

A Non-Experimental Evaluation of Curricular Effectiveness in Math

Rachana Bhatt
Georgia State University

Cory Koedel*
University of Missouri

October 2010

We use non-experimental data from a large panel of schools and districts to evaluate curricular effectiveness. Using matching methods, we estimate causal curriculum effects at a fraction of what it would cost to produce experimental estimates. Furthermore, external validity concerns that are particularly cogent in experimental curricular evaluations suggest that our non-experimental estimates may be preferred. We find large short-term differences in effectiveness across some math curricula. However, as with many other educational inputs, math-curriculum effects do not persist over time. Publishers that produce less effective math curricula in one cycle do not lose market share in the next cycle. One explanation for this result is the dearth of information available to administrators about curricular effectiveness.

* We thank Emek Basker, Julie Cullen, Gordon Dahl, Barry Hirsch, Josh Kinsler, David Mandy, Peter Mueser, Rusty Tchernis and many seminar and conference participants for useful comments and suggestions. We also thank Karen Lane and Molly Chamberlin at the Indiana Department of Education for help with data. This work was not funded or influenced by any outside entity.

I. Introduction

According to a 2002 survey sponsored by the National Education Association and the American Association of Publishers, 80 percent of teachers use textbooks in the classroom, and over half of students' in-class instructional time involves textbook use (Finn, 2004).¹ In 2006 alone, expenditures on K-12 instructional materials totaled close to \$8.1 billion.² Despite the prominent role played by curricula in schools, and the significant public expenditures devoted to curriculum purchases, we know surprisingly little about curricular effectiveness. This makes it difficult for educational administrators, who face increasing pressure from state and federal agencies to improve student outcomes, to make informed decisions regarding curriculum adoptions.

Different curricula are developed using different theories about how students learn - this results in different content, organization and structure across curricula for the same subject and grade group. While hundreds of studies have attempted to evaluate the curricular alternatives available to educational administrators, much of the literature on curricular effectiveness lacks scientific rigor, raising concerns about the reliability of the findings (for reviews of the literature see National Research Board 2004; Slavin and Lake, 2008). For example, in 2007, the What Works Clearinghouse (WWC), which was established by the Institute for Education Sciences to serve as a filter for education research, evaluated over 200 studies of curricular effectiveness in elementary mathematics and found that over 96 percent of these studies did not meet reasonable quality standards (WWC, 2007).³

¹ Textbooks are just one component of the curricula purchased by schools from publishers. Other aspects include teacher instructional support services and supplementary materials such as student workbooks and solution manuals.

² See http://www.aapschool.org/vp_funding.html

³ The WWC reviews the literature on a variety of topics in education, including the effects of curriculum adoptions, and classifies studies as either (1) meets evidence standards, (2) meets evidence standards with reservations or (3) does not meet evidence standards. Of the 237 studies on elementary math curricula reviewed by the WWC as of July, 2007, just nine were deemed to be of sufficient quality to be included in categories (1) and (2) (WWC, 2007).

Likely in response to the dearth of reliable evidence in the literature, recent research has turned to randomized controlled trials (RCTs) to evaluate curricular effectiveness (see, for example, Agodini et al., 2009; Borman et al., 2008; Resendez and Azin, 2007). RCTs randomly assign curricula across schools (and/or classrooms) and produce causal estimates of curriculum effects that are internally valid. However, a general drawback of RCTs that is particularly cogent in curricular evaluations is that the estimates may not extrapolate well outside of the experiment.

We highlight two concerns with RCTs in the context of curricular evaluation that will potentially limit their external validity.⁴ First, RCTs require voluntary participation by both *schools and curriculum publishers*. If the schools that select into the experiment differ from the general population of schools, then Manski's (1996) "experimentation on a subpopulation" concern is relevant, and the experimental results may not extend to schools that are not represented in the study.⁵ Perhaps more importantly, there is a selection problem with respect to publishers because publishers are typically actively involved in the experiments. For example, in recent experimental studies by Agodini et al. (2009), Borman et al. (2008) and Resendez and Azin (2007), publishers directly provided teacher training and support services. The active role of publishers in experimental studies means that publishers must agree to participate, and only publishers that expect their curriculum to be successful in the setting of the RCT are likely to do so. Overall, the requirements of voluntary school and publisher participation limit the extent to which experimental designs can be used to evaluate the full curricular landscape.

A second threat to the external validity of RCTs is publisher responsiveness to evaluation, commonly referred to as Hawthorne effects. In the general experimental literature, Hawthorne effects refer to the subjects of the experiment. In the case of curricular evaluation, the

⁴ See Heckman and Smith (1995) and Manski (1996) for general discussions about experimental research designs.

⁵ A common concern in educational experiments is that participating schools may differ in leadership from the average non-participating school.

active role of publishers suggests that in addition to schools and students, they *are* subjects. Given that experimental evaluations are high-stakes competitions for publishers, there is no reason to expect them to take a “business-as-usual” approach. This raises concerns about how well experimental findings will extrapolate to lower-stakes environments for publishers.⁶

In addition to these threats to external validity, the costs associated with RCTs limit the amount of information that they can provide. For example, because RCTs are expensive, they generally focus on just one or two curricula evaluated at small numbers of schools and districts.⁷ The expenses associated with RCTs also limit their usefulness in evaluating long-term impacts because it is costly to maintain the validity of the experiment over time.

Experimental evaluations are informative and offer a number of benefits; however, these issues, some of which are specific to curricular evaluation, suggest that a careful and rigorous non-experimental analysis can make a useful contribution to the literature. This is precisely what we provide in our study, using non-experimental data from the entire state of Indiana to estimate math-curriculum effects on student achievement. We evaluate the three most-used curricula in the state from 1998 - 2004, which together, accounted for 86 percent of all curriculum adoptions in the grades that we study. Indiana provides the most detailed information about curriculum adoptions over time of any of the 50 states, and also provides thorough school- and district-level data about student achievement, demographics and school finances. With the exception of the information about curriculum adoptions, similar data are available in many other states, suggesting that it would be straightforward to replicate our analysis elsewhere.

⁶ In the Agodini et al. (2009) study, the study team “provided logistical and financial support for any level of training the publishers indicated was appropriate.” Although publishers typically provide support services whenever a new curriculum is adopted, they have added incentive to provide high-quality training and support during a RCT.

⁷ In what is a relatively large-scale RCT, Agodini et al. (2009) evaluate four different curricula at four school districts and 39 schools in the first wave of their study (in the second wave they will study 100 schools). Borman et al. (2008) and Resendez and Azin (2007) each evaluate a single curriculum, at five and four schools, respectively.

