

NBER WORKING PAPER SERIES

DOES TEACHER TRAINING AFFECT
PUPIL LEARNING? EVIDENCE FROM
MATCHED COMPARISONS IN
JERUSALEM PUBLIC SCHOOLS

Joshua D. Angrist
Victor Lavy

Working Paper 6781
<http://www.nber.org/papers/w6781>

NATIONAL BUREAU OF ECONOMIC RESEARCH
1050 Massachusetts Avenue
Cambridge, MA 02138
November 1998

The Jerusalem Schools Authority funded this research and provided the data. Thanks go to Jon Guryan and Xuanhui Ng for outstanding research assistance, and to the teachers, principals, and supervisors at the schools involved for their cooperation. Special thanks also go to Mr. Avi Sela, the Project Manager in Jerusalem, and his staff for their assistance and support. The views expressed here are those of the author and do not reflect those of the National Bureau of Economic Research.

© 1998 by Joshua D. Angrist and Victor Lavy. All rights reserved. Short sections of text, not to exceed two paragraphs, may be quoted without explicit permission provided that full credit, including © notice, is given to the source.

Does Teacher Training Affect Pupil Learning?
Evidence from Matched Comparisons in
Jerusalem Public Schools
Joshua D. Angrist and Victor Lavy
NBER Working Paper No. 6781
November 1998
JEL No. I21, I28, J24

ABSTRACT

The relationship between teachers' characteristics and their pupils' achievement has been the subject of many studies. Most of this research focuses on the impact of teacher salaries, experience, and measures of teachers' pre-service training such as educational background. The effect of on-the-job or in-service training has received much less attention. In this paper, we estimate the effect of in-service teacher training on children's reading and mathematics achievement in Jerusalem elementary schools. The training was based on pedagogical methods developed in US schools. Our research uses a matched-comparison design which exploits the fact that only a few schools received extra funds for training. Differences-in-differences, regression, and nonparametric matching estimates are reported. The results suggest that the training received by teachers in the non-religious branch of the Jerusalem school system led to an improvement in their pupils' test scores. The estimates for religious schools are not clear cut, but this may be because the training program in religious schools started later and was implemented on a smaller scale. The estimates for non-religious schools suggest that, at least in this case, teacher training provided a less costly means of increasing test scores than reducing class size or adding school hours.

Joshua D. Angrist
Department of Economics
Massachusetts Institute of Technology
E52-353
50 Memorial Drive
Cambridge, MA 02139-4307
and NBER
angrist@mit.edu

Victor Lavy
Department of Economics
Hebrew University
Mt. Scopus
Jerusalem 91905
ISRAEL
msvictor@mscc.huji.ac.il

I. Introduction

The question of how teacher characteristics affect pupil learning has long been of concern to economists, educators, and parents. The interest in this question is more than academic since many parents would be willing to pay for better teacher qualifications if they knew that this would cause their children to learn more. The widespread interest in teacher characteristics has led to numerous studies that attempt to estimate the effect of teachers' characteristics on their pupils' test scores. The characteristics that have been most studied are teachers' educational background, years of teaching experience, and salaries. This research, while plentiful, has failed to produce a consensus regarding the causal impact of teacher characteristics, and the question of whether ostensibly more-qualified teachers produce better students remains open (see, e.g., Hanushek, 1986, and Hedges, Laine, and Greenwald, 1994 for contrasting assessments of the literature).¹

While many econometric and statistical studies have focused on the consequences of teachers' general skills as captured by education or experience, a second strand of research has looked at the question of how specific episodes of on-the-job or in-service training affect pupil performance. On-the-job and in-service training are probably at least as important as the more widely-studied pre-service training or general experience measures. For example, one survey of research on teacher training in developing countries concludes that "Pre-service training is essential to teach subject matter. In-service training is essential to teach teaching skills." (Farrel and Oliveira, 1993, p. 15). In spite of the importance and potential cost-effectiveness of in-service training, the literature on the effects of this sort of training is sparse. In a recent study that is similar in spirit to ours, Bressoux (1996) looks at the impact of in-service training on novice teachers in France using a quasi-experimental research design. Bressoux observes that most of the work on this topic to date falls into the category of "process evaluation" that ignores outcomes, or involves comparisons of different training programs

¹The best-known study of teacher characteristics and other school resources is probably the Coleman report (1966). Another early example is Murnane (1975). A recent US study of the effect of teachers' education and experience is Ehrenberg and Brewer (1994). Another recent study is Behrman, Khan, Ross, and Sabot (1997), who look at teacher qualifications in Pakistan. Researchers have also studied the relationship between teacher's salaries and their pupils' earnings as adults. See Welch (1966) and, more recently, Card and Krueger (1992).

with no untrained control group.²

This paper contributes to the literature on effects of in-service training by presenting an evaluation of teacher training in Jerusalem elementary schools, with the purpose of estimating the causal effect of the program on pupils' test scores. The training program we study was designed to improve the teaching of language skills and mathematics, and involves pedagogical methods developed in US schools. Most of the existing research on teacher training, like studies of the effects of other school resources, is complicated by the fact that school inputs are not exogenously determined.³ In an attempt to overcome the methodological difficulties inherent in an evaluation study of this type, our research design exploits the fact that in 1995 a handful of public schools in Jerusalem received a special infusion of funds that were primarily earmarked for teachers' on-the-job training. This program presents an unusual research opportunity because, even though the intervention was not allocated using experimental random assignment, the Jerusalem intervention can be studied with the aid of a matched group of students from schools not subject to the intervention. Considerable information is available on the students enrolled at the affected schools both before and after the intervention began. Similar information is also available for a group of comparison schools in adjacent neighborhoods and elsewhere in the city, so these schools can play the role of a control group.

Our research strategy uses a variety of statistical methods to estimate the causal effect of teachers' on-the-job training on their students' test scores. All of the methods suggest that the training received by teachers who work in the non-religious branch of the Jerusalem school system led to an improvement in their pupils'

²Perhaps surprisingly, there seems to have been more research on the impact of teacher training in developing countries than in developed countries. It should also be noted that economists have estimated the effects of teacher training in their own discipline. See, for example, Highsmith (1974) and Schober (1984). In-service training for graduate teaching assistants has been studied as well (see, e.g., Bray and Howard, 1980, or Carroll, 1980). Recent research on the productivity consequences of on-the-job and in-service training for workers other than teachers includes Bartel (1995) and Krueger and Rouse (1998).

³An exception is Dildy (1982), who reports the results of a randomized trial to evaluate the effects of teacher training in Arkansas elementary schools. The potential for bias in naive comparisons was highlighted by Lavy (1995), who noted that the apparent negative correlation between school inputs and pupil achievement in Israel at least partly due to the fact that measures of socioeconomic disadvantage are used to decide which schools get the most inputs.

test scores. In contrast with the estimates for non-religious schools, however, the estimates for religious schools are sensitive to the details of model specification. The absence of a clear effect in religious schools may be because the intervention in religious schools started later and was implemented on a smaller scale.

