savvy parents, whose children are more likely to receive extra support, will be aware of QEIA and select into a QEIA school. Although I cannot currently observe student-level characteristics, I do observe school-level averages of such things as Free and Reduced Price Lunch eligibility, whether parents have a college degree, their race, and whether the student was enrolled in the school the prior year.

Table 6 displays the results for FRPL,²⁷ percent of students whose parents have a college

²⁶Note that effects on the 2007 API scores are not calculated using both moment conditions. This is because I use 2007 API scores in the second set of moment conditions.

²⁷Estimates relying on two sets of moment conditions for FRPL in 2007 are not calculated, since FRPL in 2007 is used in that set of moment conditions.

degree, and percent of students who were in the school the prior year. Focusing on the estimates based on two sets of moment conditions, no coefficient is significant, and the magnitudes of the point estimates are quite small. Similarly, Table 7 shows no discernible impact on student enrollment, proportion black, or proportion Hispanic. This is consistent with the student population not changing in response to QEIA.

However, it is still possible that the population of test takers at schools may have changed in response to QEIA. This is particularly concerning since schools were required to improve API scores in order to remain in the program, increasing the stakes of the tests. Schools could manipulate their API scores by either encouraging more students to take alternative tests,²⁸ discouraging low-performing students from taking any test, or manipulating answer sheets.

Of these possibilities, I currently am able to observe the number of students for whom there is a valid score, and the number who take the regular standardized test. Table 8 displays the results of this analysis. There is no evidence that the number of valid scores differed between funded and unfunded schools, either before or after QEIA. Similarly, there is no evidence of a difference in the number of valid scores for low-SES students or for Hispanic students. Neither is there evidence of a change in the proportion of students taking the regular standardized test. Though this is not definitive evidence, it is at least consistent with the population of test takers not changing in response to QEIA.

Two key policy levers of QEIA, decreased class size and increased teacher experience, require changes to the teacher workforce at a school. Using the teacher-level data, I am able to observe the net changes in teacher characteristics at a school. Tables 9 and 10 list the results from this analysis. I examine differences caused by QEIA in the proportion of teachers new to the school, new to the school but not new to the district,²⁹ average

²⁸Alternative tests are included in the calculation of the API, but presumably the marginal student would find the regular California Standardized Test challenging, and the California Alternative Performance Assessment or the California Modified Assessment less so.

²⁹A teacher is new to the school but not the district if no teacher with the same characteristics is observed in the school the prior year, and the teacher has more than one year of experience in the district.

experience conditional on being new to the school, and proportion of probationary, tenured, and temporary teachers.

In 2009 QEIA appears to have caused an increase in the proportion of teachers new to the school in funded schools relative to unfunded schools. In 2009 there were 7 percentage points more ($p < 0.1$) new teachers in funded schools. The similar estimates for the change in new teachers with experience in the district suggests that nearly all teachers new to the school had experience in the district. Comparing the set of teachers new to a school in funded and unfunded schools, average experience is 0.94 years greater in funded schools in 2009 ($p < 0.1$).

Table 10 lists the change in proportion of teachers who are probationary, tenured, and long-term substitutes. The differences between funded and unfunded schools in the proportion of probationary teachers before 2008, significant at the 10% level, are presumably spurious, casting some doubt on the results in later years. Assuming that the more precisely estimated differences in proportion probationary after 2008 are not spurious, there is evidence of an increase in probationary teachers caused by QEIA. There is also marginally significant evidence of fewer tenured teachers in funded schools ($p < 0.10$).

Tables 11 and 12 show the effect of QEIA on class sizes at the grade level. Estimates based on one set of moment conditions are noisy, and are at no point statistically different from zero. Class sizes in kindergarten through 3rd grade are not affected by QEIA until 2011, at which time there are 2.7 to 3.7 fewer students in those grades in funded schools. As mentioned above, class size data are not available in 2010, and that the difference in class sizes in these earlier grades is driven by unfunded schools exiting the previous class size reduction program. Estimates based on two sets of moment conditions suggest that grades 4 and 5 decreased class sizes by about 4.4 ($p < 0.001$) and 4 ($p < 0.01$) students per class in 2009, respectively, and by 5.1 ($p < 0.001$) to 5.7 ($p < 0.001$) students in 2011.

The effect of QEIA on API scores is important, since the primary goal of the policy was to improve API scores. However, given that the API is an average across students, grades,

subjects and even test types, changes in API scores are hard to interpret or compare to other findings in the literature. Tables 13 and 14 therefore list the estimated effect of QEIA on mean scaled scores from California's Standardized Test for math and English language arts.

The effects of QEIA on math scores is greater in later years of the program and in higher grades. There's no discernible effect on 2nd grade math scores until 2011, when they are 0.28 standard deviations higher in funded schools, with respect to the population of grade-level averages ($p < 0.001$, 0.12 student-level standard deviations). The 3rd grade math scores increase one year earlier, in 2010, by 0.16 standard deviations (0.07 student-level standard deviations, $p < 0.05$), and by 0.27 standard deviations by 2011 (0.12 student-level standard deviations, $p < 0.01$).

Math scores in 4th grade improve earlier; by 2009 they show an increase of 0.32 standard deviations (0.15 student-level standard deviations, $p < 0.001$), 0.40 (0.17 student-level standard deviations, $p < 0.001$) in 2010, and level off in 2011 at 0.38 (0.16 student-level standard deviations, $p < 0.001$). Interestingly, 5th grade math scores do not begin improving until 2010, at which time they were 0.36 standard deviations higher in funded schools (0.16 student-level standard deviations, $p < 0.001$), and by 2011 they were 0.44 standard deviations higher (0.20 student-level standard deviations, $p < 0.001$).

Consistent with results from the vast majority of education reforms, the effects are smaller for English language arts. Still, the previous pattern persists: effects are larger in later years, in later grades, and there is an effect on 4th grade test scores in 2009 but not on 5th grade test scores. By 2011, ELA scores in 2nd grade are 0.18 standard deviations higher in funded schools (0.08 student-level standard deviations, $p < 0.01$), and in 5th grade they are 0.21 standard deviations higher (0.10 student-level standard deviations, $p < 0.001$).

To better understand the effects of increased exposure to QEIA, Figure 2 replicates the information in the tables, displaying the average treatment effect on class size and achievement at the grade level, but by cohort exposure. Each panel in the figure displays the change in class size and achievement that a group of students with a normal grade progression would

face. For instance, students in panel A enter kindergarten in 2005, and those who progress one grade each year are exposed to QEIA for one year, in 2009. Since class size data aren't available in 2010, I instead use the same-grade average from 2009 and 2011, e.g., the 3rd grade class size in 2010 is the average of 3rd grade class size in 2009 and 2011.

As the figure suggests, consecutive years of smaller classes does not lead to a widening of the achievement gains. Additionally, though it is not possible to empirically separate the effects of teacher training, high-stakes testing, and reduced class size, it must nonetheless be the case that if teacher training and high-stakes testing explain the improved scores, the timing would have to be correlated with changes in class size. Since both the reduced class sizes in 2nd and 3rd grade are delayed, it seems likely that the effect is driven by class size. Otherwise, there would have to be some reason that professional training was delayed, and that schools did not respond to high-stakes testing until 2011,

7 Conclusion

California's QEIA provides a unique opportunity to study the causal effects of school reform. Using district rankings and the details of the selection process, the probability of any school being selected is known. Between any two schools with the same probability of selection, being funded is uncorrelated with potential outcomes. Using this, and relying on methods described in Wooldridge (2004) and Hirano, Imbens and Ridder (2003), I am able to estimate the causal impact of QEIA by inverse probability weighting.

