

Aligning Learning Incentives of Students and Teachers: Results from a Social Experiment in Mexican High Schools

Jere R. Behrman

University of Pennsylvania

Susan W. Parker

Centro de Investigacion y Docencia Economicas

Petra E. Todd

University of Pennsylvania

Kenneth I. Wolpin

University of Pennsylvania and Rice University

This paper evaluates the impact of three different performance incentive schemes using data from a social experiment that randomized 88 Mexican high schools with over 40,000 students into three treatment groups and a control group. Treatment 1 provides individual incentives for performance on curriculum-based mathematics tests to students only, treatment 2 to teachers only, and treatment 3 gives both individual and group incentives to students, teachers, and school administrators. Program impact estimates reveal the largest average effects for treatment 3, smaller impacts for treatment 1, and no impact for treatment 2.

We worked on this study with the support of the Inter-American Development Bank and the support and collaboration of the Mexican Ministry of Education. We thank these organizations for making this study possible. We are especially indebted to Miguel Szekely, Deputy Minister of Education when the Aligning Learning Incentives program was initi-

[*Journal of Political Economy*, 2015, vol. 123, no. 2]

© 2015 by The University of Chicago. All rights reserved. 0022-3808/2015/12302-0002\$10.00

I. Introduction

This paper evaluates the impact of a large-scale social experiment (the Aligning Learning Incentives, or ALI, program) designed to promote mathematics achievement through performance-based monetary incentives. Eighty-eight Mexican federal high schools with over 40,000 students were randomly allocated to three treatment groups and a control group. Treatment 1 (T1) provides individual incentives to students only and treatment 2 (T2) to teachers only. Treatment 3 (T3) offers both individual and group incentives to students, teachers, and school administrators, thus rewarding cooperation among all the educational actors in the school. The ALI program began in the 2008–9 academic year and ran for 3 years, ending in the 2010–11 academic year. Incentives were determined by student performance on curriculum-based mathematics tests administered to students in grades 10–12 at the end of each year.

An important rationale for utilizing monetary incentives rather than intervening in the educational process directly is that policy makers are not likely to know the best means of improving education given the heterogeneous attributes of students, teachers, and administrators across schools. The production process by which school, student, and family inputs are transformed into educational outcomes is not well understood.¹ Providing monetary incentives tied to student performance allows students, teachers, and principals to choose the best means to improve performance given their circumstances. The ALI program was designed to evaluate the efficacy of alternative performance-based incentive schemes.

Research on the impact and efficacy of performance-based monetary incentives in education is relatively sparse. We review some examples from the literature of studies of teacher and student incentives in the next section, restricting attention to randomized control trials in which the measured outcomes are test scores, as in the ALI program. As far as we know, the ALI program is the first randomized control trial to incorporate incentive payments to both students and teachers.² Previous studies implement performance incentives for students or teachers only.

ated; to Marcelo Cabrol, Santiago Levy, and Ana Santiago at the Inter-American Development Bank; to Rafael De Hoyos, Martha Hernandez, Lucia Juarez, Elizabeth Monroy, Araceli Ortega, and Diana Zamora, who were responsible for the day-to-day operations and management of the project; and to George Wesolowsky, who performed the analysis of detecting nonindependent test taking.

¹ See, e.g., Hanushek (1986, 2003), Hedges, Laine, and Greenwald (1994), and Krueger (2003) for differing views about the interpretation of findings from the education production function literature. Todd and Wolpin (2003) discuss the methodological underpinnings of that literature.

² There was a recent nonexperimental program in Dallas, Texas, that paid both students and teachers for passing grades on advanced-placement exams (Jackson 2010).

A comparison of the ALI program impact estimates to those of prior studies reveals that the treatment effects associated with the ALI treatments in which students receive incentives are quite large, especially for the treatment in which both students and teachers receive incentives. However, close examination of the test book answer patterns shows that part of the reason for higher test scores in the treatment group is a higher rate of cheating, in the form of student copying, than in the control group, particularly in higher grades and in later years of the program.³ Our impact analysis therefore provides two sets of estimates of treatment effects: one that does not account for copying and one that adjusts the test scores of students identified to have been likely to have copied a part of their test answers, as determined through a comparison of multiple-choice answers of pairs of students and a statistical model for the probability of having matching answers. Even with a liberal criterion for identifying copiers and with two different ways of adjusting their scores, we find substantial program effects on student test scores, indicating that the performance incentive program significantly affected mathematics achievement.

To highlight the magnitudes of the treatment effects, consider the copying-adjusted treatment effect estimates for the 2008–9 entering tenth-grade class. The average treatment effects for the first year they were in the program, as tenth graders, were 0.17 of a standard deviation for treatment group T1, 0.01 of a standard deviation for T2, and 0.31 for T3. In the second year of the program, as eleventh graders, the treatment effects were 0.30 for T1, 0.02 for T2, and 0.44 for T3. Finally, when this cohort reached the twelfth grade, treatment effects were 0.23, 0.04, and 0.57. Treatment effects were statistically significant at conventional levels in all three years for T1 and T3 but in none of the years for T2. The pattern of positive effects for T1, even larger effects for T3, and no effect for T2 (with one exception) was also found for all the other tenth-grade entry cohorts during the years they were in the program.

A potentially important caveat to our findings is that we cannot ensure that the control students (and to a lesser extent the T2 students), for whom the ALI test is low or no stakes, applied the same level of test-taking effort as the T1 and T3 students. We develop two sets of lower bounds based on alternative assumptions that take into account the possibility of differential effort between the treatment and the control students. In the more conservative and, in our view, less plausible case, in which the entire differential between T1 and control students is attributed to differential test-taking effort in all years, the T3 effects (adjusted for copying) in year 3 were 0.31, 0.17, and 0.34 for the three grades. In the second case, where only the year 1 differential between T1 and control

³ We do not find evidence of teacher cheating.

students is attributed to test-taking effort, the year 3 effects in the three grades are 0.15, 0.12, and 0.13 for T1 and 0.47, 0.29, and 0.47 for T3.

James Heckman has a number of papers illustrating both the value and the limitations of randomized social experiments for policy evaluation purposes (see, e.g., Heckman and Smith 1995; Heckman 2000; Heckman and Urzua 2010). As he shows, randomization provides information only on mean treatment effects for a particular program design without additional assumptions. A central theme of Heckman's research is that economic models can be used to understand the mechanisms through which a program operates and to investigate effects of program designs that differ in some ways from the program that was implemented. To this end, in Section IV, we develop an economic model that can rationalize the results from the ALI experiment.

In the next section, we discuss relevant literature. In Section III, we provide details of the ALI experiment, including the overall design, the selection of the sample, a description of the tests, and the incentive schedules. In Section IV, we present the results of the experiment and briefly outline a model to interpret the results. Section V presents conclusions.

II. Related Literature

As noted, previous experimental studies of performance incentives provide incentives only to teachers or to students.

Teacher incentives.—Muralidharan and Sundararaman (2011) evaluate the effects of teacher performance incentives and school input interventions in rural India. They randomly allocated schools to four treatment groups and to a control group, with 100 schools in each group. One of the treatments was a performance incentive paid to teachers on the basis of the average improvement in their students' test scores. Tests were administered in mathematics and language at the beginning and end of the school year. The impact of the performance incentive program was a 0.28 (0.17) standard deviation increase in the average test score in mathematics (language).

Glewwe, Ilias, and Kremer (2010) conducted a randomized trial over a 2-year period that provided primary school (grades 4–8) teachers in 50 rural Kenyan schools incentives based on student performance on district-level exams in seven subjects. Students in treatment schools had higher test scores (averaged over all subjects) in the second year of about 0.14 of a standard deviation (0.07 standard error [SE]), but the gains dissipated in the third year after the program ended.

Springer et al. (2010) report on a recent 3-year teacher-incentive experiment in Nashville, Tennessee, public schools. Middle school mathematics teachers, who volunteered for the program, were randomly as-

signed to treatment and control groups. Teachers in the treatment group could earn bonuses, depending on the standardized test scores achieved by their students, of \$5,000, \$10,000, or \$15,000. The magnitudes of the treatment effect estimates in mathematics in grades 6–8 were 0.06 standard deviations or less and were not statistically significant at conventional levels. In grade 5, the average treatment effect was 0.06 of a standard deviation (0.04 SE) in year 1, 0.18 (0.06 SE) in year 2, and 0.20 (0.08 SE) in year 3.

Student incentives.—Angrist and Lavy (2009) study the effects of a student cash incentive program in Israel that offered high school students incentives for progressing from tenth to eleventh grade, for progressing from eleventh to twelfth grade, and for passing the Bagrut exam.⁴ Forty high schools were randomized in or out of the program. The school-based program increased Bagrut passing rates by 6–8 percentage points.

Kremer, Miguel, and Thornton (2009) evaluate the impact of a merit scholarship program for sixth-grade girls in Kenya. Schools were randomized to treatment and control groups, and scholarships were awarded to the top 15 percent of sixth-grade girls in treatment schools on the basis of standardized achievement test scores. On average, girls in participating schools raised their achievement by 0.2 standard deviations. They also estimate impacts by quartile of achievement and find the largest impacts in the second quartile, a group that would seem to have relatively low probabilities of winning an award. The conclusion that learning improved even for those unlikely to win an award is tempered, however, by the finding that the lowest impact is for the lowest quartile.

Fryer (2010) reports results from four different field experiments that implemented various student incentive programs, in Chicago, Dallas, New York City, and Washington, DC (in predominantly low-performing urban public schools), with substantial heterogeneity in terms of grade levels participating and incentive design. In New York City, payments were given to fourth- and seventh-grade students conditional on their performance on 10 standardized tests. The program did not yield positive impacts on final-year test scores. In Chicago, incentives were paid to ninth graders every 5 weeks for grades in five courses. The program led to higher grades but had no detectable effect on test scores. The Dallas program gave second graders \$2 per book read, with an additional requirement of passing a short quiz on the book, which led to significantly better reading test scores and also to higher grades. Finally, in Washington, incentives were given to sixth, seventh, and eighth graders on a composite index intended to capture their school attendance, behavior, and measures of inputs in educational production. The impact estimates are suggestive of

⁴ The Bagrut exam is a prerequisite for admission to university and for certain types of jobs.

substantial positive impacts but are not statistically significant. Fryer (2010) concludes from these results that children in these schools do not know what behaviors would lead to improved test score performance and are thus better served by incentives tied to inputs rather than to outcomes.

Levitt, List, and Sadoff (2010) implement a field experiment to evaluate the effects of a program that gave monthly financial incentives to ninth-grade students or their parents in two high schools in a suburb of Chicago. The incentives were given for meeting an achievement standard that is a composite of multiple measures of performance, including school attendance, behavior, grades, and test scores. The experimental design varied the award recipient (students or parents) and whether the incentive was awarded as a piece rate or as a lottery. The incentive amounts had an expected value of \$50 per month, for a total of \$400 per year. The four treatments were randomized at the student level. The study finds modest overall effects of the incentives on achieving the composite achievement standard—a 15–22 percent increase—but that the effect is not statistically significantly different across treatment types. The study does not find effects of the treatments on standardized test scores. Barrow and Rouse (2013) evaluate the effect of two performance-based scholarship programs: one for students in their last year of high school and one for postsecondary students. Incentives varied in length and magnitude. The paper finds that students eligible for incentive payments devoted more time to educational activities and less time to work and leisure activities.

