

# When Does Teacher Incentive Pay Raise Student Achievement?: Evidence from Minnesota's Q-Comp Program

Aaron Sojourner\*

Kristine West

Elton Mykerezi

University of Minnesota

University of Minnesota

University of Minnesota

Carlson School of Management

Dept. of Applied Economics

Dept. of Applied Economics

May 4, 2011

## Abstract

Since 2005, dozens of Minnesota school districts have implemented pay for performance (P4P) plans as part of the state's Quality Compensation (Q-Comp) program. This paper performs the first systematic study of Q-Comp's impact on student achievement, exploiting variation across districts in the timing of participation as well as in the design of districts' P4P plans to study effects on achievement for grades 3 through 8. Results show a consistent zero average effect of Q-Comp participation on both reading and math achievement. However, effects on reading achievement differ depending on the design of the P4P plan. Specifically, districts offering greater rewards for teacher-centered actions or outcomes evidently experienced large gains in reading ( $0.09\sigma$ /\\$1,000 bonus) while those offering rewards based on school-wide goals or formal subjective evaluations did not. Gains from specific P4P design features were not consistently evident in math. We also study effects on other outcomes, such as teacher characteristics and parent demand.

---

\*Corresponding author: [asojourn@umn.edu](mailto:asojourn@umn.edu). Thanks to Avner Ben-Ner, John Budd, David Figlio, Paul Glewwe, Michael Lovenheim, Morris Kleiner, Colleen Manchester, Jason Shaw, Chris Taber, and Joel Waldfogel for comments and to Qihui Chen, Paul Kristapovich, Qianyun Xie, and Yingchun Wang for able research assistance. All errors are ours.

## 1 Introduction

Teachers vary widely in their ability to produce student achievement gains (Hanushek 1971, Hanushek and Rivkin 2010) but this ability is not predicted by educational degrees or experience beyond the first few years of a teacher's career (Hanushek 2003, Aaronson et al 2007). This has large economic consequences (Chetty et al 2010, Hanushek 2010), which motivates policy and research interest in pay for performance (P4P).

Advocates of P4P believe that tying teacher compensation to performance will support increased efforts from incumbent teachers and attract better potential teachers to the profession (Lazear 2003). Many school districts and states are experimenting with P4P plans, which set compensation criteria beyond the conventional ones: experience and education. Plans differ on many dimensions including whether teachers are rewarded individually or in teams, whether for test scores or other measures of teacher quality, and in the magnitude of incentive pay available.

Empirical evidence on the relative and absolute merit of these programs is decidedly mixed. While reviews of the literature point to some gains from P4P (Springer and Podgursky 2007, Neal 2011), evaluations of two large-scale P4P plans that were implemented as randomized trials found null or even negative effects on student achievement (Springer et al 2010, Fryer 2011). Whether plans implemented as long-term policies rather than short-term experiments or plans with other designs would produce better results remains an open question of great interest to policymakers, educators, and economists.

Absent a clear picture of what optimal compensation reform policy should look like and given the political challenges to adoption, the U.S. Department of Education has taken a decentralized approach. Its Race to the Top Fund and Teacher Incentive Fund are allocating billions of dollars for locally-designed reforms via competitive grants. The Department of Education sets guidelines, applicants submit proposals, and federal funds are distributed to support approved plans.

The State of Minnesota took a similar approach when it implemented the Quality Com-

pensation (Q-Comp) program in 2005 as the signature education initiative of Governor Tim Pawlenty. The Minnesota Department of Education (MDE) set general guidelines for acceptable programs and invited districts to propose specific P4P program designs that they would implement. If the proposal was approved, the state authorized up to \$260 per student per year in additional funding to the district.

Districts designed plans that varied along many dimensions. Each district was required to specify the maximum incentive pay they would make available to teachers based on different types of criteria and there is great variation in what they chose. This allows us to construct continuous measures of each district's plan in terms of dollars at stake for: (1) individual teacher actions or outcomes, (2) school wide goals, or (3) through a subjective evaluation process. We exploit this variation in P4P plan designs, along with variation in when districts adopted Q-Comp, to provide evidence on the effect of plan design features on achievement scores and other outcomes.

For a number of reasons, Q-Comp provides an excellent opportunity to learn about the effects of P4P in a policy framework mirroring recent national efforts. First, the program has been in effect for over five years and was implemented as a permanent program rather than a time-limited experiment. Second, there is substantial variation in what criteria trigger P4P bonuses. Third, Minnesota has one of the longest lasting and most widely used inter-district open enrollment policies and a large number of charter schools, so parents have substantial choice in public schooling. This makes it possible to examine the effect P4P designs have on net student movements, which can reflect changes in parent demand. Overall, understanding the Minnesota P4P experience can provide valuable information to policy makers nationwide.

In addition to policy relevance, Q-Comp can provide new evidence on several issues of theoretical importance related to P4P contracts in education. For instance, it is not clear what the optimal team size for targeting bonuses should be. On one hand, incentives tied to school-level criteria may encourage efficient effort if there are positive externalities from cooperation (Weitzman and Kruse 1990) or variations in incentive strength across teachers (Ahn 2008). On the other hand, free riding may make individual or small team incentives

preferable (Kandel and Lazear 1992). Since Q-Comp districts experimented with a wide range of P4P contracts, we are able to investigate whether incentives offered at lower levels of aggregation (such as the individual teacher or grade) are more or less productive than those offered at higher levels of aggregation (such as the school or district level).

There are also important theoretical questions about how to measure teacher quality and performance. Measures based on principal or peer subjective evaluations have received some attention in the literature, especially since principals seem able to identify the best and worst teachers (Jacob and Lefgren 2008). However, high-stakes subjective evaluation processes may be captured and converted into de facto salary augmentations (Neal 2011). Minnesota's Q-comp offers a valuable opportunity to examine if a high-stakes P4P plan based on subjective evaluations affects educational outcomes.

Subjective evaluations have also received attention because they can provide a solution to multitasking problems (Holstrom and Milgrom 1991). Attaching high stakes to standardized test scores may induce teachers to narrow the curriculum, teach to the test, or worse. Subjective evaluations may alleviate this problem by harnessing richer information about a teacher's performance (Baker 1992, Rockoff et al 2010, Tyler et al 2010). Q-Comp's diverse P4P plans, along with the availability of net student movements in an environment of open enrollment as well as student attendance data, used as a proxy for family engagement with learning, allows us to examine if there are significant issues related to multitasking. If multitasking issues are important, incentives tied to test scores may have positive effects on test scores but adverse effects on other educational outcomes and incentives tied to broad subjective evaluation criteria may have adverse impacts on test scores but positive impacts on other measures of quality such as parent demand.

This study finds that districts that put higher stakes on individual teacher (or small team) level actions or outcomes experience bigger increases in student reading scores in standardized tests, on the order of 0.09 standard deviations per \$1,000 of bonus offered. This is a shockingly large estimated effect for a relatively low price.<sup>1</sup> The finding appears quite robust within the

---

<sup>1</sup>Based on Hanushek (2010) and Chetty et al (2010), a lower bound on the social value of a  $0.2\sigma$  achievement gain for a teacher's class each year is conservatively \$200,000.

limits of our study design. Districts that link more rewards to outcomes measured at the school or district level do not see increases in reading achievement and districts that link rewards to subjective evaluations see small declines in reading scores. For math, there are no apparent effects of the incentives tied to teacher incentives or actions but there is evidence on schoolwide incentive effects and negative effects of high-stakes subjective evaluation are similar in some specifications.

The paper finds no evidence that evaluation as a solution to multitasking provides an explanation for the small negative effect of dollars tied to subjective evaluations on test scores. These dollars are also found to reduce enrollment and do not have a statistically significant positive impact on net pupil movements. School P4P dollars, on the other hand, have null or marginally positive effects on student test scores, to which they are almost always explicitly tied, but there is some suggestive evidence that they may increase net pupil inflows.

In addition to the inquiry into the heterogeneous effects of program design features, we investigate the average effect of the Q-Comp program overall. We find no average effects on reading or math achievement. While many districts chose programs that created gains, others chose programs that offset these to produce no effect on program participation on average.

The paper proceeds as follows. Section 2 provides more detail on the Q-Comp program and some of the challenges and opportunities it presents for study. Section 3 briefly reviews relevant theoretical and empirical literature. Section 4 introduces a model and discusses identification. In section 5 we present results on the relative success of different plan designs along with robustness checks and results on the effects of Q-Comp participation in the aggregate. Section 6 concludes with a discussion of how our results add to the existing literature on P4P in education as well as plans for future research.

## 2 Design of and selection into Q-Comp

### 2.1 Q-Comp participation

Q-Comp is sizable. Since its inception in 2005, over one million student-years have been taught in dozens of participating districts and charters and over \$200 million of state funds have been distributed to districts. As one of the nation's largest teacher P4P programs, Q-Comp has attracted significant policy and political attention. Yet little is known about the designs of the local P4P plans it funds or their effects.<sup>2</sup>

Selection into Q-Comp works as follows. The state promised additional annual funding to districts that would implement Q-Comp plans and defined guidelines regarding the content of these plans. Districts (and charters) decided whether to apply and proposed specific program designs. The state then decided whether to accept the proposal. Where teachers were unionized, teachers voted on whether to accept the proposed changes.<sup>3</sup> Districts that clear all these hurdles then participate in Q-Comp.

Since 2005, new districts joined the program each year. Table 1 on page 7 describes the number of districts, schools and students participating and not participating in Q-Comp each year. The population is all Minnesota public schools, including charters each constituting its own "district." In 2005, only eight of the state's 504 districts participated (1.6%). These included 59 of the 2,256 schools with 33,674 of the 838,997 students (4.0%). By the 2009-10 academic year, 14.1% of districts with 28.6% of students participated.<sup>4</sup> Most analysis will focus on grades 3 to 8 because, in these grades, all students took both math and reading tests. Participation statistics are provided for schools in this sample in the bottom panel.

---

<sup>2</sup>Neal (2011) summarizes U.S. and international empirical evaluations of performance pay plans and notes that there has been no previous independent study of Q-Comp. Nadler and Wiswall (2011) study selection into Q-Comp and do not address whether Q-Comp impacts student achievement.

<sup>3</sup>Almost all Minnesota districts are unionized though many charter schools are not. They usually informally negotiated the proposals with administrators in advance and officially ratified after state approval.

<sup>4</sup>A few participating districts dropped out of Q-Comp. These tables reflect stock given exit and entry flow.

Table 1: District and School Q-Comp Participation by Year

Year	Participants			Non-Participants		
	Districts	Schools	Students	Districts	Schools	Students
<b>All schools</b>						
2005-06	8	59	33,674	496	2,197	805,323
2006-07	50	322	183,216	458	1,922	657,346
2007-08	60	397	231,465	456	1,856	606,113
2008-09	70	429	252,716	457	1,786	583,218
2009-10	74	411	239,489	451	1,796	597,141
<b>Schools including at least one grade in 3 to 8</b>						
2005-06	7	52	23,131	404	1,511	567,202
2006-07	36	255	129,754	379	1,338	463,862
2007-08	43	309	162,499	379	1,278	462,980
2008-09	52	328	176,870	381	1,258	413,023
2009-10	56	315	166,697	375	1,256	427,549

## 2.2 District P4P Design Features

Q-Comp is a package of reforms with P4P at its center. We focus our attention on the performance pay component because these are the most interesting theoretically, the best measured in the available data, and the most likely to constitute a real policy change from the pre-adoption period. Data on performance pay are collected primarily from letters sent by the MDE to each district upon approval of its Q-Comp application. The letters detail agreed-upon features of the plan.<sup>5</sup> From these, we create three variables for each district measuring the maximum performance pay available to teachers for the following types of criteria:

1. Teacher P4P\$: anything under a teacher's primary influence. This includes inputs

<sup>5</sup>The stated policy, taken from the application guidelines available online, requires that each district enumerate the maximum incentive pay available to teachers for meeting various goals. These include the amounts of performance pay each teacher

1. is eligible to earn for meeting specified goals for student achievement measured at the teacher, team, or grade level by formative, summative or standardized tests
2. is eligible to earn for meeting specified goals for student achievement measured school-wide or district-wide
3. will earn through the teacher evaluation/observation process.