We use school-level matching estimators in our study, adopting the pairwise-comparison approach suggested by Lechner (2002) to evaluate the three curricula. Drawing on the extensive methodological literature on matching, we show that the data conditions in Indiana are generally favorable to our approach, particularly in the comparison of the two most popular curricula in the state. A key feature of our study is our extended data panel of Indiana schools containing information from multiple cohorts of students who were never exposed to the curricula that we evaluate. We use these cohorts to perform a series of falsification tests for our primary estimates, which suggest that our findings are unlikely to be driven by selection into the different curricula.

We highlight three primary results from our study: (1) differences across some math curricula have large short-term effects on student achievement, (2) as has been found with other educational inputs (see, e.g., Jacob et al., 2008; U.S. Department of Health and Human Services, 2010), math-curriculum effects do not persist over time, and (3) curriculum publishers that are relatively less effective in one adoption cycle do not lose market share in future adoption cycles. This latter result shares a common theme with prior research suggesting that educational administrators may not make optimal choices (Ballou, 1996). In this case, one explanation is the limited availability of reliable evidence on curricular effectiveness.

II. The Curriculum Selection Process

Curriculum adoptions in Indiana occur annually across the entire state, and rotate in six-year cycles by subject. For example, Indiana's districts adopted new math curricula in 1998, 2004, and 2010. Similarly, recent reading adoptions occurred in 1994, 2000 and 2006. We focus our attention primarily on the math-curriculum adoption that occurred in 1998.⁸

The curriculum selection process begins with the state's Department of Education (DOE) providing a list of approved curricula to school districts. Upon receiving this list, districts have

⁸ We focus on this adoption in order to maximize the number of grade cohorts whose achievement data we observe.

three choices. First, and most commonly, they can adopt one or more of the state-approved curricula. Second, districts may apply to use alternative curricula that are not on the list, but this rarely happens in practice (e.g., no more than one out of the roughly 300 districts chooses this option in any grade in our data). Third, districts can apply for “continued use” where they quite literally continue to use the old textbooks from the prior adoption cycle. Overall, over 98 percent of the districts in Indiana adopted new math curricula from the approved list during the 1998 adoption cycle in each grade.⁹

III. Data

We construct a 17-year data panel of schools and districts in Indiana to evaluate the effects of math-curriculum adoptions in grades one, two and three on grade-3 test scores in math (grade-3 is the first time that students are tested in Indiana). Indiana is the only state where curriculum-adoption information is available at the district level for multiple statewide adoption cycles. Upon request, Indiana provides detailed school- and district-level information on test scores (from the Indiana state test, the ISTEP), attendance rates, enrollment demographics, and district-level financial information.

We evaluate the three curricula that dominated the market during the adoption cycle of interest (1998-2004). These curricula were published by Saxon, Silver-Burdett Ginn and Scott-Foresman, and they accounted for roughly 48, 23 and 15 percent of observed curriculum adoptions in the state, respectively. We denote the Saxon curriculum as curriculum *A*, the Silver-Burdett Ginn curriculum as curriculum *B*, and the Scott-Foresman curriculum as curriculum *C*.

⁹ Indiana is one of 22 states that have a centralized component to the adoption process. The state’s role in approving the curricula adds a constraint to the adoption environment. However, it is not clear that the constraint is meaningfully binding given that the alternative-curriculum option is rarely exercised. Perhaps more telling, the majority of the curriculum market belongs to just a handful of publishers, with 86 percent of all curriculum adoptions in the grades that we study involving just three of the ten state-approved curricula. Also, note that Tulley (1989) finds that review processes and adoption lengths are similar in states that don’t have a state-level component.

Because we first observe student outcomes in grade 3, our estimated curriculum effects characterize the impacts of *sequences of treatments*. That is, grade-3 test scores are presumably a function of the curricula to which students are exposed in grades one, two and three. To allow for cleanly identified curriculum effects, we exclude districts that adopted more than one curriculum in grades one, two and/or three from our analysis. To illustrate the assignment problem for these districts, consider a district that adopted curriculum *A* in grade one and curriculum *B* in grades two and three. In identifying the effect of curriculum *A* relative to curriculum *B*, the schools in this district are not well-defined as either treatments or controls.¹⁰

We refer to districts that used the same curriculum in all three grades as “uniform curriculum adopters.” Restricting our analysis to these districts reduces our district sample size by eight percent and our analogous school sample size by seven percent (see Appendix Table A.1 for details). That is, most districts are “uniform adopters.” Overall, our analysis includes data from 213 districts and 716 schools. Contrasted with the experimental literature, our non-experimental design allows for a much broader evaluation of curricular effectiveness.

In Table 1 we report differences in means across the schools and districts that adopted the different curricula, using pre-adoption information from 1997. There are only small differences in test score performance and attendance outcomes across adopters of the different curricula, suggesting that selection into the curricula may be limited. However, there are noticeable differences in terms of school demographics, district size, and to some extent, median household income (measured at the district level from the US Census). Among other things, Table 1 indicates that Saxon adopters are disproportionately rural districts, as evidenced by their much smaller district sizes (and corresponding revenues) and their larger shares of white students.

¹⁰ Although we want to distinguish our estimates from estimates of single-year curriculum effects, our analysis is not related to the literature on sequences of treatments that also involve sequential decisions (see, for example, Lechner, 2004; Lechner and Miquel, 2010). In our study, districts make a treatment decision at a single point in time. Thus, methodologically, our evaluation procedure is the same as in the typical one-shot treatment case.

IV. Curriculum Descriptions

In 1998, Mathematically Correct (MC), a national organization of mathematicians, scientists and engineers, qualitatively evaluated eight grade-2 math curricula, including the three curricula that we evaluate here. We briefly highlight the key differences between the curricula as indicated by MC, and report the MC rating for each curriculum, which was based on a 5-point scale. It is important to note that while the MC reviews provide useful insights, they are not based on analyses of actual implementation, let alone student outcomes. We present the descriptions simply to highlight some of the differences that exist across the curricula.

Curriculum A: Saxon Math (overall rating: 3.6)

The MC evaluation indicates that the program design is “easily implemented by teachers,” and instructions to teachers are “clear and direct.” The teacher’s manual includes scripted statements, worksheet problems are assigned for homework, and students are given periodic written and oral assessments. Saxon Math is very thorough in the topics that are covered, but more advanced topics are generally not covered. As one example, of the three curricula of interest here, Saxon is the only one that does not cover addition and subtraction with three-digit numbers in the second grade. Overall, the MC evaluation suggests that Saxon Math may be the most effective curriculum for low-achieving students given its thorough coverage of the topics it covers, but will be less effective for high-achieving students.

Curriculum B: Silver-Burdett Ginn Math (overall rating: 3.4)

The teacher’s manual provides guidance to teachers, although the guidance is not as direct as in Saxon Math. Student worksheets are tied to the daily lesson, but no information is given about the regularity of assessments or homework assignments. The MC review identifies this curriculum’s heavy reliance on graphics to aid in calculations as a weakness (however,

notably, MC still rated curriculum *B* similarly to the other curricula in our study). The level of this curriculum appears to be higher than that of Saxon Math – MC reports that students using this program have a “reasonable chance of moderate achievement levels” but also that the program is “not seen as supporting high achievement levels.”