Doing so, I find that QEIA caused a decrease in class size, and had no discernible effect on teacher experience. Two of the other QEIA requirements applied to all QEIA eligible schools, and are therefore not considered part of the treatment here. The remaining components of treatment, professional training for teachers, which is unobserved, and added incentive to increase achievement to maintain funding, may also contribute to the improvement in test scores.

Test scores improved significantly, albeit unevenly across grades and years. Grades 4 and 5, in which class sizes were first reduced, experienced the largest and earliest increase in test scores. In the first fully-funded year of the program, math scores in 4th grade increased by 0.32 standard deviations in the population of school-grade averages, and by the second fully-funded year 5th grade math scores improved by 0.36 standard deviations. The improvement in test scores in 2nd and 3rd grade, like the reduction in class sizes in those grades, occurred later, and was more modest. By the third fully-funded year of the program, math scores in 2nd grade were 0.28 standard deviations higher in the distribution of school-grade averages, and 0.27 standard deviations higher 3rd grade. For teacher professional training and added test pressure to explain the improvement, they would have to exhibit a similar pattern of implementation across grades and years. Gains in English language arts were modest, but exhibit the same pattern across grades and years.

Though the design of QEIA precludes separate identification of effects of its constituent reforms, it is nonetheless a remarkable policy, unprecedented in education for being a large-scale policy intervention with random assignment. In terms of students affected, it dwarfs Tennessee's Project STAR by an order of magnitude. Though certainly imperfect, more policies like QEIA could vastly improve our understanding of the effectiveness of school reform.

References

- Angrist, J., Lavy, V., 1999. Using Maimonides' rule to estimate the effect of class size on scholastic achievement. *The Quarterly Journal of Economics* 114, 533–575.
- Balcom, F., 2007. Quality Education Investment Act (QEIA) of 2006. <http://www.cde.ca.gov/fg/fo/r16/documents/qeia07present.ppt>.
- Barrett, N., Butler, J., Toma, E.F., 2012. Do less effective teachers choose professional development does it matter? *Evaluation Review* 36, 346–374.
- Bluth, A.H., 2005. Lawsuit seeking cash for schools; Governor broke his word, say teachers and schools chief. *Sacramento Bee* <http://www.mikemcmahon.info/ctasuit.htm>.
- CDE, 2010. Report to the Legislature and the Governor; Quality Education Investment Act First Progress Report. <http://www.cde.ca.gov/ta/lp/qe/documents/qeialegrpt.doc>.
- Chetty, R., Friedman, J.N., Hilger, N., Saez, E., Schanzenbach, D.W., Yagan, D., 2011a. How does your kindergarten classroom affect your earnings? Evidence from Project STAR. *The Quarterly Journal of Economics* 126, 1593–1660.
- Chetty, R., Friedman, J.N., Rockoff, J.E., 2011b. The long-term impacts of teachers: Teacher value-added and student outcomes in adulthood. Technical Report. National Bureau of Economic Research.
- Chiang, H., 2009. How accountability pressure on failing schools affects student achievement. *Journal of Public Economics* 93, 1045–1057.
- Chingos, M.M., 2012. The impact of a universal class-size reduction policy: Evidence from Florida's statewide mandate. *Economics of Education Review* 31, 543–562.
- Cohodes, S., Deming, D., Jennings, J., Jencks, C., 2013. School accountability, postsecondary attainment and earnings. NBER Working Paper .

- Crump, R.K., Hotz, V.J., Imbens, G.W., Mitnik, O.A., 2009. Dealing with limited overlap in estimation of average treatment effects. *Biometrika* 96, 187–199.
- Dieterle, S., 2013. Class-size reduction policies and the quality of entering teachers.
- Figlio, D.N., Getzler, L.S., 2006. Accountability, ability and disability: Gaming the system? *Advances in applied microeconomics* 14, 35–49.
- Greenwald, R., Hedges, L., Laine, R., 1996. The effect of school resources on student achievement. *Review of Educational Research* 66, 361–396.
- Hahn, J., 1998. On the role of the propensity score in efficient semiparametric estimation of average treatment effects. *Econometrica* , 315–331.
- Han, C., Kim, B., 2011. A GMM interpretation of the paradox in the inverse probability weighting estimation of the average treatment effect on the treated. *Economics Letters* 110, 163–165.
- Hanushek, E., 1979. Conceptual and empirical issues in the estimation of educational production functions. *Journal of Human Resources* , 351–388.
- Hanushek, E., 1986. The economics of schooling: Production and efficiency in public schools. *Journal of Economic Literature* 24, 1141–1177.
- Hanushek, E., 1996. A more complete picture of school resource policies. *Review of Educational Research* 66, 397–409.
- Hanushek, E., 1997. Assessing the effects of school resources on student performance: An update. *Educational Evaluation and Policy Analysis* 19, 141.
- Hanushek, E., 1999. Some findings from an independent investigation of the Tennessee STAR experiment and from other investigations of class size effects. *Educational Evaluation and Policy Analysis* 21, 143.