In creating variables for use in our analysis some judgments were necessary in interpreting the letters. Our results are robust to different interpretations of vague letters. These additional results are available upon request.

related to professional development (e.g., attending meetings, taking classes, completing professional development plans and self-assessments) and outputs close to the teacher (e.g., student performance on teacher-created assessments, a teacher's own students' standardized test scores). It also includes analogous small team or grade-wide outcomes. The descriptions in the approval letters do not allow us to consistently distinguish between various elements within this domain.

2. School P4P\$: anything linked to school-wide or district-wide outcomes. These primarily involve hitting standardized test score targets.
3. Evaluation P4P\$: anything linked to classroom observations or subjective evaluations performed by peers, administrators, or a district-sanctioned mentors including a formal annual review process.

Districts vary in the total levels of pay available across these three dimensions as well as the shares available through each dimension. The value of these variables is shared by all school-grades within a district-year. Table 2 summarizes the distribution of these measures across participating districts.<sup>6</sup> We prefer the summary statistics where district observations are weighted by the number of students tested because this approximates the average at the teacher level and because the regression results use these weights. Participating teachers can earn an average maximum of \$872 a year in incentive pay through actions or outcomes closely connected to them (Teacher P4P\$). They can earn an average maximum of \$1,100 by meeting criteria tied to annual evaluations and classroom observations (Evaluation P4P\$). School or district level goals – usually based on student achievement test scores – can earn them an average maximum of around \$247 (School P4P\$). The triggers for paying out on these dimensions are set according to various locally-designed criteria within and across districts. The marginal distributions of the three variables are graphed in the Figure 1 histograms.

---

<sup>6</sup>If only the share assigned to each dimension, and no dollar values, were listed in the approval letter nor in any available program documents, we assumed the modal total amount among observed districts: \$2,000. Five letters were so ambiguous as to be impossible to code on these dimensions. Although districts and schools may change their designs over time, at this point we do not have measures of this, so we assume they stay constant at the initial levels. White noise measurement error would bias estimates to zero.



Table 2: Summary statistics for district Q-Comp program design variables measuring maximum pay available through each dimension, in thousands of dollars

	Unweighted		Weighted by students		Min.	Max.
	Mean	Std. Dev.	Mean	Std. Dev.		
Teacher P4P\$	0.809	0.581	0.872	0.692	0	2.5
School P4P\$	0.379	0.347	0.247	0.214	0	2.5
Evaluation P4P\$	0.813	0.566	1.1	0.694	0	2.5
Number of participating districts	77					

Note: the 2010 cohort included additional districts but their plans are not coded.

Figure 1: Marginal frequencies of P4P design variables across Q-Comp districts, in \$1,000

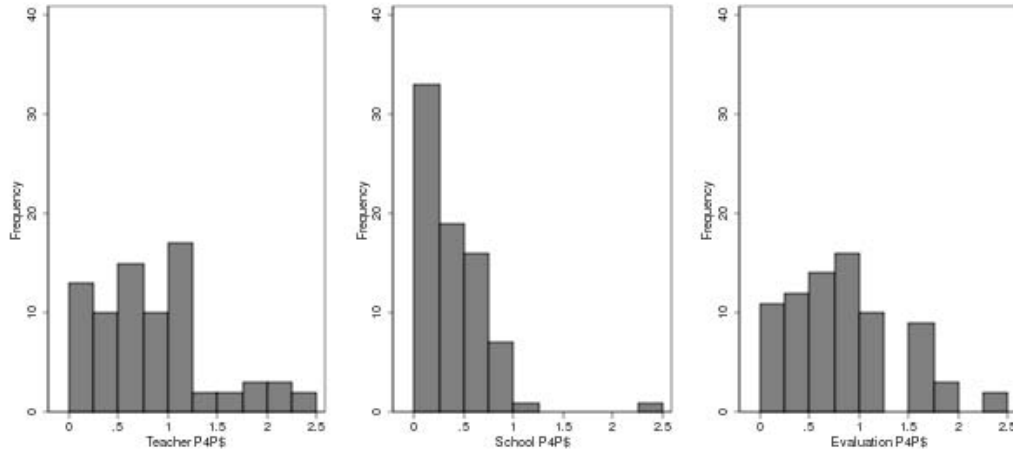


Table 3 describes how the maximums in the three dimensions are correlated. For Teacher P4P\$, we observe a 0.12 correlation with School P4P\$ and a  $-0.80$  correlation with Evaluation P4P\$. We also observe a  $-0.15$  correlation between School P4P\$ and Evaluation P4P\$. So programs that reward teacher-centered actions or outcomes also tend to reward school- or district-level goals at the expense of rewards based on subjective evaluations.

Table 3: Correlation of districts' maximum pay available by dimension, weighted

	Teacher P4P\$	School P4P\$	Evaluation P4P\$
Teacher P4P\$	1.00		
School P4P\$	0.12	1.00	
Evaluation P4P\$	-0.80	-0.15	1.00

Figure 2 on page 10 displays the joint distribution in participating districts across program designs. Each point represents a district's (or a charter school's) Q-Comp P4P design. The size of each point represents the maximum total bonus available to teachers in that district, the sum of Teacher, School and Evaluation P4P\$. Each district's share of awards tied to Teacher P4P\$ criteria is graphed horizontally. The share tied to School P4P\$ criteria is graphed vertically. The remaining share, tied to Evaluation P4P\$ criteria, is represented by the distance to the frontier. For instance, the large dot appearing on the frontier represents a district with a plan offering each teacher over \$4,000 in additional pay annually. Being on the frontier line means that none of the bonus is tied to subjective evaluation. Half the bonus is tied to Teacher P4P\$ criteria and the other half to School P4P\$ criteria. The small dot at the origin represents a district with a plan that awards between \$1,000 and \$2,000 based solely on subjective evaluations.

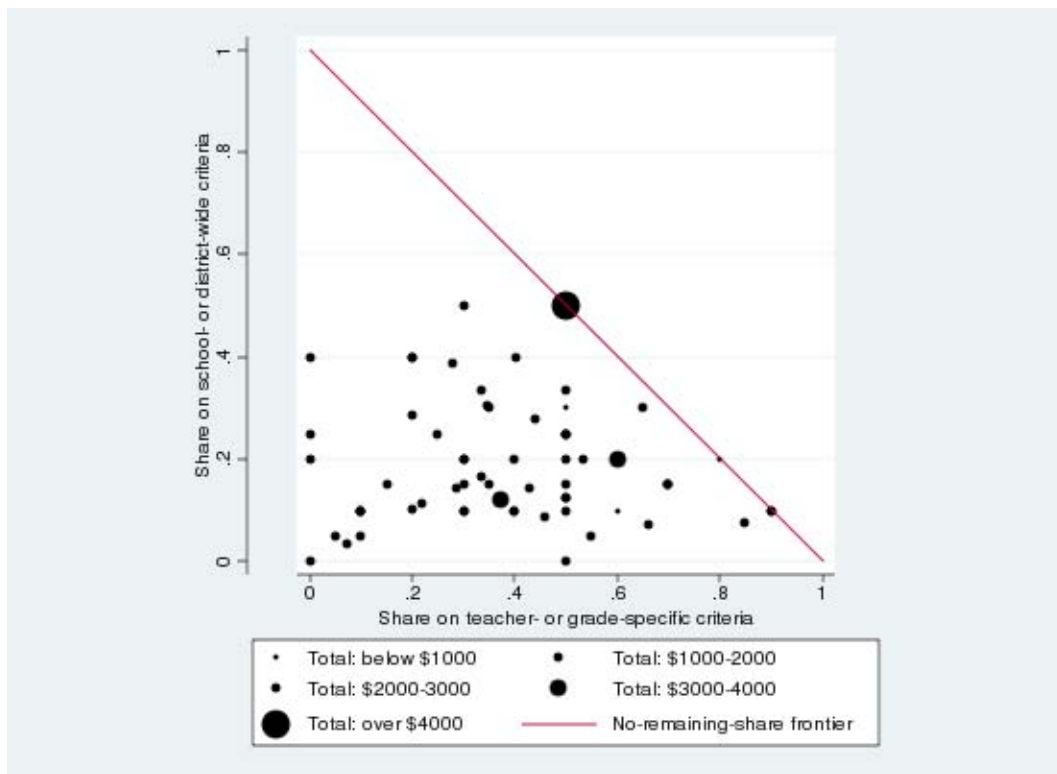


Figure 2: Joint distribution of P4P designs across Q-Comp districts

This graph makes some important points about the program designs clear. First, there is a lot of variation across districts in these four dimensions. They do not cluster around some general “optimal contract.” Second, few districts offer less than \$1,000 or more than \$4,000 in total and most districts offer a mix across all three dimensions. Few lie on the edges of the triangle. Third, none offer more than half tied to School P4P\$ criteria, although there is a lot variation in shares below half. Because of this variation in program design, analyzing Q-Comp in the aggregate likely masks important differences across districts. Accounting for this heterogeneity in program design provides a unique opportunity to understand how P4P plan specifics impact educational outcomes.

One additional element of P4P plan designs is worthy of note. Many Q-Comp districts’ approval letters also specified the academic subjects to which districts elected to tie their school-wide performance bonuses (School P4P\$). These often varied by school within the same district. We coded whether these goals focus on reading or math. Schools were more likely to tie School P4P\$ to reading rather than to math achievement. Three times more school-grades (15.6%) chose to focus exclusively on reading rather than exclusively on math (4.5%). The remainder divided their attention between math, reading, or other subjects. This provides an additional opportunity to examine a number of issues.

### **3 Literature Review**

The substantial variation in the plans’ designs allows us to address at least two major issues where economic theory provides ambiguous guidance and existing evidence is thin. First, as described in the organization design literature on team production, there is a trade-off between offering rewards based on individual versus team outcomes. Given complementarities in production, individual-level incentives can discourage productive cooperation. This pressure pushes for incentives tied to criteria defined at the team level. On the other hand, team incentives also open the door for free riding, a problem that worsens in team size. In schools, grade levels are a natural grouping and additionally, many middle and high schools are organized

into even smaller teams of core subject teachers. These groups may be small enough to exert sufficient peer pressure to overcome the free-rider problem (Kandel and Lazear 1992).

In our analysis of Q-Comp, we group together rewards based on criteria defined at the teacher level with those at the small-group level (i.e. subject teams or grade-level groups). We distinguish these Teacher P4P\$ from School P4P\$ tied to criteria defined at higher levels, such as whole schools or districts, where the free rider problem may be especially severe. Previous empirical literature on the effect of incentives at various levels is mixed.