The most plausible estimates for non-religious schools suggest effect sizes in the range of .2-.4 standard deviations, which is noteworthy given the modest cost of the intervention. Given the recent concern with the level of teacher training in subject matter (see, e.g., Boston Globe, 1998; New York Times, 1998), it is also noteworthy that the Jerusalem training was in pedagogy and not subject content. In an attempt to assess the economic value of the training program, we compare the treatment effect and costs of training to the effect size and costs of alternative school improvement strategies involving reductions in class size and lengthening the school day. This analysis suggests that teacher training may provide a less expensive strategy for raising test scores than reducing class size or adding school hours.

II. The Intervention and the Data

A. Teacher training in Jerusalem

Beginning in 1995, the *30 Towns* intervention led to a major increase in resources for schools in two neighborhoods in North Jerusalem (Neve Yaakov and Pisgat Zeev) as well as in other towns and cities in Israel. The intervention essentially came in the form of budget increases, though in some cities the money was earmarked for certain uses. In Jerusalem, the *30 Towns* money was spent mainly on teacher training, though some resources were used for reorganization and for programs for children considered to have special needs, including immigrants.⁴ The extra training was provided in the school on a weekly basis by outside instructors who focused on improving instruction techniques for Hebrew language skills (which we call “reading”), mathematics, and English. The goals of the intervention were to increase the pass rate from grade to grade, to

⁴The total amount involved was about \$960,000 per year for elementary schools.

increase scores on achievement tests, and (more vaguely) to improve the school climate. Our study is limited to the effect of the program on reading and math achievement since there is no test score data for English.⁵

Table 1 summarizes the program impact on school inputs for our sample of treatment and control schools using data from the JSA (described in greater detail below). Public schools in Jerusalem and elsewhere in Israel are separated into religious and non-religious systems, so we discuss the effect of the program on the two types of schools separately. In Jerusalem, a total of 10 elementary schools were affected, 7 non-religious and 3 religious, though we drop one of the religious schools from the analysis since it opened after the intervention began. Panel A shows the data for non-religious schools (7 treated, 6 control) and panel B shows the corresponding data for religious schools (2 treated, 5 control).⁶ The information is shown for three dates: 1994, which is before the training program started, 1995, which is the year training began, and 1996, which is a year in which training inputs continued to be at the new higher level in treated schools. Training is measured as the total hours of instruction received by all teachers in each school per week, on average over the course of the relevant school year, and is recorded separately for reading and math.⁷

Panel A of the table shows a marked increase in hours of training provided to teachers in the non-religious treated schools between 1994 and 1995, and again between 1995 and 1996. For example, in 1994, treatment schools received an average of 1.2 training hours per week in math. This increased to 6.7 hours of training in 1995 and to 12.5 hours of training in 1996. In contrast, training hours changed little in the control schools. Thus, by 1996, the treatment schools received 10.5 more math training hours per week than the

⁵The training interventions varied somewhat from school to school, but all involved a mixture of similar counseling and feedback sessions for teachers, changes in the organization of class time, and the use of instructional aids. There was generally an emphasis on developing strategies for helping children who are doing poorly, including immigrants and those with learning disabilities. The Math teachers received training based on a modern “Humanistic Mathematics” philosophy of teaching (see White, 1987 and 1993). The reading teachers received training based on the “Individualized Instruction” approach to schooling. See, e.g., Bishop (1971), for a description, and Bangert, Kulik, and Kulik (1983) and Romberg (1985) for research summaries.

⁶The selection of control schools is discussed in the next section.

⁷The total school allocation of training hours is divided among the teachers in the school. On average, each school in our data has about 450 pupils and 12 classes (including all grades).

control schools. Similarly, the treatment schools actually received fewer hours of training in reading in 1994, but by 1996 the gap between treatment and control schools was 7.7 hours in favor of the treated schools.

After the data on training, the next row of the table shows expenditure on special projects per pupil by type of school (in nominal shekels), including the additional *30 Towns* money but excluding regular teachers' salaries. The changes in expenditure reflect the additional expenditure on training. In contrast with the data on training and expenditure, instruction hours per pupil (which reflects the number of teachers and their work hours) and class size were essentially unchanged in both treated and control schools.

Panel B repeats the analysis of inputs for religious schools, though this analysis is limited by the fact that we were unable to obtain data on teacher training in the religious control schools. The data on training inputs for religious schools actually show a decline in training hours between 1994 and 1995, with a modest increase in 1996. This accords with the fact that in religious schools, the *30 Towns* money was not used until September of 1995, i.e., in the 1996 school year (see Angrist and Lavy, 1998a). Thus, for religious schools, the 1995 school year is a pre-treatment year. In addition to starting later, the intervention was also less intense in the religious schools than in the non-religious schools. This can be seen in the expenditure data, which show no differential increase in the religious treated schools.

B. Test score data and the choice of control schools

Table 2 summarizes information on test scores and other variables available in each of three years: in 1994, before the intervention began; in 1995, when the intervention first started; and in 1996. Although the intervention was school-wide, the only pupils for whom we have data on test scores both before and after the intervention began is the cohort of fourth graders enrolled in the 1993-94 school year. Test score data are available for these children when they were finishing fourth grade in June 1994, as fifth graders in June 1995, and as sixth graders in June 1996. The 1994 and 1995 data were collected in a routine JSA testing program in which all schools in the treated area as well as 20 other schools participated. The study is limited to the

1994 fourth graders in these schools because test score data for 1996 are available only for sixth graders.

The training effort in non-religious schools began in January 1995, so the June 1995 test scores provide an early indicator of program effectiveness in these schools. But training in religious schools did not begin until September 1995, so the June 1995 data provide an additional pre-treatment observation for religious schools. To obtain information on test scores in 1996, which provides follow-up information 1 ½ years after training began in non-religious schools and 1 year after training began in religious schools, we used data from two sources. The Hebrew test scores come from a middle-school placement exam given to all sixth graders in Jerusalem. Since the placement test does not cover math, at our request, special math tests were developed and administered in treatment and control schools for the purposes of evaluating *30 Towns*. These tests were designed to be similar in style to the placement test in reading. To minimize differences between the tests across years, we analyze standardized scores with mean zero and unit standard deviation (except in the presentation of descriptive statistics in Table 3)

The evaluation strategy compares the population of fourth grade pupils who were enrolled in treated schools in 1994 with the population of fourth graders who attended a sample of control schools. The control schools were selected for comparability and for practical reasons related to data. First, the control schools had to be in the group of 20 non-treated schools that were tested by the JSA in 1994 and 1995. Second, to maximize comparability and to encourage the cooperation of school principals with our 1996 testing effort, we chose control schools that report to the same district supervisors as the treatment schools. A total of 11 out of the 20 candidate control schools met the criteria for inclusion in the study -- 6 non-religious schools, and 5 religious schools.⁸