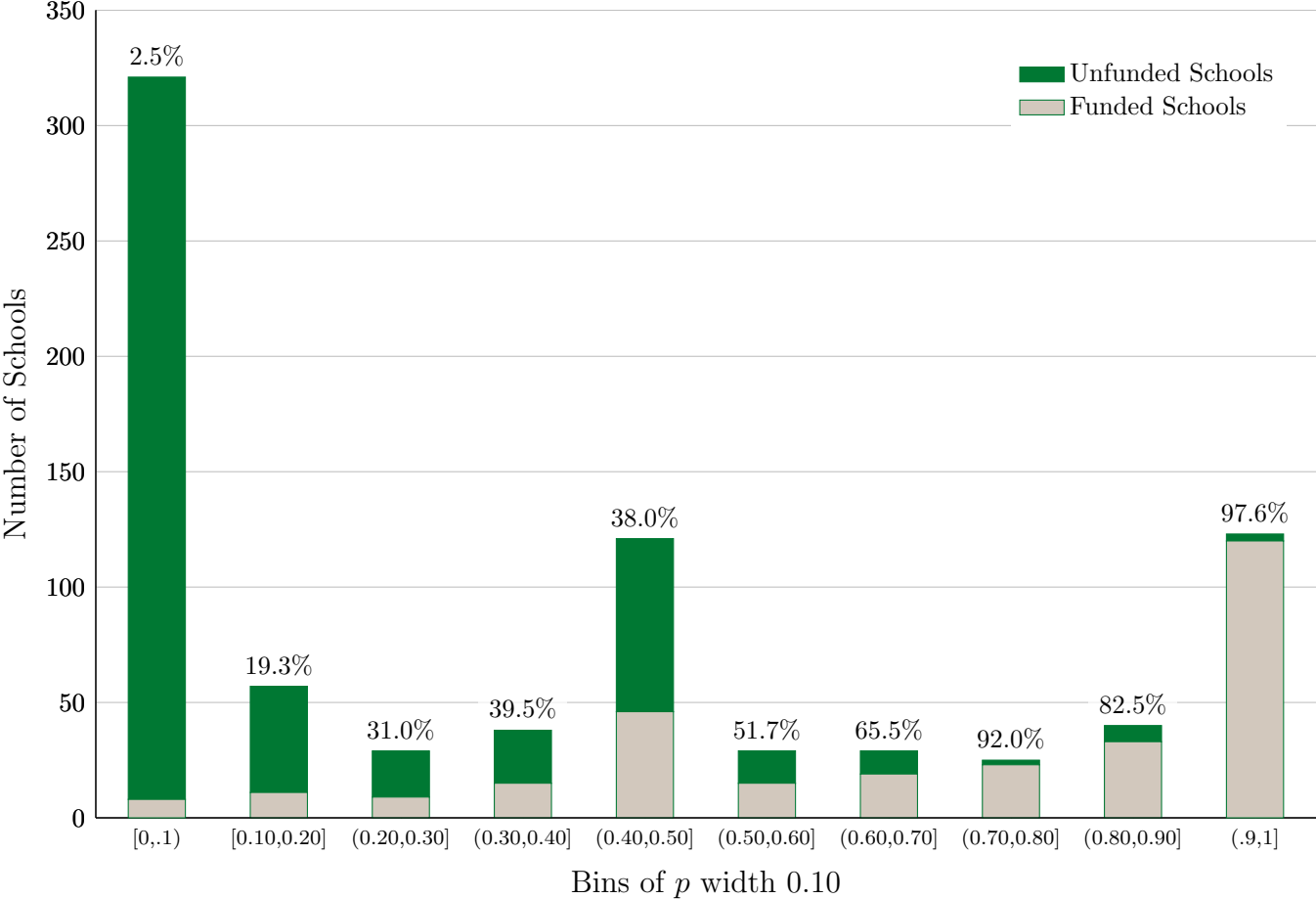
- Henmi, M., Eguchi, S., 2004. A paradox concerning nuisance parameters and projected estimating functions. *Biometrika* 91, 929–941.
- Hirano, K., Imbens, G.W., Ridder, G., 2003. Efficient estimation of average treatment effects using the estimated propensity score. *Econometrica* 71, 1161–1189.
- Hitomi, K., Nishiyama, Y., Okui, R., 2008. A puzzling phenomenon in semiparametric estimation problems with infinite-dimensional nuisance parameters. *Econometric Theory* 24, 1717–1728.
- Horvitz, D.G., Thompson, D.J., 1952. A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association* 47, 663–685.
- Hoxby, C., 2000. The effects of class size on student achievement: New evidence from population variation. *Quarterly Journal of Economics* 115, 1239–1285.
- Imbens, G.W., Wooldridge, J.M., 2009. Recent developments in the econometrics of program evaluation. *Journal of Economic Literature* 47, 5–86.
- Jacob, B.A., 2005. Accountability, incentives and behavior: The impact of high-stakes testing in the Chicago public schools. *Journal of Public Economics* 89, 761–796.
- Jacob, B.A., Levitt, S.D., 2003. Rotten apples: An investigation of the prevalence and predictors of teacher cheating. *The Quarterly Journal of Economics* 118, 843–877.
- Jepsen, C., Rivkin, S., 2009. Class size reduction and student achievement: The potential tradeoff between teacher quality and class size. *Journal of Human Resources* 44, 223–250.
- Krueger, A., 1999. Experimental estimates of education production functions. *The Quarterly Journal of Economics* 114, 497–532.
- Krueger, A., 2002. Understanding the magnitude and effect of class size on student achievement. *The Class Size Debate* , 7–35.

- Krueger, A., Whitmore, D., 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.
- Neal, D., Schanzenbach, D.W., 2010. Left behind by design: Proficiency counts and test-based accountability. *The Review of Economics and Statistics* 92, 263–283.
- Nye, B., Hedges, L., Konstantopoulos, S., 1999. The long-term effects of small classes: A five-year follow-up of the Tennessee class size experiment. *Educational Evaluation and Policy Analysis* 21, 127–142.
- Prokhorov, A., Schmidt, P., 2009. GMM redundancy results for general missing data problems. *Journal of Econometrics* 151, 47–55.
- Rice, J., Schwartz, A., 2008. *Handbook of Research in Education Finance and Policy*. Routledge New York. chapter 8. pp. 131–165.
- Rockoff, J., 2009. Field experiments in class size from the early twentieth century. *The Journal of Economic Perspectives* 23, 211–230.
- Santa Rosa City Schools, 2007. School board minutes, Quality Education Investment Act. <http://www.srcs.k12.ca.us/board/agendas/attachments/032807-BR-F7.pdf>.
- Todd, P., Wolpin, K., 2003. On the specification and estimation of the production function for cognitive achievement. *The Economic Journal* 113, F3–F33.
- Wooldridge, J.M., 2004. Estimating average partial effects under conditional moment independence assumptions. CeMMAP Working Paper Number CWP03/04 .
- Wooldridge, J.M., 2010. *Econometric analysis of cross section and panel data*. Second ed., MIT Press.
- Yoon, K.S., Duncan, T., Lee, S.W.Y., Scarloss, B., Shapley, K.L., 2007. Reviewing the evidence on how teacher professional development affects student achievement. *National*

Center for Educational Evaluation and Regional Assistance, Institute of Education Sciences, US Department of Education.