A randomized trial in Nashville, Tennessee, in which teachers assigned to the treatment group could earn up to \$15,000 based on their individual students' gains on state mathematics tests, found no significant treatment effect on student achievement (Springer et al 2010). We find that Q-Comp districts see increases in reading achievement on the order of 0.09 standard deviations per \$1,000 Teacher P4P\$ bonus available.<sup>7</sup>

Most Q-Comp plans include some school- or even district-wide measures of success in determining bonuses. P4P plans in North Carolina and New York reward exclusively at the school-level. Analysis of the North Carolina program finds no significant impact on student achievement resulting from school wide bonuses (Vigdor 2008). Fryer (2011) presents evidence from a randomized trial in New York City also showing no positive impact on student achievement. In fact, he writes that the school-level bonuses may have decreased student achievement, especially in larger schools. We do not find a significant impact on school- or district-level rewards.<sup>8</sup>

Second, theory provides mixed predictions on the role of subjective versus objective measures of teacher quality. Most P4P plans adjust pay based on objective measures of teaching quality such as test scores. These may not fully capture teachers' ability to affect critical

---

<sup>7</sup>In a randomized trial in Andhra Pradesh, India, Muralidharan and Sundararaman (2011) finds that individual rewards do impact student achievement. Their evidence suggests that individual and small group rewards both had a positive impact on language and math exams with effect sizes between 0.12 to 0.27 standard deviations. In the first year, individual and small group rewards were equally effective. In the second year, individual rewards were more effective. However, in their study, 92% of the treatment schools had between two and five teachers. This is a similar size to small groups of teachers defined at the team or grade level in Minnesota. Also, average teacher effort levels are much lower in India, as measured by high absenteeism.

<sup>8</sup>Contrasting evidence from Israel indicates that school level incentives based on measures such as the average number of credits per student and the dropout rate had a positive and significant impact on student test scores (Lavy 2002).

thinking, non-cognitive skills, or other unobserved yet valuable aspects of learning. This is a multitasking problem (Holstrom and Milgrom 1991). High powered incentives tied to test scores may induce teachers to spend too much time on tested skills at the expense of other socially valuable skills, leading to “teaching to the test” or a “narrowing of the curriculum.”

Adding subjective evaluation criteria may mitigate this problem (Baker, Gibbons and Murphy 1994). Subjective evaluations are especially attractive in schools because they can be used in non-tested subjects. Principals are able to distinguish effective from ineffective teachers as measured by value-added statistics (Jacob and Lefgren 2008, Rockoff et al 2010; Tyler et al 2010). Interestingly, the principal’s overall rating of the teacher is a better predictor of future parent requests for the teacher than value-added statistics (Jacob and Lefgren 2008). This suggests principals’ judgments about which teachers contribute to overall student learning agree with parents’ judgements better than test scores do.

However, principals may be reluctant to use their knowledge of teacher effectiveness when making high-stakes decisions (Jacob 2010, Neal 2011). Neal (2011) attributes the failure of P4P programs in England (Atkinson et al 2009) and Portugal (Martins 2009) to the fact that they were largely based on subjective evaluations done by local staff. Such plans may not improve student achievement because evaluators lack incentives to assess teachers accurately. Neal asserts, with specific mention of Q-Comp, that plans which base pay on locally-defined goals and locally-conducted evaluations can become a “vehicle for raising base pay of most or all teachers whether or not these teachers improve their performance.”<sup>9</sup>

We distinguish between Q-Comp plans that focus on subjective evaluations and those that tie bonuses to specific actions or student outcomes and find some support for Neal’s prediction. Almost all Q-Comp contracts reward subjective evaluations but they differ in the amount of

---

<sup>9</sup>There is evidence that almost all teachers in Q-Comp districts earn at least *some* performance-based pay, often through the subjective evaluation portion. A Minneapolis Star-Tribune investigation found that, in the 22 Q-Comp districts they researched, only 27 teachers got absolutely no performance payment out of the roughly 4,200 teachers eligible. (2/1/2009, Minneapolis Star Tribune). The article reports that “...many educators say [Q-Comp] is strengthening teacher evaluations and training. But others are questioning whether Q-Comp has just become a cash handout.” However, not everyone earns the maximum evaluation payout nor meets the teacher-centered or school- or district-level standards based on student achievement. There are incentives unclaimed so this is not strictly a cash transfer program. In order to get teachers’ assent for Q-Comp adoption, which creates uncertainty and variance in their future pay, teachers may demand some augmentation to base pay. Dollars tied to subjective evaluation may be the vehicle for it.

money at stake for a successful evaluation and in the details of the process or payout criteria.<sup>10</sup> While we are unable to say whether the evaluators correctly identify effective teachers, we can test how making the subjective evaluation process central to a P4P plan affects student achievement. Our evidence suggests that plans which tie bonuses primarily to evaluations may even be detrimental to student achievement on standardized tests. Districts with larger rewards linked to evaluation appear to score slightly worse in reading than they did before Q-Comp adoption.

Neal’s concern that subjective evaluations may be made into de facto salary augmentations is premised on the assumption that schools do not face the competitive pressures that are at work in the private sector. In Minnesota, however, there are two mechanisms for subjecting districts to the discipline of the market. Minnesota has the nation’s longest standing open enrollment and charter school laws. Parents can enroll their children in any public district or charter as long as there is space. Because of this, we also test the impact of Q-Comp on net pupil movements between regular districts and from districts to charters.

## 4 Model and Data

To learn about the impact of Q-Comp on student achievement, we analyze a panel of student achievement, demographic, and school characteristic data defined at the year-school-grade level using generalized difference-in-difference methods. Our primary achievement measure is the Minnesota Comprehensive Assessments Series II (MCA-II) average scores. Since 2005-06 (coincidentally the first year of Q-Comp), these have been mandated for every student in third to eighth grade in both reading and math.<sup>11</sup>

We study how schools’ student achievement changes as their Q-Comp participation changes. The main outcome is average student achievement on MCA-II tests each academic year indexed  $t = 2005, 2006, \dots, 2009$ , in each school indexed  $s = 1, 2, \dots, S$ , in each tested grade indexed

---

<sup>10</sup>The wording in many Q-Comp plans indicates rewards given simply based on the completion of a subjective evaluation cycle (i.e. a set number of observations coupled with pre- and post-observation meetings with the evaluator). Not all plans differentiate pay based on the evaluator’s assessment of teacher effectiveness.

<sup>11</sup>Unfortunately, Minnesota switched testing regimes in 2005. Prior to 2005, only grades 3, 5 and 7 were tested and on a different test.

$g = 3, 4, \dots, 8$ , and in either math or reading indexed  $b \in \{M, R\}$ .<sup>12</sup>

In explaining average student achievement, we use variants of this model:

$$y_{tsgb} = \beta_{gb}Q_{tsgb} + \alpha_{gb}w_{tsg} + \gamma_{sgb} + \delta_{tgb} + \epsilon_{tsgb} \quad (1)$$

Interest centers on the effects of Q-Comp participation and of features of the P4P designs adopted. These are captured by  $\beta_{gb}$ . In the simplest form, we will characterize Q-Comp participation as either

- (A) a simple participation dummy, 1(Post-adoption) or
- (B) the participation dummy paired with a dummy indicating academic years two or more years prior to adoption, 1(2+ pre-adoption) (reference category: the single year immediately prior to adoption).

Moving across specifications (A) and (B) the definition of  $Q$  changes but the model is otherwise the same. Specification (A) treats the whole pre-adoption period as the reference category. Specifications (B) conditions on and measures pre-adoption differences in achievement levels between adopters and non-adopters using a dummy to indicate observations that come from years more than one year pre-adoption. The fact that our data start in the year that the first cohort adopted and that different-sized cohorts adopted during each year of our data generate imbalances in what data are available to identify various parameters. All observations from more than one year pre-adoption come from the smaller cohorts of districts that adopted in 2007, 2008 or 2009. Defining  $Q$  as a simple participation dummy allows  $\beta$  to measure the effect of Q-Comp participation on average.

To measure the effects of various aspects of Q-Comp P4P program design, we use different definitions of  $Q$ . We define  $Q$  to be a vector measuring P4P design features in years that they are participating and zeros otherwise. In particular, as discussed in section 2.2 we measure maximum P4P bonus available to teachers in a district for three kinds of criteria, measured

---

<sup>12</sup>Before third grade, students are not tested. After eighth, tenth graders are tested only in reading and eleventh graders only in math. Estimated effects for these two series are also available on request.

as Teacher P4P\$, School P4P\$, and Evaluation P4P\$.

Additionally, we include a variable to indicate district-years where the district once participated in Q-Comp but has since dropped out. This only affects a small number of districts. If the estimated coefficient on this dummy is negative it indicates districts do worse after leaving Q-Comp than they did in the year(s) prior to adoption.

Since Q-Comp participation is not randomly assigned, there may be systematic unobserved differences between districts that influence both Q-Comp adoption and our outcomes, which would bias estimates of program effects. We use four main strategies to guard against this threat.

First, since within any given school and grade, average student achievement may vary over time due to differences in student cohorts, we condition on a vector of year-school-grade student demographic characteristics and school-level variables ( $w_{tsg}$ ).<sup>13</sup> These are listed in the top panel of Table 6 on page 34, which also provides summary statistics. These characteristics do not vary across subject, although their coefficients  $\alpha_{gb}$  can.

Second, school-grade-subject fixed effects ( $1_{sgb}$ ) are included to remove time-invariant, additive unobserved differences in achievement levels  $\gamma_{sgb}$  between schools. The model is identified from within-school-grade-subject, across-time variation.

Fixed effects for each year-grade-subject  $1_{tgb}$  are also included.<sup>14</sup> These terms identify counter-factual year effects for each grade and subject ( $\delta_{tgb}$ ). The comparison group matters here because their experience across years defines these time effects. This is a generalization of difference-in-difference analysis that relies on differences in the timing of adoption across districts to separate time effects from program effects.<sup>15</sup>

---

<sup>13</sup>This alone improves over available evidence on Q-Comp's effects. Neither a legislative auditor's report (Nobles, 2009) nor a state-commissioned external report (Hezel Assoc., 2009) dealt with selection or covariates.

<sup>14</sup>Also, each year-grade-subject distribution of student scores is standardized to have mean zero and standard deviation one in order to facilitate pooling across grades.

<sup>15</sup>The first difference is the within-school comparison across time periods. The second difference is between the first-differences at adopting schools and those at non-adopting schools across the same time period. A within-school change between any two points in time is evaluated against changes across those same two years among other schools. With a simple participation dummy,  $\beta_{gb}$  measures the difference in average grade- $g$ , subject- $b$  achievement within adopting-schools in the years after adoption compared to the years prior to adoption conditional on changes in  $w_{tsg}$  and the average change experienced across these years by other schools. Lovenheim (2010) uses a similar approach to study the effects of teachers unionizing in different districts over time.



The model is identified by assuming that program variables  $Q_{ts}$  are uncorrelated with unobserved influences  $\epsilon_{tsgb}$  conditional on other observables, school-grade fixed effects, and year-grade fixed effects,

$$Cov[Q_{ts}, \epsilon_{tsgb} | (w_{ts}, 1_{sg}, 1_{tgb})] \equiv 0 \quad (2)$$

Within the restrictions of functional form, this model yields unbiased estimates of program effects even if selection into Q-Comp is based on stable differences in achievement levels. If, for instance, schools with higher achievement levels are more likely to adopt or to adopt earlier than schools with lower achievement levels, that is not a problem. The crucial assumption is that within-school, time-varying, unobserved influences on achievement levels are not systematically related to whether or when a school adopted Q-Comp or the features of the design it adopted. The estimates of  $\beta$  may be biased if districts select into participation or design based on fluctuations in achievement levels. For example, if a school is more likely to adopt in a year when levels would rise for other reasons than in a year when they would fall (perhaps, districts experimenting with Q-Comp are also experimenting with other reforms), this violates the identifying condition and would bias the estimated program effect upwards.