In June 1994, 406 fourth-grade pupils took the reading test in the non-religious treatment schools and 405 pupils took the reading test in the non-religious control schools. The corresponding figures for the math

⁸In addition to test score data in 1994, 1995, and 1996, the JSA also provided information on the characteristics of pupils and schools in our sample in 1994. School characteristics were assigned based on school identifiers at baseline. Two schools later merged, another later split, and a new school opened, but there was no crossover between treatment and control groups.

test are 428 treated test-takers and 420 control test-takers. This can be seen in Table 3, which reports average test scores and sample sizes in reading and math by treatment/control status, for 1994-96. About 90 percent of all pupils who were enrolled in the treated and control schools in 1994 took the tests. Of the original 406 pupils tested in reading in 1994 in treated schools, 364 were tested in 1995 and 301 were tested in 1996. Of the original 405 tested in reading in 1994 in control schools, 284 were tested in 1995 and 290 were tested in 1996. The sample sizes for math tests show a similar pattern of attrition. Note also that some pupils who were not tested in 1994 were tested in 1995 and 1996 since the tests covered whoever was enrolled at the time. The number of those tested in 1995 and 1996 who were also tested in earlier years is reported in the table, along with average scores for subsamples of pupils where we have repeated observations. A similar analysis for religious schools appears in Panel B of Table 3.

The proportion of pupils tested in 1994 who were also tested in 1996 ranges from .62 for math tests in control schools to .87 for math tests in treatment schools. The higher attrition rates in the control group are potentially a cause for concern. It should be noted, however, that most of the non-participation in follow-up waves results from entire classes not taking the test rather than from individual pupils missing the test. Attrition at the class level probably generates less bias than attrition by pupils. This claim is supported by the fact that in most cases, average scores in the entire 1994 cross section are similar to those in the longitudinal samples with information on scores in 1995 and 1996. Similarly, average scores in the 1996 cross section are similar to those in the longitudinal samples with scores in 1995 and 1994. The scores for different years and samples are shown in columns 1, 3, 5, and 7 in Table 3.

C. Pupil characteristics in treatment and control schools

Pupils in treated and control schools are similar along many dimensions, but pupils in treated schools have less educated parents. This can be seen in Table 4, which compares the characteristics of pupils in the treatment and control schools in 1994. The difference in parents' education is over 1.5 years in religious

schools, and almost 1.5 years for fathers' education in non-religious schools. This suggests that it may be important to control for pupils' family background when estimating the effect of teacher training. We use a variety of strategies for doing this: differencing test scores to eliminate time-invariant characteristics, control for pupils' characteristics using regression, and regression and matching methods to control for 1994 scores.

III. Evaluation Strategies

The ideal evaluation strategy for the 30 Towns intervention (or any similar intervention) would involve the random assignment of pupils to treatment and control groups, with those in the treatment group being taught by teachers who had received training. Random assignment would insure that pupils in the control group are indeed comparable to pupils in the treatment group, so that any difference between pupils in the two groups could be confidently attributed to the intervention. In the absence of random assignment, statistical methods must be used to control for differences between pupils in treated and non-treated schools. We use a variety of models for this purpose.

One simple model that has been used in numerous evaluation studies (see, e.g., Ashenfelter, 1978; Ashenfelter and Card, 1985) is based on the assumption that any differences between pupils in the treatment and control groups are fixed over time. In this case, repeated observations on the same pupils can be used to make the treatment and control groups comparable. Let D_{it} be a dummy variable that indicates treatment status and let $Y_{it}(0)$ denote the potential test score of any pupil if he or she were not exposed to the treatment. The fixed-effects model says that in the absence of the intervention, the potential test score of pupil i at time t can be written:

$$Y_{it}(0) = \phi_i + \delta_t + \epsilon_{it}, \tag{1a}$$

where D_{it} is assumed to be independent of the time-varying component, ϵ_{it} , although it may be correlated with the pupil-specific intercept, ϕ_i . The term δ_t is a period effect common to all pupils. Thus, while pupils in the treatment group may have lower scores than pupils in the control group, this difference is assumed to be due

to differences in family background or neighborhood characteristics that can be viewed as permanent *and* that have time-invariant effects on test scores.

We begin with a simple model that also has a constant treatment effect, so treated pupils are assumed to have test scores equal to $Y_{it}(0)$ plus the treatment effect:

$$Y_{it}(1) = Y_{it}(0) + \alpha. \tag{1b}$$

Equations (1a) and (1b) can therefore be used to write the observed test score of pupil i at time t as

$$Y_{it} = Y_{it}(0)[1 - D_{it}] + Y_{it}(1)D_{it} = \phi_i + \delta_t + \alpha D_{it} + \epsilon_{it}, \tag{2}$$

where ϵ_{it} is the error term in (1a), assumed to be uncorrelated with D_{it} . In this model, simple post-treatment comparisons of test scores by treatment status do not estimate the causal effect of treatment because the time-invariant characteristics of pupils in the treatment and control groups differ. Formally, this means that $E[\phi_i | D_{it}=1] \neq E[\phi_i | D_{it}=0]$. On the other hand, the assumptions of the model imply that the *change in test scores* in the treatment group can be compared to the *change in test scores* in the control group. Let $t=a$ to denote post-treatment scores (“after”) and let $t=b$ to denote the pre-treatment test scores (“before”). Note that $D_{it}=0$ for both treatment and control pupils. Then we have,

$$E[Y_{ia} - Y_{ib} | D_{ia}=1] - E[Y_{ia} - Y_{ib} | D_{ia}=0] = \alpha \tag{3}$$

The sample analog of equation (3) is called a differences-in-differences estimate of the program effect because it contrasts the change in test scores between treatment and control pupils.

The most important difference between pupils in treatment and control schools is that pupils in the treatment schools have lower pre-treatment test scores, on average, than pupils in the control schools. The differences-in-differences method controls for these pre-treatment differences by subtracting them from the post-treatment difference in scores. To see this, note that equation (3) can be re-arranged to read

$$(E[Y_{ia} | D_{ia}=1] - E[Y_{ia} | D_{ia}=0]) - (E[Y_{ib} | D_{ia}=1] - E[Y_{ib} | D_{ia}=0]) = \alpha.$$

The term $(E[Y_{ib} | D_{ia}=1] - E[Y_{ib} | D_{ia}=0])$ is the pre-treatment difference in scores by treatment status.