Curriculum C: Scott-Foresman Addison Wesley Math (overall rating: 3.8)

The teacher’s edition received mixed reviews from the MC evaluation. Like the Silver-Burdett Ginn curriculum, the lessons also involve some discretion for teachers (although there appear to be fewer teacher choices). Vocabulary development is an important part of this curriculum – new vocabulary words are introduced at the beginning of each lesson, and a verbal skills assessment occurs after each lesson. A one page homework sheet is also attached to each lesson. The level of this program appears to be somewhere in between the levels in the prior two curricula. On the one hand, the MC review indicates that “the level is low in a few topics” and “at the top level of students...some topics should be augmented.” On the other hand, the review also notes that “some areas are very well taught and at an excellent level.”

V. Methodology

We use school-level matching estimators to identify curriculum effects. Matching is an increasingly common empirical technique, and the conditions under which matching will identify causal estimates of treatment effects have been well-documented (see, for example, Rosenbaum and Rubin, 1983; Heckman et al., 1997). The key benefits of matching relative to simple regression analysis are (1) matching imposes weaker functional form restrictions and (2) matching resolves any “extrapolation” problems that may arise in regression analysis by limiting the influence of non-comparable treatment and control units in the data (Black and Smith, 2004).

Briefly, the key assumption under which matching will return causal estimates of treatment effects is the conditional independence assumption (CIA). The CIA requires that potential outcomes are independent of the curriculum uptake decision conditional on observable information. Denoting potential outcomes by $\{Y_0, Y_1, \dots, Y_K\}$, curriculum treatment options by $D \in \{0, 1, \dots, K\}$, and X as a vector of (pre-treatment) observable school- and district-level information, the CIA in our multi-treatment context can be written as:¹¹

$$Y_0, Y_1, \dots, Y_K \perp D \mid X \quad (1)$$

Conditional independence will not be satisfied if there is unobserved information that influences both treatment and outcomes. For example, if districts have access to information that is unobserved to the econometrician, Z , such that $P(D = k \mid X, Z) \neq P(D = k \mid X)$, and the additional information in Z influences outcomes, matching estimates will be biased.

We estimate average treatment effects (*ATEs*) for the three curricula using the pairwise-comparison approach suggested by Lechner (2002), and match schools using an estimated propensity score (Rosenbaum and Rubin, 1983). For example, for the comparison between curricula j and m , where Y_j and Y_m are outcomes for treated and control schools, respectively, we estimate $ATE_{j,m} \equiv E(Y_j - Y_m \mid D \in \{j, m\})$. Defining P_j as the probability of choosing j , we

match schools by $\rho_{jm} \equiv \left(\frac{P_j}{P_j + P_m} \right)$, where P_j and P_m are estimated using a multinomial probit

that simultaneously models the three treatment options (Lechner, 2002).

We use kernel and local-linear-regression (LLR) matching estimators. These estimators construct the match for each “treated” school using a weighted average over multiple “control”

¹¹ The CIA is actually a stronger assumption than is required to identify causal treatment effects, although it is difficult to imagine an environment where only the weaker but necessary condition of conditional *mean* independence is satisfied (Heckman et al., 1997; Imbens, 2004).

schools, and vice versa. We estimate $ATE_{j,m}$ by:

$$\hat{\theta}_{j,m} = \frac{1}{N^S} \left[\sum_{j \in N_j \cap S_p} \{Y_j - \sum_{m \in I_{0j} \cap S_p} W(j,m)Y_m\} - \sum_{m \in N_m \cap S_p} \{Y_m - \sum_{j \in I_{0m} \cap S_p} W(m,j)Y_j\} \right] \quad (2)$$

In (2), N^S is the number of schools using j or m on the common support, S_p . I_{0j} indicates the set of schools that chose m in the neighborhood of observation j , and I_{0m} indicates the set of schools that chose j in the neighborhood of observation m (neighborhoods are defined based on propensity scores using a bandwidth parameter – see Appendix B). $W(j,m)$ and $W(m,j)$ weight each comparison school outcome depending on its distance, in terms of estimated propensity scores, from the observation of interest. We omit a detailed discussion of these matching estimators for brevity. For more information, see Heckman et al. (1997, 1998), and Fan (1993).¹²

Our matching estimators condition on all of the observable information detailed in Table 1. *Ex ante*, it is unclear how unobserved selection might bias our estimates. For example, we might be worried that adopters of different curricula have student populations that differ in unobservable ways, or that differences in administrator quality that are correlated with curriculum adoptions may bias our results.¹³ Although the CIA is not a testable assumption, in Section IX we provide some insight about the likely role of unobserved selection in our analysis using a series of falsification tests. In these tests, we estimate curriculum “effects” for multiple cohorts of students who were never actually exposed to the curricula of interest. If our matching procedure is producing estimates that are not biased by unobserved selection, we should estimate effects of *zero* for these cohorts. We present 80 different falsification estimates along these lines, which show that our primary findings are unlikely to be driven by selection on unobservables.

Finally, average treatment-on-the-treated effects (*ATTs*) may also be of interest. *ATTs* can provide important information if the curricula differentially affect different subgroups of

¹² Our results are robust to alternative matching estimators, and weighting estimators (see Section VIII).

¹³ Additionally, students may move across districts in response to curriculum adoptions. We find no evidence of such movement in the data.

schools. For example, consider a case where $\theta_{j,m} = 0$. This could occur even if schools that chose j were better off for having chosen j , and schools that chose m were also better off for having chosen m . In addition to our *ATE* estimates, we estimate *ATT*'s for all of the curriculum comparisons in both directions (i.e., we estimate $ATT_{j,m}$ and $ATT_{m,j}$). We briefly discuss our findings in Section VIII, but in general, we gain little additional insight by estimating the *ATT*s.

VI. Timing and Treatment Definition

Timing is an important issue in our analysis. Our data panel spans 17 years, starting with the 1991-1992 school year and ending with the 2007-2008 school year. The curricula of interest were adopted in the fall of 1998, and replaced with new curricula in the fall of 2004. We observe seven cohorts of grade-3 students who were never exposed to the curricula of interest during the pre-period (1991-1992 through 1997-1998), one cohort that was exposed to the curricula in grade three only (1998-1999), one cohort that was exposed in grades two and three only (1999-2000), four cohorts that used the curricula in grades one, two and three and were thus “fully exposed” (2000-2001 through 2003-2004), one cohort that was exposed in grades one and two only (2004-2005), one cohort that was exposed in grade one only (2005-2006), and two additional cohorts that were never exposed to the curricula in the post-period (2006-2007 and 2007-2008).