Third, we estimate these models with three different comparison groups. We compare the experience of participants to that of either (1) all other schools in the state, (2) districts that applied to Q-Comp but failed to adopt, due either to the state rejecting the proposal or their teachers voting against it,<sup>16</sup> and (3) just Q-Comp adopters who have not yet adopted. We refer to comparison group (2) as “interested” nonadopters and to the these three samples as the full, interested-only, and adopters-only samples, respectively. Excluding never-applicants from the analysis reduces precision because they contain information about the effect of observable characteristics  $w$  and the time effects  $\delta$ . However, excluding them can reduce bias if they are fundamentally different from adopters or applicants in unobservable, time-varying ways. Also, unlike never-applying districts, interested nonadopters passed the first hurdle to participation; they choose to apply. Some even cleared the second hurdle (state approval). In

---

<sup>16</sup>Failed applications had to be obtained through a Freedom of Information Act to the state Dept. of Education.

this sense, interested nonadopters are more similar to adopters than the never-apppliers are. Parameter estimates across all samples are provided for comparison and results turn out to be very stable.

Figures 3 and 4 on pages 18 and 19 present trends in average reading and math achievement among each adoption cohort, the cohort of never-adopters, and among interested non-participants. In each grade-year-subject, test scores are normalized to be mean zero and standard deviation one across schools weighted by students tested.

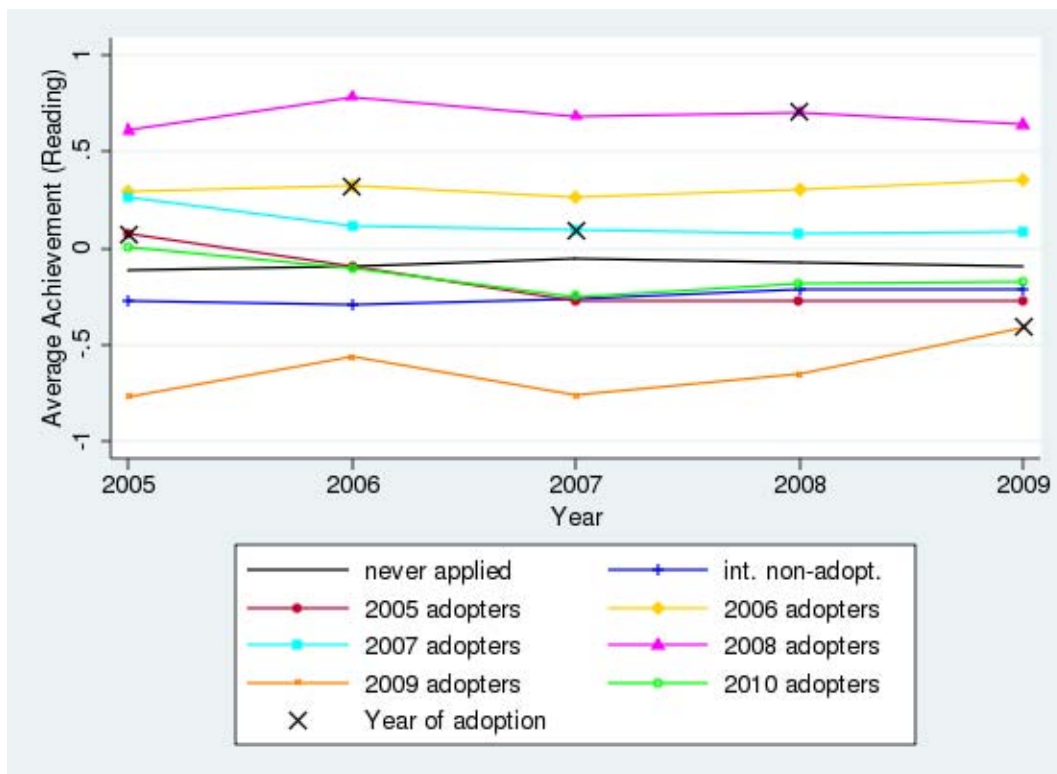


Figure 3: Trend in average reading achievement by Q-Comp adoption cohort

There are three points to make about these trends. First, there are differences in average achievement levels between cohorts. The never-applied cohort is the largest and hovers just below state mean achievement throughout the period. The interested non-adopters' trends are just below the never appliers generally. Among Q-Comp adopters, the 2005 and 2010 cohorts are most similar to the never-adopters and the interested non-adopters on average.<sup>17</sup> However,

<sup>17</sup>Achievement data for the 2010-11 year is not yet available.

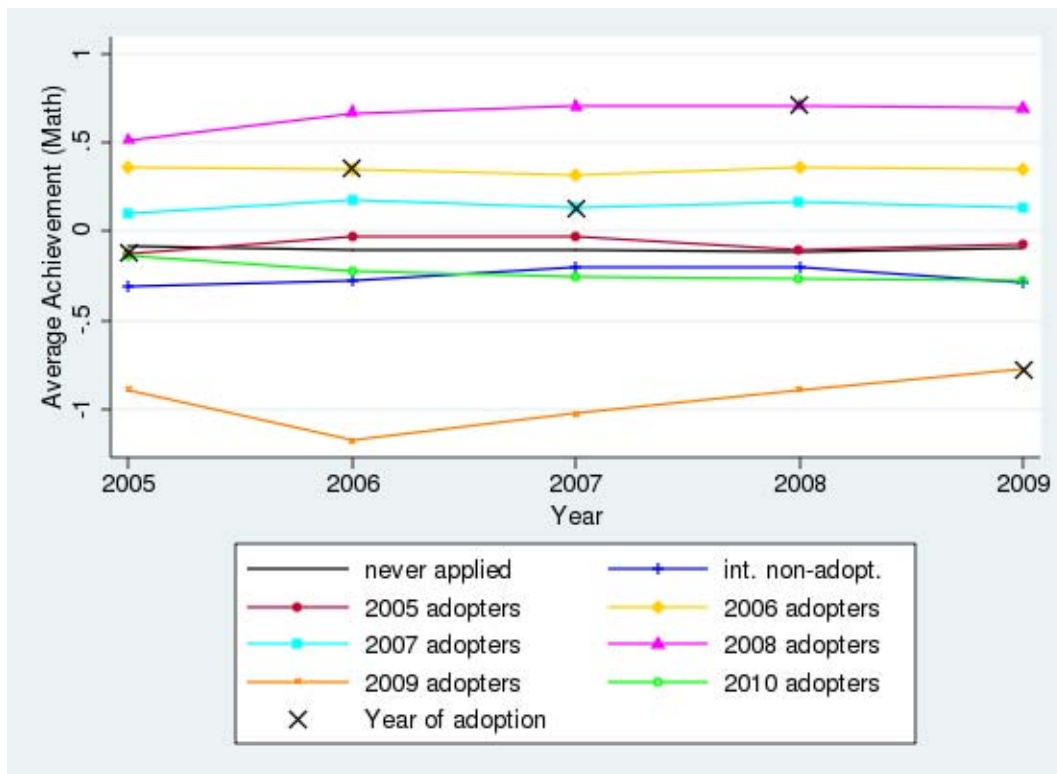


Figure 4: Trend in average math achievement by Q-Comp adoption cohort

the 2006, 2007, and 2008 adoption cohorts were higher achieving than average. The 2009 adopters are lower achieving on average, around a half to a full standard deviation below the mean in math and reading. Second, there do not seem to be large differences in achievement trends between cohorts, aside from the fluctuations in the very small 2009 cohort. Third, this foreshadows one of our main conclusions: the effects of Q-Comp participation appear to be null on average. Increases in achievement do not seem to follow Q-Comp adoption in the aggregate.

The above school-grade-subject level models can, however, still be susceptible to endogenous plan design based on time varying unobservables. If administrators were able to forecast future achievement and successfully designed plans that differed based on these forecasts the difference in difference estimators above could produce biased results.<sup>18</sup>

<sup>18</sup>We say “successfully” because if unobservable were simply correlated with application the “interested only” sample still produces unbiased results. So time-varying unobservables have to be correlated with actual implementation and/or design of P4P plans upon implementation.

So, lastly, we estimate growth models that condition on lagged achievement. We do this at the district level in order to avoid losing students as they transition between schools across grades. At the district level we can obtain lagged scores based on all students, adding  $\vec{y}_{(t-1)d(g-1)b}$  as a covariate to explain  $y_{(t-1)d(g-1)b}$ . These specifications do not use all the variation across grades and schools that the above models do but are more robust to omitted time varying variables that affect P4P plans and achievement growth. The results turn out to be quite similar qualitatively.

In order to boost power, our primary results pool across grades 3 to 8 and restrict program effects to be the same across grades within subject,  $\beta_{gb} \equiv \beta_b$  but some results by grade level are also provided.

Although most participation decisions are made at the district level, a few large districts allowed individual schools to participate in the program, so the participation decision was coded accordingly at the school level. Use of school-grade level data increases precision. However, standard errors adjusted for correlation at the district level are provided throughout.<sup>19</sup>

## 5 Results

### 5.1 Did Q-Comp Change Teacher Pay and Incentives?

Before examining the impact of Q-Comp P4P design features we first present three pieces of evidence on whether Q-Comp actually changed teacher incentives as advertised. It is often the case that innovative grant recipients find ways to elicit funds to accomplish tasks that they would perform even in the absence of the grant. Null effects would result. Was this the case with Minnesota school districts? By supplementing our data with information from the National Center for Educational Statistics' Schools and Staffing Survey (SASS) and by conducting a survey of district human resource professionals, we provide evidence that Q-Comp actually did change the way teachers were paid.

The SASS asks districts whether they use any pay incentives to “reward excellence in

---

<sup>19</sup>Clustering at the school or school-grade level provides identical estimates and smaller standard errors. However, we present the most conservative estimates.

teaching.” Q-Comp participation is significantly associated with switches from “No” before Q-Comp adoption to “Yes” after adoption among Minnesota districts. Table 4 reports on the 55 Minnesota districts sampled in both waves. Among districts not participating in Q-Comp in 2007-08, 96% report no pay for excellence both before and after Q-Comp started. Among districts participating in Q-Comp in 2007-2008, none reported paying for excellence before Q-Comp, in 2003-04. However, in contrast to the nonparticipants, 58% of participants report paying for excellence in the post-adoption SASS survey wave. This suggests that, many districts perceived something programmatic to have changed. The fact that 42% of surveyed Q-Comp districts still reported no pay for excellence also suggests that not all districts experienced deep changes or conceptualized Q-Comp in this way.

Table 4: Evidence on change in “Pay for Excellence” among Minnesota districts by Q-Comp participation status from the Schools and Staffing Survey (SASS)

District in Q-Comp in 2007-08	Districts in both waves	In 2003-04 In 2007-08	Can teachers earn extra pay “for excellence”?				Total
			No	Yes	No	Yes	
Yes	12		42%	0%	58%	0%	100%
No	43		96%	2%	0%	2%	100%

Note: only these 55 districts appear in both the 2003-04 and 2007-08 waves of SASS. Q-Comp began in 2005.

In order to get more detail on the particular aspects of the P4P programs implemented in Q-Comp schools, we conducted an independent phone survey of district human resource professionals about their district pay practices and without mentioning Q-Comp. It found that participating districts are significantly different from nonparticipants in how they compensate teachers. We obtained data from 92 districts (38% response rate). Twenty-one of these, we know from administrative data, participate in Q-Comp. Table 5 summarizes our findings. Among Q-Comp participants, 86% report paying for student performance and 90% report paying for subjective evaluations. In stark contrast, none of the non Q-Comp districts report paying on either of these dimensions. Unsurprisingly, participating districts are just as likely to pay for years of experience and educational credentials as are non-participants. Q-Comp

P4P is a supplement to, rather than a replacement of, traditional compensation criteria.