We also explore alternative strategies that control directly for pre-treatment score differences between

the treatment and control groups using regression or matching. The regression approach to controlling for pre-treatment score differences (and other covariates) is based on a model where lagged test scores determine future test scores and treatment is independent of potential scores conditional on lagged scores and other covariates, X_i . Formally, the model is

$$Y_{ia}(0) = X_i' \beta + \gamma Y_{ib} + \delta_a + \epsilon_{ia}, \quad (4)$$

where ϵ_{ia} is assumed to be independent of D_{it} . Maintaining the constant treatment effects assumption, we have

$$Y_{ia} = X_i' \beta + \gamma Y_{ib} + \delta_a + \alpha D_{ia} + \epsilon_{ia}, \quad (5)$$

This amounts to replacing ϕ_i in (2) with $X_i' \beta + \gamma Y_{ib}$. Estimates of α are produced by regressing post-treatment test scores on a constant, pre-treatment scores, X_i , and D_{ia} . It is important to emphasize that the causal interpretation of α in (5) turns on the assumption that once X_i and pre-treatment scores are held constant, the only reason for a post-treatment difference between treated and control pupils is, in fact, the treatment.

Equation (5) relies on a linear model to control for the covariates X_i and Y_{ib} . In addition to regression estimates, we also use a matching strategy to nonparametrically control for pre-treatment score differences. This is accomplished by dividing the pre-treatment (1994) test score distribution into four equal-sized groups (quartiles). We then compare the post-treatment (1996) scores of treatment and control pupils in each group. The final step in this method is to produce a single treatment-control contrast from the four groups by averaging across groups using the number of treated pupils in each group as weights. Assuming that the treatment is independent of potential outcomes conditional on past test scores, this matching strategy produces an estimate of the average “effect of treatment on the treated” in a model with heterogeneous treatment effects (see, e.g., Rubin, 1977, Card and Sullivan, 1988, or Angrist, 1998).

Finally, note that the evaluation strategies discussed in this section involve matching treatment and control pupils and not treatment and control schools. This is because pupils are the unit of observation; we would like pupils in the control group to be comparable to pupils in the treatment group. For some of the regression estimates, however, we limit our sample to treatment and control schools where school-level

average test scores are similar in the two groups. This is useful if there is something special about the treatment schools (“random school effects”) that cannot be captured by matching pupils on personal characteristics and previous test scores. The cost of school-matching is that information on many pupils is discarded.

IV. Empirical Results

A. Differences-in-differences estimates

Table 5 presents differences-in-differences estimates of the effect of the intervention. The row labeled “control difference” is an estimate of the change in test scores for pupils in control schools, i.e., the sample analog of $E[Y_{it} - Y_{it'} | D_{it} = 0]$. The row labeled “difference-in-differences” is the sample analog of $E[Y_{it} - Y_{it'} | D_{it} = 1] - E[Y_{it} - Y_{it'} | D_{it} = 0]$. Test score growth is measured between 1995 and 1994, between 1996 and 1994, and between 1996 and 1994. The estimates in each column of the table were computed by regressing the change in scores on a constant and a dummy for pupils in the treatment group. Two sets of standard errors are presented: those based on the usual regression assumptions, and standard errors that allow for within-school correlation in test scores using the formula in Moulton (1986).

The results for non-religious schools show no significant change in test scores in the control schools between 1994 and 1996. For example, reading scores fell by .06 and math scores fell by .05, but neither of these changes is significantly different from zero. In contrast, test scores increased considerably and significantly in the treatment schools, both in absolute terms and relative to the control schools. Thus, the differences-in-differences estimate based on 1994-to-1996 changes is .62 for reading scores and .46 for math scores. This implies the program increased test scores by roughly half a standard deviation for both reading and math. The estimates are smaller when 1994-to-1995 changes are used, especially for math scores. On the other hand, math scores in treatment schools increased between 1994 and 1995 and between 1995 and 1996, while reading scores increased only between 1994 and 1995.

The differences-in-differences estimates for non-religious schools are significant whether conventional standard errors or Moulton standard errors are used to measure precision. It is interesting to note, however, that the Moulton standard errors are 2-3 times larger than the conventional standard errors. This reflects the fact that even though the within-school correlation between pupils' test scores is not very large (on the order of .1), the model being estimated has features that cause a small amount of within-group correlation to have a large effect on regression standard errors. First, the regressor (treatment status) is fixed within groups. Second, there are not many groups (schools) so the group size is large relative to the sample. Moulton shows that when the regressor of interest is fixed within groups and the groups are large, conventional standard errors may seriously overestimate the precision of regression estimates.

Estimates for religious schools are reported in Panel B of Table 5. The data for religious schools do not include reading scores for 1996.⁹ Between 1994 and 1995, test scores in the religious treated schools fell sharply, both in absolute terms and relative to control schools. This can be seen in columns 1 and 4 of the table. As noted earlier, the treatment in the religious schools started only in September of 1995 (i.e., the beginning of the 1996 school year), so the 1994-95 score decline in treated schools is probably unrelated to the intervention. A consequence of the decline in pre-treatment years is that the differences-in-differences estimates of effects on math scores are positive when 1995 is taken as the base year and negative when 1994 is taken as the base year (see columns 5 and 6). Sensitivity to the choice of base year implies that the assumptions underlying the differences-in-differences strategy are not satisfied (Ashenfelter and Card, 1985). In fact, the results for religious schools seem generally less reliable and harder to interpret than those for non-religious schools. The treated sample includes less than 100 pupils, from only two schools. Moreover, as noted earlier, we do not have data on inputs in religious control schools, while the data on treated schools actually show a decline in teacher training between 1994 and 1995.

⁹Since most religious schools combine primary and middle grades, they do not use the reading test that is used in secular schools to place students in the transition from primary to middle school.

B. Regression estimates

Comparisons of means by treatment status show that pupils in treated schools had lower test scores than pupils in control schools in 1994, before treatment began. Thus, pupils in the two sets of schools are not directly comparable. The differences-in-differences strategy is motivated by a model where the difference between pupils in treatment and control schools is captured by the fixed effect in equation (1a), ϕ_i . Differences-in-differences estimates do not have a causal interpretation, however, if the lower pre-treatment scores in treatment schools are due to *temporarily* low scores as opposed to permanent differences. In that case, test-score *growth* in the treatment schools will tend to be larger than test-score growth in the control schools regardless of the program effect, generating a spurious positive estimate of the treatment effect.¹⁰

To avoid this sort of bias, the regression approach discussed in Section III replaces the fixed-effects assumption with the assumption that pupils in treatment and control schools are comparable conditional on past test scores and other observed covariates. The resulting estimates are reported in Table 6. The first column in the table shows a specification that includes only the treatment dummy as an explanatory variable. The estimates reported in column 2 are from a model where a limited set of student characteristics (year of birth, sex, immigration status) were added as controls. The estimates in column 3 are from a model where additional controls (parental schooling, family status, and dummy variables for continent of origin) were added to the list of regressors. Column 4 shows the results of adding lagged test scores and column 5 shows the results from a model where lagged test scores are the only control variable. Columns 1-5 are for reading scores while columns 6-10 repeat this sequence for math scores. As before, conventional standard errors are reported in parentheses and standard errors corrected for within-school correlation are reported in brackets.