It is difficult to generalize from these studies (or from a larger set that includes additional nonexperimental studies). First, there are not many evaluation studies of performance-based incentive programs. Second, the existing studies differ in their designs and in the populations studied. And, third, in contrast to the ALI program, the agent receiving the incentive is not varied within the same population.⁵ However, the literature generally finds small measured effects of performance-based incentive programs, regardless of whether the student or teacher receives the incentive; an average treatment effect above 0.3 standard deviations in test scores appears to be unusual.

III. The ALI Experiment

The ALI experiment began with the 2008–9 academic year and ended with the 2010–11 academic year. There were a total of 88 high schools in the experiment, consisting of three groups of 20 treatment schools, each

⁵ Levitt et al.'s (2010) study is an exception, as noted, in which there is randomized variation in whether the student or his or her parents received the incentive payment.

subject to a different incentive design, and a fourth group of 28 control schools with no incentives. In the first program year, there were approximately 12,000 students per treatment and 16,800 students in the control schools. All students in each high school, grades 10–12, participated in the program in each of the three years.

Incentive payments were based on standardized curriculum-based mathematics examinations in grades 10, 11, and 12 given at the end of each academic year. Specifically, the four groups were as follows:

1. Treatment group 1 (T1): Payments to students were based on their own performance.
2. Treatment group 2 (T2): Payments to mathematics teachers were based on the performance of the students in their classes.
3. Treatment group 3 (T3): Payments to students were based on their own performance and on the performance of the other students in their class. Payments to mathematics teachers were based on the performance of the students in their classes and on the performance of the students in all other mathematics classes. Payments to nonmathematics teachers and school administrators were based on the performance of all the students in the school.
4. Control group (C): There were no payments.

A. *Sample Selection and Characteristics*

There are 706 federal upper-secondary schools in Mexico.⁶ From that set, we identified 357 schools that were not in their first year of operation and had only one session per day, in the morning.⁷ There are four types of federal upper-secondary schools in Mexico: academically oriented schools, technically oriented schools with a marine focus, technically oriented schools with an agricultural focus, and technically oriented schools with an industrial focus. The technical orientation of the latter three school types has diminished over time, and all the schools are now considered college preparatory. Of the 357 “morning-only” schools, 14

⁶ Upper-secondary schools (high schools) in Mexico encompass grades 10–12 and lower-secondary schools grades 7–9.

⁷ The selection of schools was based on data supplied by the Ministry of Education. There were seven schools in their first year of operation, four schools for which the number of sessions could not be determined, and 262 schools with multiple sessions. We dropped multiple-session schools because each session essentially constitutes an autonomous school, having a separate principal (when this is accounted for, the number of schools is thus about 1,000). Clearly, it would have been problematic to have only one session in a given multi-session school as part of the program or to have different sessions with different treatments. Moreover, given the likely similarities in the student bodies and some overlap in teachers, having both sessions within the same treatment group would reduce the effective number of schools.

fall into the first category, 26 into the second, 183 into the third, and 134 into the fourth. Academically oriented schools have a different mathematics curriculum and so were not included in the ALI program. The schools with a marine specialization also were not included because, after further selection criteria described below, there were too few to randomize across the treatment and control groups.

In addition to federally administered upper-secondary schools, there are also state-administered schools that are publicly funded and private schools.⁸ Students successfully completing lower-secondary school (ninth grade) may apply to any of these schools. Admissions to public (federal and state) and private schools are determined on a competitive basis. To minimize the impact of the ALI program on students' application decisions, and thus on the composition of entering students in years 2 and 3 of the program across the treatment and control groups, schools that were located within 10 miles of another federal upper-secondary school were eliminated.⁹ Very small schools (fewer than 200 students) and very large schools (over 2,000) were also eliminated, as were schools that had satellite divisions. Finally, schools located in the states of Oaxaca and Michoacan were eliminated because of feasibility constraints. With these restrictions, 135 schools remained out of which 88 were chosen for the experiment.

Randomization was performed using a school-based block randomization design. Schools were first grouped into nine blocks, where the block definitions were based on school size and the previous year's graduation rate.¹⁰ Within each block, schools were allocated at random to treatment regimes.¹¹ The block definitions (cutoffs on school size and graduation rates) were chosen to have roughly similar numbers of schools within each block.

⁸ In 2008, about 25 percent of upper-secondary school students attended federal public schools, 42 percent attended state public schools, and 33 percent attended private and other schools.

⁹ Most of the schools are located in rural areas because of the distance criterion. We do not know if there are state schools, autonomous schools, or private schools located closer than 10 miles to the ALI schools.

¹⁰ Blocking is a widely used method for improving the precision of estimated treatment effects by increasing the comparability of the variables used to define the blocks across treatment/control groups. As described in Cox and Reid (2000), the rationale for blocking is to improve precision by using prior knowledge on which baseline characteristics of the units being randomized are likely to be associated with the treatment responses. For maximal efficiency, units should be grouped into blocks so that all units within a block might be expected to give similar responses in the absence of treatment differences.

¹¹ Following the recommended procedure of Cox and Reid (2000), we randomized six times and chose the randomization in which the groups (T1, T2, T3, and C) are most comparable in terms of baseline observed characteristics. Cox and Reid discuss that if there is any imbalance in observed characteristics that may be related to treatment response, it is better to rerandomize to achieve a better balance in the covariates than to do ex post regression adjustment to adjust for imbalance, which would entail a loss in degrees of freedom.

Tables 1 and 2 present evidence on the quality of the randomization among the 88 schools. Table 1 compares the treatment and control schools and the federal schools not in the experiment on aggregate school-level data supplied by the Ministry of Education. The first two variables, the student population and the graduation rate, as noted, were used for blocking. The other variables in the table were used as additional evidence on the quality of the randomization. They included the following baseline characteristics: percentage of Oportunidades recipients within the school, mean class size, percentages of teachers with university degrees, percentages of new principals, mean distance to the nearest federal upper-secondary school, school type distribution (DGETI or DGETA),

TABLE 1
COMPARISON OF TREATMENT, CONTROL, AND OTHER FEDERAL NON-ALI SCHOOLS:
2007–8 ACADEMIC YEAR

	C	T1	T2	T3	NON-ALI	
	(1)	(2)	(3)	(4)	(5)	(6)
Number of schools	28	20	20	20	269	408
Blocking variables:						
Mean number of students	582 (.77)	632 (.55)	609 (.72)	550 (.63)	656 (.12)	691 (.02)
Mean graduation rate (%)	58.3 (.74)	60.4 (.53)	56.2 (.61)	57.9 (.94)	55.3 (.24)	54.2 (.15)
Nonblocking variables:						
Oportunidades (%)	40.3 (.99)	39.5 (.90)	40.6 (.97)	40.1 (.97)	37.6 (.42)	18.9 (.00)
Mean class size	35.8 (.42)	41.0 (.15)	39.0 (.41)	35.7 (.97)	34.7 (.56)	42.1 (.00)
Teachers with university degree (%)	82.3 (.99)	79.4 (.46)	81.7 (.87)	84.8 (.94)	81.5 (.74)	81.0 (.60)
New directors (%)	25.0 (.72)	25.0 (1.00)	30.0 (.71)	40.0 (.29)	29.4 (.62)	25.2 (.98)
Mean distance to nearest federal upper-secondary school (km)	32.9 (.99)	32.8 (.97)	31.4 (.81)	32.4 (.91)	23.9 (.00)	15.9 (.00)
DGETI (%)	46.4 (.92)	50.0 (.81)	55.0 (.57)	45.0 (.92)	33.8 (.20)	80.1 (.00)
Region 1 (%)	35.7	35.0	50.0	50.0	30.8	36.0
Region 2 (%)	39.3	45.0	40.0	35.0	47.0	50.1
Region 3 (%)	17.9	10.0	5.0	10.0	15.4	10.4
Region 4 (%)	7.1 (.94)	10.0 (.88)	5.0 (.58)	5.0 (.76)	6.8 (.89)	3.5 (.43)

NOTE.—Numbers in parentheses are as follows: col. 1: p -value for test $C = T1 = T2 = T3$; col. 2: p -value for test $C = T1$; col. 3: p -value for test $C = T2$; col. 4: p -value for test $C = T3$; cols. 5 and 6: p -value for test $C = \text{non-ALI schools}$. For col. 5, like the ALI schools, these are schools that have one session per day (morning); for col. 6, these are schools that have two sessions per day (morning and afternoon). The figures pertain only to the morning session because data for afternoon sessions were often missing.

TABLE 2
NINTH-GRADE ENLACE: TREATMENT AND CONTROL SCHOOLS AT BASELINE

Variable	C (1)	T1 (2)	T2 (3)	T3 (4)
9th-grade ENLACE mean test score in mathematics—fall term enrollees:*				
10th-grade class	515.9 (.86)	519.6 (.81)	512.6 (.68)	522.6 (.57)
11th-grade class	516.0 (.91)	516.6 (.96)	517.4 (.86)	524.7 (.47)
% with ENLACE score:				
10th-grade class	90.6 (.30)	88.7 (.23)	88.8 (.44)	86.8 (.08)
11th-grade class	78.3 (.62)	74.0 (.25)	75.2 (.37)	75.3 (.39)

NOTE.—Numbers in parentheses are as follows: col. 1: p -value for test $C = T1 = T2 = T3$; col. 2: p -value for test $C = T1$; col. 3: p -value for test $C = T2$; col. 4: p -value for test $C = T3$. All account for school-level clustering.

* The national mean is 500 and the standard deviation is 100.

and regional distribution.¹² All the variables in table 1 do not differ from the control group at conventional levels of statistical significance, an indication that the randomization procedure was successful. Table 2 compares treatment and control schools on the basis of the mean score of students on the ninth-grade mathematics ENLACE.¹³ This variable was unavailable at the time of the randomization. As seen in table 2, the mean score on the ninth-grade ENLACE does not differ between treatment and controls at conventional levels of significance, and the largest difference from the mean score of the control group is less than 7 (9) standardized points (0.07 [0.09] of a standard deviation) for the tenth- (eleventh-) grade class.

B. ALI Mathematics Tests

Incentive payments (see below) for all treatment regimes depended on performance on standardized mathematics tests administered at the end of the school year. The tests were designed by CENEVAL on the basis of the input of Mexican experts on upper-secondary school mathematics.¹⁴

¹² Oportunidades is a conditional cash transfer program that provides payments for school attendance to low-income families. DGETI (Dirección General de Educación Tecnológica Industrial) schools have an industrial focus and DGETA (Dirección General de Educación Tecnológica Agropecuaria) schools have an agricultural focus.

¹³ The ENLACE is a national test with separate mathematics and language components that began in 2007 with only two grades and is now administered each year to students in all grades between grade 3 and grade 9 and in grade 12.

¹⁴ CENEVAL is a nongovernmental organization similar to the Educational Testing Service in the United States.