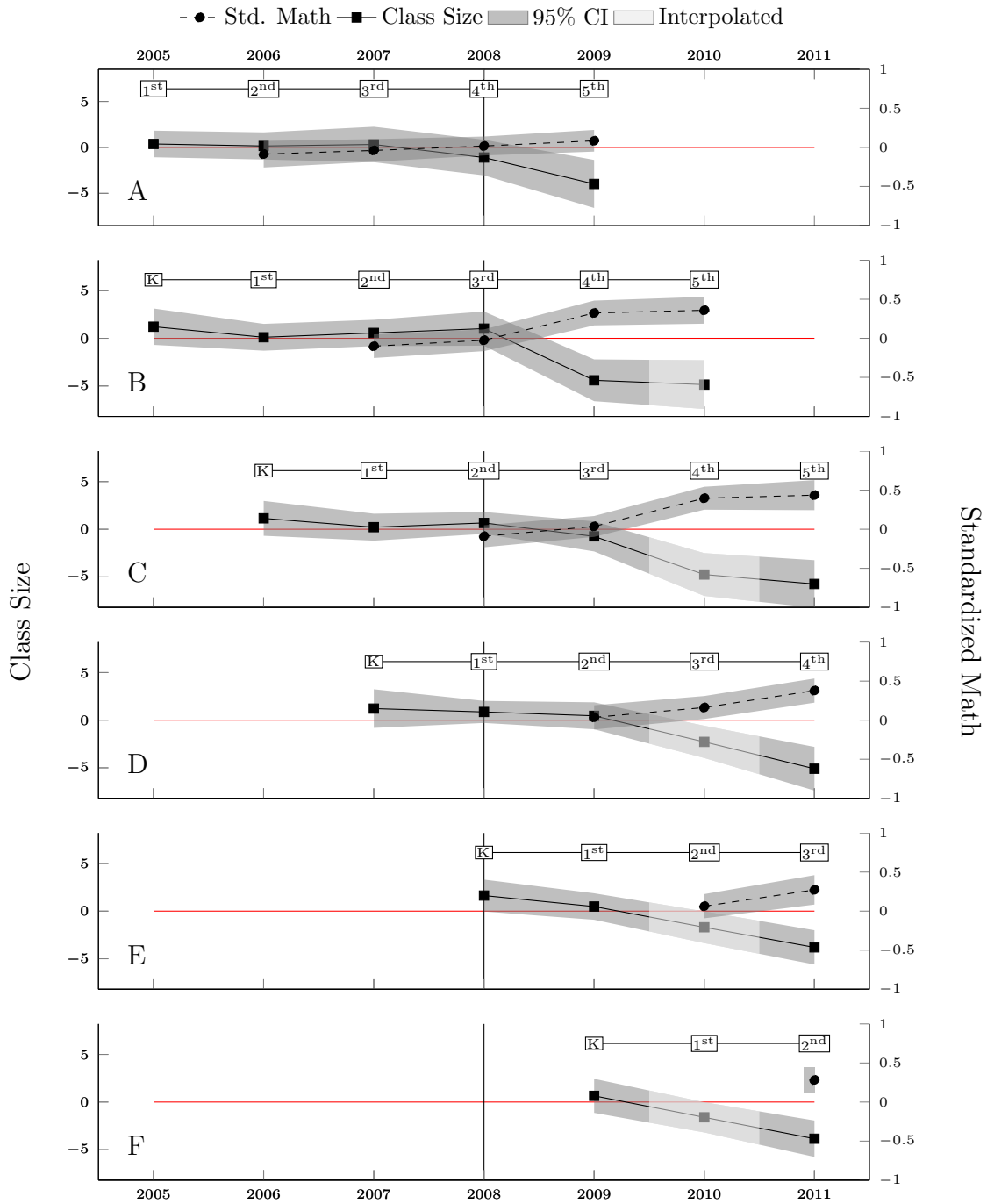
A Figures

Figure 1: Support Over p , All Elementary Schools



Note: Numbers above bars refer to the percentage of schools funded in that bin. Sample includes all elementary schools participating in the regular QEIA program.

Figure 2: Cohort-Level Class Size and Math Achievement Comparison



Note: Estimates are of average treatment effect using two moment conditions. Left axis refers to class size, right axis refers to standardized math scores on California's CST. Class size data are missing for 2010. Shaded regions indicate use of average of 2009 and 2011 same-grade class size. For example 2010 4th grade class size is average of 2009 4th grade class size and 2011 4th grade class size. Only grades 2 and above are tested.

B Tables

Table 1: Descriptives, Elementary Regular QEIA Schools 2007

| | All Elem. | $p_i \in [0, 1], 2007$ | | $p_i \in [.1, .9], 2007$ | | $p_i \in [.1, .9], 2011$ | |
|------------------------------------|--------------------|------------------------|--------------------|--------------------------|--------------------|--------------------------|--------------------|
| | 2007 | Unfunded | Funded | Unfunded | Funded | Unfunded | Funded |
| Average Class Size | 21.86 (3.54) | 21.93 (1.87) | 21.96 (2.04) | 22.14 (1.75) | 22.12 (2.05) | 25.02 (3.17) | 20.36 (2.20) |
| Class Size Kindergarten | 20.64 (4.14) | 20.24 (3.39) | 20.41 (3.78) | 20.30 (3.47) | 20.74 (4.21) | 23.97 (4.29) | 20.23 (3.48) |
| Class Size 1 st Grade | 19.32 (1.85) | 19.29 (1.40) | 19.13 (1.31) | 19.33 (1.47) | 19.21 (1.29) | 23.84 (4.00) | 19.54 (2.27) |
| Class Size 2 nd Grade | 19.14 (1.87) | 18.84 (1.52) | 18.97 (1.40) | 18.79 (1.56) | 19.00 (1.34) | 23.88 (4.18) | 19.35 (2.53) |
| Class Size 3 rd Grade | 19.89 (3.17) | 19.42 (2.71) | 19.60 (3.01) | 19.76 (3.19) | 19.68 (3.22) | 23.92 (4.48) | 19.56 (2.64) |
| Class Size 4 th Grade | 28.47 (4.23) | 28.01 (3.82) | 28.18 (3.80) | 28.27 (3.56) | 28.58 (3.80) | 28.15 (4.31) | 22.40 (3.99) |
| Class Size 5 th Grade | 28.89 (4.15) | 28.28 (3.72) | 28.51 (3.77) | 28.68 (3.56) | 28.66 (3.97) | 28.42 (4.51) | 22.24 (3.35) |
| Average Experience | 13.01 (3.95) | 11.92 (3.17) | 11.61 (3.35) | 12.37 (3.12) | 11.68 (3.38) | 13.68 (3.57) | 12.94 (3.17) |
| TEI Relative | -0.04 (1.05) | -0.17 (1.02) | -0.28 (1.00) | -0.06 (0.89) | -0.26 (0.98) | -0.12 (0.80) | -0.15 (0.69) |
| Highly Qualified Teachers | 0.96 (0.11) | 0.94 (0.10) | 0.94 (0.14) | 0.96 (0.09) | 0.94 (0.15) | 0.99 (0.04) | 0.99 (0.11) |
| <i>Williams</i> Settlement Applies | 0.24 (0.43) | 0.94 (0.24) | 0.96 (0.20) | 0.95 (0.21) | 0.95 (0.22) | 0.95 (0.21) | 0.95 (0.22) |
| Standardized API | 0.00 (1.00) | -1.10 (0.48) | -1.25 (0.48) | -1.14 (0.44) | -1.15 (0.47) | -1.10 (0.56) | -0.77 (0.60) |
| API Percentile Rank | 0.50 (0.29) | 0.16 (0.11) | 0.12 (0.10) | 0.15 (0.11) | 0.15 (0.11) | 0.17 (0.14) | 0.26 (0.17) |
| Met Growth Target | 0.70 (0.46) | 0.68 (0.47) | 0.65 (0.48) | 0.66 (0.48) | 0.60 (0.49) | 0.60 (0.49) | 0.71 (0.46) |
| Proportion Black | 0.08 (0.12) | 0.09 (0.14) | 0.11 (0.16) | 0.07 (0.12) | 0.08 (0.13) | 0.06 (0.11) | 0.08 (0.12) |
| Proportion Hispanic | 0.46 (0.30) | 0.79 (0.21) | 0.74 (0.25) | 0.78 (0.21) | 0.76 (0.25) | 0.80 (0.20) | 0.78 (0.24) |
| Proportion White | 0.33 (0.28) | 0.06 (0.10) | 0.07 (0.10) | 0.09 (0.13) | 0.07 (0.10) | 0.08 (0.12) | 0.06 (0.10) |
| English Language Learners | 0.29 (0.23) | 0.55 (0.18) | 0.55 (0.20) | 0.55 (0.18) | 0.56 (0.20) | 0.52 (0.18) | 0.53 (0.20) |
| Proportion FRPL | 0.55 (0.31) | 0.89 (0.11) | 0.88 (0.11) | 0.87 (0.12) | 0.86 (0.12) | 0.89 (0.14) | 0.85 (0.17) |
| Parent College Grad | 0.18 (0.14) | 0.06 (0.06) | 0.06 (0.07) | 0.06 (0.06) | 0.07 (0.09) | 0.07 (0.04) | 0.07 (0.05) |
| Student Enrollment | 376.60 (192.56) | 472.24 (195.68) | 434.61 (179.70) | 446.60 (164.72) | 418.85 (151.15) | 404.83 (139.09) | 377.09 (133.58) |
| Proportion Same School | 0.92 (0.08) | 0.91 (0.04) | 0.91 (0.04) | 0.91 (0.04) | 0.91 (0.04) | 0.92 (0.05) | 0.93 (0.08) |
| Los Angeles | 0.09 (0.29) | 0.23 (0.42) | 0.09 (0.29) | 0.06 (0.24) | 0.03 (0.17) | 0.06 (0.24) | 0.03 (0.17) |
| Proportion Teachers New to School | 0.33 (0.27) | 0.30 (0.26) | 0.35 (0.25) | 0.34 (0.26) | 0.36 (0.26) | 0.50 (0.31) | 0.48 (0.29) |
| New to School, Not District | 0.24 (0.24) | 0.22 (0.23) | 0.25 (0.23) | 0.25 (0.24) | 0.25 (0.23) | 0.47 (0.30) | 0.43 (0.29) |
| Average Experience New Teachers | 3.21 (3.64) | 2.76 (3.18) | 3.13 (3.12) | 3.16 (3.47) | 3.13 (3.07) | 6.43 (4.88) | 5.83 (4.47) |
| Proportion Temp Teachers | 0.06 (0.10) | 0.05 (0.08) | 0.05 (0.10) | 0.06 (0.09) | 0.06 (0.11) | 0.04 (0.08) | 0.04 (0.07) |
| Proportion Probationary | 0.14 (0.18) | 0.13 (0.14) | 0.16 (0.16) | 0.15 (0.15) | 0.18 (0.17) | 0.06 (0.10) | 0.09 (0.13) |
| N | 6476 | 546 | 307 | 197 | 171 | 197 | 171 |