Table 5: Evidence on Q-Comp’s impact on compensation from author survey in 2010

Districts in Q-Comp in 2010-11?	Percent of districts paying for:				<i>N</i>
	Student Perform.	Subjective Evaluation	Years of Experience	Education Credentials	
Yes	86%	90%	95%	95%	21
No	0 %	0%	100%	100%	71

Lastly, using mean teacher pay as a dependent variable in our analysis shows that the introduction of Q-Comp is associated with a 2.5% increase in average teacher salaries. This is consistent with an average salary of \$55,000 and an average Q-Comp bonus paid of \$1,375 per year per teacher in participating districts.

## 5.2 Impact of P4P Design

We begin by estimating the impact of program design parameters on standardized test scores. Table 7 presents estimates for the effects on reading pooled across grades 3-8. As noted, specification (A) compares scores in post-adoption years to all pre-adoption years. Specification (B) compares post-adoption years to the single year prior to adopting. All specifications condition on time-varying student demographics, school-grade effects and grade-year effects. The full sample includes 4,677 school-grades with multiple observations across years for each. Together they include 1,749,818 tested student-years. Each school-grade-year-subject observation is weighted by the number of students tested.

Schools which offer more performance pay based on factors over which teachers exert closer control produce somewhat large achievement gains in reading. This result is consistent across alternative comparison groups. Columns 1 and 2 present estimates using the full sample, columns 3 and 4 present estimates using only the sample of interested districts (those that ever applied for Q-Comp), and columns 5 and 6 present estimates using only districts that ever participate in the program at some point. The parameter estimates are positive and significant across specifications and samples, ranging from 0.087 (0.025) to 0.112 (0.026) per \$1,000 at stake.

The parameter estimates on School P4P\$ are positive and estimated very imprecisely. So there is no evidence that school or district-level incentives increase test scores. This could be related to the fact that the average maximum bonus for School P4P\$ is quite low, \$247, and only 1 small charter school exceeds \$1,000. Evaluation P4P\$ in the full sample have a negative and statistically significant impact on achievement. Parameter estimates are slightly smaller in the interested only and adopters only samples, and standard errors are higher so the results are not statistically significant at conventional levels. The results suggest that districts adding high stakes subjective evaluations did worse on reading student achievement.

Because plan design is endogenously chosen by the districts rather than randomly assigned, these design feature “effects” really measure the combination of selection and the effect of the plan feature. If this were purely selection, however, to get oppositely signed effects on Teacher P4P\$ and Evaluation P4P\$ would require that districts anticipating abnormally *positive* (compared to non-adopters) time-varying growth trends be more likely to adopt incentives tied to teacher-based targets and those anticipating abnormally *negative* trends be more likely to adopt incentives tied to evaluation-based targets. That this would be the case is far from not obvious, especially since Teacher P4P\$ are not generally tied to standardized achievement test scores. They are commonly tied to professional growth plans or professional learning communities, through which teachers are engaged in a process of reflection, planning, goal-setting, commitment and accountability to their peers and administrators.

We also estimated the effect of P4P on achievement growth rather than on achievement levels. The results are robust to this alternative specification. These models include lagged measures of achievement as predictors and, in this specification, district-grade fixed effects pick up differences in stable growth trends for each grade across districts rather than differences in levels. Parameter estimates in Table 8 continue to indicate a significant impact of Teacher P4P\$ on reading scores with magnitudes similar to those in the levels models (Table 7). The estimated impact of Teacher P4P\$ is now 0.081 (0.037) in the interested only sample, specification (B). The estimated impact of Evaluation P4P\$ is still negative, and of somewhat smaller magnitude. The effect is also estimated less precisely absent within-district variation.

The estimated impact of these same incentives on math scores is less clear.<sup>20</sup> Estimates presented in Table 9 indicate no effect of Teacher or Evaluation P4P\$ on achievement *levels*. School P4P\$, on the other hand, show a marginally significant positive effect, but only on the specifications that use all pre-adoption years as a reference and not in those that use the single pre-adoption year and condition on differences in prior years' achievement. However, the district level *growth* models presented in Table 10 indicate patterns more similar to those in reading for the Evaluation P4P\$, but not for teacher-level incentives. Specifically, estimates imply a negative Evaluation P4P\$ impact on math scores of similar magnitude to the impact on reading scores. For instance, the estimated effect is -0.044 (0.019) in the interested only sample specification (B).

The estimated negative impact of Evaluation P4P\$ is consistent with an environment where multitasking is of concern and evaluators can observe a richer measure of teacher quality than is reflected on test scores. If so, rewarding this broad measure may focus teacher effort away from tested material resulting in lower test scores, but still benefit students and society through improved learning on other valuable dimensions.

In a state with a long standing and widely used open enrollment policy and many charter schools, net pupil movements between districts as well as total enrollment could reflect parent satisfaction with schools. If parents observe a richer measure of the overall educational experience, an examination of the impact that P4P plans had on pupil movements would offer some evidence on whether substantial multitasking issues are at play.

Table 11 shows the effect of different P4P plan dimensions on several teacher attributes and student outcomes.<sup>21</sup> The only statistically significant effect of high stakes subjective evaluations appears to be on enrollments and it is negative. These results do not suggest an environment where subjective evaluations promote unobserved quality valued by parents. Of course, this may also reflect the fact that our alternative educational outcome measures are imperfect. Parent demand may respond to changes in quality only slowly. Teacher P4P\$

---

<sup>20</sup>The number of observations is slightly different for reading and math because year-school-grade-subject scores are not released by the state when there are fewer than ten students tested and this varies across subject.

<sup>21</sup>This analysis uses data back to 2003. Summary statistics are in Table 6.



appears to have no effect on enrollment or net student flow, suggesting that parents do not respond to the induced achievement gains either because they do not value them highly or do not know about them. School P4P\$ has a marginally significant positive effect on net pupil movements.

In Table 11 we also examine the impact P4P design features had on teacher average salaries to examine if the incentives actually had an impact on pay, as well as experience and education levels in order to examine whether P4P plans induced teacher sorting. There is evidence of an impact of both Teacher and Evaluation P4P\$ on salaries, but not of School P4P\$. The estimated parameters are consistent with average salaries and average bonuses paid out by Q-Comp and, thus, provide good evidence that Q-Comp is associated with real changes to compensation.

We find no evidence of teacher sorting on two observable dimensions, percentage of teachers with masters degrees and years of experience. Additionally, teacher sorting with respect to teacher education or experience is not an important channel through which these P4P contracts generated changes in student test scores. When the percentage of teachers with a masters and average teacher experience are entered into the student achievement models as control variables, they are generally not significant and do not mediate the relationship between Q-Comp and student achievement. This is consistent with prior research showing these variables to have a weak relationship with value-added.

We also investigate possible multitasking by using data on what subject school P4P\$ are tied to. Table 12 presents additional models using an indicator of whether the subject in question is the only subject that School P4P\$ are tied to:  $1(\text{only high stakes goal})_{tsgb}$ . For each subject, we estimate the main effect of the three P4P\$ dimensions as well as interacting them with the only-goal indicator. None of the P4P\$ dimensions are significantly more effective when applied in a high-stakes year-school-grade-subject. While imprecise, the magnitude of the School P4P\$ interactions are positive and large, consistent with what one would expect from theory. We interpret this as weak evidence against multitasking as a big concern. Further, these results suggest that any differences in the P4P\$ effects between math and reading

are not primarily due to differences in the incidence of goals set across subjects.

As a robustness check, we also estimate the model with alternative sets of conditioning variables ( $w_{tsg}$ ). This generates further evidence about the lack of a mediating role played by student, teacher and district changes. Our primary analysis in Tables 7 and 9 uses student demographics and total enrollment at the school-grade-year level. Table 13 shows that the results are robust to alternative conditioning sets. The first column shows the effect of P4P\$ excluding student demographics and grade enrollment. Only fixed effects and a pre-trend are included. The top panel shows reading and the bottom math. The second column reproduces the results from Tables 7 and 9 for comparison. The third column adds two teacher variables: average experience and percent with a masters. The fourth column adds three district administrative variables: general reserve fund balance as a percent of previous year expenditures, net pupil inter-district movements, and  $\log(\text{average teacher salaries})$ . All the results for reading and math are quite stable. Table 14 performs the same exercise for the growth models in Tables 8 and 10, with similarly stable results.

As an additional robustness check, we also test whether omitting any one cohort dramatically changes the results. This is especially important for two reasons. First, our identification depends on variation in the timing of adoption and the assumption that timing is not correlated with unobserved achievement trends. Therefore, dropping cohorts will help clarify if different cohorts are getting different effects from their designs. Second, because both the Q-Comp program and the outcome data start in 2005, there are no pre-adoption trends for the first two cohorts. Dropping these cohorts can reveal whether they are driving the results or whether the results generalize to the later cohorts where pre-trends are available. Table 15 gives the results. Generally, the results are stable. When the biggest cohort, 2006, is dropped, the reading results in the top panel weaken somewhat and become less precise. However, the results are qualitatively very similar. For math, the results are qualitatively stable. Teacher P4P\$ becomes larger and significant when either the 2006 or 2007 cohorts are dropped. Table 16 performs the same exercise for growth models with similar results.

Finally, we take a look at whether the different P4P designs examined in this paper

have disproportionate impacts by grade. Table 17 indicates that the positive impact of P4P incentives at the teacher level on reading is present in most grades, with largest magnitudes in higher grades (seven and eight). The negative impact of high stakes peer evaluations is spread across most grades. In math, Evaluation P4P\$ show up as negative and significant in many grades, though this does show up in the pooled results. School P4P\$ are usually positive and not significant. Against pattern, the estimated impact of teacher-level incentives is significant and negative in grade five.

### 5.3 The Overall Effect of Q-Comp Participation

Recent national efforts to reform teacher compensation follows a similar general approach as Q-Comp in that they set guidelines and accept proposals from districts. How did Minnesota's program fare overall with this flexible approach? What was the average effect of the program after six years?

Table 18 presents estimated effects of program participation on reading (math) achievement pooled across grades 3 to 8 in the upper (lower) panel. Across all samples and in both subjects, we see evidence of a null effect. In math, specification (B) reveals evidence that participating districts may have been already improving in the years prior to adoption. The omitted category here is the year immediately prior to adoption. Therefore, the  $-0.074$  ( $0.038$ ) estimated coefficient on 1(Pre-adoption) implies that adopting districts were doing worse between four and two years prior to adoption than they were in the year immediately prior to adoption. However, once they adopted, the progress did not continue.

## 6 Conclusion

Incentives tied to criteria defined at the teacher- or small team-level had a large, robust, positive impact on reading achievement. It is worth recalling that the Teacher P4P\$ were, sometimes but not usually, tied to measure of student achievement analyzed here. Often they were tied to completion of professional development plans or student achievement on teacher-

developed assessments. Getting such large effects from incentives of this size is somewhat surprising. Springer et al (2010) report experimental results with much larger incentives and null effects on math achievement. This difference could be due to many possible differences: prospective time-horizon for the programs, research designs, the substance and attainability of the payoff targets, or a difference between math and reading. We also found no effects of Teacher P4P\$ on math. We developed evidence that the different effects of Q-Comp Teacher P4P\$ between reading and math is not due to the difference in the incidence of school-level math versus reading goals.

We find no evidence that rewards tied to larger groups (i.e. school- or district-level) led to achievement gains. This is consistent with evidence North Carolina (Vigdor, 2008) and New York City (Fryer, 2011). Taken with the prior result, this evidence is consistent with the idea that free riding may be an important problem for incentives defined at high levels of aggregation in education.

Subjective evaluations have been proposed as a potentially important component of performance pay for teachers. This is largely based on studies such as Jacob and Lefgren (2008), Rockoff et al (2010) and Tyler et al (2010) which show that evaluations are correlated with value-added measures of teacher effectiveness. These evaluations, however, are conducted in low-stakes environments rather than being attached to bonuses.