The tests were based on the curriculum in each grade.¹⁵ The tests were administered by the Ministry of Education with procedures designed to minimize possible testing abuses. In particular, the tests were not administered or monitored by school personnel, but by representatives of the Ministry of Education state offices, with one monitor assigned to each class and an overall supervisor assigned to the school. The same administrators collected the answer sheets and were required to account for all copies of the tests after test administration to reduce the possibility of teaching to the test based on past tests.¹⁶

For the purpose of determining incentive payments, performance on each test was categorized, as in the ninth-grade ENLACE, into four levels: Pre-Basic, Basic, Proficient, and Advanced. A popular method for determining cutoffs for the categories is the bookmark method, which is used for the ENLACE (see Cizek, Bunch, and Koons 2004). Using that method with the ALI test scores, in the first program year the percentage of students scoring in the top two categories was zero for the tenth grade, 4.0 percent for the eleventh grade, and zero for the twelfth grade for the treatment and control groups combined. The corresponding percentages for the Pre-Basic category were 76.5, 92.3, and 92.8. This performance reflects the fact that the test design faithfully adhered to the curriculum content, which is quite advanced, especially in light of the low level of pre-high school mathematics skills. Using the bookmark cutoffs would have resulted in few students or teachers receiving incentive payments and almost none receiving the larger payments associated with performance in the top two categories. This result would likely have had a deleterious impact on student and teacher effort in subsequent years of the program. For this reason, the bookmark procedure was not used to establish cut scores for determining incentive payments. Instead, the cut scores for the ALI tests were chosen to mimic the ninth-grade ENLACE distribution of the control schools for the tenth and eleventh grades and to mimic the twelfth-grade ENLACE distribution of the control schools for the twelfth grade.¹⁷

¹⁵ The standardized curriculum for each grade is as follows: grade 10: algebra, geometry, and trigonometry (class hours, 4 hours per week); grade 11: analytical geometry and differential calculus (class hours, 4 hours per week); grade 12: probability and statistics and applied statistics (class hours, 5 hours per week). In 2010–11, applied statistics was replaced with integral calculus.

¹⁶ Barlevy and Neal (2011) develop a model of teacher effort choice under an incentive scheme based on the ordinal ranking of students. They show that such a scheme would not require the equating of assessments over time and would thus eliminate the incentive for teachers to “teach to the test.” Their framework does not incorporate student effort as a joint determinant of achievement.

¹⁷ The twelfth-grade mathematics ENLACE was administered for the first time in the 2007–8 academic year.

C. Structure of Incentive Payments

1. Treatment 1 (Student Incentives Only)

Table 3 shows the incentive payment schedule for students at each grade level that serves as the basis for treatments 1 and 3. The amount in each cell represents the payment in pesos for a student with a given level of performance at the start of the grade (the baseline test score as defined above) and at the end of the grade.¹⁸ Payment levels were intended to be large enough to be expected to induce behavioral changes. The payments are similar in magnitude to the attendance incentives given by the Oportunidades program and to a scholarship program offered by the Ministry of Education.

As seen in the table, in the tenth and eleventh grades, there was no payment for performance on the ALI tests at the Pre-Basic level. In those two grades, students who scored at the Pre-Basic level on the baseline test received a payment of Mex\$4,000 (pesos) if they improved to the Basic level, \$9,000 if they improved to the Proficient level, and \$15,000 if they improved to the Advanced level. The increments become progressively larger (\$4,000 for the first increment, \$5,000 for the second, and \$6,000 for the third), recognizing that the effort necessary to improve from Pre-Basic to Proficient, for example, is likely to be greater than twice the effort in going from Pre-Basic to Basic.

In the tenth grade, students who originally scored at the Basic level received a payment of \$2,500 for maintaining their achievement level, a smaller amount than those at the Pre-Basic level who improved to the Basic level. The smaller payment reflects the presumably greater effort associated with improvement than with maintenance, a premise that is reflected throughout the incentive schedule. Considerably larger payments were given for improvement beyond the original Basic level: \$7,500 for achieving the Proficient level and \$13,500 for achieving the Advanced level. As before, the increments were increasing (\$5,000 for the first and \$6,000 for the second).

Students in the tenth grade who began at the Proficient level received no payment if they fell back to the Basic or Pre-Basic level. If they remained at the Proficient level, they received \$6,000, less than that received by students who improved to that level, while if they improved to the Advanced level, they received \$12,000. Students who originally scored at the Advanced level received \$4,500 if they fell back to the Proficient level and \$10,500 if they remained at the Advanced level.

Bonus amounts were the same for students in the eleventh grade with the exception that there was no payment for remaining at the Basic level.

¹⁸ A dollar was equivalent to about 11 pesos at the time.

TABLE 3
SCHEDULE OF INCENTIVE PAYMENTS (Pesos) FOR STUDENT ACHIEVEMENT

START OF GRADE	END OF GRADE			
	Pre-Basic	Basic	Proficient	Advanced
10th grade:				
Pre-Basic	0	4,000	9,000	15,000
Basic	0	2,500	7,500	13,500
Proficient	0	0	6,000	12,000
Advanced	0	0	4,500	10,500
11th grade:				
Pre-Basic	0	4,000	9,000	15,000
Basic	0	0	7,500	13,500
Proficient	0	0	6,000	12,000
Advanced	0	0	4,500	10,500
12th grade:				
Pre-Basic	0	0	5,000	10,000
Basic	0	0	5,000	10,000
Proficient	0	0	5,000	10,000
Advanced	0	0	5,000	10,000

The lack of any reward to those in the eleventh grade who remained at the Basic level reflected the ALI program emphasis on making progress toward achieving proficiency by the end of the twelfth grade. This formulation does, however, lead to an incentive for tenth graders who would perform at the Basic level on the tenth- and eleventh-grade tests to score at the Pre-Basic level on the tenth-grade test instead. In that case, the student would receive \$4,000 in total (\$0 in the tenth grade and \$4,000 in the eleventh grade) as opposed to \$2,500 in total (\$2,500 in the tenth grade and \$0 in the eleventh grade). Of course, students would be uncertain of their eleventh-grade score and indeed would have been better off scoring at the Basic level in the tenth grade if their score was at the Proficient or Advanced level in the eleventh grade. Although the potential incentive incompatibility could have been avoided by giving a bonus of \$1,500 for remaining at the Basic level in the eleventh grade, our view was that, given the newness of the ALI test and thus the inherent uncertainties students would have about how well they would perform on the tests (as well as their uncertainty about the cutoffs), the saliency of a zero bonus in emphasizing the goal of attaining proficiency outweighed the potential incentive problem.¹⁹

¹⁹ A more serious issue would arise, particularly if the program were universally adopted, with respect to performance on the ninth-grade ENLACE, where it would be unequivocally better for a student faced with the prospect of tenth-grade ALI incentives to have performed at the Pre-Basic level. However, this incentive would be mitigated by the fact that the ninth-grade ENLACE is a high-stakes test used by high schools as part of their competitive admissions criteria.

The twelfth-grade payment schedule provided a bonus only for achieving the Proficient or Advanced levels of performance, reflecting the goal that students reach at least the Proficient level by the time they graduate. There was a significantly higher payment for performance at the Advanced level, \$10,000, as opposed to the Proficient level, \$5,000.

2. Treatment 2 (Teacher Incentives Only)

In treatment 2, mathematics teachers were rewarded for the performance of the particular students they taught during the year. The reward was based on the sum of the rewards earned by the students as described in table 4. The per-student bonus was 5 percent of the bonus payments in the student schedules, except for the modification that teachers were penalized for students in the tenth and eleventh grades who were not at the Proficient or Advanced level at the end-of-grade test and who performed more poorly on the end-of-grade test than on the baseline test.

Consider, as an example, a teacher who had a tenth-grade mathematics class. For each student who improved from the Pre-Basic to the Basic level, the teacher would receive \$200 (5 percent of the \$4,000 that such a student would earn for himself or herself). If such a student instead improved to the Proficient category, the teacher would earn \$450 (again 5 percent of the \$9,000 the student would receive). If, however, a student lost ground, for example, moving from the Basic to the Pre-Basic level, the total of the student payments used to calculate the teacher reward would be reduced by \$125. A teacher's total payment was bounded below by zero.

A teacher with an eleventh-grade mathematics class faced the same incentive schedule as for a tenth-grade class except that, in conformity with the student incentive schedule, there was no payment for a student whose starting and ending test scores were at the Basic level. A teacher with a twelfth-grade class received \$250 for each student who reached the Proficient level and \$500 for each student who reached the Advanced level (5 percent of the student payment).

Schools operate on a semester basis. To obtain the academic year (two-semester) sum, each student whom the teacher had in his class during a semester was counted as one-half a student. The total earned by a teacher was the sum of the earnings from all the students in that teacher's classes over both semesters. To get an idea of the magnitude of the payments, consider a full-time teacher, one teaching five classes in each semester, who had 200 (year-equivalent) students (an average class size each semester of 40 students). Suppose that the average student payment was \$2,500 and that no student fell back. Such a teacher would earn \$25,000, which is a bonus of between 10 and 15 percent of the annual salary of a teacher in a federal high school.

TABLE 4
SCHEDULE OF INCENTIVE PAYMENTS (Pesos) PER STUDENT
FOR MATHEMATICS TEACHERS

START OF GRADE	END OF GRADE			
	Pre-Basic	Basic	Proficient	Advanced
10th grade:				
Pre-Basic	0	200	450	750
Basic	-125	125	375	675
Proficient	-125	-125	300	600
Advanced	-125	-125	225	525
11th grade:				
Pre-Basic	0	200	450	750
Basic	-125	0	375	675
Proficient	-125	-125	300	600
Advanced	-125	-125	225	525
12th grade:				
Pre-Basic	0	0	250	500
Basic	0	0	250	500
Proficient	0	0	250	500
Advanced	0	0	250	500

A specific aim of the teacher incentives design was to avoid teachers concentrating their effort on high-performing students.²⁰ It did so in three ways. First, teachers gained more from a lower-performing student achieving a given level than from a higher-performing student achieving that same level; for example, a teacher with a tenth-grade class earned \$200 if a student who scored initially at the Pre-Basic level improved to a Basic level but only \$125 for a student who scored initially at the Basic level and remained there. As with the students, the teacher effort required to elicit an improvement in a student initially scoring at the Pre-Basic level is presumably greater than the effort to maintain the score of the Basic-level student. As seen in table 4, the \$75 differential between initial Pre-Basic and Basic test scores carried over to the other final test categories. It also carried over between any two adjacent initial test score categories, that is, between Basic and Proficient and between Proficient and Advanced. This pattern resulted in a doubling of the differential between a student initially scoring at the Pre-Basic level and improving to the Proficient level and a student initially scoring at the Advanced level and falling back to the Proficient level (\$450 vs. \$225).

Second, and opposite to the “carrot” that compensated teachers for the extra effort associated with improving scores of low-performing students, the \$125 penalty incurred if a student regressed acts as a “stick”

²⁰ Neal and Schanzenbach (2010) use data from No Child Left Behind to analyze how teachers may have incentives to concentrate on subsets of students, in their case, students near cutoff values that determine whether school-level goals are met.

aimed at maintaining students at least at their initial performance level. Third, for students initially at the Pre-Basic level, the possible payments to the teachers (and to the students) were strictly nonnegative, with relatively larger payments if the students improved a great deal. For example, the mathematics teacher received \$200 if a student who started at Pre-Basic in the tenth grade advanced to Basic, \$450 if a student who started at Pre-Basic advanced to Proficient, and \$750 if a student who started at Pre-Basic advanced to Advanced.

In contrast to the ALI design, in an incentive system that depended on, say, the average performance of teachers' students, teachers would receive a greater reward if low-performing students were encouraged to not take the examination. Under the ALI design, that was not the case for the lowest-performing students; those who scored Pre-Basic in the baseline test could not subtract from the teacher's reward and, as noted, could contribute considerably if their performance improved in the end-of-year test. For students with Basic and Proficient baseline scores, there was some potential for teachers losing if the students dropped back because of the "stick," but the "carrot" was intended to be sufficient to offset that potential loss.