Note: Funded and unfunded includes all elementary schools participating in the regular QEIA program

Table 2: Sample Moment Conditions

| Variable | Mean | S.D. |
|--|---------|------------|
| (Funded _{<i>i</i>} - <i>p_i</i>) | 0.0032 | (0.4540) |
| (2007 Met Growth Target)(Funded _{<i>i</i>} - <i>p_i</i>) | -0.0195 | (0.3607) |
| (2007 Proportion FRPL)(Funded _{<i>i</i>} - <i>p_i</i>) | -0.0029 | (0.3950) |
| (2007 Student Enrollment)(Funded _{<i>i</i>} - <i>p_i</i>) | -3.4513 | (203.2044) |
| (2007 Std. API)(Funded _{<i>i</i>} - <i>p_i</i>) | 0.0009 | (0.5296) |
| (2007 Proportion Hispanic)(Funded _{<i>i</i>} - <i>p_i</i>) | -0.0021 | (0.3596) |
| (2007 English Language Learners)(Funded _{<i>i</i>} - <i>p_i</i>) | 0.0029 | (0.2611) |
| (2007 Migrant)(Funded _{<i>i</i>} - <i>p_i</i>) | -0.0074 | (0.0623) |

Note: Sample analogs of moments in condition $E[\mathbf{X}(\text{Funded}_i - p_i)] = \mathbf{0}$.

Table 3: Select Regression Results

| | (1) | (2) | (3) | (4) | (5) |
|------------------------------------|----------------------|----------------------|----------------------|----------------------|----------------------|
| | Full Sample | | | | |
| Avg. Class Size | -4.454*** (0.559) | -4.761*** (0.564) | -4.538*** (0.401) | -4.616*** (0.462) | -4.578*** (0.419) |
| Experience | -0.831** (0.256) | -0.730† (0.431) | -0.903** (0.271) | -0.932** (0.280) | -0.977*** (0.250) |
| Std. API | 0.143† (0.075) | 0.407*** (0.071) | 0.079 (0.069) | 0.260*** (0.063) | 0.221*** (0.053) |
| 5 th Grade Math | 0.273** (0.086) | 0.462*** (0.094) | 0.248** (0.078) | 0.365*** (0.083) | 0.341*** (0.078) |
| 5 th Grade ELA | 0.116* (0.059) | 0.301*** (0.069) | 0.056 (0.056) | 0.173** (0.054) | 0.157** (0.048) |
| Prop. Enrolled Since Previous Year | 0.0008 (0.005) | 0.010 (0.008) | 0.006* (0.003) | 0.004 (0.005) | 0.007* (0.004) |
| | $p_i \in [.10, .90]$ | | | | |
| Avg. Class Size | -4.661*** (0.439) | -4.587*** (0.471) | -4.544*** (0.400) | -4.839*** (0.429) | -4.749*** (0.437) |
| Experience | -0.738† (0.381) | -0.659 (0.450) | -0.741† (0.398) | -0.573† (0.345) | -0.649† (0.340) |
| Std. API | 0.334*** (0.063) | 0.400*** (0.070) | 0.256*** (0.060) | 0.365*** (0.063) | 0.333*** (0.062) |
| 5 th Grade Math | 0.439*** (0.086) | 0.453*** (0.095) | 0.408*** (0.086) | 0.502*** (0.097) | 0.482*** (0.100) |
| 5 th Grade ELA | 0.256*** (0.061) | 0.305*** (0.070) | 0.192** (0.057) | 0.262*** (0.060) | 0.244*** (0.061) |
| Prop. Enrolled Since Previous Year | 0.006 (0.007) | 0.009 (0.009) | 0.012* (0.005) | 0.006 (0.007) | 0.011* (0.005) |
| Covariates | No | No | Yes | No | Yes |
| Propensity Score | No | Yes | No | No | No |
| Fixed Effects | No | No | No | Yes | Yes |

Note: † indicates $p < 0.10$, * indicates $p < 0.05$, ** indicates $p < 0.01$, *** indicates $p < 0.001$. Standard errors robust and clustered at district level. Omitted variable is 2005 unfunded. Covariates include dummy for whether school met growth target in 2007, percent of students eligible for Free and Reduced Price Lunches in 2007, an indicator for being in LA, and total enrollment in 2007. Results are from regression with time dummies, and time dummies interacted with treatment indicator. When propensity score is included, it is interacted with time dummies. Reported coefficients are from interaction of dummy for 2011 and treatment indicator.