The fully-exposed cohorts provide our cleanest estimates of curriculum effects because they were exposed to the curricula of interest in all three grades. Treatment effects can still be estimated for the partially-exposed cohorts (those that were exposed to the curricula for at least one year, but less than three), however, because we do not observe curriculum treatments outside the adoption cycle of interest, inference will be somewhat limited. A similar concern is relevant for our falsification tests; this issue will be addressed in more detail in Sections VIII and IX.

An additional concern related to timing is that “curriculum familiarity” in schools may be

important. For example, the cohorts of students who used the curricula when the curricula were first introduced to teachers may have had a different experience than the cohorts of students who used the curricula toward the end of the adoption cycle. Unfortunately, familiarity effects cannot be identified using the partially-exposed cohorts in our data because differences in familiarity are inseparable from differences in the length of time cohorts were exposed to the curricula, as well as the grade level(s) in which they were exposed. We can only identify familiarity effects across the four fully exposed cohorts who used the curricula in all three grades.

Finally, a third timing issue involves district restructuring over the course of our 17-year data panel, where there is a pattern of school consolidations in the data. As we discuss in the next section, we match schools based on their static characteristics from the 1996-1997 and 1997-1998 school years. School consolidations may alter the populations of students served by the schools that remain in our data over time, reducing the quality of our matches and potentially introducing bias into our estimates.

In order for the school consolidations to bias our estimates they must be correlated with curriculum adoptions. However, this does not appear to be the case. Using a χ^2 test for independence, we fail to reject the null hypothesis that curriculum adoptions are independent of whether a district experiences a school closing (p-value ≈ 0.40). Additional evidence that our results are unlikely to be biased by school consolidations is provided in Section VII, where we evaluate covariate balance across matched treatment and control schools over the entire course of the data panel (see Table 2). If the schools that drop out of our sample over time systematically adopted specific curricula, we should find that our treatment and control samples become less balanced as we move away from the matching years. We find little evidence of this, supporting

our contention that school closings are not correlated with curriculum adoptions.¹⁴

Although we do not expect the school consolidations to bias our results, they will reduce the quality of our matches as we move away from the 1996-1997 and 1997-1998 school years in the data panel. This will add noise to our estimates. Ultimately, we simply report this issue as a caveat, and caution the reader to interpret results that are estimated far away from the matching years more liberally.¹⁵

VII. The Propensity Score

We use a multinomial probit (MNP) to estimate the propensity scores for schools. The covariates from the MNP are documented in Table 1, and include both school- and district-level information. At the school level, we include controls for enrollment, demographics (race, free and reduced-price lunch status, language status) and outcomes (grade-3 test scores in math and language arts, and attendance) from the 1996-1997 school year, and controls for enrollment and demographics from the 1997-1998 school year. At the district level, we include enrollment, outcome and finance controls from 1996-1997, and enrollment and finance controls from 1997-1998. We also use district-level zip codes to assign year-2000 Census measures of local-area socioeconomic status to each school; namely, median household income and the share of the adult population without a high-school diploma. We treat the census variables as fixed area characteristics.

The covariates in the MNP were selected to represent the relevant information set

¹⁴ Of course, this use of balancing to test for non-random attrition will only catch non-random attrition if it is correlated with observables.

¹⁵ In an omitted analysis, we also considered a more direct solution to this problem – at any point where a school closing was observed in a district, we dropped all school-level observations from that district for the remainder of the data panel (an analogous procedure was done for schools that came into existence between 1991-1992 and 1996-1997). This alternative approach produces estimates that are qualitatively similar to what we report in the text, but comes at the cost of reduced efficiency. Another limitation of this district-dropping strategy is that a school closing may not only re-shuffle students to other schools within that district, but also across district lines. Consequently, this approach may still retain schools which have been altered by the closings of other schools.

available to schools and districts at the time of the adoption decision.¹⁶ For instance, because the curriculum adoption decision was made by the summer of 1998, it is unlikely that decision makers would have had access to spring 1998 test scores, and consequently we do not include them in the model (we also omit annual attendance figures from 1997-1998 for the same reason). That said, our findings are not qualitatively sensitive to reasonable adjustments to the MNP specification, including the addition of the 1997-1998 outcome variables. Similarly, our findings do not depend on whether we include additional years of lagged test scores in the MNP. An important reason for limiting the number of lagged years of achievement in the model is that we want to use as many years of data as possible for the falsification tests. Each year of data that we use to match schools is one less year that we can use in the falsification exercise.

In each comparison we match treatment and control schools based on the estimated pairwise propensity scores, and test for balance in the covariates among the treated and control samples used for estimation.¹⁷ Balancing tests are motivated by Rosenbaum and Rubin (1983). The tests determine whether $X \perp D | P(D = K | X)$, a necessary condition if the propensity score is to be used to reduce the dimensionality of the matching problem to one.

Although achieving covariate balance is important for any matching analysis that relies on a propensity score, there is no clearly preferred test for balance. Furthermore, in some cases, different balancing tests return different results (Smith and Todd, 2005). Given this limitation we consider two different tests. The first is a regression-based test suggested by Smith and Todd (2005), estimated separately for each pairwise comparison, and for each covariate in each year of our analysis. In the comparison between curricula j and m we estimate:

¹⁶ The timeline for the current math-curriculum adoption cycle is available at <http://www.doe.in.gov/olr/docs/CHRONOLOGYFORTHE2009MATHEMATICSADOPTIONApr09.pdf>.

¹⁷ For brevity we do not report the results from the propensity-score model, but they are available upon request. To provide a sense of the predictive power of the covariates in the model, we estimate separate linear-regression models for each curriculum comparison where the dependent variable indicates the adoption of one of the curricula, and the independent variables are the covariates from the MNP. Within comparison pairs, the covariates explain 23 to 42 percent of the variability in curriculum adoptions.

$$\begin{aligned}
X_k = & \beta_0 + \beta_1 \rho_{jm} + \beta_2 \rho_{jm}^2 + \beta_3 \rho_{jm}^3 + \beta_4 \rho_{jm}^4 \\
& + \beta_5 D + \beta_6 * D * \rho_{jm} + \beta_7 * D * \rho_{jm}^2 + \beta_8 D * \rho_{jm}^3 + \beta_9 * D * \rho_{jm}^4 + \varepsilon
\end{aligned} \tag{3}$$

In (3), X_k represents a covariate from the propensity-score specification, ρ_{jm} is the estimated pairwise propensity score, and D indicates treatment. We test whether the coefficients β_5 - β_9 are jointly equal to zero in each regression – that is, we test whether treatment predicts the X 's conditional on a quartic of the propensity score.