The regression results are reported in separate rows for test scores in 1994, 1995, and 1996. The 1994 results are shown because it is of interest to know whether the pre-treatment differences between pupils in treatment and control schools are explained by differences in observed covariates. The raw difference in 1994

¹⁰See, e.g., Ashenfelter and Card (1985).

reading scores by treatment status is $-.74$, while the regression-adjusted difference in column 3 is $-.563$. The raw difference in 1994 math scores is $-.37$, while the regression-adjusted difference in column 8 is $-.28$. Thus, most of the difference in 1994 scores between treatment and control pupils remains even after accounting for treatment-control differences in the observed covariates.

Both 1995 and 1996 are post-treatment years for non-religious schools. The estimated treated effects on 1995 and 1996 scores in models that do not control for pre-treatment test scores, reported in columns 1-3, are negative. These negative estimates are not surprising since pupils in treated schools had lower test scores before the treatment began, and since this difference is not accounted for by the included covariates. Controlling for 1994 scores, however, the treatment effects become positive in both 1995 and 1996, and at least marginally significant in 1996, whether or not covariates other than the 1994 scores are included. For example, the effect of the treatment on 1996 reading scores is $.32$ with a (corrected) standard error of $.13$ in a model with lagged scores and the extended set of control variables (reported in column 4). Dropping covariates reduces the estimate to $.25$. The estimate for math scores controlling for lagged test scores and the full set of covariates is $.26$ with a standard error of $.15$. For both math and reading scores, this is about half as large as the corresponding differences-in-differences estimate, suggesting that the differences-in-differences estimates are biased upwards.

In addition to results for the full sample of non-religious schools, we also computed regression results for a matched sample of schools with similar 1994 school-level average scores in treatment and control schools. As noted earlier, this is a non-parametric procedure for making the treatment and control schools as comparable as possible. An attractive feature of this approach is that it reduces the likelihood of bias from school effects that are correlated with treatment status. The matched sample was constructed by selecting groups or pairs of treatment and control schools where the difference in average math scores in 1994 is less than $.1$ (in standard deviation units). Treatment schools for which no control school could be found with a comparable average score were discarded and vice versa. The matched subsample includes 5 treatment and

3 control schools. The results of analyzing this subsample are reported in Panel B of Table 6. Column 6-8 of the table show that matching effectively eliminates treatment-control differences in 1994 math scores though not in reading scores (we were not able to match schools very closely on 1994 reading scores, and therefore use the same matched sample for analysis of both math and reading scores).

Although they are less precise, the results for math scores in the matched subsample are generally similar to those in the full sample. An important difference, however, is that since the schools have already been chosen to have similar math scores in 1994, the estimated effects on 1996 math scores are positive even without controlling for 1994 scores. For example, the estimate in column 8, from a model with the extended set of covariates but no lagged test score, is .241 with a corrected standard error of .19. For reading scores, the results in the matched sample show smaller treatment gaps in 1994 than in the full sample, but the difference by treatment status for 1994 reading scores is not completely eliminated. Controlling for 1994 reading scores, the estimated effect on 1996 reading scores is positive as before, though not significantly different from zero and smaller than in the full sample.

The last set of regression estimates, for religious schools, is reported in Panel C of Table 6. These results are for math scores only since there are no reading scores in 1996 (the only post-treatment year in religious schools). Columns 6-8 show that test scores in treated schools were higher in 1994 but much lower in 1995, regardless of which covariates are included in the regression. In contrast with the estimates for non-religious schools, adding the 1994 score to the equation for 1996 scores actually makes the estimated treatment effect more negative. On the other hand, models that include the 1995 score as a control variable, the results of which are reported in column 10, lead to a positive treatment effect (though not significantly different from zero). This sensitivity to the choice of lagged control variable is similar to the sensitivity to the choice of base year observed in the differences-in-differences estimates for religious schools. Finally, column 11 shows that adding both pre-treatment scores to the regression for religious schools leads to a very small negative estimate that is not significantly different from zero. Since religious treatment and control schools clearly differed in

both pre-treatment years, this last estimate seems most credible. The absence of an effect in religious schools may be due to the fact that the training program began later than in the non-religious schools, and because the scale of the intervention in religious schools was smaller.

C. Additional matching estimates

As a final check on the estimates for non-religious schools, we matched individual pupils on the basis of their 1994 test scores instead of schools as in Panel B of Table 6. An advantage of pupil matching over school matching is that for both reading and math it is possible to find pupils with similar 1994 scores in all treatment and control schools even though their school averages may differ. The pupil-matching strategy was implemented by dividing the distribution of 1994 test scores (including treatment and control observations) into quartiles, comparing treatment and control scores in each quartile, and then summing the quartile-by-quartile treatment effects into a single weighted average with weights given by the distribution of treated observations across quartiles. Assuming the only difference between treatment and control pupils besides the treatment is their 1994 scores, this procedure produces a non-parametric estimate of the effect of treatment on the treated. We chose to match pupils based on quartiles of the 1994 score distribution instead of a finer breakdown because it turns out that division into four groups is enough to eliminate almost all treatment-control differences in 1994 test scores.

The only significant within-quartile treatment-control difference in 1994 scores is for math in the first quartile. This can be seen in column 3 of Table 7, which reports the within-quartile difference in 1994 scores for math and reading. The overall average effect of treatment status on 1994 scores (i.e., the average across quartiles weighted by the number treated in each quartile) is .031 with a standard error of .031 for math and -.052 with a standard error of .04 for reading. This suggests that the quartile-matching strategy does indeed balance the treatment and control groups.¹¹

¹¹ The standard errors reported in this table were corrected for within-school correlation in test scores.

The matching estimates of treatment effects on 1996 scores are reported in column 6 of Table 7. Except for the upper quartile of the reading score distribution, the within-quartile contrasts for 1996 are all positive and, in some cases, individually significant. The overall matching estimate is .25 with a standard error of .16 for math scores and .4 with a standard error of .16 for reading scores. The matching estimates for math are virtually identical to the corresponding regression estimates (reported in column 10 of Table 6). The matching estimates for reading are somewhat larger than the corresponding regression estimates (reported in column 5 of Table 6). Thus, the matching results reinforce the finding that pupils with same 1994 test scores did better in 1996 if they were in schools where teachers received additional training than if they were enrolled in control schools.

V. The Cost-effectiveness of Training Inputs

The estimates in Tables 6 and 7 suggest that a conservative estimate of the effect of the training program in non-religious schools is $.25\sigma$, where σ denotes the standard deviation of pupils' test scores. The overall cost of the program in Jerusalem was about \$12,000 per class. To determine whether this expenditure was worthwhile, we would need to estimate the value of an increase in children's test scores, something that is hard to quantify. On the other hand, we can compare the cost of attaining a given increase in test scores obtained by increasing teacher training with estimates of the cost of achieving a similar increase in test scores obtained by reducing class size or lengthening the school day. This tells us which of these three important inputs is most likely to be worth increasing.