Table 4: Class Size and HQT

| | Avg. Class Size | | HQT | | Experience | | TEI | |
|------|-------------------|----------------------|------------------|-------------------|--------------------|-------------------|--------------------|--------------------------------|
| | ATE 1 | ATE 2 | ATE 1 | ATE 2 | ATE 1 | ATE 2 | ATE 1 | ATE 2 |
| 2005 | 1.143 (4.119) | 0.528 (0.880) | 0.055 (0.187) | 0.019 (0.033) | 0.347 (2.083) | -0.135 (0.576) | -0.044 (0.102) | -0.013 (0.102) |
| 2006 | 1.237 (3.963) | 0.567 (0.854) | 0.054 (0.165) | 0.015 (0.040) | 0.223 (2.094) | -0.207 (0.636) | -0.102 (0.106) | -0.053 (0.092) |
| 2007 | 1.191 (3.813) | 0.493 (0.838) | 0.036 (0.171) | -0.004 (0.038) | -0.229 (2.193) | -0.655 (0.674) | -0.304* (0.130) | -0.231 [†] (0.121) |
| 2008 | 0.977 (3.861) | 0.484 (0.739) | 0.058 (0.177) | 0.016 (0.040) | 0.060 (2.207) | -0.431 (0.641) | -0.289* (0.135) | -0.229 [†] (0.131) |
| 2009 | -0.537 (3.863) | -1.116 (0.817) | 0.052 (0.173) | 0.016 (0.036) | -0.0008 (2.244) | -0.355 (0.755) | -0.269* (0.131) | -0.182 (0.126) |
| 2010 | | | | | -0.158 (2.423) | -0.532 (0.706) | | |
| 2011 | -3.356 (3.984) | -3.952*** (0.923) | 0.043 (0.183) | 0.010 (0.037) | -0.037 (2.480) | -0.455 (0.653) | -0.017 (0.100) | 0.008 (0.118) |

Note: [†] indicates $p < 0.10$, * indicates $p < 0.05$, ** indicates $p < 0.01$, *** indicates $p < 0.001$. ATE1 is estimate of average treatment effect using one moment condition; ATE2 is estimate of average treatment effect using two moment conditions. Standard errors bootstrapped and clustered at district level. Class size data, which is used to calculate TEI and includes HQT data, is not available in 2010.

Table 5: Average Performance Index

| | Std. API | | Std. API Hispanic | | Std. API Low SES | |
|------|-------------------|-------------------------------|--------------------|-------------------------------|--------------------|---------------------|
| | ATE 1 | ATE 2 | ATE 1 | ATE 2 | ATE 1 | ATE 2 |
| 2005 | -0.069 (0.219) | -0.050 (0.059) | 0.035 (0.213) | 0.020 (0.075) | -0.120 (0.206) | -0.088 (0.062) |
| 2006 | -0.039 (0.219) | -0.047 (0.050) | 0.039 (0.205) | 0.010 (0.069) | -0.077 (0.192) | -0.080 (0.052) |
| 2007 | -0.017 (0.200) | | 0.040 (0.182) | 0.048 (0.058) | -0.046 (0.183) | -0.006 (0.039) |
| 2008 | 0.027 (0.194) | 0.012 (0.043) | 0.119 (0.161) | 0.085 (0.060) | 0.034 (0.165) | 0.028 (0.055) |
| 2009 | 0.123 (0.179) | 0.091 [†] (0.054) | 0.201 (0.165) | 0.149 [†] (0.083) | 0.161 (0.159) | 0.122 (0.078) |
| 2010 | 0.248 (0.182) | 0.208*** (0.061) | 0.414* (0.168) | 0.325*** (0.085) | 0.332* (0.159) | 0.290*** (0.078) |
| 2011 | 0.383* (0.169) | 0.330*** (0.065) | 0.414** (0.152) | 0.394*** (0.082) | 0.445** (0.165) | 0.404*** (0.082) |

Note: [†] indicates $p < 0.10$, * indicates $p < 0.05$, ** indicates $p < 0.01$, *** indicates $p < 0.001$. ATE1 is estimate of average treatment effect using one moment condition; ATE2 is estimate of average treatment effect using two moment conditions. Standard errors bootstrapped and clustered at district level. 2007 API effect not calculated since that covariate is used in the second moment condition.

Table 6: Demographics

| | FRPL | | Parents College Degree | | Enrolled Since Previous Year | |
|------|-------------------|-------------------|------------------------|--------------------|------------------------------|------------------|
| | ATE 1 | ATE 2 | ATE 1 | ATE 2 | ATE 1 | ATE 2 |
| 2005 | 0.030 (0.162) | 0.023 (0.034) | 0.014 (0.015) | 0.005 (0.006) | 0.056 (0.162) | 0.021 (0.031) |
| 2006 | 0.018 (0.152) | 0.008 (0.034) | 0.014 (0.014) | 0.005 (0.007) | 0.060 (0.157) | 0.027 (0.030) |
| 2007 | 0.017 (0.154) | | 0.013 (0.017) | 0.005 (0.008) | 0.052 (0.160) | 0.017 (0.030) |
| 2008 | 0.024 (0.153) | 0.021 (0.030) | 0.006 (0.014) | -0.001 (0.006) | 0.059 (0.157) | 0.023 (0.029) |
| 2009 | 0.013 (0.158) | 0.005 (0.034) | 0.009 (0.014) | 0.0001 (0.006) | 0.061 (0.158) | 0.026 (0.032) |
| 2010 | 0.013 (0.162) | 0.002 (0.032) | 0.015 (0.015) | 0.005 (0.005) | 0.067 (0.168) | 0.025 (0.031) |
| 2011 | -0.006 (0.155) | -0.018 (0.039) | 0.010 (0.016) | -0.0002 (0.005) | 0.064 (0.170) | 0.029 (0.031) |

Note: † indicates $p < 0.10$, * indicates $p < 0.05$, ** indicates $p < 0.01$, *** indicates $p < 0.001$. ATE1 is estimate of average treatment effect using one moment condition; ATE2 is estimate of average treatment effect using two moment conditions. Standard errors bootstrapped and clustered at district level. FRPL, parents with college degree, and enrolled previous year are expressed in proportions.

Table 7: Demographics

| | Enrollment | | Black | | Hispanic | |
|------|--------------------|--------------------|------------------|------------------|------------------|------------------|
| | ATE 1 | ATE 2 | ATE 1 | ATE 2 | ATE 1 | ATE 2 |
| 2005 | -6.541 (89.279) | 1.826 (22.459) | 0.026 (0.027) | 0.018 (0.014) | 0.007 (0.130) | 0.006 (0.028) |
| 2006 | -1.332 (82.240) | 5.158 (18.640) | 0.025 (0.027) | 0.016 (0.013) | 0.011 (0.138) | 0.008 (0.027) |
| 2007 | 0.522 (82.737) | | 0.020 (0.025) | 0.014 (0.013) | 0.018 (0.137) | |
| 2008 | -9.580 (78.651) | -4.517 (18.204) | 0.018 (0.024) | 0.012 (0.013) | 0.020 (0.140) | 0.015 (0.028) |
| 2009 | -5.770 (74.625) | 0.289 (16.940) | 0.020 (0.025) | 0.014 (0.012) | 0.017 (0.131) | 0.013 (0.029) |
| 2010 | 1.656 (76.704) | 2.082 (18.336) | 0.015 (0.024) | 0.008 (0.012) | 0.027 (0.138) | 0.015 (0.029) |
| 2011 | -2.720 (72.262) | 2.720 (18.119) | 0.015 (0.023) | 0.009 (0.011) | 0.019 (0.145) | 0.015 (0.029) |

Note: † indicates $p < 0.10$, * indicates $p < 0.05$, ** indicates $p < 0.01$, *** indicates $p < 0.001$. ATE1 is estimate of average treatment effect using one moment condition; ATE2 is estimate of average treatment effect using two moment conditions. Standard errors bootstrapped and clustered at district level. Black and Hispanic are expressed in proportions.