3. Treatment 3 (Aligned Student, Teacher, and Administrator Incentives)

Students.—In treatment 3, in the tenth and eleventh grades, each student received a reward based on individual performance as in treatment 1 (according to the schedule in table 4) and also on the performance of the other students in his or her mathematics class. In the twelfth grade, the student received a reward based only on individual performance. The first component was calculated in exactly the same way as in treatment 1. The second component was calculated as a fixed proportion, 1 percent, of the total payments earned by classmates.²¹

The rationale for paying students for the performance of their classmates rests on possible synergies of two different kinds, both of which depend on there being a fundamental complementarity between a student's own effort and classmates' efforts. In one case, the effort of one's classmates is a pure externality in which a positive climate or culture of learning can be created by the overall effort within the classroom. This climate affects the amount of effort each student puts into his own learning. A second synergy arises when students actively help other students, for example, if higher-performing students tutor lower-performing stu-

²¹ In calculating the class sum over a year's time, the class sums in each of the two semesters are multiplied by one-half and then added together. This procedure accounts for compositional changes in classes, i.e., that a student may not be in classes with the same students in both semesters.

dents. The component of the bonus payment based on class performance provided an (extra) incentive for this behavior. Not only may such activities improve learning among low performers; they may also improve learning among high performers as teaching itself can lead to a deeper understanding.

Mathematics teachers.—The reward to full-time mathematics teachers was the sum of the total performance payments earned by the students in their classes calculated as in treatment 2 (according to the schedule in table 4) and a fixed proportion, 25 percent, of the average full-time-equivalent adjusted performance payments earned by the other mathematics teachers (across all grade levels).²² The rationale for the second component of the reward was to stimulate cooperation among the mathematics teachers. Such cooperation may take the form of formal or informal discussions of teaching methods and subject matter, mentoring less experienced teachers, and directly sharing lesson plans or other class materials.

Nonmathematics teachers.—Nonmathematics teachers received a cash payment that is 25 percent of the schoolwide average (full-time-equivalent) mathematics teacher performance payment. Payments for part-time teachers were adjusted for their own full-time equivalence status. The rationale for this payment recognizes the potential importance of the overall learning environment in the school and of the potential value of interactions among teachers from different disciplines in sharing ideas about pedagogy (and perhaps subject matter in the case of allied fields like physics) and about students that they have in common.

Principals and associate principals.—Principals received a cash payment that is 50 percent of the average full-time-equivalent mathematics teacher performance payment. Associate principals received a cash payment that is 25 percent of the schoolwide average full-time-equivalent mathematics teacher performance payment, adjusted for their own full-time equivalence status. These payments recognize the importance of support services provided by administrative personnel in fostering learning within the school.²³

IV. Results

A. ALI Test Completion Rates

Most students in Mexico who complete ninth grade enter high school (tenth grade). Of the cohort that entered kindergarten in 1996, about 65 percent completed ninth grade and 95 percent of those continued to

²² Rewards to part-time teachers were prorated by their own full-time equivalence status.

²³ The formulas used for calculating the bonuses for students, mathematics teachers, nonmathematics teachers, principals, and associate principals are given in App. B.

tenth grade. However, only 76 percent of those who entered tenth grade actually completed the grade, and of those, 81 percent graduated from high school. Given high dropout rates at the national level, one might expect to see a significant proportion of students in the ALI schools who begin the school year but do not take the ALI examination at the end of the year. To the extent that attrition is not uniform across the treatment and control schools, treatment effect estimates based on simple mean comparisons could be biased.

Table 5 provides attrition figures for the 2008–9 tenth-grade cohort, both for enrollment between the first two (year 1) semesters and for the completion of the ALI tests over the three years. The first row of figures shows the percentage of students in treatment and control schools enrolled in the fall semester who were also enrolled in the spring semester. Attrition between the fall and spring semesters is 9.9 percent for the controls and 12.0, 9.5, and 13.2 percent for T1, T2, and T3; none of the differences are statistically significant at conventional levels. The fact that continuation rates do not differ significantly either in magnitude or statistically across the treatment and control groups may be surprising, particularly for T1 and T3, where students receive direct monetary incentives. Indeed, one might have expected the ALI program to have reduced dropout rates. However, as already noted, ALI is not the only program providing attendance incentives. Almost 40 percent of the students in ALI schools receive a substantial attendance subsidy as part of the Oportunidades conditional cash transfer program. In addition, as part of another scholarship program, students whose family income is below the poverty line and who successfully progress from one grade to

TABLE 5
CONTINUATION RATES AND ALI TEST COMPLETION RATES
FOR TENTH GRADE, YEAR 1 COHORT

	Control Schools (1)	Treatment 1 Schools (2)	Treatment 2 Schools (3)	Treatment 3 Schools (4)
% enrolled in spring of year 1 given enrollment in the fall	90.1	88.0 (.384)	90.5 (.851)	86.8 (.203)
% taking ALI exam in year 1 given enrollment in both semesters	85.9	88.9 (.255)	85.6 (.926)	88.6 (.388)
% taking ALI exam in year 2 given test taken in year 1	7.74	8.19 (.053)	7.69 (.820)	8.08 (.172)
% taking ALI exam in year 3 given test taken in year 1	6.78	7.09 (.317)	6.64 (.567)	7.03 (.375)

NOTE.—Numbers in parentheses are as follows: col. 2: p -value for test C = T1; col. 3: p -value for test C = T2; col. 4: p -value for test C = T3. All account for school-level clustering.

another receive a scholarship payment; the baseline scholarship amount is considerable, with increments for achieving a high grade point average. Given these subsidies already in place in all federal high schools, it is less surprising that there is no discernible additional effect of the ALI program on the dropout rate between semesters. The existence of these additional programs does not pose a problem for the estimation of the marginal effect of the ALI program, although it could affect the generalizability of the ALI program impact results to other settings.²⁴

As also seen in table 5, conditional on enrollment in both semesters, 85.9 percent of the controls took the ALI exam, about the same as for T2, but about 3 percentage points lower than for T1 and T3; again, the differences from the control group are not statistically significant. This pattern of test taking is maintained in years 2 and 3. For example, 67.8 percent of the controls who took the ALI exam in year 1 as tenth graders also took the exam in year 3 as twelfth graders. The comparable figures for T1, T2, and T3 are 70.9, 66.4, and 70.3 percent, which do not differ statistically from the controls. We previously noted that there is potentially an incentive for teachers in T2 and T3 schools to try to identify students in the tenth and eleventh grades who would do worse on the ALI test than on the baseline ninth-grade ENLACE test and in some way have them not take the ALI test. However, T3 students cannot lose anything by taking the test, which acts to counterbalance the teacher incentive. Moreover, the poorest-performing students on the ENLACE test, those who scored at the Pre-Basic level, cannot fall back, so there is no incentive for the teachers in T2 or T3 to discourage them from taking the test.²⁵ Given the data in table 5, it does not appear that teachers manipulated the test-taking sample in any appreciable way.

B. Treatment Effects

As we discussed, the protocol for the administration of the ALI test, which called for an external monitor in each classroom, was intended to minimize the possibility of cheating. To determine whether and to what extent cheating had occurred, an independent statistical analysis of student answer sheets was conducted by George Wesolowsky, based on

²⁴ Recall from table 2 that the percentage of Oportunidades recipients did not differ between the treatment and control schools. Student self-reports in year 3 indicate that the percentage of students receiving a Minister of Public Education scholarship ranges between 11 and 13 percent across the treatment and control groups.

²⁵ In addition, teachers are not officially informed of student ninth-grade ENLACE scores, which is relevant for determining teacher bonuses for the tenth grade (and for the eleventh grade in the first year), nor are they informed about the ALI test score results of individual students, which is relevant for determining teacher bonuses for the eleventh grade.

the methods described in his 2000 study.²⁶ The method specifies a statistical model for the probability that student i correctly answers multiple-choice question j , incorporating a parametric function of the “difficulty” of the question and the “ability” of the student. Using that model, it is possible to determine, for every pair of students and for each question, the probability that two students will have the same answer assuming that all wrong answers are equally likely. With a critical value for the number of observed matches chosen, the null hypothesis of no copying was tested for each student pair. The critical value was based on a Bonferroni correction such that the probability was one that at least one pair of students would be falsely accused. Given that criterion, we interpret the amount of cheating that is identified as an upper bound.²⁷

In each pair, we assigned the student with the higher ninth-grade ENLACE score as the source of the answers and the other student as a “copier.”²⁸ If a student is ever a copier in any pair, even if the student was a source in another pair, the student’s final designation is as a copier. Table 6 provides the results from the cheating analysis. The table reports, for each grade and programmatic year, the percentage of students who were members of a cheating pair and the percentage who were copiers.²⁹ The estimated (upper-bound) percentage of copiers in the control group varied between 2 and 6 percent, depending on grade and year, with no obvious pattern across grades or years. The extent of copying was similar for students in T2, who, like the controls, had no direct ALI incentive. However, the percentage of copiers was considerably higher and of similar magnitudes in both T1 and T3, the treatments for which there were direct student incentives. The percentage of copiers was espe-

²⁶ To our knowledge, student copying has not been studied in relation to student incentive performance programs. However, the education literature demonstrates that a significant fraction of students admit to cheating during regular school examinations (Cisek 1999).

²⁷ The analysis compared pairs of students within a grade independently of treatment group, school, or class. The number of pairs identified who were not in the same treatment group or school was always negligible, supporting the validity of the method. However, that was not always the case across classes, indicative of the relatively low critical value used to determine cheating pairs. In particular, 6.0 percent of cheating pairs crossed class boundaries in year 2 for grade 12, 9.8 percent in year 1 for grade 12, 14.4 percent in year 3 for grade 11, and 29.2 percent in year 3 for grade 12. Unfortunately, we cannot rule out that in some cases classes were combined for the test administration, though no such occurrences were reported by the overall external supervisor assigned to the school.

²⁸ To determine the accuracy of this classification, we compared the difference between the ninth-grade ENLACE and the ALI test scores for three groups: noncheaters, sources, and copiers. In essentially all years, grades, and treatment groups, the difference in scores for copiers was far greater than for the other two groups, which were themselves similar. For example, in year 2 among the controls, the difference between the scores was -34 for noncheaters, -40 for sources, and 53 for the copiers. The similar figures for T3 were 8.9, 21.2, and 190.8. In defining copiers, for cases in which the ninth-grade ENLACE was missing, the assignment was based on the ALI test score.

²⁹ The average percentage of copiers is more than half the percentage of all cheaters because some students were sources for more than one copier.

TABLE 6
 PERCENTAGE OF STUDENTS WITH NONINDEPENDENT TEST SCORES
 BY YEAR, GRADE, AND TREATMENT

	GRADE 10		GRADE 11		GRADE 12	
	Copiers (%)	Cheaters (%)	Copiers (%)	Cheaters (%)	Copiers (%)	Cheaters (%)
Year 1:						
C	3.7	6.4	4.5	7.8	5.7	9.3
T1	5.1	9.1	10.9	14.9	5.2	8.4
T2	3.4	5.8	3.9	6.5	3.7	6.5
T3	3.7	6.7	10.1	14.9	2.7	4.7
Year 2:						
C	3.5	6.1	3.6	6.2	2.4	4.5
T1	6.4	11.0	19.1	27.6	12.7	17.3
T2	4.3	7.4	6.2	9.8	3.4	5.5
T3	6.6	10.6	17.2	23.9	10.6	16.0
Year 3:						
C	3.1	5.7	4.6	7.8	2.5	4.7
T1	8.1	13.2	19.8	28.2	17.5	24.7
T2	4.2	7.3	4.1	7.1	4.0	6.8
T3	10.3	16.2	23.8	31.3	15.4	21.3

cially high for the eleventh grade in years 2 and 3 and for the twelfth grade in year 3; the highest percentage of copiers was for the eleventh grade in year 3, 19.8 percent for T1 and 23.8 percent for T3.³⁰

As noted, teachers did not administer the test or handle the test booklets, and all test copies were to be returned to the state ministry offices. The finding that the extent of cheating was no greater in T2 than in C nor greater in T3 than in T1 is consistent with that protocol being followed.³¹ In addition, we analyzed the difference in scores between treatment and control schools based on the “anchor” questions, that is, the 30 percent of questions that were repeated each year.³² Larger treatment effects based on those questions alone might indicate that teachers somehow gained access to the tests and used them to teach to the test in subsequent years. However, we found the treatment effects (either adjusted for cheating or unadjusted) to be no larger for the anchor questions.