We test whether attaching bonuses to these evaluations benefits achievement, while remaining agnostic on whether evaluators are able or are choosing to distinguish teachers by quality. There is no evidence here that high stakes evaluations result in improvements in student achievement. If anything, we find that reading test scores decrease in districts that attach high stakes to subjective evaluation. The fact that high stakes evaluations may decrease test scores does not necessarily mean that they are ineffective. Subjective evaluations may be doing their job and solving the multitasking problem, discouraging teachers from teaching to the test. However, our results using measures of family demand for education — namely attendance, enrollment and net pupil movements — do not support this hypothesis. If tying dollars to evaluations led teachers to produce more engaging and desirable lessons

then we might expect to see increases in these alternative measures of educational quality. We do not.

A pessimistic interpretation of these results would be that high stakes evaluations do not elicit productive effort, perhaps because of the capture issues discussed in Neal (2011), and may even divert effort from productive activities. There may be a *dog-and-pony show* effect, where teachers divert effort towards developing observational experiences evaluators value but that do not benefit measured student achievement or parent-assessed education quality.

The experience in Minnesota adds to our understanding of locally-designed P4P plans. The grantor-grantee relationship between education authorities and districts has advantages because it allows use of local information and experimentation in finding appropriate, feasible designs. Our findings suggest that if a granting authority proposes a range of reforms and allows districts to design plans locally, many districts (in cooperation with local teachers' unions) will design plans that base rewards largely on subjective evaluations and this does not seem to benefit student achievement. On the other hand, some districts (in cooperation with their local teachers' unions) will weight rewards to more specific teacher-centered criteria and this appears beneficial for reading achievement.

The fact that, despite large gains in some areas of the program, Minnesota spent \$200,000,000 to get a net effect of zero also points out risks associated with too much local control over the plans. Some plans will operate to extract rents from the state more than to improve education. State and federal governments can, however, use the experiences of early adopters, such as Minnesota, to choose more appropriate program guidelines. These tradeoffs have been explored in other contexts through principal-agent models and there is great potential to apply them in this setting.

## 7 Works Cited

1. Aaronson, Daniel, Lisa Barrow, and William Sander. 2007. Teachers and Student Achievement in Chicago Public High Schools. *Journal of Labor Economics*, 25(1).
2. Ahn, Tom. 2008. The Missing Link: Estimating the Impact of Incentives on Effort and Effort on Production Using Teacher Accountability Legislation. Duke University Working Paper.
3. Atkinson, Adele, Burgess, Simon, Croxson, Bronwyn, Gregg, Paul, Propper, Carol, Slater, Helen and Deborah Wilson. 2009. Evaluating the Impact of Performance-related Pay for Teachers in England. *Labour Economics*, 16(3): 251–61.
4. Baker, George. 1992. Incentive Contracts and Performance Measurement. *Journal of Political Economy*, 100: 598–614.
5. Baker, George, Robert Gibbons, and Kevin J. Murphy. 1994. Subjective Performance Measures in Optimal Incentive Contracts. *Quarterly Journal of Economics*, 109, 1125–56.
6. Chetty, Raj, John N. Friedman, Nathaniel Hilger, Emmanuel Saez, Diane Whitmore Schanzenbach, and Danny Yagan. 2010. How Does Your Kindergarten Classroom Affect Your Earnings? Evidence from Project STAR. National Bureau of Economic Research Working Paper No. 16381.
7. Fryer, Roland. 2011. Teacher Incentives and Student Achievement: Evidence from New York City Public Schools. NBER Working Paper 16850.
8. Glazerman, Steven and Allison Siefullah. 2010. An Evaluation of the Teacher Advancement Program (TAP) in Chicago: Year Two Impact Report. Mathematica Policy Research, Inc.
9. Hanushek, Eric. 1971. Teacher Characteristics and Gains in Student Achievement: Estimation Using Micro Data. *The American Economic Review*. 61(2): 280-288.

10. Hanushek, Eric. 2003. The Failure of Input-Based Resource Policies. *The Economic Journal*. 113(485): F64-F98.
11. Hanushek, Eric. 2010. The Economic Value of Higher Teacher Quality. National Bureau of Economic Research Working Paper 16606.
12. Hanushek, Eric. and Steven Rivkin. 2010. Generalizations about Using Value-Added Measures of Teacher Quality. *American Economic Review*. 100(May 2010):267-271.
13. Hezel Associates, LLC. 2009. Quality Compensation for Teachers Summative Evaluation.
14. Holmstrom, Bengt, and Paul Milgrom. 1991. Multitask Principal-Agent Analyses: Incentive Contracts, Asset Ownership, and Job Design. *Journal of Law, Economics, and Organization*, 7: 24-52.
15. Jacob, Brian and Lars Lefgren. 2008. Can Principals Identify Effective Teachers? Evidence on Subjective Performance Evaluation in Education. *Journal of Labor Economics*, 26(1):101-36.
16. Jacob, Brian. 2010. Do Principals Fire the Worst Teachers? National Bureau of Economic Research Working Paper 15715.
17. Johns, Emily. 2009. Is It “Merit Pay” If Nearly All Teachers Get It? *Minneapolis Star Tribune*.
18. Kandel, Eugene and Edward Lazear. 1992. Peer Pressure and Partnerships. *Journal of Political Economy*, 100(4): 801–17.
19. Lavy, Victor. 2002. Evaluating the Effect of Teachers’ Group Performance Incentives on Pupil Achievement. *Journal of Political Economy*, 110(6): 1286–1317.
20. Lazear, Edward P. 2003. Teacher Incentives. *Swedish Economic Policy Review*, 10.

21. Lovenheim, Michael. 2009. The Effect of Teachers' Unions on Education Production: Evidence from Union Election Certifications in Three Midwestern States. *Journal of Labor Economics*, 27(4):525-87.
22. Martins, Pedro. 2009. Individual Teacher Incentives, Student Achievement and Grade Inflation. IZA Discussion Paper No. 4051.
23. Muralidharan, Karthik and Venkatesh Sundararaman. 2011. Teacher Performance Pay: Experimental Evidence from India. *Journal of Political Economy*. 119(1).
24. Nadler, Carl and Matthew Wiswall. 2011. Risk Aversion and Support for Merit Pay: Theory and Evidence from Minnesota's Q Comp Program. *Education Finance and Policy*, 6(1):75-104.
25. Neal, Derek. 2011. The Design of Performance Pay in Education. NBER Working Paper No. 16710.
26. Nobels, James. 2009. Evaluation Report: Q Comp Quality Compensation. Minnesota Office of the Legislative Auditor.
27. Podgursky, Michael and Matthew Springer. 2007. Teacher Performance Pay: A Review. *Journal of Policy Analysis and Management*, 26(4):909-50.
28. Rockoff, Jonah E., Brian A. Jacob, Thomas J. Kane, and Douglas O. Staiger. 2008. Can You Recognize An Effective Teacher When You Recruit One? NBER Working Paper 14485.
29. Springer, M., D. Ballou and Peng. 2008. Impact of the Teacher Advancement Program on Student Test Score Gains: Findings for an Independent Appraisal. National Center on Performance Incentives Working Paper 2008-19.
30. Springer, Matthew, Laura Hamilton, Daniel McCaffrey, Dale Ballou, Vi-Nhuan Le, Matthew Pepper, J.R. Lockwood and Brian Stecher. 2010. Teacher Pay for Perfor-



- mance: Experimental Evidence from the Project on Incentives in Teaching. National Center on Performance Incentives.
31. Taylor, Eric and John H. Tyler. 2011. The Effect of Evaluation on Performance: Evidence from Longitudinal Student Achievement Data of Mid-Career Teachers. National Bureau of Economic Research Working paper 16877.
  32. Taylor, Lori and Matthew Springer. 2008. Optimal Incentives for Public Sector Workers: The Case of Teacher-Designed Incentive Pay in Texas. National Center on Performance Incentives Working Paper 2009-05.
  33. Tyler, John, Eric Taylor, Thomas Kane, and Amy Wooten. 2010. Using Student Performance Data to Identify Effective Classroom Practices. *American Economic Review*, 100(2): 256-60.
  34. Vigdor, Jacob L. 2008. Teacher Salary Bonuses in North Carolina. Conference paper, National Center on Performance Incentives.
  35. Weitzman, Martin L. and Douglas Kruse. 1990. Profit sharing and productivity. In *Paying for Productivity*. Edited by A. Blinder. Brookings Institution: Washington, D.C.

## 8 Tables

Table 6: Descriptive statistics

Variable	Mean	Std. Dev.	Min.	Max.	N
<b>Student: school-grade-year, weighted by enrollment</b>					
Total enrollment	167.8	139.3	1	826	1,826,036
Share male	0.513	0.064	0	1	1,826,036
Share free lunch	0.256	0.202	0	1	1,826,036
Share reduced price	0.083	0.052	0	1	1,826,036
Share special educ.	0.139	0.08	0	1	1,826,036
Share Afr.-American	0.092	0.146	0	1	1,826,036
Share Hispanic	0.063	0.092	0	1	1,826,036
Share Asian-American	0.061	0.099	0	1	1,826,036
Share Native American	0.021	0.075	0	1	1,826,036
<b>Teacher</b>					
% teachers with masters	17.9	9.0	0	40.1	3,142
Mean years of experience	14.7	2.9	1	27.9	3,142
<b>District: district-year</b>					
Inter-district flow	-0.36	498.6	-11,037	2599	3,244
General Reserve Fund/Expend.	12.4	10.7	-54.7	174.0	3,199

Table 7: Program design effects on student achievement *levels* - reading - pooled across grades 3 to 8 and academic years 2005-06 to 2009-10

Sample Specification	Full		Interested Only		Adopters Only	
	(A)	(B)	(A)	(B)	(A)	(B)
Teacher P4P\$	0.087*** (0.025)	0.087*** (0.025)	0.096*** (0.03)	0.097*** (0.03)	0.108*** (0.027)	0.112*** (0.026)
School P4P\$	0.036 (0.08)	0.034 (0.08)	0.037 (0.075)	0.033 (0.074)	0.032 (0.077)	0.024 (0.076)
Evaluation P4P\$	-.051** (0.021)	-.051** (0.023)	-.044 (0.029)	-.045 (0.029)	-.035 (0.031)	-.034 (0.031)
2+ pre-adoption		-.007 (0.046)		-.013 (0.047)		-.024 (0.042)
1(Dropped Q-Comp)	-.030 (0.068)	-.032 (0.069)	-.022 (0.09)	-.022 (0.091)	0.0005 (0.094)	0.005 (0.096)
Enrollment, 1,000s	-.174 (0.222)	-.175 (0.224)	-.242 (0.345)	-.249 (0.353)	-.270 (0.355)	-.288 (0.361)
Share free lunch	-1.211*** (0.118)	-1.211*** (0.118)	-1.304*** (0.156)	-1.305*** (0.156)	-1.366*** (0.17)	-1.367*** (0.169)
Share red. price	-.763*** (0.132)	-.763*** (0.131)	-.693** (0.285)	-.695** (0.283)	-1.052*** (0.325)	-1.055*** (0.323)
Share special Ed.	-1.855*** (0.099)	-1.855*** (0.099)	-1.837*** (0.17)	-1.837*** (0.17)	-1.698*** (0.207)	-1.700*** (0.207)
Share Male	-.484*** (0.064)	-.483*** (0.064)	-.391*** (0.106)	-.390*** (0.106)	-.439*** (0.132)	-.437*** (0.132)
Share Afr.-American	-1.589*** (0.279)	-1.588*** (0.279)	-1.815*** (0.166)	-1.809*** (0.167)	-1.690*** (0.245)	-1.671*** (0.242)
Share Hispanic	-1.311*** (0.188)	-1.311*** (0.188)	-1.129*** (0.29)	-1.124*** (0.285)	-1.191*** (0.315)	-1.181*** (0.312)
Share Asian-American	-.723** (0.291)	-.721** (0.293)	-.460* (0.255)	-.451* (0.266)	-.456* (0.259)	-.432 (0.266)
Share Native American	-.738*** (0.261)	-.738*** (0.261)	-1.267*** (0.38)	-1.267*** (0.38)	-.738* (0.396)	-.741* (0.399)
School-grade FE	Yes	Yes	Yes	Yes	Yes	Yes
Year-grade FE	Yes	Yes	Yes	Yes	Yes	Yes
<i>N</i> districts	471	471	134	134	101	101
<i>N</i> school-years	4677	4677	1785	1785	1335	1335
<i>N</i> tested students	1749818	1749818	755801	755801	607067	607067
Adj. $R^2$	0.886	0.886	0.916	0.916	0.91	0.91

Coefficient (within-district-correlation corrected SE). Significance: \* : 10% \*\* : 5% \*\*\* : 1%.