Angrist and Lavy (1998) estimated that reducing the maximum class size in Israel from 40 to 30 would raise test scores by $.15\sigma$, and require a 28 percent increase in the number of classes. An unpublished memorandum from the ministry of Education suggests that the annual operating cost of each new class would be about \$75,000, so the proposed reduction in class size would cost about \$21,500 per existing class. Assuming the relationship between class size and performance is linear, we estimate that it would cost \$35,000

per existing class to achieve a $.25\sigma$ test-score gain. Using the same data, Lavy (1998) estimated that lengthening the school week by 3.8 hours would raise test scores by $.15\sigma$. Since the annual cost of an extra weekly hour of instruction is \$2,000 for each class, this would cost \$7,600 per class. Again, assuming a linear relationship between inputs and performance, the cost of a $.25\sigma$ increase in scores is estimated to be about \$12,600. These calculations suggest that if the objective is improving pupil achievement, teacher training may be at least as cost-effective a strategy as lengthening the school day, and considerably cheaper than reducing class size.

VI. Summary and Conclusions

The relationship between teachers' characteristics and their pupils' achievement has been the subject of many studies, but few have looked at the impact of in-service training. Our estimates suggest that an in-service training program run in Jerusalem's non-religious elementary schools raised children's achievement in reading and mathematics. These findings appear using a variety of statistical methods, including differences-in-differences, regression, and matching. The estimates for religious schools are not clear cut, but this is possibly because the training program in religious schools started later and was implemented on a smaller scale. The estimates for non-religious schools suggest that teacher training may provide a less costly means of increasing test scores than reducing class size or adding school hours.

Of course, the question remains whether the particular training program studied here is similar to training programs that might be used in other settings. Discussions with school officials lead us to believe that the approach taken in Jerusalem is not unusual in the Israeli context. Moreover, the type of training given to reading and math teachers using *30 Towns* money was based on widely used pedagogical strategies ("Humanistic Mathematics" and "Individualized Instruction") originally developed in U.S. schools. Given the recent interest in training in subject matter, it is also noteworthy that the Jerusalem training was in pedagogy and not subject content. At a minimum, the results here suggests that the impact of relatively inexpensive teacher training of this type warrants further study.

References

- Angrist, J. (1998) "Using Social Security Data on Military Applicants to Estimate the Effect of Military Service Earnings." *Econometrica* 66 (2): 249-288.
- Angrist, J. and Lavy, V. (1998a) *Evaluation of the 30 Towns Project in Jerusalem: Pisgat Zeev-Neve Yaakov* (Hebrew). Jerusalem: Jerusalem Schools Authority, January.
- Angrist, J. and Lavy, V. (1998b) "Using Maimonides' Rule to Estimate the Effect of class size on Scholastic Achievement." *Quarterly Journal of Economics*: forthcoming.
- Ashenfelter, O. A. (1978) "Estimating the Effect of Training Programs on Earnings." *The Review of Economics and Statistics* 60 (1): 47-57.
- Ashenfelter, O.A. and Card, D. (1985) "Using the Longitudinal Structure of Earnings to Estimate the Effect of Training Programs on Earnings." *The Review of Economics and Statistics* 67 (4): 648-660.
- Bangert, R., Kulik, J. and Kulik, C. (1983) "Individualized Systems of Instruction in Secondary Schools." *Review of Educational Research* 53 (2), 143-158.
- Bartel, A.P. (1995) "Training, Wage Growth and Job Performance: Evidence from a Company Database." *Journal of Labor Economics* 13 (3): 401-425.
- Behrman, J. R., Khan, S., Ross, D. and Sabot, R. (1997) "School Quality and Cognitive Achievement Production: A Case Study for Rural Pakistan." *Economics of Education Review* 16 (2): 127-42.
- Bishop, Lloyd K. (1971) *Individualizing Educational Systems: The Elementary and Secondary School: Implications for Curriculum, Professional Staff and Students*. New York: Harper and Row.
- Boston Globe (1998), "Towards Better Teachers," July 23, 1998, page A14.
- Bray, J. H. and Howard, G. S. (1980) "Methodological Considerations in the Evaluation of a Teacher-Training Program." *Journal of Educational Psychology* 72 (1): 62-70.
- Bressoux, P. (1996) "The effect of Teachers' Training on Pupils' Achievement: The Case of Elementary Schools in France." *School Effectiveness and School Improvement* 7 (3): 252-79.
- Card, D. and Krueger, A. (1992) "Does School Quality Matter? Returns to Education and the Characteristics of Public Schooling in the United States." *Journal of Political Economy* 100 (1): 1-40.
- Card, D. and Sullivan, D. (1988) "Estimating the Effect of Subsidized Training on Movements In and Out of Employment." *Econometrica* 56 (3): 497-530.
- Carroll, J. G. (1980) "Effects of Training Programs for University Teaching Assistants." *The Journal of Higher Education* 51 (2): 167-183.
- Coleman, J. S., et al. (1966) *Equality of Educational Opportunity*. Washington, D.C., US GPO.
- Dildy, P. (1982) "Improving Student Achievement by Appropriate Teacher In-Service Training: Utilizing Program for Effective Teaching (PET)." *Education* 103 (2): 132-38.
- Ehrenberg, R. G. and Brewer, D. J. (1994) "Do School and Teacher Characteristics Matter? Evidence from *High School and Beyond*." *Economics of Education Review* 13 (1): 1-17.
- Farrell, J. P. and Oliveira, J. (1993) "Teacher Costs and Teacher Effectiveness in Developing Countries." In *Teachers in Developing Countries: Improving Effectiveness and Managing Costs*, eds. Farrell, J.P. and Oliveira, J. B.:175-86. EDI Seminar Series, World Bank, Washington, D.C.
- Hanushek, Eric A. (1986) "The Economics of Schooling: Production and Efficiency in Public Schools." *Journal of Economic Literature* 24 (3): 1141-57.
- Hedges, L. V., Laine, R.D. and Greenwald, R. (1994) "Does Money Matter? A Meta-analysis of Studies of the Effects of Differential School Inputs on Student Outcomes." *Educational Researcher* 23 (3), 5-14.
- Highsmith, R. (1974) "A Study to Measure the Impact of In-Service Institutes on the Students of Teachers who Have Participated." *Journal of Economic Education* 5 (2): 77-81.
- Krueger, A.B., C. Rouse. (1998) "The Effect of Workplace Education on Earnings, Turnover, and Job Performance." *Journal of Labor Economics* 16 (1), 61-94.
- Lavy, V. (1995) "Endogenous School Resources and Cognitive Achievement in Primary Schools in Israel." Falk Institute for Economic Research in Israel, Discussion Paper No. 95.03.
- Lavy, V. (1998) "Using Quasi-Experimental Designs to Evaluate the Effect of School Hours and Class size