Given the results of the cheating analysis, we report average treatment effects in table 7 both with and without an adjustment to the ALI scores

³⁰ Cheating was highly concentrated among a few schools. In many cases, the top three schools accounted for three times as many copiers in the treatment groups in a given year and grade as the schools accounted for the total number of students. At the extreme, e.g., the top three T1 schools in year 2 for grade 11 accounted for 73.3 percent of the copiers and only 21.7 percent of all of the students.

³¹ For an analysis of teacher cheating, see Jacob and Levitt (2003).

³² Anchor questions are not supposed to be more or less difficult than other questions.

TABLE 7
AVERAGE TREATMENT EFFECTS (ATE) WITH AND WITHOUT ADJUSTMENTS FOR COPIERS: ALL PROGRAM YEARS

GRADE	YEAR 1: ACADEMIC YEAR 2008/9			YEAR 2: ACADEMIC YEAR 2009/10			YEAR 3: ACADEMIC YEAR 2010/11		
	T1	T2	T3	T1	T2	T3	T1	T2	T3
With copying adjustment:									
10th grade:									
ATE	16.9 (4.90)	1.27 (5.74)	31.4 (5.79)	29.1 (4.57)	.11 (5.34)	46.6 (7.61)	32.3 (4.77)	13.5 (5.54)	63.4 (10.4)
p -value: Tj = T3	.010	<.001040	<.001002	<.001	...
11th grade:									
ATE	13.6 (5.40)	-4.84 (5.50)	18.6 (7.39)	29.7 (4.89)	2.11 (6.05)	43.7 (8.33)	25.2 (4.24)	-2.00 (4.31)	42.1 (5.64)
p -value: Tj = T3	.545	.004098	<.001011	<.001	...
12th grade:									
ATE	9.63 (6.85)	4.71 (6.58)	28.8 (6.36)	21.9 (5.04)	-4.46 (6.10)	34.8 (6.46)	22.7 (7.49)	3.99 (7.54)	56.7 (15.1)
p -value: Tj = T3	.010	<.001078	<.001015	<.001	...
No copying adjustment:									
10th grade:									
ATE	18.5 (5.02)	1.11 (5.35)	32.3 (6.18)	32.4 (5.24)	.31 (5.74)	54.7 (11.1)	41.5 (6.25)	15.9 (6.16)	83.4 (16.9)
p -value: Tj = T3	.025	<.001073	<.001014	<.001	...
11th grade:									
ATE	22.4 (7.22)	-2.98 (6.74)	27.8 (9.93)	55.5 (7.51)	6.17 (6.91)	67.4 (12.7)	51.3 (9.05)	-1.36 (8.89)	106.4 (25.6)
p -value: Tj = T3	.639	.006382	<.001037	<.001	...
12 grade:									
ATE	9.73 (7.04)	4.73 (6.62)	29.3 (6.67)	36.0 (7.32)	-1.81 (6.30)	44.6 (7.99)	42.3 (8.15)	7.33 (7.98)	90.2 (21.3)
p -value: Tj = T3	.011	<.001400	<.001022	<.001	...

NOTE.—Standard errors are in parentheses. Tj refers to either T1 or T2.

of the copiers. The treatment effects are based on standardized test scores, normalized (as in the ENLACE) to have a mean of 500 and a standard deviation of 100 for the controls. The predicted ALI test scores for copiers in grades 10 and 11 are based on regressions of the ALI scores of the control students (who were not copiers) on their ninth-grade ENLACE score and a dummy variable for a missing ENLACE score; for the twelfth grade, the regression included in addition the student's twelfth-grade ENLACE score and a missing dummy.

There are five tenth-grade entry cohorts that are observed in different grades over the three years of the ALI program. For example, the 2006–7 tenth-grade entry cohort was observed only once, in the twelfth grade in year 1, and the 2007–8 cohort was observed twice, in the eleventh grade in year 1 and in the twelfth grade in year 2. The samples for each cohort shown in table 7 are based on nonattriters over the years each cohort is observed.³³ The average treatment effect is the difference in the mean test score of each treatment and the control group conditioning on the following preprogram school-level variables to increase precision: the schoolwide average ninth-grade mathematics ENLACE scores for the current tenth-, eleventh-, and twelfth-grade students; the 2007–8 academic year average twelfth-grade mathematics and language ENLACE scores; regional dummies; and state dummies (for the five states in which there was at least one school in each ALI group). Standard errors account for clustering at the school level. The estimates with the adjustment represent, in our view, a reasonable lower bound on treatment effects.³⁴

In year 1, as seen, the adjusted and unadjusted treatment effects are very close for grade 10. On the basis of the adjusted estimates, T3 students in the tenth grade scored 31.4 standardized points higher, on average, than did C students and T1 students 16.9 standardized points higher.³⁵ A 95 percent confidence interval for T3 ranges from 19.8 to 43.0 points and for T1 from 7.1 to 26.7. The score for T2 students is almost identical to that of C students (the point estimate is 1.27 points [SE 5.74]). The *p*-value for the (two-sided) test that the T3 treatment effect differs from the T1 effect is .010 and that for the test of T3 against T2, less than .001. The adjusted and unadjusted estimates diverge slightly for grade 10 in year 2. With the more conservative adjusted estimates, the T3 treatment increases test scores by 46.6 standardized points over the controls and the

³³ The results are not sensitive to the exclusion of attriters. As discussed, attrition was similar across the treatment and control groups.

³⁴ In App. table C1, we report results in which the ALI scores of copiers were set equal to the 25th percentile score of the control group students who were noncopiers. The results do not differ much for tenth grade in all years and for eleventh and twelfth grades in year 1. The magnitudes of the treatment effects for the upper two grades in years 2 and 3 are somewhat smaller than those reported in table 7.

³⁵ Each standardized point represents 0.01 of a standard deviation.

T1 treatment by 29.1 points. As in year 1, the T2 treatment has no effect on test scores. Both the T3 and T1 effects have relatively narrow 95 percent confidence intervals, and the difference between T3 and T1 is statistically significant. A larger difference between the adjusted and unadjusted treatment effects emerges in year 3, with the adjusted effects being about 25 percent lower. The adjusted treatment effects are 63.4 points for T3, 32.3 for T1, and 13.5 for T2 (the only time when there is a statistically significant T2 effect across all grades and years).

There are a number of reasons why the program effect, as seen, would be expected to grow over time. First, the schools were not informed of the ALI program until well into the first semester of the first year, and it took some time after that for the students to be apprised of the program details. Second, the only information given to the students and teachers about the test was that it would adhere strictly to the curriculum of each grade. In the first year, there would be considerable uncertainty about the format and difficulty of the test for both students and teachers. Having experience with the test in the first year would resolve some of that uncertainty, even though the tests themselves were not left behind to minimize the possibility of teaching to the test.

The eleventh- and twelfth-grade results in many ways mimic those of the tenth grade. The adjusted treatment effects are larger in years 2 and 3 than in year 1 and of a magnitude similar to those for the tenth grade. Estimated treatment effects are larger for T3 than for T1, each is precisely estimated, and they are statistically distinguishable. There is no discernible T2 effect. In year 2 (3), on the basis of the adjusted estimates, T3 students in the eleventh grade scored 43.7 (42.1) standardized points higher, on average, than did the C students, and twelfth graders 34.8 (56.7) points higher.

1. Heterogeneous Treatment Effects

Gender.—Table 8 reports both adjusted and unadjusted average treatment effects for the 2008–9 tenth-grade cohort for each of the three years of the program by gender. Recall that this sample consists of students who took the ALI test in all three years. There is little gender difference in average treatment effects at any grade.

Ninth-grade ENLACE.—Table 8 also reports average treatment effects for this cohort distinguished by performance on the ninth-grade ENLACE. Recall that the incentive schedule was designed so that low-performing students would have a larger incentive to improve their mathematics knowledge, owing to the presumed greater effort necessary to achieve any standard, and teachers would not have an incentive to specialize their efforts on higher-performing students. Table 8 provides evidence on whether the incentive design accomplished these goals.

TABLE 8
AVERAGE TREATMENT EFFECTS BY GENDER AND BY NINTH-GRADE ENLACE: 2008-9 TENTH-GRADE COHORT

	10TH GRADE (Year 1)			11TH GRADE (Year 2)			12TH GRADE (Year 3)		
	T1-C	T2-C	T3-C	T1-C	T2-C	T3-C	T1-C	T2-C	T3-C
Adjusted score:									
Gender:									
Female	18.7 (5.65)	1.51 (6.39)	35.8 (5.30)	33.8 (5.62)	4.71 (6.40)	51.0 (7.43)	28.8 (7.85)	6.72 (7.57)	63.9 (15.8)
Male	15.0 (5.91)	1.32 (6.42)	33.0 (7.48)	25.3 (5.79)	-.32 (6.64)	45.5 (9.98)	14.7 (7.84)	-1.10 (9.03)	63.7 (14.9)
9th-grade ENLACE:									
Pre-Basic	15.0 (4.07)	1.95 (4.49)	26.8 (4.84)	24.4 (3.59)	2.11 (4.79)	33.4 (5.98)	23.6 (6.28)	4.75 (6.32)	50.7 (12.7)
Basic	18.2 (5.92)	-1.70 (7.43)	30.8 (7.71)	35.3 (5.95)	-.15 (7.32)	48.9 (9.54)	22.5 (8.87)	2.72 (8.76)	57.4 (16.6)
Proficient or Advanced	28.0 (12.5)	1.19 (16.1)	45.3 (17.5)	47.3 (13.5)	-2.12 (16.1)	58.1 (19.8)	45.6 (16.1)	17.9 (17.7)	70.2 (23.7)
Unadjusted score:									
Gender:									
Female	19.9 (5.87)	.62 (5.73)	37.2 (5.76)	54.7 (7.54)	5.91 (6.54)	77.6 (9.94)	44.8 (8.93)	9.04 (8.32)	101 (22.9)
Male	17.1 (6.01)	2.09 (6.42)	34.6 (7.62)	55.2 (8.11)	5.90 (8.33)	76.2 (12.8)	37.2 (7.75)	2.21 (9.76)	99.4 (18.7)
9th-grade ENLACE:									
Pre-Basic	16.4 (4.46)	1.44 (4.87)	28.8 (6.26)	45.4 (7.21)	2.94 (5.96)	65.3 (11.5)	43.9 (7.59)	6.89 (7.92)	97.0 (22.4)
Basic	21.0 (6.10)	-.67 (6.89)	31.8 (7.90)	64.0 (8.85)	5.06 (8.38)	69.7 (14.8)	39.9 (9.32)	7.27 (8.79)	80.6 (21.2)
Proficient or Advanced	27.8 (12.4)	-.26 (15.9)	44.3 (17.5)	73.6 (16.0)	5.86 (17.0)	72.4 (22.1)	61.4 (17.6)	25.2 (18.5)	86.3 (26.7)

NOTE.—There are 5,388 males, 5,665 females, 4,926 observations in the Pre-Basic category, 4,294 observations in the Basic category, and 1,223 observations in the Proficient or Advanced category. Standard errors, in parentheses, account for clustering at the school level.