The single year immediately prior to adoption is always omitted.

Table 8: Program design effects on student achievement *growth* - reading

DV: Reading average achievement for district-grade-year						
Sample Specification	Full		Interested Only		Adopters Only	
	(A)	(B)	(A)	(B)	(A)	(B)
Teacher P4P\$	0.074** (0.036)	0.073** (0.035)	0.083** (0.038)	0.081** (0.037)	0.081** (0.04)	0.081** (0.037)
School P4P\$	-.128 (0.106)	-.121 (0.103)	-.163 (0.108)	-.153 (0.104)	-.126 (0.107)	-.153 (0.104)
Evaluation P4P\$	-.030 (0.033)	-.028 (0.037)	-.029 (0.039)	-.027 (0.041)	-.028 (0.038)	-.027 (0.041)
Lagged reading	0.309*** (0.02)	0.309*** (0.02)	0.281*** (0.037)	0.282*** (0.038)	0.291*** (0.035)	0.282*** (0.038)
Lagged math	0.143*** (0.016)	0.143*** (0.016)	0.161*** (0.031)	0.16*** (0.03)	0.149*** (0.035)	0.16*** (0.03)
2+ pre-adoption		0.02 (0.051)		0.031 (0.046)		0.031 (0.046)
Student observables	Yes	Yes	Yes	Yes	Yes	Yes
District-grade FE	Yes	Yes	Yes	Yes	Yes	Yes
Year-grade FE	Yes	Yes	Yes	Yes	Yes	Yes
<i>N</i> districts	446	446	132	132	98	132
<i>N</i> school-years	1989	1989	584	584	442	584
<i>N</i> students	1339042	1339042	578414	578414	446951	578414
Adjusted $R^2$	0.914	0.914	0.947	0.947	0.932	0.947

Coefficient (within-district SE). Significance: \*: 10% \*\*: 5% \*\*\*: 1%.

Variables are year-district-grade averages.

Lags are prior year, prior grade  $(t - 1)d(g - 1)b$ .

Table 9: Program design effects on student achievement *levels* - math - pooled across grades 3 to 8 and academic years 2005-06 to 2009-10

Sample Specification	Full		Interested Only		Adopters Only	
	(A)	(B)	(A)	(B)	(A)	(B)
Teacher P4P\$	-.031 (0.031)	-.028 (0.031)	-.040 (0.032)	-.033 (0.032)	-.049 (0.031)	-.039 (0.032)
School P4P\$	0.182* (0.107)	0.157 (0.11)	0.195* (0.107)	0.169 (0.11)	0.187* (0.108)	0.166 (0.109)
Evaluation P4P\$	-.005 (0.022)	-.011 (0.02)	-.006 (0.026)	-.009 (0.024)	-.015 (0.025)	-.015 (0.024)
2+ pre-adoption		-.065* (0.039)		-.070* (0.038)		-.061 (0.041)
1(Dropped Q-Comp)	0.046 (0.111)	0.033 (0.111)	0.052 (0.125)	0.052 (0.122)	0.016 (0.115)	0.027 (0.113)
Enrollment, 1,000s	-.952** (0.386)	-.964** (0.385)	-.958* (0.573)	-1.000* (0.571)	-1.060* (0.607)	-1.106* (0.605)
Share free lunch	-1.077*** (0.135)	-1.079*** (0.135)	-1.166*** (0.207)	-1.168*** (0.208)	-1.262*** (0.249)	-1.263*** (0.249)
Share red. price	-.547*** (0.139)	-.549*** (0.139)	-.675** (0.317)	-.683** (0.316)	-1.039*** (0.365)	-1.044*** (0.367)
Share special Ed.	-1.907*** (0.122)	-1.909*** (0.122)	-2.041*** (0.222)	-2.046*** (0.221)	-1.840*** (0.235)	-1.848*** (0.233)
Share Male	-.008 (0.078)	-.007 (0.078)	0.146 (0.146)	0.15 (0.145)	0.118 (0.175)	0.122 (0.174)
Share Afr.-American	-1.653*** (0.346)	-1.643*** (0.347)	-2.051*** (0.206)	-2.021*** (0.212)	-1.940*** (0.299)	-1.892*** (0.308)
Share Hispanic	-.892*** (0.17)	-.887*** (0.17)	-.715*** (0.248)	-.686*** (0.251)	-.556* (0.285)	-.530* (0.288)
Share Asian-American	0.16 (0.226)	0.179 (0.228)	0.265 (0.303)	0.313 (0.304)	0.16 (0.285)	0.224 (0.284)
Share Native American	-.691*** (0.263)	-.690*** (0.263)	-1.198*** (0.399)	-1.198*** (0.397)	-.613 (0.383)	-.625 (0.386)
School-grade FE	Yes	Yes	Yes	Yes	Yes	Yes
Year-grade FE	Yes	Yes	Yes	Yes	Yes	Yes
<i>N</i> districts	469	469	134	134	101	101
<i>N</i> school-years	4666	4666	1779	1779	1329	1329
<i>N</i> tested students	1698331	1698331	729520	729520	586667	586667
Adj. $R^2$	0.86	0.86	0.89	0.89	0.884	0.884

Coefficient (within-district-correlation corrected SE). Significance: \* : 10% \*\* : 5% \*\*\* : 1%.

The single year immediately prior to adoption is always omitted.

Table 10: Program design effects on student achievement *growth* - math

DV: Math average achievement for district-grade-year						
Sample Specification	Full		Interested Only		Adopters Only	
	(A)	(B)	(A)	(B)	(A)	(B)
Teacher P4P\$	0.019 (0.039)	0.02 (0.04)	0.004 (0.038)	0.007 (0.039)	-.006 (0.039)	0.007 (0.039)
School P4P\$	0.042 (0.125)	0.028 (0.129)	0.059 (0.122)	0.046 (0.126)	0.066 (0.123)	0.046 (0.126)
Evaluation P4P\$	-.032* (0.016)	-.036** (0.016)	-.041** (0.02)	-.044** (0.019)	-.049** (0.02)	-.044** (0.019)
Lagged reading	0.164*** (0.017)	0.164*** (0.017)	0.167*** (0.035)	0.166*** (0.035)	0.163*** (0.042)	0.166*** (0.035)
Lagged math	0.344*** (0.016)	0.345*** (0.016)	0.36*** (0.03)	0.361*** (0.03)	0.368*** (0.035)	0.361*** (0.03)
2+ pre-adoption		-.041 (0.032)		-.040 (0.034)		-.040 (0.034)
Student observables	Yes	Yes	Yes	Yes	Yes	Yes
District-grade FE	Yes	Yes	Yes	Yes	Yes	Yes
Year-grade FE	Yes	Yes	Yes	Yes	Yes	Yes
<i>N</i> districts	445	445	132	132	98	132
<i>N</i> school-years	1985	1985	584	584	442	584
<i>N</i> students	1295202	1295202	556746	556746	433988	556746
Adjusted $R^2$ 0.9	0.9	0.936	0.936	0.924	0.936	

Coefficient (within-district SE). Significance: \*: 10% \*\*: 5% \*\*\*: 1%.

Variables are year-district-grade averages.

Lags are prior year, prior grade  $(t - 1)d(g - 1)b$ .

Table 11: Program design effects on alternative outcomes using academic years from 2003-2004 to 2009-2010 and all available grades (K-12)

	Student		Teacher		
	Log (Enrolmt.)	Inter-dist. movement	Log (mean pay)	Mean yrs. exper.	% M.A.
Teacher P4P\$	-.0009 (0.032)	-30.62 (45.83)	0.029*** (0.010)	-.123 (0.224)	0.48 (0.655)
School P4P\$	0.149 (0.094)	133.70* (69.62)	-.036 (0.039)	-.944 (0.853)	-1.184 (2.361)
Evaluation P4P\$	-.060*** (0.021)	-15.59 (21.49)	0.027** (0.012)	0.18 (0.141)	0.227 (0.359)
2+ pre-adoption	-.096*** (0.028)	-51.99*** (14.83)	0.019* (0.011)	0.193 (0.178)	-.523 (0.388)
<i>N</i> districts	558	516	498	500	500
Weighted by		Years	Years	Years-FTE	
<i>N</i>	3974	3244	3120	357307	357307
Adj R <sup>2</sup>	0.986	0.934	0.804	0.887	0.947

Significance: \*: 10% \*\*: 5% \*\*\*: 1%.

Coefficient (within-district SE). Year effects and district effects included. All use district-level variables, except enrollment (district-grade).

Table 12: Effects of design features on achievement allowing for differential effects by whether School P4P\$ are tied to student achievement exclusively in a single subject - pooled across grades 3 to 8 and academic years 2005-06 to 2009-10

Dep. Variable:	Reading	Math
Teacher P4P \$	0.077*** (0.029)	-.029 (0.033)
1(Goal for this subject only) · Teacher P4P\$	0.026 (0.063)	-.046 (0.073)
School P4P\$	0.01 (0.08)	0.158 (0.112)
1(Goal for this subject only) · School P4P\$	0.101 (0.248)	0.326 (0.506)
Evaluation P4P\$	-.041 (0.025)	-.018 (0.022)
1(Goal for this subject only) · Evaluation P4P\$	-.047 (0.031)	0.074 (0.084)
2+ pre-adoption	-.003 (0.047)	-.066* (0.04)
<i>N</i> districts	471	469
<i>N</i> school-grades	4,677	4,666
<i>N</i> tested students	1749818	1698331
Adj. $R^2$	0.887	0.86

Coefficient (within-district SE). Significance: \* : 10% \*\* : 5% \*\*\* : 1%.

Each column is a separate regression of specification B in full sample as in second column of Tables 7 and 9.