The second test measures the absolute standardized difference in observables after matching, and was originally suggested by Rosenbaum and Rubin (1985). The formula for the absolute standardized difference for covariate X_k is given by:

$$SDIFF(X_k) = \frac{|\frac{1}{N^S} [\sum_{j \in N_j \cap S_p} \{X_{kj} - \sum_{m \in I_{0j} \cap S_p} W(j, m) X_{km}\} - \sum_{m \in N_m \cap S_p} \{X_{km} - \sum_{j \in I_{0m} \cap S_p} W(m, j) X_{kj}\}]|}{\sqrt{\frac{Var(X_{kj}) + Var(X_{km})}{2}}} * 100 \tag{4}$$

The numerator in (4) is analogous to the formula for our matching estimators in (2) where we replace Y with X_k and take the absolute value (note the denominator is calculated using the full sample). A weakness of using standardized differences is that there is not a clear rule by which to judge the results, although Rosenbaum and Rubin (1985) suggest that a value of 20 is large.

Our MNP specification uses 32 school- and district-level covariates. The results from the balancing tests are reported in Table 2 by comparison and year. From the regression tests we report the number of covariates where the F-tests reject the null hypothesis at the 5- and 10-percent levels (the former group is a subset of the latter), and the average p-values across all F-tests. We also report the average absolute standardized difference across all covariates.

Table 2 shows that our comparison between B and A achieves better balance than our other comparisons. For this comparison, both the regression tests and the standardized-difference

results suggest that schools are well-matched. For our comparisons between C and A , and C and B , the covariates appear to be less balanced, although it is not clear that the levels of imbalance in these comparisons are cause for concern. For example, the average p-values from the F-tests in both comparisons are fairly close to 0.50 in all years, which suggests good balance, despite there being more unbalanced covariates than would be expected by chance in both cases. Similarly, although the average absolute standardized difference is larger in these comparisons than in our comparison of B and A , by some standards it is still quite reasonable.¹⁸

We also calculate the divergence between the densities of the estimated propensity scores for treated and control units in each comparison. Intuitively, density divergence will affect the precision of the estimates obtained from matching. Frölich (2004) measures density divergence using the Kullback-Leibler (KL) information criterion; we follow his approach here, using kernel-density plots based on the Epanechnikov kernel. We estimate the divergence between the densities of ρ_{jm} for treatment and control schools as:

$$KL = \int \ln \left(\frac{f_{p|D=j}(\rho_{jm})}{f_{p|D=m}(\rho_{jm})} \right) f_{p|D=j}(\rho_{jm}) d\rho_{jm} + \int \ln \left(\frac{f_{p|D=m}(\rho_{jm})}{f_{p|D=j}(\rho_{jm})} \right) f_{p|D=m}(\rho_{jm}) d\rho_{jm} \quad (5)$$

In (5), $f_{p|D=j}(\rho_{jm})$ is the density function of ρ_{jm} among schools treated with j , and $f_{p|D=m}(\rho_{jm})$ is the analogous density function for schools that used m . A KL-information-criterion measure of zero suggests that the densities are identical, and the measure increases with density divergence. Note that when the parameters of interest are average effects of treatment on the treated, researchers use a unidirectional version of the KL information criterion (Frölich, 2004). In our

¹⁸ Although it is not obvious that the level of imbalance in any of our comparisons is large enough to be problematic, in unreported results we considered many alternative propensity score specifications where we added higher-order and interaction terms in an effort to improve covariate balance. Likely due in part to our relatively small samples, these alternative models generated only modest improvements in covariate balance, and did not affect our findings qualitatively. Thus, we proceed using the MNP model described in the text, noting that the balancing results are less compelling in our comparisons between C and A , and C and B .

case, where average treatment effects are the parameters of interest, we use the bidirectional information criterion originally suggested by Kullback and Leibler (1951).

Figure 1 plots the estimated density functions of the propensity scores for treatment and control schools for each pairwise comparison, and Table 3 reports the corresponding KL information criteria. Similarly to the balancing tests, the density-divergence measures suggest that the data conditions are most favorable in our comparison between B and A . Density divergence is largest in our comparison between C and A .¹⁹

Both the balancing tests and the density-divergence measures indicate that our data are best-suited to compare curricula B and A , which combined, accounted for over 70 percent of the curriculum market in Indiana during the 1998 adoption cycle. In the other two comparisons the data conditions are generally less favorable; however, even in these comparisons, it is not clear that they are cause for concern. We consider the reliability of our results further in Section IX.

VIII. Estimates of Curricular Effectiveness in Math

Rather than overwhelm the reader with estimates using the numerous matching algorithms available in the literature, we instead present estimates using kernel and LLR matching only (for details on these and other matching estimators, see, for example, Mueser et al., 2007). Frölich's (2004) analysis indicates that kernel matching in particular should perform well in our context. As for LLR matching, the evidence in the literature is mixed.²⁰ Although our estimates using LLR matching are less precise than the kernel-matching estimates, they are generally very similar. We present results using the Epanechnikov kernel for both types of

¹⁹ Frölich (2004) uses unidirectional density divergence measures in his study. Although the one and two-sided measures are not directly comparable; roughly speaking, our comparison of B and A corresponds to his most favorable design, C and B to his middle, and C and A to his least favorable design. This is purely by coincidence.

²⁰For instance, Caliendo and Kopeinig (2005) suggest LLR is useful when controls are distributed asymmetrically around treated observations. Frölich (2004) notes that LLR will perform worse in regions of sparse data, which is consistent with the large standard errors that we estimate using LLR in our comparisons with less density overlap.

matching estimators. In unreported results available upon request we show that our results are robust to alternative estimators, including kernel and LLR matching estimators that use the Gaussian kernel, other matching estimators based on simple pair matching or radius matching using various radii, and regression-adjusted and weighting estimators (see Imbens (2004) and Millimet and Tchernis (2009) for discussions of weighting estimators).

Table 4 presents results for all grade-3 cohorts who were ever exposed to the curricula of interest using fixed-bandwidth matching estimators where the bandwidths are obtained via conventional cross-validation (see Appendix B). All of our matching estimators impose the common support condition. We also report OLS estimates where we regress test score outcomes on the covariates used in the propensity score model and indicator variables for curriculum adoptions, retaining the pairwise comparisons (i.e., when we compare B to A , we drop all schools at districts that adopted C). The standard errors for the matching and OLS estimates are clustered at the district level and the matching-estimator standard errors are bootstrapped with 250 repetitions. We obtain the optimal number of bootstrap repetitions to use for our standard errors following Ham et al. (2006), who use a special case of Andrews and Buchinsky (2001).²¹

Each cohort is labeled in the tables according to the year of its spring test score (e.g., the 1998-1999 cohort is labeled “1999”). All of the effects in the table are standardized using the distribution of school-level test scores. For example, the estimate in Table 4 for $ATE_{B,A}$ in 2002 indicates that, among the sample of schools that chose B or A , the average effect of using B instead of A was 0.40 standard deviations of the distribution of school-level test scores. More typically, researchers report effects that are standardized based on the distribution of *individual-*

²¹ Our bootstrapping procedure re-samples entire districts. Abadie and Imbens (2006) show that bootstrapping cannot be used to obtain standard errors for nearest-neighbor matching estimators, but their study does not apply to smoother estimators like those used here. For our estimators, the optimal number of bootstrap repetitions is consistently near 200, so we use 250 repetitions to ensure that we meet or exceed the optimal count in each year.

level scores, but we do not have access to the distributions of individual-level scores over the entire course of the data panel. In Appendix Table A.2, for the years where we have access to the distribution of individual-level test scores, we show the scaling factors that convert the estimates in Table 4 into the more common metric. Roughly speaking, dividing the coefficients by three returns estimates in units of standard deviations of the individual-level distribution of scores.