ALI incentives induced significant increases in test scores across all ninth-grade ENLACE performance categories: Pre-Basic, Basic, and Proficient and Advanced combined. Concentrating on the adjusted estimates, the T3 treatment effect for those students in the Pre-Basic ENLACE category was 26.8 points in the tenth grade, 33.4 points in the eleventh grade, and 50.7 points in the twelfth grade. The comparable effects for T1 were 15.0, 24.4, and 23.6. Those students scoring in the Basic category on the ENLACE scored higher than students in the Pre-Basic category in the tenth and eleventh grades, but about the same in the twelfth grade. Students scoring in the two highest ENLACE categories have the highest treatment effects: 28.0, 47.3, and 45.6 for T1 in the three grades and 45.3, 58.1, and 70.2 in T3 for the three grades. Treatment 2 effects for all ENLACE categories and in all grades are small, in a few cases negative, and not statistically significant.

B. Potential Caveats

Test-taking effort.—The estimated treatment effects for T1 and T3, even those that are adjusted for copying, appear to be large relative to the literature, especially for T3 in year 3. An implicit assumption is that the test-taking effort of the students in the control schools (and to some extent the T2 schools) is the same as that of the students in the T1 and T3 schools. Given that the ALI test is a low-stakes (if not a no-stakes) test for the C students, it is possible that the average treatment effects are exaggerated. Some evidence can be brought to bear on the issue.

We use two methods to obtain a lower bound on the treatment effects that accounts for differential test-taking effort, in one case a lower bound only for T3 and in the other for both T1 and T3. The first method assumes that the test-taking effort of T1 students is no less than that of T3 students. Although there are greater monetary incentives for T3 students and teachers due to group incentives and administrators are also rewarded, it seems reasonable to assume that the incentives for the students in T1 are substantial enough for students to want to maximize their effort on the test. Under that assumption and also assuming that the entire difference between the test outcome of T1 and C students in each year is due to a lack of test-taking effort on the part of the controls, the difference in the treatment effect between T3 and T1 would be a lower-bound estimate of the treatment effect for T3. The results in table 7 (with the copying adjustment) imply that, for year 3, the lower bound for the treatment effect for T3 is 31.1 (63.4 – 32.3) standardized points for grade 10, 16.9 for grade 11, and 33.9 for grade 12.

The second method assumes that test-taking effort of the students in the control schools is at least as great in year 3 as in year 1 and, further, that the treatment effect for T1 in year 1 was entirely due to the lack of

test-taking effort of the controls, that is, that the true treatment effect in year 1 was zero for T1. In that case, a lower-bound estimate of the treatment effect is the difference between the treatment effect in year 3 and the treatment effect of T1 students in year 1. This assumption leads to the following lower-bound estimates of treatment effects (adjusted for copying): in year 3 for T1, 15.4 for the tenth grade, 11.6 for the eleventh grade, and 13.1 for the twelfth grade. Similarly, the lower-bound estimates for T3 are 46.5 for the tenth grade, 28.5 for the eleventh grade, and 47.0 for the twelfth grade.

Clearly, the first bounding approach is more severe in terms of the influence of the lack of test-taking effort on treatment effects, implying as it does a zero treatment effect for T1 in all three years and thus, given the increased performance of T1 students over time, a sizable reduction in test-taking effort of the controls in year 3 relative to year 1. Perhaps one could make an argument that the twelfth graders, who, in year 3, were taking the ALI test for the third time, and the eleventh graders for the second time, had diminished test effort, but such an argument is less compelling for the tenth-grade students in year 3 who were taking the ALI test for the first time. Indeed, raw scores (the percentage of questions answered correctly) for the control school students were similar in all three years in any of the grades, which is consistent with a constant level of test-taking effort as assumed in the second set of lower-bound estimates.

Twelfth-grade ENLACE.—In the middle of the spring semester of the twelfth grade, students are administered a national mathematics “competency” test. The scores of the students in the ALI schools, which were made available to us, can be used to see whether the effects of the treatments on the curriculum-based ALI tests carry over to the twelfth-grade ENLACE. On the basis of the same regression specification as in table 7, the treatment effect estimates for the twelfth-grade students in year 3 of the ALI program were (standard errors in parentheses) -19.1 (9.12) for T1, -15.6 (9.12) for T2, and 18.3 (13.7) for T3. Not only are these results quite different from those for the ALI exam, but the negative and statistically significant effect for T1 (given the randomization) is anomalous. Whether these results are due to the very different content of the ALI and ENLACE tests or perhaps to diminished test-taking effort, particularly for the T1 and T3 students who are concentrating their effort and attention on the end-of-year ALI tests, is unclear.

Translation to raw percentage scores.—Measured in standard deviation units, the treatment effects estimated under the ALI experiment are large relative to the range of estimates reported in the experimental incentives literature. However, as measured by raw percentage scores, the performance of the treatment groups is less striking. For example, for the 2008/9 cohort, students in T3 (T1) answer only 45.3 (43.6) percent

of the questions correctly on the tenth-grade test as compared to 41.7 percent for the controls. Similarly, on the eleventh-grade test, the average score for the T3 (T1) students is 42.8 (41.3) versus 38.2 for the control students, and on the twelfth-grade test, the T3 (T1) score was 48.4 (45.1) as opposed to the control average score of 42.9. As previously noted, the test covered the curriculum for each grade. Although students in T3 and T1 gain relative to controls, in absolute terms they master less than half of the curriculum.

3. Student and Teacher Effort

As part of the ALI project, students and teachers participating in the ALI program were given (self-administered paper) questionnaires that attempted to measure learning and teaching effort in each semester of each year. Table 9 compares the survey answers in year 3 across the control and treatment groups for the students who took the ALI test and responded to the survey and whose teachers also responded to the survey (which is true for about 90 percent of the students who took the exam). The estimated treatment-control differences are broadly consistent with students in T1 and T3 groups having higher levels of effort than those in the control group and with students in the T2 group behaving very similarly to the controls. For example, students in T1 and T3 spent more time studying mathematics, were significantly less likely to text or watch television while doing homework, were significantly more likely to give help to classmates, and to put "much effort" into their school work. Students in T1 and T3 also reported spending no less time studying subjects other than mathematics, which implies that the ALI incentives did not appear to have shifted study effort away from other subjects.

Table 9 also compares measures of teacher effort. The evidence is mixed on whether teachers in the treatment groups exerted more effort. A higher fraction of teachers in the treatment groups reported preparing their students for the ALI exams and helping their students outside of class to prepare for the exam, with the highest fractions observed in the T3 group. However, there were a number of measures related to time spent preparing for class (not shown in the table) in which no difference was observed between the different groups of teachers.

C. Rationalizing the Results with a Model of Student and Teacher Effort Choice

A perhaps puzzling feature of the results is the large differential between T3 and T1 effects and a zero effect of T2. One potential explanation is that T3 is not simply a combination of T1 and T2 given that T3 includes bonuses for students and mathematics teacher based on the performance of peers and for nonmathematics teachers and administrators. It

TABLE 9
STUDENT AND TEACHER EFFORT MEASURES FOR CONTROLS AND TREATMENT/CONTROL DIFFERENCE: YEAR 3

	C			T1 - C			T2 - C			T3 - C		
	10th	11th	12th	10th	11th	12th	10th	11th	12th	10th	11th	12th
Student:												
Average hours/week study math	4.68	4.45	4.53	.199 (.095)	.408 (.135)	.385 (.124)	-.138 (.091)	-.070 (.182)	-.097 (.165)	.397 (.112)	.301 (.135)	.370 (.127)
Average hours/week study nonmath subjects	5.56	5.48	5.32	.109 (.122)	.189 (.122)	.250 (.156)	-.161 (.129)	-.134 (.168)	-.040 (.153)	.152 (.122)	.074 (.134)	.168 (.127)
Fraction pay attention >75% of time	.473	.479	.481	.070 (.022)	.048 (.021)	.042 (.024)	.015 (.028)	.007 (.030)	-.006 (.026)	.101 (.028)	.070 (.023)	.050 (.032)
Fraction never or almost never text while doing homework	.423	.429	.415	.109 (.023)	.093 (.028)	.056 (.027)	.023 (.026)	.004 (.028)	-.007 (.028)	.126 (.024)	.097 (.022)	.061 (.022)
Fraction never or almost never watch TV while doing homework	.493	.517	.498	.077 (.028)	.075 (.018)	.066 (.024)	-.021 (.025)	-.010 (.022)	-.010 (.020)	.088 (.026)	.093 (.022)	.060 (.027)
Fraction gave help to classmates	.599	.608	.643	.055 (.020)	.058 (.022)	.026 (.023)	-.017 (.020)	-.014 (.019)	-.041 (.028)	.086 (.020)	.087 (.022)	.026 (.028)
Fraction report putting much effort	.466	.489	.486	.077 (.022)	.090 (.026)	.087 (.028)	-.039 (.021)	-.029 (.030)	-.017 (.025)	.114 (.022)	.093 (.021)	.092 (.037)
Teacher:												
Fraction prepared students for ALI test	.167	.260	.241	.202 (.103)	.181 (.121)	.211 (.107)	.182 (.091)	.155 (.106)	.111 (.114)	.412 (.106)	.256 (.110)	.176 (.098)
Fraction helped students outside of class to prepare for ALI test	.241	.220	.204	.338 (.104)	.339 (.126)	.453 (.102)	.341 (.103)	.390 (.111)	.391 (.122)	.435 (.098)	.554 (.092)	.482 (.103)

NOTE.—Standard errors, in parentheses, are corrected for clustering at the school level.

should be noted, however, that the bonuses based on peer performance account for only about 25 percent of the total bonus for both students and teachers. Given the large differential between T1 and T3 treatment effects in many cases, it may seem implausible that it is mostly due to the additional elements in T3 beyond the incentives for student and teacher own performance.

It is useful to ask whether this pattern of results can be generated by a model of the determination of student performance. To that end, we consider a strategic model of the effort choices of the students in a class and of the teacher. Each student, s , begins a grade with a given level of knowledge, denoted by k_{0s} , and the teacher, t , with a given stock of instructional capital, k_{0t} . During the school year, each student supplies learning effort, ε_s , and the teacher supplies instructional effort, ε_t , which is a public input. End-of-year knowledge, K_s , is given by the production function

$$K_s = F(\varepsilon_s, \varepsilon_t; k_{0s}, k_{0t}, S), \quad (1)$$

where S is the size of the class. Students care about their end-of-year knowledge and teachers about the knowledge of the students in their class. Students and the teacher face effort cost functions $c_s(\varepsilon_s)$ and $c_t(\varepsilon_t)$. Each student in the class maximizes

$$V_s = U_s(K_s) - c_s(\varepsilon_s), \quad (2)$$

where $U_s(K_s)$ is the utility the student receives from end-of-year knowledge, and the teacher maximizes

$$V_t = U_t(\mathbf{K}_s) - c_t(\varepsilon_t), \quad (3)$$

where \mathbf{K}_s is the class vector of student end-of-year knowledge (determined by [1] for each student).³⁶ Student and/or teacher bonuses can be accommodated in the model by augmenting the student and teacher utility functions to include the ALI bonus schedules.