Table 8: Test Taking

| | Valid Scores | | Number Low SES Scores | | Number Hispanic Scores | | CST | |
|------|--------------------|--------------------|-----------------------|-------------------|------------------------|--------------------|------------------|------------------|
| | ATE 1 | ATE 2 | ATE 1 | ATE 2 | ATE 1 | ATE 2 | ATE 1 | ATE 2 |
| 2005 | -2.275 (77.676) | 4.076 (19.553) | -5.931 (72.290) | 6.224 (18.422) | -5.118 (64.237) | 6.959 (16.237) | 0.075 (0.184) | 0.022 (0.037) |
| 2006 | 3.588 (75.051) | 9.148 (15.011) | -1.484 (67.519) | 8.239 (15.449) | 0.820 (59.916) | 12.317 (12.972) | 0.080 (0.176) | 0.027 (0.035) |
| 2007 | -1.175 (69.466) | 4.731 (13.761) | -7.220 (61.021) | 4.915 (12.123) | -2.778 (53.774) | 8.814 (11.992) | 0.072 (0.187) | 0.029 (0.035) |
| 2008 | -0.971 (69.978) | -2.421 (15.299) | -5.290 (58.471) | 0.826 (13.492) | -2.617 (54.816) | 2.248 (13.093) | 0.056 (0.174) | 0.021 (0.033) |
| 2009 | 6.147 (66.625) | 4.356 (14.462) | -1.070 (62.504) | 3.495 (14.762) | 3.515 (54.005) | 8.631 (13.615) | 0.079 (0.176) | 0.013 (0.032) |
| 2010 | 5.079 (68.400) | 6.011 (15.908) | -4.223 (64.110) | 3.950 (15.185) | 5.240 (56.355) | 12.440 (13.371) | 0.083 (0.172) | 0.012 (0.034) |
| 2011 | 4.168 (68.329) | 8.715 (15.966) | -7.012 (60.272) | 2.553 (17.070) | 1.128 (56.573) | 12.162 (14.543) | 0.076 (0.161) | 0.016 (0.035) |

Note: † indicates $p < 0.10$, * indicates $p < 0.05$, ** indicates $p < 0.01$, *** indicates $p < 0.001$. ATE1 is estimate of average treatment effect using one moment condition; ATE2 is estimate of average treatment effect using two moment conditions. Standard errors bootstrapped and clustered at district level. Valid, low SES, and Hispanic scores refers to all standardized tests used in API. CST is proportion of students taking the California Standardized Test, a subset of the API.

Table 9: Teacher Mobility

| | New to School | | New to School, Not New to District | | Average Experience of New Teachers | |
|------|-------------------|-------------------------------|---------------------------------------|-------------------|---------------------------------------|-------------------------------|
| | ATE 1 | ATE 2 | ATE 1 | ATE 2 | ATE 1 | ATE 2 |
| 2005 | 0.080 (0.076) | 0.061 (0.045) | 0.086 (0.066) | 0.068 (0.043) | 0.853 (0.800) | 0.632 (0.534) |
| 2006 | 0.058 (0.068) | 0.042 (0.041) | 0.046 (0.055) | 0.036 (0.037) | 0.485 (0.713) | 0.274 (0.486) |
| 2007 | 0.023 (0.073) | 0.005 (0.041) | -0.008 (0.064) | -0.024 (0.040) | -0.221 (0.826) | -0.490 (0.650) |
| 2008 | 0.081 (0.073) | 0.066 (0.042) | 0.050 (0.060) | 0.042 (0.040) | 0.663 (0.841) | 0.405 (0.642) |
| 2009 | 0.101 (0.073) | 0.077 [†] (0.040) | 0.089 (0.056) | 0.073* (0.036) | 1.141 (0.755) | 0.942 [†] (0.552) |
| 2010 | 0.060 (0.077) | 0.056 (0.040) | 0.062 (0.071) | 0.060 (0.039) | 0.434 (0.894) | 0.469 (0.524) |
| 2011 | -0.013 (0.094) | 0.010 (0.047) | -0.026 (0.087) | -0.002 (0.046) | -0.333 (1.386) | 0.011 (0.777) |

Note: [†] indicates $p < 0.10$, * indicates $p < 0.05$, ** indicates $p < 0.01$, *** indicates $p < 0.001$. ATE1 is estimate of average treatment effect using one moment condition; ATE2 is estimate of average treatment effect using two moment conditions. Standard errors bootstrapped and clustered at district level.

Table 10: Teacher Composition

| | Probationary | | Tenured | | Long-term Substitute | |
|------|-------------------|-------------------------------|--------------------|--------------------------------|-------------------------------|-------------------|
| | ATE 1 | ATE 2 | ATE 1 | ATE 2 | ATE 1 | ATE 2 |
| 2005 | 0.048 (0.039) | 0.041 [†] (0.023) | -0.015 (0.124) | -0.044 (0.053) | 0.004 (0.014) | 0.003 (0.013) |
| 2006 | 0.035 (0.036) | 0.029 (0.023) | 0.0004 (0.130) | -0.028 (0.052) | 0.0003 (0.012) | -0.001 (0.013) |
| 2007 | 0.044 (0.034) | 0.037 [†] (0.020) | 0.0006 (0.128) | -0.021 (0.056) | 0.005 (0.015) | 0.004 (0.014) |
| 2008 | 0.045 (0.031) | 0.041* (0.018) | -0.006 (0.136) | -0.038 (0.037) | 0.013 (0.018) | 0.012 (0.015) |
| 2009 | 0.033 (0.027) | 0.034* (0.015) | -0.035 (0.140) | -0.064 [†] (0.038) | 0.030 [†] (0.017) | 0.025 (0.017) |
| 2010 | 0.038* (0.019) | 0.040** (0.013) | -0.003 (0.156) | -0.033 (0.037) | 0.004 (0.015) | -0.003 (0.014) |
| 2011 | 0.028 (0.019) | 0.027* (0.013) | -0.0002 (0.163) | -0.032 (0.039) | 0.008 (0.013) | 0.004 (0.011) |

Note: [†] indicates $p < 0.10$, * indicates $p < 0.05$, ** indicates $p < 0.01$, *** indicates $p < 0.001$. ATE1 is estimate of average treatment effect using one moment condition; ATE2 is estimate of average treatment effect using two moment conditions. Standard errors bootstrapped and clustered at district level. Probationary, tenured, and long-term substitute are expressed in proportions.

Table 11: Class Size Grades K-2

| | Class Size Kindergarten | | Class Size 1 st grade | | Class Size 2 nd grade | |
|------|-------------------------|-------------------------------|----------------------------------|----------------------|----------------------------------|----------------------|
| | ATE 1 | ATE 2 | ATE 1 | ATE 2 | ATE 1 | ATE 2 |
| 2005 | 1.330 (3.716) | 1.218 (0.975) | 0.792 (3.603) | 0.377 (0.740) | 0.978 (3.623) | 0.513 (0.733) |
| 2006 | 1.328 (3.890) | 1.143 (0.934) | 0.018 (3.542) | 0.107 (0.718) | 0.674 (3.483) | 0.158 (0.755) |
| 2007 | 1.306 (3.930) | 1.218 (1.035) | 0.587 (3.584) | 0.216 (0.722) | 0.600 (3.457) | 0.563 (0.706) |
| 2008 | 1.578 (3.808) | 1.633 [†] (0.855) | 0.929 (3.422) | 0.865 (0.591) | 0.769 (3.366) | 0.662 (0.587) |
| 2009 | 1.154 (3.778) | 0.651 (0.913) | 0.480 (3.620) | 0.482 (0.709) | 0.798 (3.537) | 0.459 (0.722) |
| 2010 | | | | | | |
| 2011 | -1.988 (4.151) | -2.692* (1.059) | -3.457 (4.187) | -3.693*** (0.909) | -3.128 (4.001) | -3.851*** (0.975) |

Note: [†] indicates $p < 0.10$, * indicates $p < 0.05$, ** indicates $p < 0.01$, *** indicates $p < 0.001$. ATE1 is estimate of average treatment effect using one moment condition; ATE2 is estimate of average treatment effect using two moment conditions. Standard errors bootstrapped and clustered at district level. Class size data are not available in 2010.