Focusing first on our largest comparison between *B* and *A*, and the estimates for the fully-exposed cohorts (2001 – 2004), we find that curriculum *B* meaningfully outperformed curriculum *A*. Averaging the kernel-matching estimates across all four fully-exposed cohorts, and using the appropriate scaling factors in Appendix Table A.2, the effect of using curriculum *B* instead of *A* was approximately 0.12 standard deviations of the test. Our estimates are also consistent with *C* outperforming *A*. There we estimate an average effect of roughly 0.06 standard deviations of the student-level distribution of scores, although only two of the four estimates are statistically significant and the estimate from 2004 is particularly small. Our results also suggest, at least weakly, that *B* outperformed *C*, although inference from this comparison is limited because the estimates are imprecise. Finally, we do not observe any pronounced trends in curricular effectiveness across the four fully-exposed cohorts, providing no evidence to suggest a role for curriculum-familiarity effects.²²

The magnitudes of the estimated curriculum effects are economically meaningful, particularly when weighed against the marginal costs associated with choosing one curriculum over another. Fryer and Levitt (2006) show that between grades one and three, the black-white

²² Although we note that our data are not well-suited to definitively evaluate curriculum-familiarity effects. For example, it may be that familiarity is quite important during the first few years after implementation, but we cannot distinguish familiarity effects for the partially-exposed cohorts per the discussion in Section VI.

achievement gap grows at a rate of approximately 0.10 standard deviations per year.²³ In our most-compelling curriculum comparison, we estimate that the effect of choosing curriculum *B* over curriculum *A* is equivalent to roughly one year's worth of expansion of the black-white achievement gap. Given that the curricula are very similarly priced (the texts from *A*, *B* and *C* were, averaged over grades 1-3, \$23.08, \$24.80 and \$25.34 each, respectively, in 1998 dollars), selecting a better curriculum appears to be a cost-effective way to improve student achievement.

Our results for the partially-exposed cohorts differ by comparison. One common theme is that the point estimates for the 2005 and 2006 cohorts are generally larger than for the 1999 and 2000 cohorts. In fact, in our comparison between *B* and *A*, the estimates for the 2005 and 2006 cohorts are large and statistically distinguishable from zero. One explanation for this finding is that although we cannot distinguish any curriculum-familiarity effects using the fully-exposed cohorts, there may be familiarity issues upon immediate adoption, which would affect the 1999 and 2000 cohorts but not the 2005 and 2006 cohorts. Also, per Section IV, curriculum *B* is distinguished from the other curricula by its reliance on models to teach mathematics. Although it would be entirely speculative to link the benefits of curriculum *B* to any specific attribute, this pedagogical difference may have the potential to stay with teachers beyond the 1998 adoption cycle. Equally interesting is that MC rated curriculum *B* similarly to *A* and *C* despite its use of models, which MC does not favor. That curriculum *B* was likely downgraded for using models, but still received a rating similar to *A* and *C*, suggests that its overall quality may be high.

It is important to note that the students in the partially exposed cohorts were exposed to other curricula in other adoption cycles, and this may attenuate the partial-exposure estimates. The degree of attenuation will depend on the extent to which curriculum quality is correlated across adoption cycles for treatment and control schools. We explore this issue to the extent

²³ Fryer and Levitt (2006) analyze a different testing instrument; however, similar estimates of the black-white achievement gap spread are available elsewhere (see, for example, Chubb and Loveless, 2002).

possible in Table 5, where we compare curriculum adoptions in the 2004 adoption cycle across uniform adopters from 1998 (curriculum data are unavailable prior to the 1998 adoption cycle, meaning that we can only directly evaluate the across-cycle curriculum adoptions that are relevant for the 2005 and 2006 cohorts).

For brevity, Table 5 shows adoption shares in 2004 only for the four most popular curricula from that adoption cycle (published by Saxon, Harcourt, Houghton-Mifflin and Scott-Foresman). The table shows that while Saxon adopters (curriculum *A*) in 1998 were much more likely to adopt Saxon in 2004, adopters of the other two curricula are quite dispersed across alternative options. Without knowing the respective qualities of the different curricula adopted outside of the 1998 adoption cycle, including those from the same publishers (there is no evidence that we are aware of on the persistence of publisher quality), it is difficult to form expectations based on the patterns in Table 5. Ultimately, given the potential for attenuation in the estimates for the partially exposed cohorts, and the sizes of our standard errors, we cannot make strong inference about partial-exposure curriculum effects.²⁴

Table 5 is also informative about the changing market shares of curriculum publishers over time. It shows that the publisher of curriculum *A*, despite its relative underperformance, maintained its near-50-percent market share from the 1998 adoption cycle to the 2004 cycle. Although curriculum *B* was the most effective curriculum during the 1998 adoption cycle, it did not appear in 2004. The publisher of curriculum *B* was bought by Pearson Publishing during the 1998 cycle and Pearson phased out curriculum *B* in favor of *C*, which it also publishes. Curriculum *C*'s market share fell from roughly fifteen to nine percent across adoption cycles.

Finally, in an omitted analysis we also consider the possibility that the treatment effects

²⁴ Evidence on the persistence of publisher quality would be difficult to obtain without the availability of consistent comparisons over time. For example, because Silver-Burdett Ginn did not offer a curriculum in Indiana during the 2004 adoption, our most reliable comparisons (per Section VII) cannot be replicated in the later adoption cycle. Even more, we cannot reliably compare Saxon and Scott-Foresman in 2004 because of the large decline in Scott-Foresman's market share across adoption cycles. Even in cases where curriculum publishers are consistently represented across adoption cycles, long cycle durations imply that long data panels will be required to evaluate the persistence of publisher quality.

depend on treatment status. For example, despite our finding that curriculum *B* outperformed *A* on average, it could be that curriculum *A* was still better for schools that actually chose *A*, while curriculum *B* performed better for schools that chose *B*. To investigate the extent to which the curriculum effects might depend on treatment status, for each of our comparisons we estimated average treatment-on-the-treated effects (*ATT*) in each direction. Our findings provide few insights. In our comparison between *B* and *A*, the treatment effects do not depend on treatment status. Similarly, the *ATT*'s in our comparison between *C* and *B* do not suggest differential effects (although again, these estimates are noisy). Only in our comparison between *C* and *A* do we find any evidence of differential effects. There, curriculum *A* appears to perform less poorly relative to *C* at schools that actually chose *A*. Nonetheless, even our estimates of $ATT_{A,C}$ suggest that schools that actually chose *A* would have been better off had they instead chosen *C*.