Todd and Wolpin (2012) develop sufficient conditions, in which students and the teacher play a Nash game, that can generate the ALI results and an estimation method to empirically implement the model. The critical assumptions are (1) the student can supply a minimal level of effort at zero cost but must pay a fixed (and variable) cost for supplying effort above that minimal level, (2) the marginal product of the teacher's effort is zero if the student chooses to supply only minimum effort, and (3) student and teacher effort are complementary inputs in equation (1). Given these

³⁶ One also could allow the teacher and/or student utility to depend on initial knowledge or value added.

assumptions, a teacher bonus alone may not be sufficient to induce enough students within a class, who without the bonus were supplying minimum effort, to supply above-minimum effort in response to an increase in teacher effort (given complementarity of student and teacher effort). A student bonus alone, given that it directly affects student incentives, can induce such a response and will also increase teacher effort (given complementarity). Once a student bonus is in place and students are supplying above-minimum effort, an additional teacher bonus can further augment both teacher and student effort. It is thus possible to observe that the $T2 - C$ difference is zero, and both the $T1 - C$ effect and the $T3 - T1$ effect are positive. Alternatively, it is possible to generate the ALI results through restrictions on the shape of the production function. Such an explanation would seem to require a significant degree of non-convexity (akin to a fixed cost).

Within any class, the solution of the model depends on the teacher's preference over each student's end-of-year knowledge, the teacher's instructional capital and effort cost, and the distribution within the class of student initial knowledge, preferences over knowledge, and effort costs. In the ALI project, as noted, we collected information from surveys of teachers and students that provide measures and determinants of these characteristics across all the schools and classes. Todd and Wolpin (2012) exploit those data to estimate a version of the model that allows for heterogeneity across schools, teachers, and students in these attributes. Extrapolating the ALI results to other populations in which the distributions of these attributes differ would require specifying how they differ and related data to account for those differences. The model provides the framework for the data collection that would be necessary.

D. The Cost of the Treatments

Table 10 summarizes the program cost in terms of the incentive payments for the second year and for each of the three treatments.

Treatment 3.—Tenth-grade students in T3 on average earned \$2,991 based on their own performance and an additional \$1,108 based on the performance of the students in their class. Thus, the class contribution was about 25 percent of the total. Eleventh-grade students earned \$2,679, on average, from their own performance, lower than what tenth-grade students received primarily because there was no payment for students who performed at the Basic level on the baseline ENLACE test and on the ALI test. They earned an additional \$861 from the performance of their classmates. Twelfth graders earned, on average, \$991, less than the other grades because their payment was contingent on scoring at least at the Proficient level; there was no additional payment for the performance of classmates.

TABLE 10
 PERCENTAGE RECEIVING PAYMENT AND INCENTIVE PAYMENT COST (Pesos): YEAR 2

	Treatment 3	Treatment 1	Treatment 2
% of students receiving payment:			
Grade 10:			
For own performance	64.6	58.8	
For class performance	100.0	...	
Grade 11:			
For own performance	41.3	38.8	
For class performance	99.4	...	
Grade 12:			
For own performance	17.3	15.3	
Mean student payment:			
Grade 10:			
For own performance	2,991	2,515	
For class performance	1,108	...	
Total	4,099	2,515	
Grade 11:			
For own performance	2,679	2,541	
For class performance	861	...	
Total	3,540	2,541	
Grade 12:			
For own performance	991	915	
% of teachers receiving payment:			
For own performance	97.2		93.5
For class performance	100.0		...
Mean math teacher payment (FTE):			
For own performance	15,330		6,332
For other teacher performance	3,779		...
Total	19,109		6,332
Mean nonmath teacher and assistant director payments:			
Payment per FTE	3,872		...
Mean director payments:			
Payment per director	7,744		...
Incentive payment cost per student	3,303	2,080	43
Amount controls would receive	1,643	1,163	44
% of total	49.7	55.9	100

Full-time mathematics teachers received, on average, \$15,330 for the performance of the students in their classes and an additional \$3,779 for the performance of the students taught by the other mathematics teachers in the school. Each full-time-equivalent nonmathematics teacher or assistant principal received \$3,872 and the principal received \$7,744.

Taken together, the cost per student was \$3,303 (about US\$275), about 15 percent of the current per-student expenditure in federal upper-secondary schools.

Treatment 1.—Tenth-grade students in T1 received \$2,515, eleventh-grade students \$2,541, and twelfth-grade students \$915. The cost per student was \$2,080, about 10 percent of the current per-student annual expenditure.

Treatment 2.—Full-time mathematics teachers received, on average, \$6,332 for the performance of the students in their classes. The cost per student was \$43.

A part of the incentive payments of the ALI program is “wasted” in the sense that some students were paid for results that would have been achieved without the incentive. It is possible to compute the magnitude of the waste on the basis of the payments the C students would have received under the different treatments using the transition rates between the ALI score categories for the C students. For example, of those C students who scored at the Pre-Basic level on the ninth-grade ENLACE, 24.0 percent jumped at least one category on the tenth-grade ALI test, thus earning a reward in the hypothetical case that they had been in the T1 or T3 treatment. Of course, had they been in one of those treatments, our estimates indicate that many more would have jumped to a higher category, and some of those who jumped one category would have instead jumped two or three. Nevertheless, these students would have been rewarded, in at least some part, for an outcome that would have been achieved without the treatment. As seen in table 10, the waste is substantial, amounting to 49.7 percent of the cost of T3 and 55.9 percent of the cost of T1. Although one would prefer a program in which the waste was small, what matters for assessing the efficacy of the program is a comparison of the program’s benefits relative to its costs.

V. Conclusions

This paper evaluates the effect of the ALI pilot program that randomly assigned 88 Mexican high schools to three treatment groups and to a control group to measure the effectiveness of three alternative performance incentive schemes on mathematics tests scores: (T1) incentives for students only, (T2) incentives for teachers only, and (T3) individual and group incentives for students, teachers, and administrators. Previous studies used randomized trials to analyze effects of student-only or teacher-only incentive schemes, but the ALI program is the first experimental study to combine student and teacher incentives.

An analysis of student test scores finds very large treatment impacts for treatments T1 and T3 that pay incentives to students. Further examination of the test score answer booklets revealed that the student incentive payments also induced a higher rate of cheating in the form of student copying. We proposed an adjustment procedure to account for the higher rates of copying estimated to have occurred in these treatment groups versus the control group. Even after making adjustments for copying, we find substantial program impact estimates on mathematics achievement.

Our general findings, based on the copying adjustment, can be summarized as follows: (i) Providing the ALI incentives to students alone increased mathematics test scores by 0.2–0.3 of a standard deviation, depending on the grade and year. (ii) Providing ALI incentives to teachers alone did not affect test scores. (iii) Providing ALI incentives to students and mathematics teachers (both for their own performance and for that of their peers and for other teachers and school administrators) led to the largest treatment effect estimates, increasing test scores by 0.3–0.6 of a standard deviation. (iv) Analysis of treatment effects conditional on initial test score performance categories shows that there are positive impacts across the entire baseline test score distribution. Finally, our sensitivity analysis explored the robustness of the impact estimates to allowing for differential test-taking effort between the treatment and control groups. Assuming, for example, that all of the measured T1 impact in the first year is due to test-taking effort, we still find large impacts of treatment T3.

The fact that incentive payments in the ALI program are performance based means that the cost of any treatment increases with its success. The per-student cost of T3, the most successful treatment, is 50 percent more than that of T1. However, T3 and T1 differed in many ways, with T1 providing rewards based on individual student performance only and T3 rewarding both individual and group performance for students, teachers, and administrators. We can isolate the cost of each of the components of T3, but the experimental design does not identify the relative contribution of each component to T3's greater success. Similarly, the cost of T2 was less than 1.5 percent of the cost of T3, but T2 was found to be ineffective.

A limitation of the ALI experiment, as of all experiments, is that we can learn the impact only of the treatments that were tried and only in the population studied. Extrapolations to other populations would have to account for the fact that the subset of Mexican federal high schools studied was selective and that there were preexisting programs that provided school attendance subsidies. To learn about effects of other treatments, such as variations in the performance incentive schedules, requires either additional experiments or the development and estimation of behavioral models that can be used to extrapolate to other hypothetical treatments or to other populations.³⁷

³⁷ See Todd and Wolpin (2006) for an example of the use of experiments and modeling for the purpose of performing counterfactual experiments.

Appendix A

As discussed in Bloom (2005), carrying out power calculations before implementing an experiment requires some guesswork, because one needs an idea of the variance of the outcome and some estimate of the intraclass correlation. Once the experiment has been implemented, however, these preliminary estimates can be refined using the experimental data. As an easy way of understanding the power of an experimental design, Bloom introduces the notion of a minimum detectable effect size (MDES) for a given level of power (k), significance level (α), group size (n), number of groups (J), intraclass correlation (ρ), and a proportion of subjects allocated to the treatment group (P). The MDES is the program impact (divided by the standard deviation of the outcome for the target population) that can be detected with the specified parameters (power, significance level) under a particular experiment. The formula for the MDES is

$$\text{MDES} = \frac{M_{J-2}}{J^{1/2}} \left(\rho + \frac{1 - \rho}{n} \right)^{1/2} \left[\frac{1}{P(1 - P)} \right]^{1/2}, \quad (\text{A1})$$

where M_{J-2} varies depending on whether a one-tailed or two-tailed t -test is conducted; $M_{J-2} = t_{1-k} + t_{\alpha/2}$ for a two-tailed test and $M_{J-2} = t_{1-k} + t_{\alpha}$ for a one-tailed test. In the case of ALI, the average group size (n) is 200 students per school per grade. The number of groups is 28 for the control group and 20 for the treatment group, for a total (J) of 48 and a proportion (P) equal to 0.4167. For the calculations below, we require a power of 0.8 and assume an intraclass correlation (ρ) equal to .12, which is the estimated correlation using the full sample (all 88 schools) for the baseline ninth-grade mathematics ENLACE (for each grade). With these parameters, the MDESs in a one-tailed test are 0.34, 0.26, and 0.22 for the three critical values $\alpha = .01$, .05, and .10. Similarly, at those same critical values, the MDESs are 0.37, 0.30, and 0.26 for a two-tailed test.

As noted above, the required sample size can be substantially smaller if the analysis is conditioned on baseline school-level covariates. In the case of ALI, lagged school-level baseline student test scores are powerful predictors of current test scores, and including them in a treatment impact regression substantially reduces the estimated intraclass correlation. Therefore, the MDESs reported above are conservative.

Appendix B

A. Calculation of the Student Bonus

For tenth-grade students, s is the semester, $s = 1, 2$; c_{ks} is the mathematics class designation k in semester s ; $n_{\{ij\}}^{c_{ks}}$ is the number of students in class k in semester s with test scores that place them in the ij th cell of the tenth-grade performance bonus schedule; $b_{\{ij\}}$ is the bonus payment to a student in the ij th cell of the tenth-grade performance bonus schedule; and $B_{\{ij\}}^m$ is the total bonus payment to a student with a test score in the $\{ij\}$ th cell of the tenth-grade performance bonus schedule who was in class m in semester 1 ($s = 1$) and class n in semester 2 ($s = 2$).

Bonus formula:

$$B_{\{ij\}}^{mn} = b_{\{ij\}} + \left[\sum_{\{ij\}} \frac{1}{2} (n_{\{ij\}}^{c_{m1}} b_{\{ij\}} + n_{\{ij\}}^{c_{m2}} b_{\{ij\}}) \right] \times .01.$$

For eleventh-grade students, s is the semester, $s = 3, 4$; c_{ks} is the mathematics class designation k in semester s ; $n_{\{ij\}}^{c_{ks}}$ is the number of students in class k in semester s with test scores that place them in the ij th cell of the eleventh-grade performance bonus schedule; $b_{\{ij\}}$ is the bonus payment to a student in the ij th cell of the eleventh-grade performance bonus schedule; and $B_{\{ij\}}^{mn}$ is the total bonus payment to a student with a test score in the $\{ij\}$ th cell of the eleventh-grade performance bonus schedule who was in class m in semester 1 ($s = 3$) and class n in semester 2 ($s = 4$).