Table 12: Class Size Grades 3-5

| | Class Size 3 rd grade | | Class Size 4 th grade | | Class Size 5 th grade | |
|------|----------------------------------|----------------------|----------------------------------|----------------------|----------------------------------|----------------------|
| | ATE 1 | ATE 2 | ATE 1 | ATE 2 | ATE 1 | ATE 2 |
| 2005 | 1.340 (3.563) | 0.282 (0.937) | 2.138 (5.325) | 0.635 (1.166) | 3.222 (5.609) | -0.177 (1.252) |
| 2006 | 0.844 (3.644) | 0.533 (0.896) | 1.586 (5.248) | 0.174 (1.142) | 1.890 (5.237) | -0.292 (1.176) |
| 2007 | 1.130 (3.689) | 0.322 (0.985) | 0.815 (5.082) | 0.852 (1.216) | 1.874 (5.450) | 0.425 (1.283) |
| 2008 | 1.265 (3.705) | 1.017 (0.917) | -0.435 (4.668) | -1.103 (0.991) | -0.383 (4.989) | -0.870 (1.007) |
| 2009 | -0.294 (3.574) | -0.748 (0.807) | -3.451 (4.538) | -4.408*** (1.120) | -1.964 (4.875) | -3.979** (1.333) |
| 2010 | | | | | | |
| 2011 | -2.795 (4.045) | -3.807*** (0.918) | -3.404 (4.673) | -5.098*** (1.170) | -4.173 (4.752) | -5.744*** (1.275) |

Note: † indicates $p < 0.10$, * indicates $p < 0.05$, ** indicates $p < 0.01$, *** indicates $p < 0.001$. ATE1 is estimate of average treatment effect using one moment condition; ATE2 is estimate of average treatment effect using two moment conditions. Standard errors bootstrapped and clustered at district level. Class size data are not available in 2010.

Table 13: Math Standardized Test

| | 2 nd Grade | | 3 rd Grade | | 4 th Grade | | 5 th Grade | |
|------|-----------------------|---------------------|-----------------------|--------------------|-----------------------|---------------------|-----------------------|---------------------|
| | ATE 1 | ATE 2 | ATE 1 | ATE 2 | ATE 1 | ATE 2 | ATE 1 | ATE 2 |
| 2005 | -0.069 (0.222) | -0.045 (0.081) | -0.114 (0.196) | -0.089 (0.064) | -0.088 (0.201) | -0.060 (0.062) | -0.154 (0.205) | -0.114 (0.076) |
| 2006 | -0.108 (0.213) | -0.087 (0.087) | -0.033 (0.209) | -0.025 (0.070) | -0.090 (0.190) | -0.082 (0.068) | -0.076 (0.192) | -0.039 (0.061) |
| 2007 | -0.122 (0.189) | -0.101 (0.077) | -0.078 (0.192) | -0.039 (0.074) | -0.029 (0.178) | -0.018 (0.053) | -0.024 (0.168) | 0.003 (0.055) |
| 2008 | -0.086 (0.172) | -0.090 (0.072) | -0.051 (0.187) | -0.025 (0.071) | 0.015 (0.179) | 0.018 (0.062) | 0.131 (0.160) | 0.096 (0.070) |
| 2009 | 0.062 (0.168) | 0.036 (0.079) | 0.038 (0.174) | 0.036 (0.068) | 0.286* (0.145) | 0.324*** (0.081) | 0.117 (0.153) | 0.085 (0.071) |
| 2010 | 0.091 (0.163) | 0.063 (0.079) | 0.177 (0.171) | 0.162* (0.075) | 0.383** (0.146) | 0.397*** (0.075) | 0.409** (0.154) | 0.359*** (0.088) |
| 2011 | 0.323* (0.143) | 0.282*** (0.084) | 0.277† (0.167) | 0.272** (0.096) | 0.414** (0.132) | 0.379*** (0.079) | 0.472** (0.148) | 0.436*** (0.098) |

Note: † indicates $p < 0.10$, * indicates $p < 0.05$, ** indicates $p < 0.01$, *** indicates $p < 0.001$. ATE1 is estimate of average treatment effect using one moment condition; ATE2 is estimate of average treatment effect using two moment conditions. Standard errors bootstrapped and clustered at district level.

Table 14: ELL Standardized Test

| | 2 nd Grade | | 3 rd Grade | | 4 th Grade | | 5 th Grade | |
|------|-----------------------|--------------------|-----------------------|-------------------|-----------------------|---------------------|-----------------------|---------------------|
| | ATE 1 | ATE 2 | ATE 1 | ATE 2 | ATE 1 | ATE 2 | ATE 1 | ATE 2 |
| 2005 | 0.035 (0.207) | 0.068 (0.066) | -0.060 (0.214) | -0.036 (0.050) | -0.019 (0.203) | -0.007 (0.052) | -0.073 (0.222) | -0.021 (0.075) |
| 2006 | -0.060 (0.209) | -0.044 (0.062) | -0.003 (0.213) | 0.021 (0.062) | -0.041 (0.199) | -0.040 (0.054) | -0.051 (0.200) | -0.019 (0.057) |
| 2007 | -0.103 (0.190) | -0.068 (0.051) | -0.027 (0.209) | 0.006 (0.065) | 0.044 (0.196) | 0.047 (0.055) | -0.047 (0.197) | -0.013 (0.044) |
| 2008 | -0.060 (0.189) | -0.062 (0.066) | -0.046 (0.186) | -0.033 (0.059) | 0.031 (0.200) | 0.012 (0.053) | 0.099 (0.182) | 0.064 (0.048) |
| 2009 | 0.053 (0.175) | 0.024 (0.058) | -0.0010 (0.193) | -0.026 (0.066) | 0.135 (0.183) | 0.153* (0.062) | 0.038 (0.189) | 0.025 (0.055) |
| 2010 | 0.037 (0.169) | 0.042 (0.071) | 0.125 (0.187) | 0.074 (0.076) | 0.195 (0.182) | 0.186** (0.062) | 0.207 (0.189) | 0.186** (0.064) |
| 2011 | 0.227 (0.151) | 0.188** (0.066) | 0.161 (0.184) | 0.147 (0.090) | 0.278† (0.163) | 0.219*** (0.056) | 0.276 (0.173) | 0.217*** (0.066) |

Note: † indicates $p < 0.10$, * indicates $p < 0.05$, ** indicates $p < 0.01$, *** indicates $p < 0.001$. ATE1 is estimate of average treatment effect using one moment condition; ATE2 is estimate of average treatment effect using two moment conditions. Standard errors bootstrapped and clustered at district level.