- on Student Achievement." Mimeo, The Hebrew University of Jerusalem, Department of Economics.
- Murnane, Richard J. (1975) *The Impact of School Resources on the Learning of Inner City Children*. Cambridge, Mass.: Ballinger.
- New York Times (1998), "The Flaw in Student-Centered Learning," July 20, 1998, p. 15.
- Romberg, Thomas A. (1985) *Toward Effective Schooling - The IGE Experience*. Lanham: University Press of America.
- Rubin, D.B. (1977) "Assignment to a Treatment Group on the Basis of a Covariate." *Journal of Educational Statistics* 2 (1), 1-26.
- Schober, H. M. (1984) "The Effects of In-service Training on Participating Teachers and Students in their Economics Classes." *Journal of Economic Education* 15 (4): 282-95.
- Welch, F. (1966) "Measurement of the Quality of Schooling." *The American Economic Review* 56 (1): 379-92.
- White, Alvin. (1987) "Mathematics as a Humanistic Discipline." In *The International Encyclopedia of Education*. Oxford: Pergamon Press
- White, Alvin. (1993) *Essays in Humanistic Mathematics*. Washington, DC: The Mathematical Association of America.

Table 1: Inputs and School Resources

Inputs	1994 (Before training)			1995 (Training begins)			1996 (Training continues)		
	Treatment (1)	Control (2)	Difference (3)	Treatment (4)	Control (5)	Difference (6)	Treatment (7)	Control (8)	Difference (9)
Teachers training, math (hours per week)	1.17 (.75)	0.50 (.34)	.67 (.83)	6.67 (2.11)	1.17 (.83)	5.50 (2.27)	12.50 (3.10)	2.00 (1.00)	10.50 (3.37)
Teacher training, reading (hours per week)	.00 (.00)	5.67 (1.91)	-5.67 (1.89)	11.00 (2.29)	4.17 (1.60)	6.83 (2.78)	16.00 (2.45)	8.33 (3.78)	7.67 (2.91)
Expenditure per pupil (current NIS)	515.6 (49.0)	580.5 (47.7)	-64.9 (68.5)	759.6 (50.3)	568.4 (40.2)	191.2 (63.7)	780.5 (55.5)	663.9 (35.8)	116.6 (64.8)
Class-size	30.6 (2.59)	32.3 (2.68)	-1.70 (3.79)	30.5 (2.89)	32.5 (2.57)	-2.00 (3.85)	28.62 (2.78)	31.69 (2.96)	-3.07 (4.06)
Total instruction hours per pupil	1.44 (.12)	1.31 (.10)	.13 (.14)	1.48 (.10)	1.44 (.10)	.04 (.13)	1.54 (.10)	1.48 (.10)	.06 (.14)
Instruction hours per pupil for special projects	.37 (.04)	.33 (.06)	.04 (.07)	.23 (.05)	.24 (.05)	-.01 (.09)	.44 (.05)	.45 (.07)	-.01 (.09)

Notes: There are 7 treated and 6 control schools. Standard errors are reported in parentheses. Years refer to school years; for example, 1994 is the school year beginning in September 1993. Training in non-religious schools began in January 1995, during the 1995 school year.

Table 1: Inputs and School Resources (continued)

Inputs	1994 (Before training)			1995 (Before training)			1996 (Training begins)		
	Treatment (1)	Control (2)	Difference (3)	Treatment (4)	Control (5)	Difference (6)	Treatment (7)	Control (8)	Difference (9)
Teachers training, math (hours per week)	3.5 (3.5)	-	-	0.0 (.0)	-	-	7.0 (.0)	-	-
Teacher training, reading (hours per week)	3.5 (3.5)	-	-	0.0 (.0)	-	-	7.0 (.0)	-	-
Expenditure per pupil (current NIS)	568.5 (134.1)	718.9 (87.3)	-150.4 (160.0)	715.9 (291.1)	710.8 (206.7)	5.1 (364.3)	903.5 (350.6)	984.7 (128.4)	-81.2 (406.0)
Class-size	33.7 (4.8)	26.4 (.9)	7.3 (4.8)	30.3 (2.8)	29.6 (1.3)	.7 (3.3)	28.3 (3.6)	27.2 (1.1)	1.1 (4.5)
Total instruction hours per pupil	1.44 (.20)	1.61 (.07)	-.17 (.21)	1.72 (.15)	1.65 (.10)	.07 (.18)	1.83 (.38)	1.80 (.05)	.03 (.33)
Instruction hours per pupil for special projects	.41 (.05)	.44 (.04)	-.03 (.08)	.36 (.09)	.32 (.12)	.04 (.10)	.60 (.20)	.56 (.06)	.04 (.20)

Notes: There are 2 treated and 5 control schools. Training in religious schools began in September 1995, at the beginning of the 1996 school year.

Table 2: Information Availability

	June 1994	January 1995	June 1995	September 1995	June 1996
A. Pupil Status:					
Non-religious	Baseline	Treatment begins	First follow-up	--	Second follow-up
Religious	Baseline	--	Pre-treatment	Treatment begins	First follow-up
B. Information available for/on:					
Grade level	4 th Graders in treated schools	--	5 th Graders	--	6 th Graders
Test scores	JSA feedback tests in Math and Hebrew	--	JSA feedback tests in Math and Hebrew	--	Special tests in Math (similar to JSA feedback tests); JSA middle school placement tests in Hebrew
Control variables	Pupil and school characteristics	--	--	--	Pupil and school characteristics

Table 3: Basic Test Score Data

	Reading						Math					
	Treated schools			Control schools			Treated schools			Control schools		
	Mean (1)	N (2)		Mean (3)	N (4)		Mean (5)	N (6)		Mean (7)	N (8)	
1994												
All	60.48 (17.22)	406		73.18 (15.68)	405		65.18 (17.27)	428		71.20 (16.52)	420	
Tested in 95	60.33 (17.08)	364		75.82 (15.08)	284		65.24 (17.14)	387		70.33 (16.02)	307	
Tested in 96	59.38 (17.17)	301		73.19 (15.10)	290		64.79 (17.07)	372		71.15 (16.74)	262	
1995												
All	76.49 (15.35)	449		79.16 (12.87)	347		76.25 (19.32)	456		77.48 (16.72)	354	
Tested in 94	77.02 (14.83)	364		80.44 (11.64)	284		75.95 (19.07)	387		78.20 (16.40)	307	
1996												
All	66.39 (17.72)	395		68.68 (16.29)	357		75.88 (18.36)	452		73.10 (19.37)	307	
Tested in 94	66.69 (17.05)	301		68.87 (16.43)	290		75.51 (18.37)	372		73.64 (19.40)	262	
Number of schools		7			6			7			6	

Notes: Standard deviations are reported in parentheses. The statistics presented in the table are based on the distribution of raw scores.