Bonus formula:

$$B_{\{ij\}}^{mn} = b_{\{ij\}} + \left[\sum_{\{ij\}} \frac{1}{2} (n_{\{ij\}}^{c_{m3}} b_{\{ij\}} + n_{\{ij\}}^{c_{m4}} b_{\{ij\}}) \right] \times .01.$$

For twelfth-grade students, B is the bonus payment.

Bonus formula:

$$B = \begin{cases} 5,000 & \text{if test score is Proficient} \\ 10,000 & \text{if test score is Advanced.} \end{cases}$$

B. Calculation of Mathematics Teacher Bonus

With s as the semester, $s = 1, 2$ for tenth grade, $s = 3, 4$ for eleventh grade, and $s = 5, 6$ for twelfth grade; $n_{\{ij\}k}^s$ is the number of teacher k 's students in semester s with test scores that place them in the ij th cell of the performance bonus schedule; $b_{\{ij\}}^s$ is the bonus payment per student in the ij th cell of the performance bonus schedule for semester s , where $b_{\{ij\}}^1 = b_{\{ij\}}^2$, $b_{\{ij\}}^3 = b_{\{ij\}}^4$, and $b_{\{ij\}}^5 = b_{\{ij\}}^6$; f_k is teacher k 's full-time equivalence (FTE) status, which is the number of math classes taught over an entire year (both semesters) divided by 10;³⁸ M is the total number of mathematics teachers over the entire year regardless of the number of classes taught; $F = \sum_{m=1}^M f_m$ is the total number of full-time-equivalent mathematics teachers in the school; B_k^1 is the bonus payment to teacher k for the performance of teacher k 's students; B_k^2 is the bonus payment to teacher k for the performance of students schoolwide; B_k is the total bonus payment to teacher k ; and $F - f_k$ is the number of full-time-equivalent teachers subtracting off teacher k 's FTE status.

Bonus formula:

$$B_k^1 = \frac{1}{2} \sum_{s=1}^{s=6} \sum_{\{ij\}} n_{\{ij\}k}^s b_{\{ij\}}^s,$$

$$B_k^2 = \left(\frac{1}{F - f_k} \sum_{m=1, m \neq k}^M B_m^1 \right) \times f_k \times .25,$$

$$B_k = B_k^1 + B_k^2.$$

³⁸ For example, a teacher who taught five classes in the fall term and five classes in the spring term will have an FTE status of one ($f_k = 1$), while a teacher who taught four and three classes in the two terms will have an FTE status of .7 ($f_k = .7$). It is possible to have an FTE status that is greater than one.

C. Calculation of Nonmathematics Teacher, Director, and Associate Director Bonus

With s as the semester, $s = 1, 2$ for tenth grade, $s = 3, 4$ for eleventh grade, and $s = 5, 6$ for twelfth grade; f_k is the nonmathematics teacher or administrator k 's FTE status; M is the total number of mathematics teachers over the entire year regardless of the number of classes taught; $F = \sum_{m=1}^M f_m$ is the total number of full-time-equivalent mathematics teachers in the school; B_m^1 is the bonus payment to mathematics teacher m for the performance of teacher m 's students; B_k is the total bonus payment to nonmathematics teacher or administrator k ; and F is the number of full-time-equivalent mathematics teachers.

Bonus formulas: For a director:

$$B_k = \left(\frac{1}{F} \sum_{m=1}^M B_m^1 \right) \times f_k \times .50.$$

For a nonmathematics teacher and associate administrator:

$$B_k = \left(\frac{1}{F} \sum_{m=1}^M B_m^1 \right) \times f_k \times .25.$$

Appendix C

TABLE C1
AVERAGE TREATMENT EFFECTS (ATE) WITH AND WITHOUT ADJUSTMENTS FOR COPIERS: ALL PROGRAM YEARS

ATE	YEAR 1: ACADEMIC YEAR 2008/9			YEAR 2: ACADEMIC YEAR: 2009/10			YEAR 3: ACADEMIC YEAR: 2010/11			
	T1	T2	T3	T1	T2	T3	T1	T2	T3	
With copying adjustment:*										
10th grade	15.3 (5.02)	2008/9 1.20 (6.32) 2007/8	31.4 (5.84)	27.2 (4.63)	2009/10 -0.41 (5.68) 2008/9	43.9 (6.62)	28.2 (4.67)	2010/11 12.1 (5.62) 2009/10	57.0 (9.49)	
11th grade	9.21 (5.62)	-5.18 (5.36) 2006/7	14.5 (6.93)	15.7 (5.72)	-1.40 (6.42) 2007/8	33.3 (8.20)	14.2 (4.63)	-2.55 (4.37) 2008/9	25.2 (5.90)	
12th grade	10.9 (6.86)	7.13 (6.47)	31.5 (6.48)	12.3 (6.06)	-6.19 (6.56)	25.1 (6.68)	13.0 (8.08)	0.84 (7.94)	41.2 (14.9)	
No copying adjustment:										
10th grade	18.5 (5.02)	2008/9 1.11 (5.35) 2007/8	32.3 (6.18)	32.4 (5.24)	2009/10 0.31 (5.74) 2008/9	54.7 (11.1)	41.5 (6.25)	2010/11 15.9 (6.16) 2009/10	83.4 (16.9)	
11th grade	22.4 (7.22)	-2.98 (6.74) 2006/7	27.8 (9.93)	55.5 (7.51)	6.17 (6.91) 2007/8	67.4 (12.7)	51.3 (9.05)	-1.36 (8.89) 2008/9	106.4 (25.6)	
12th grade	9.73 (7.04)	4.73 (6.62)	29.3 (6.67)	36.0 (7.32)	-1.81 (6.30)	44.6 (7.99)	42.3 (8.15)	7.33 (7.98)	90.2 (21.3)	

NOTE.—The table is based on regressions that include school averages of ninth-grade ENLACE in each of the three grades in the program year in which the test was taken, the school average twelfth-grade ENLACE for the academic year 2007/8, and region dummies and state dummies for states in which there was at least one school in all treatment groups. Standard errors (in parentheses) account for clustering at the school level. The sample is based on students who took an ALI test in each year of possible enrollment in program years. There are 11,896 observations in the 2006/7 cohort, 11,385 in the 2007/8 cohort, 11,314 in the 2008/9 cohort, 13,778 in the 2009/10 cohort, and 17,813 in the 2010/11 cohort.

* Copiers are assigned the 25th percentile score of control group noncheaters.

References

- Angrist, Joshua, and Victor Lavy. 2009. "The Effects of High Stakes High School Achievement Awards: Evidence from a Randomized Trial." *A.E.R.* 99:1384–1414.
- Barlevy, Gadi, and Derek Neal. 2011. "Pay for Percentile." Working Paper no. 17194, NBER, Cambridge, MA.
- Barrow, Lisa, and Cecilia E. Rouse. 2013. "Financial Incentives and Educational Investment: The Impact of Performance-Based Scholarships on Student Time Use." Working Paper no. 19351, NBER, Cambridge, MA.
- Bloom, Howard S. 2005. "Randomizing Groups to Evaluate Place-Based Programs." In *Learning More from Social Experiments: Evolving Analytic Approaches*, edited by Howard S. Bloom. New York: Russell Sage Found.
- Cizek, Gregory J. 1999. *Cheating on Tests: How to Do It, Detect It and Prevent It*. Mahwah, NJ: Erlbaum.
- Cizek, Gregory J., Michael B. Bunch, and Heather Koons. 2004. "Setting Performance Standards: Contemporary Methods." *Educ. Measurement: Issues and Practice* 23 (4): 31–50.
- Cox, D. R., and N. Reid. 2000. *The Theory of the Design of Experiments*. Boca Raton, FL: Chapman & Hall.
- Fryer, Roland. 2010. "Financial Incentives and Student Achievement: Evidence from Randomized Trials." Working Paper no. 15898, NBER, Cambridge, MA.
- Glewwe, Paul, Nauman Ilias, and Michael Kremer. 2010. "Teacher Incentives." *American Econ. J.: Appl. Econ.* 2:205–27.
- Hanushek, Eric A. 1986. "The Economics of Schooling: Production and Efficiency in Public Schools." *J. Econ. Literature* 24:1141–77.
- . 2003. "The Failure of Input-Based Schooling Policies." *Econ. J.* 113:F64–F98.
- Heckman, James J. 2000. "Causal Parameters and Policy Analysis in Economics: A Twentieth Century Retrospective." *Q.J.E.* 115 (1): 45–97.
- Heckman, James J., and Jeffrey A. Smith. 1995. "Assessing the Case for Social Experiments." *J. Econ. Perspectives* 9 (2): 85–110.
- Heckman, James J., and Sergio Urzua. 2010. "Comparing IV with Structural Models: What Simple IV Can and Cannot Identify." *J. Econometrics* 156 (1): 27–37.
- Hedges, L. V., R. Laine, and R. Greenwald. 1994. "Does Money Matter: A Meta-Analysis of Studies of the Effects of Differential School Inputs on Student Outcomes." *Educ. Researcher* 23:5–14.
- Jackson, C. Kirabo. 2010. "A Little Now for a Lot Later: A Look at a Texas Advanced Placement Incentive Program." *J. Human Resources* 45:591–639.
- Jacob, Brian, and Steven D. Levitt. 2003. "Rotten Apples: An Investigation of the Prevalence and Predictors of Teacher Cheating." *Q.J.E.* 118:843–77.
- Kremer, Michael, Edward Miguel, and Rebecca Thornton. 2009. "Incentives to Learn." *Rev. Econ. and Statis.* 91:437–56.
- Krueger, Alan B. 2003. "Economic Considerations and Class Size." *Econ. J.* 113: F34–F63.
- Levitt, Steven D., John A. List, and Sally Sadoff. 2010. "The Effect of Performance-Based Incentives on Educational Achievement: Evidence from a Randomized Experiment." Working paper, Univ. Chicago.
- Muralidharan, Karthik, and Venkatesh Sundararaman. 2011. "Teacher Incentives in Developing Countries: Experimental Evidence from India." *J.P.E.* 119:39–77.

- Neal, Derek, and Diane Whitmore Schanzenbach. 2010. "Left Behind by Design: Proficiency Counts and Test-Based Accountability." *Rev. Econ. and Statis.* 92:263–83.
- Springer, Matthew G., et al. 2010. "Teacher Pay for Performance: Experimental Evidence from the Project on Incentives in Teaching." Manuscript, Nat. Center Performance Incentives, Vanderbilt Univ.
- Todd, Petra E., and Kenneth I. Wolpin. 2003. "On the Specification and Estimation of the Production Function for Cognitive Achievement." *Econ. J.* 113: F3–F33.
- . 2006. "Assessing the Impact of a School Subsidy Program in Mexico: Using a Social Experiment to Validate a Dynamic Behavioral Model of Child Schooling and Fertility." *A.E.R.* 96:1384–1417.
- . 2012. "Estimating a Coordination Game within the Classroom." Manuscript, Univ. Pennsylvania.
- Wesolowsky, George. 2000. "Detecting Excessive Similarity in Answers on Multiple Choice Exams." *J. Appl. Statis.* 22:908–21.