Overall, our most reliable estimates of curriculum effects come from the four fully-exposed cohorts of students. Our estimates from these cohorts indicate that curriculum *B* outperformed curriculum *A* by a substantial margin. We also find that *C* outperformed *A*, although the differential effect was smaller. The statistical imprecision associated with our comparison between *C* and *B* limits inference, but if anything, our estimates suggest that *B* outperformed *C*. The relative underperformance of *A* did not adversely impact the publisher's market share in the next adoption cycle in Indiana.

IX. Falsification Tests

Matching estimators will not return causal estimates if conditional independence is violated. Although conditional independence is not a testable assumption, we provide some evidence on its plausibility using a series of falsification tests. We present falsification tests based on data from students who were never actually exposed to the curricula of interest, and

from students who were exposed, but we estimate curriculum effects on reading test scores. For the former, we expect to estimate “effects” that are statistically indistinguishable from zero, whereas for the latter, timing does not rule out the possibility of causal effects for some cohorts. However, at most, we would expect only small across-subject spillover effects.

Potentially confounding both types of falsification estimates are correlations in curriculum adoptions across grades, subjects, and adoption cycles. Recall from Table 5 that there are non-zero correlations in math-curriculum adoptions across adoption cycles. Not surprisingly, in unreported results (omitted for brevity and available upon request) we also find that math-curriculum adoptions are correlated across grades within adoption cycles, and to a lesser extent, with adoptions in other subjects. However, in practice, the correlations do not seem to be strong enough to limit inference from our falsification exercise – as we show below, almost all of the falsification estimates are statistically indistinguishable from zero.²⁵

For brevity, we only report falsification estimates using kernel matching with the Epanechnikov kernel. We present 80 falsification estimates in all (but note that the tests are not independent).²⁶ Summarizing the results, the tests do not uncover any consistent evidence of selection bias in any of our comparisons, although similarly to Table 4, the falsification estimates are noisy in our comparison between *C* and *B*, limiting inference.

We first estimate curriculum “effects” on math test scores for cohorts of grade-3 students from 1992 through 1996, and 2007 and 2008. The results are reported in Table 6. Our most-convincing falsification estimates are from the 1992-1996 cohorts, who passed through Indiana schools prior to the 1998 adoption cycle. For these cohorts, all of the estimates are small and statistically indistinguishable from zero with the exception of the 1992 estimate in our

²⁵ Also note that we match schools based on 1997 achievement, which will include curriculum effects from the prior adoption cycle. For relevant cohorts, this should reduce any bias from across-cycle correlations in curriculum quality.

²⁶ If the falsification tests were independent we would expect roughly eight “false positives” in total. However, treatment and control schools are uniformly defined over time, making it unclear how many false positives to expect.

comparison between *C* and *A*. Although the 2007 and 2008 cohorts were not exposed to the curricula of interest, their outcomes are observed after the curriculum-adoption cycle we study. This leaves open the possibility of non-zero treatment effects, limiting inference to some degree, but even so, none of the estimates from these cohorts are statistically significant.

Next we estimate curriculum “effects” using cohorts of grade-6 students who were never exposed to the curricula of interest (cohorts from 1993-2001). For these falsification tests we use the same matching procedure to predict the same treatments (the uniform adoption of curriculum *A*, *B* or *C* in grades one, two and three), only we match schools that have grade-6 classrooms and estimate the “effects” of the curricula on grade-6 achievement. Unfortunately, our grade-6 analysis is limited by sample-size issues. Specifically, because many districts teach grade six in middle school, and multiple elementary schools generally feed into a single middle school, our grade-6 sample of schools is much smaller than our grade-3 sample. In the data, the number of grade-6 curriculum-*C* schools is particularly small (roughly 80, on average, across the data panel), and we cannot balance treatment and control schools in either of our comparisons involving this curriculum. Because the unbalanced comparisons will not be informative, we present grade-6 falsification estimates only for our comparison between *B* and *A*.²⁷ These estimates are reported in Table 7, where we estimate one non-zero “effect” in 1993, but otherwise, the point estimates are generally small and statistically indistinguishable from zero.

In Table 8 we return to our well-balanced grade-3 samples and estimate math-curriculum effects on reading scores for all cohorts in the data. Students in the cohorts from 1992 through 1996, and 2007 and 2008, were never exposed to the curricula of interest. The other cohorts were exposed, and it is unclear *a priori* whether we should expect any across-subject spillover

²⁷For example, taking the average p-values from the Smith and Todd (2005) balancing regressions across years for the comparisons involving *C*, they fall from roughly 0.50 in the grade-3 analysis (Table 2), to roughly 0.20 in the grade-6 analysis. In contrast, the average p-value falls to just below 0.50 in our grade-6 *B* to *A* comparison.

effects.²⁸ Although we do not have a strong prior about whether math curricula affect reading outcomes, one straightforward expectation is that their effects on math test scores should be larger than their effects on reading test scores. Thus, at its most basic level, this final test should confirm this result. Table 8 presents estimates for the effects of the math curricula on reading test scores throughout the data panel. The point estimates are generally small and there is only one statistically significant estimate (in the comparison between *B* and *A* in 2002).²⁹

X. Persistence

Finally, we use our extended data panel to evaluate the persistence of curriculum effects over time (which would be quite difficult, and expensive, with an experimental research design). Specifically, in our comparison between *B* and *A*, we ask whether the cohorts of students who were exposed to curriculum *B* in grades one, two and three still performed better by grade six.³⁰ We measure persistence using test scores for cohorts of grade-6 students between 2002 and 2008, where the fully-exposed cohorts were in grade six between 2004 and 2007.

Two issues merit attention in our persistence analysis. First, if there are test-score ceilings in higher grades on the Indiana test, it will be difficult to detect persistence effects because the tests in later grades may not adequately differentiate student learning. We test for ceiling effects

²⁸We suggest five possible mechanisms that may generate spillover effects. First, math curricula may directly affect reading performance (i.e. via exposure to word problems). Second, the training for teachers associated with each math curriculum could affect teacher performance in other subjects. Third, a better math curriculum may afford teachers more time to spend on reading instruction. Fourth, a better math curriculum may increase the return to math instruction and encourage teachers to substitute out of reading and into math instruction. Fifth, students could become more interested in school in general if they do well in math or perceive it to be more fun.