Table 4: Pupil Characteristics in 1994

Characteristics	Non-religious schools			Religious schools		
	Treatment (1)	Control (2)	Differences (3)	Treatment (4)	Control (5)	Differences (6)
Number of observations	428	420		114	183	
Percent male	.520 (.024)	.533 (.024)	-.013 (.026)	.53 (.024)	.52 (.024)	.01 (.02)
Percent immigrants	.133 (.016)	.095 (.014)	.038 (.022)	.15 (.03)	.14 (.03)	.01 (.06)
Family size	4.89 (.057)	4.85 (.055)	.04 (.08)	5.44 (.13)	6.34 (.15)	-.90 (.20)
Percent status: married	.84 (.018)	.84 (.018)	-.00 (.01)	.87 (.03)	.87 (.03)	-.00 (.04)
Father's education	12.36 (.146)	13.83 (.179)	-1.47 (.23)	12.12 (.28)	13.82 (.32)	-1.70 (.43)
Mother's education	12.47 (.129)	13.95 (.154)	-1.48 (.20)	12.26 (.25)	13.77 (.25)	-1.51 (.36)
Child age	10.16 (.019)	10.10 (0.15)	.06 (.02)	10.14 (.04)	10.14 (.03)	.01 (.05)

Notes: Standard errors are reported in parentheses. The sample in this table includes all the pupils that were tested in 1994.

Table 5: Difference-in-differences in Test Scores

Estimate	Reading			Math		
	1995-94 (1)	1996-95 (2)	1996-94 (3)	1995-94 (4)	1996-95 (5)	1996-94 (6)
A. Non-religious schools						
Control differences	-.128 (.055) [.196]	.115 (.055) [.119]	-.061 (.057) [.163]	.045 (.055) [.184]	-.110 (.057) [.141]	-.052 (.058) [.147]
Differences-in-differences	.602 (.073) [.264]	-.021 (.073) [.162]	.616 (.080) [.230]	.173 (.073) [.251]	.249 (.072) [.186]	.460 (.075) [.203]
N	648	624	591	694	549	634
B. Religious schools						
Control differences	-.343 (.071) [.173]	-	-	.054 (.066) [.133]	-.241 (.088) [.116]	-.192 (.097) [.208]
Differences-in-differences	-.490 (.133) [.312]	-	-	-.867 (.110) [.236]	.246 (.137) [.193]	-.719 (.147) [.370]
N	234	-	-	255	187	196

Notes: Control differences are the change in test scores for pupils in control schools. Differences -in-differences are the difference in changes between treatment and control schools. Conventional standard errors reported in parenthesis. Standard errors corrected for within-school correlation are reported in brackets. Standardized scores were used for this analysis.

Table 6: Regression Estimates of training effects

Dependent Variable	Reading					Math				
	No control (1)	Basic control (2)	Extended control (3)	Extended and 1994 score (4)	1994 score only (5)	No control (6)	Basic control (7)	Extended control (8)	Extended and 1994 score (9)	1994 score only (10)
A. All non-religious schools										
1994 Score	-.744 (.072) [.248]	-.709 (.072) [.242]	-.563 (.074) [.204]	-	-	-.370 (.080) [.150]	-.362 (.079) [.150]	-.284 (.086) [.162]	-	-
1995 Score	-.231 (.072) [.220]	-.195 (.071) [.197]	-.049 (.076) [.201]	.221 (.076) [.204]	.173 (.071) [.217]	-.124 (.076) [.243]	-.091 (.075) [.222]	-.002 (.088) [.256]	.057 (.076) [.243]	.034 (.065) [.238]
1996 Score	-.128 (.081) [.144]	-.116 (.082) [.139]	.069 (.084) [.095]	.315 (.082) [.130]	.251 (.078) [.152]	.089 (.072) [.158]	.082 (.072) [.151]	.138 (.076) [.145]	.260 (.067) [.146]	.262 (.063) [.165]
B. Non-religious schools, matched sample										
1994 Score	-.468 (.087) [.321]	-.399 (.088) [.299]	-.373 (.092) [.275]	-	-	-.086 (.097) [.122]	-.083 (.096) [.122]	.035 (.105) [.153]	-	-
1995 Score	-.411 (.096) [.215]	-.341 (.096) [.178]	-.284 (.103) [.152]	-.045 (.096) [.139]	-.102 (.087) [.175]	-.016 (.084) [.298]	-.012 (.083) [.284]	.053 (.096) [.310]	.030 (.082) [.290]	.025 (.072) [.297]
1996 Score	-.063 (.098) [.192]	-.063 (.100) [.183]	.008 (.102) [.111]	.178 (.096) [.162]	.173 (.091) [.206]	.147 (.085) [.195]	.148 (.084) [.189]	.241 (.091) [.188]	.228 (.082) [.200]	.182 (.075) [.219]

Notes: Conventional standard errors are reported in parentheses. Standard errors corrected for within-school correlation are reported in square brackets. Standardized scores were used in this analysis.

Table 6: Regression Estimates (continued)

Dependent Variable	Reading					Math					
	No control (1)	Basic control (2)	Extended control and 1994 score (3)	Extended control and 1994 and 1995 score (4)	Extended control and 1995 score (5)	No control (6)	Basic control (7)	Extended control (8)	Extended control and 1994 score (9)	Extended control and 1995 score (10)	Extended control and 1995 and 94+ 95 score (11)
C. Religious schools											
1994 Score	-	-	-	-	-	.171 (.131) [.271]	.170 (.134) [.261]	.357 (.152) [.178]	-	-	-
1995 Score	-	-	-	-	-	-.732 (.122) [.181]	-.729 (.124) [.166]	-.542 (.145) [.110]	-	-	-
1996 Score	-	-	-	-	-	-.478 (.167) [.206]	-.492 (.169) [.212]	-.304 (.213) [.300]	-.696 (.187) [.324]	.190 (.177) [.238]	-.059 (.182) [.263]

Table 7: Matching on 1994 Scores – Estimates for non-religious schools

Subject	1994 Quartile	1994 Scores			1996 Scores		
		Treatment (1)	Control (2)	Differences (3)	Treatment (4)	Control (5)	Differences (6)
Math	I	-1.483	-1.611	.128 (.062)	-.484	-.636	.152 (.201)
	II	-.350	-.364	.014 (.055)	-.092	-.149	.241 (.188)
	III	.339	.370	-.032 (.055)	.482	.131	.350 (.192)
	IV	1.045	1.070	-.025 (.057)	.737	.441	.296 (.210)
Average effect		.031 (.031)			.250 (.158)		
Reading	I	-1.433	-1.355	-.078 (.069)	-.402	-1.108	.706 (.235)
	II	-.589	-.538	-.051 (.058)	-.020	-.193	.173 (.201)
	III	.164	.200	-.036 (.055)	.395	.062	.333 (.189)
	IV	.969	.957	.012 (.063)	.402	.601	-.199 (.214)
Average effect		-.052 (.040)			.399 (.157)		

Notes: Standard errors are reported in parenthesis. Standardized test scores were used in this analysis. The standard errors are corrected for within-school correlation in test scores.