Table 13: Robustness of levels model to alternative conditioning sets - pooled across grades 3 to 8 and academic years 2005-06 to 2009-10

	(1)	(2)	(3)	(4)
<b>Reading</b>				
Teacher P4P\$	0.087*** (0.03)	0.087*** (0.025)	0.087*** (0.024)	0.092*** (0.027)
School P4P\$	0.03 (0.074)	0.034 (0.08)	0.035 (0.078)	0.001 (0.104)
Evaluation P4P\$	-.067*** (0.023)	-.051** (0.023)	-.051** (0.023)	-.059** (0.028)
2+ pre-adoption	-.015 (0.048)	-.007 (0.046)	-.006 (0.046)	0.003 (0.055)
<i>N</i> districts	471	471	471	436
<i>N</i> district grades	4677	4677	4670	4439
<i>N</i> tested students	1749818	1749818	1749080	1384099
Adj. R <sup>2</sup>	0.873	0.886	0.886	0.893
<b>Math</b>				
Teacher P4P\$	-.025 (0.038)	-.028 (0.031)	-.028 (0.031)	-.013 (0.032)
School P4P\$	0.131 (0.106)	0.157 (0.11)	0.159 (0.109)	0.122 (0.119)
Evaluation P4P\$	-.017 (0.02)	-.011 (0.02)	-.011 (0.02)	-.002 (0.02)
2+ pre-adoption	-.063 (0.041)	-.065* (0.039)	-.062 (0.04)	-.056 (0.046)
<i>N</i> districts	469	469	469	434
<i>N</i> district grades	4666	4666	4659	4420
<i>N</i> tested students	1698331	1698331	1697597	1347064
Adj. R <sup>2</sup>	0.848	0.86	0.859	0.873
Student observables	No	Yes	Yes	Yes
Teacher observables	No	No	Yes	Yes
District observable	No	No	No	Yes
District-grade FE	Yes	Yes	Yes	Yes
Year-grade FE	Yes	Yes	Yes	Yes

Coefficient (within-district SE). Significance: \*: 10% \*\*: 5% \*\*\*: 1%.

Reading (math) analogous to column 2 of Table 7 (9), except for changes in covariate sets.

Table 14: Robustness of growth model to alternative conditioning sets - pooled across grades 3 to 8 and academic years 2005-06 to 2009-10

	(1)	(2)	(3)	(4)
<b>Reading</b>				
Teacher P4P\$	0.073** (0.036)	0.073** (0.035)	0.073** (0.035)	0.098** (0.042)
School P4P\$	-.092 (0.093)	-.121 (0.103)	-.114 (0.102)	-.243* (0.132)
Evaluation P4P\$	-.032 (0.032)	-.028 (0.037)	-.030 (0.037)	-.033 (0.037)
2+ pre-adoption	0.013 (0.054)	0.02 (0.051)	0.021 (0.051)	0.03 (0.07)
<i>N</i> districts	446	446	446	415
<i>N</i> district grades	1989	1989	1987	1890
<i>N</i> tested students	1339042	1339042	1338696	1038698
Adj. R <sup>2</sup>	0.91	0.914	0.914	0.92
<b>Math</b>				
Teacher P4P\$	0.021 (0.04)	0.02 (0.04)	0.02 (0.04)	0.07** (0.031)
School P4P\$	0.024 (0.125)	0.028 (0.129)	0.027 (0.13)	-.113 (0.115)
Evaluation P4P\$	-.037** (0.016)	-.036** (0.016)	-.037** (0.016)	-.037** (0.016)
2+ pre-adoption	-.043 (0.03)	-.041 (0.032)	-.042 (0.031)	-.035 (0.032)
<i>N</i> districts	445	445	445	415
<i>N</i> district grades	1985	1985	1983	1886
<i>N</i> tested students	1295202	1295202	1294863	1005407
Adj. R <sup>2</sup>	0.899	0.9	0.9	0.911
Student observables	No	Yes	Yes	Yes
Teacher observables	No	No	Yes	Yes
District observable	No	No	No	Yes
District-grade FE	Yes	Yes	Yes	Yes
Year-grade FE	Yes	Yes	Yes	Yes

Coefficient (within-district SE). Significance: \*: 10% \*\*: 5% \*\*\*: 1%.

Reading (math) analogous to column 2 of Table 8 (10), except for changes in covariate sets.

Table 15: Robustness of levels model to dropping any Q-Comp adoption cohort - pooled across grades 3 to 8 and academic years 2005-06 to 2009-10

	Adoption cohort excluded from analysis:					
	2005	2006	2007	2008	2009	2010
<b>Reading</b>						
Teacher P4P\$	0.108*** (0.03)	0.054 (0.046)	0.08*** (0.024)	0.088*** (0.026)	0.087*** (0.025)	0.087*** (0.025)
School P4P\$	-.056 (0.09)	0.085 (0.083)	0.12 (0.078)	-.005 (0.088)	0.035 (0.08)	0.039 (0.08)
Evaluation P4P\$	-.053** (0.023)	-.038 (0.031)	-.074*** (0.024)	-.041* (0.023)	-.052** (0.023)	-.049** (0.024)
2+ pre-adoption	-.010 (0.047)	0.0003 (0.056)	-.042 (0.046)	0.002 (0.052)	-.006 (0.047)	0.007 (0.054)
<i>N</i> districts	464	432	462	461	465	443
<i>N</i> district grades	4509	4041	4475	4591	4637	4474
<i>N</i> tested students	1680075	1428053	1638260	1700066	1743176	1702211
Adj. R <sup>2</sup>	0.886	0.881	0.883	0.884	0.886	0.887
<b>Math</b>						
Teacher P4P\$	0.002 (0.032)	-.080 (0.063)	-.037 (0.03)	-.026 (0.032)	-.029 (0.031)	-.028 (0.031)
School P4P\$	0.048 (0.099)	0.217* (0.131)	0.259* (0.139)	0.137 (0.12)	0.16 (0.111)	0.16 (0.111)
Evaluation P4P\$	-.007 (0.02)	-.010 (0.021)	-.029 (0.025)	-.003 (0.021)	-.011 (0.02)	-.010 (0.02)
2+ pre-adoption	-.071* (0.04)	-.067* (0.038)	-.068 (0.05)	-.064 (0.046)	-.065 (0.041)	-.061 (0.044)
<i>N</i> districts	462	430	460	459	463	441
<i>N</i> district grades	4498	4034	4466	4580	4626	4463
<i>N</i> tested students	1631582	1386350	1591116	1650008	1691809	1652454
Adj. R <sup>2</sup>	0.86	0.853	0.856	0.857	0.858	0.86

Coefficient (within-district SE). Significance: \*: 10% \*\*: 5% \*\*\*: 1%.

Reading (math) analogous to column 2 of Table 7 (9), except for exclusion of adoption cohorts.

Table 16: Robustness of growth model to dropping any adoption cohort - pooled across grades 3 to 8 and academic years 2005-06 to 2009-10

	Adoption cohort excluded from analysis:					
	2005	2006	2007	2008	2009	2010
	<b>Reading</b>					
Teacher P4P\$	0.121*** (0.043)	0.03 (0.049)	0.068** (0.033)	0.071* (0.036)	0.074** (0.035)	0.072** (0.035)
School P4P\$	-.302** (0.124)	-.096 (0.074)	-.008 (0.108)	-.125 (0.115)	-.118 (0.106)	-.111 (0.102)
Evaluation P4P\$	-.025 (0.037)	0.019 (0.033)	-.071** (0.032)	-.024 (0.04)	-.028 (0.038)	-.026 (0.039)
2+ pre-adoption	0.011 (0.053)	0.04 (0.059)	-.0009 (0.039)	0.003 (0.063)	0.023 (0.052)	0.04 (0.066)
<i>N</i> districts	439	407	438	436	440	419
<i>N</i> district grades	1954	1808	1951	1942	1962	1877
<i>N</i> tested students	1292480	1094541	1257031	1301639	1335331	1306279
Adj. R <sup>2</sup>	0.914	0.907	0.911	0.91	0.914	0.914
	<b>Math</b>					
Teacher P4P\$	0.085** (0.035)	-.055 (0.06)	0.01 (0.039)	0.02 (0.04)	0.019 (0.04)	0.019 (0.04)
School P4P\$	-.194 (0.122)	0.134 (0.112)	0.057 (0.159)	0.072 (0.138)	0.041 (0.13)	0.032 (0.129)
Evaluation P4P\$	-.030* (0.016)	-.016 (0.019)	-.045* (0.023)	-.044** (0.017)	-.036** (0.016)	-.036** (0.017)
2+ pre-adoption	-.054* (0.032)	-.032 (0.036)	-.031 (0.043)	-.063* (0.036)	-.039 (0.032)	-.037 (0.034)
<i>N</i> districts	438	406	437	435	439	418
<i>N</i> district grades	1950	1804	1947	1938	1958	1873
<i>N</i> tested students	1249991	1058062	1215480	1258601	1291574	1263516
Adj. R <sup>2</sup>	0.901	0.89	0.897	0.897	0.9	0.901

Coefficient (within-district SE). Significance: \*: 10% \*\*: 5% \*\*\*: 1%.

Reading (math) analogous to column 2 of Table 8 (10), except for exclusion of adoption cohorts.

Table 17: Design features on achievement with heterogeneous effects by grade-subject

Grade:	3	4	5	6	7	8
<b>Reading</b>						
Teacher P4P\$	0.093** (0.047)	-.020 (0.03)	0.038* (0.023)	0.061 (0.072)	0.165*** (0.039)	0.186*** (0.06)
School P4P\$	0.108 (0.138)	0.217** (0.087)	0.015 (0.102)	-.083 (0.226)	0.113 (0.11)	-.153 (0.205)
Evaluation P4P\$	-.067* (0.039)	-.069*** (0.025)	-.049 (0.039)	-.050 (0.04)	-.113*** (0.04)	0.032 (0.044)
2+ pre-adoption	0.02 (0.075)	0.026 (0.051)	-.042 (0.06)	0.112 (0.1)	-.133* (0.068)	-.038 (0.078)
<b>Math</b>						
Teacher P4P\$	-.035 (0.037)	-.045 (0.049)	-.103*** (0.034)	-.024 (0.073)	-.023 (0.048)	0.066 (0.053)
School P4P\$	0.226* (0.12)	0.2 (0.149)	0.052 (0.139)	-.020 (0.193)	0.288 (0.202)	0.211 (0.224)
Evaluation P4P\$	-.067* (0.039)	-.069*** (0.025)	-.049 (0.039)	-.050 (0.04)	-.113*** (0.04)	0.032 (0.044)
2+ pre-adoption	-.123** (0.055)	-.037 (0.053)	-.034 (0.095)	-.006 (0.067)	-.146** (0.063)	-.053 (0.093)

Coefficient (within-district SE). Significance: \*: 10% \*\*: 5% \*\*\*: 1%.

In each sample, estimates from a single regression with separate effects by grade from specification B in full sample as in second column of Tables 7 and 9.

Table 18: Participation effects on achievement - pooled across grades 3 to 8 and academic years 2005-06 to 2009-10

Sample Specification	Full		Interested Only		Adopters Only	
	(A)	(B)	(A)	(B)	(A)	(B)
<b>Reading</b>						
Post-adoption	0.016 (0.024)	0.001 (0.024)	0.014 (0.033)	0.006 (0.032)	-.004 (0.033)	-.004 (0.035)
2+ yrs. pre-adoption		0.013 (0.054)		0.011 (0.057)		0.004 (0.055)
<i>N</i> districts	471	471	134	134	101	101
<i>N</i> school-grade-years	4677	4677	1785	1785	1335	1335
<i>N</i> tested students	1749818	1749818	755801	755801	607067	607067
Adj. $R^2$	0.886	0.886	0.915	0.915	0.909	0.909
<b>Math</b>						
Post-adoption	0.016 (0.024)	0.001 (0.024)	0.014 (0.033)	0.006 (0.032)	-.004 (0.033)	-.004 (0.035)
2+ yrs. pre-adoption		-.074** (0.038)		-.081** (0.037)		-.073* (0.04)
<i>N</i> districts	469	469	134	134	101	101
<i>N</i> school-grade-years	4666	4666	1779	1779	1329	1329
<i>N</i> tested students	1698331	1698331	729520	729520	586667	586667
Adj. $R^2$	0.859	0.859	0.89	0.89	0.883	0.884

Coefficient (within-district SE). Significance: \*: 10% \*\*: 5% \*\*\*: 1%.

Specification as in Table 7 and Table 9, except 1(Post-adoption)<sub>*t*<sub>sg</sub></sub> substituted for the three P4P\$ variables.

The single year immediately prior to adoption is always omitted.