²⁹We also consider how a pure-bias interpretation of the reading estimates would impact our results by assuming that across-subject spillover effects are zero. To do this, we estimate math-curriculum effects on schools' de-trended math test scores, where we de-trend each school's math score by separately standardizing its math and reading test scores, and subtracting the reading score from the math score. We omit the estimates for brevity, but note that they are in line with what would be expected by subtracting the stand-alone reading estimates from the stand-alone math estimates. The estimates that are statistically significant in Table 4 for our comparisons between *B* and *A*, and *C* and *A* remain statistically significant in the de-trended analysis. In the comparison between *C* and *B*, the curricula are not statistically distinguishable in any year using the de-trended estimates.

³⁰As discussed in the previous section, we are unable to construct observationally equivalent comparisons of treated and control schools from the grade-6 sample in our evaluations involving curriculum *C*. Therefore, we only examine persistence in our comparison between *B* and *A*.

following Koedel and Betts (2010) and find that the testing instruments should be sufficient to detect any persistence effects should such effects exist. A second concern is that we cannot track individual students over time, and must assume that the district a student attends during grade six is the same district they attended during grades one, two and three. Student movement across districts will add noise to our treatment classifications, attenuating the persistence estimates.³¹

Table 9 presents our persistence findings (using kernel matching with the Epanechnikov kernel). The estimates provide little indication that curriculum effects persist over time, despite the likelihood that there is some attenuation bias. Put differently, for the estimates in the table to be driven by downward bias from student movement across districts, the amount of movement would need to be inordinately large. Our persistence findings are consistent with a large body of evidence pointing to a general lack of persistence in the effects of educational inputs (see, for example, Jacob et al., 2008; U.S. Department of Health and Human Services, 2010), and raise doubts about the extent to which administrators can improve student performance in the long run by choosing more effective curricula.

XI. Conclusion

We identify causal curriculum effects using non-experimental data from the state of Indiana. A key component of our study is our falsification exercise, where we use data from multiple cohorts of students who were never exposed to the curricula of interest, and students' out-of-subject test scores, to show that our findings are unlikely to be driven by selection into the different curricula. In cases where data conditions are favorable, and some form of confirmation exercise is possible (like our falsification tests), much can be learned from careful, non-experimental work. A caveat relating to curricular evaluation is that, somewhat surprisingly,

³¹ Student churning across districts is also an issue in our primary analysis, to a lesser extent, and implies some attenuation bias in the estimates in Table 4 as well.

most states do not centrally track curriculum adoptions. Given that it would be relatively inexpensive to track this information, and that curricula play such a large role in students' everyday learning experiences, this seems peculiar.

Currently, the bulk of the curricular-effectiveness debate is not based on rigorous evidence from analyses of implementation. For example, in addition to the general lack of rigor in comparative curricular evaluations (WWC, 2007), much of the literature relies on case studies or content studies, where curriculum impacts on student outcomes are not measured (National Research Board, 2004). Rigorous scientific evidence about how different curricula actually affect achievement is needed in order for administrators and educators to make informed decisions. Our study provides such evidence on a scale not yet seen in the curriculum-evaluation literature.

That our study is non-experimental allows us to bypass some of the limitations inherent to experimental analyses of curricular effectiveness. These limitations include the experimentation on a subpopulation problem (Manski, 1996), and the possibility of publisher Hawthorne effects. The latter concern seems particularly important given publishers' active roles in curriculum experiments. Another benefit of our non-experimental approach is that it is feasible to replicate in other environments both methodologically *and fiscally*. In contrast to the ongoing project by Agodini et al. (2009), a particularly well-designed RCT that is funded by the Institute for Education Sciences for roughly 21 million dollars over five years, our study was performed using publicly available data at only a small fraction of the cost.

We also note several limitations of our study. For one, we do not have enough data, or the right kind of data (i.e., student level), to evaluate the extent to which curricula differentially affect different types of students (e.g., high and low-achieving, English-proficient and English as second language, etc.). This deficiency in our analysis is likely to be less problematic in the

future because districts and states continue to develop longitudinal databases that track individual students. These data could be linked to curriculum data, if such data were available, quite easily. We also depend on the standardized test administered by the state of Indiana as our outcome measure (the ISTEP). While we expect our results to extrapolate well to other states or districts that use similar tests, they may not carry over to states or districts where the testing instrument differs greatly in content. Our results also may not extrapolate well to states or regions where the population differs greatly from the population in Indiana, which is a fairly rural state.

Our findings indicate that students in Indiana who used curricula *B* or *C* outperformed students who used curriculum *A*. In our most compelling comparison, between *B* and *A*, the effect of exposure to the better curriculum for three consecutive years is roughly 0.12 standard deviations of the grade-3 ISTEP test. This effect is similar in magnitude to one year's growth of the black-white achievement gap over these grades (Fryer and Levitt, 2006). Interestingly, despite the consistent underperformance of curriculum *A* in our analysis, the publisher of curriculum *A* did not lose market share in the next adoption cycle in Indiana. There are many possible explanations for this finding, ranging from a lack of reliable information available to administrators about curriculum quality (WWC, 2007), to poor decision making by educational administrators (also see Ballou, 1996).

Overall, our results are encouraging because choosing a better curriculum can non-negligibly improve student performance. Further, the near-zero marginal cost of choosing one curriculum over another suggests that implementing a better curriculum will be quite cost-effective. However, our finding that curriculum effects do not persist over time, although not unique to curricula in education, dampens enthusiasm about the potential benefits of improved curricula. By grade six, the benefits of the most-effective curriculum in our study are no longer

distinguishable.

References

Abadie, Alberto and Guido W. Imbens. 2006. On the Failure of the Bootstrap for Matching Estimators. NBER Technical Working Paper No. 325.

Agodini, Robert and Barbara Harris and Sally Atkins-Burnett and Sheila Heaviside, and Timothy Novak. 2009. *Achievement Effects of Four Early Elementary School Math Curricula*. National Center for Education Evaluation and Regional Assistance, U.S. Department of Education, Institute of Education Sciences. NCEE 2009-4052.

Association of American Publishers. 2001. *Less than a Penny: The Instructional Materials Shortage & How It Shortchanges Students, Teachers, & Schools*. Association of American Publishers Report.

Andrews, Donald W. K. and Moshe Buchinsky. 2001. "Evaluation of a Three-Step Method for Choosing the Number of Bootstrap Repetitions," *Journal of Econometrics*, 103, 345-386

Ballou, Dale. 1996. "Do Public Schools Hire the Best Applicants?" *Quarterly Journal of Economics* 111(1), pp. 97-133.

Black, Dan and Jeffrey Smith. 2004. "How Robust is the Evidence on the Effects of College Quality? Evidence From Matching," *Journal of Econometrics* 121 (2), 99-124.

Borman, Geoffrey D. and N. Maritza Dowling and Carrie Schneck. 2008. "A Multisite Cluster Randomized Field Trial of Open Court Reading," *Education Evaluation and Policy Analysis* 30 (4), 389-407.

Caliendo, Marco and Sabine Kopeinig. 2005. "Some Practical Guidance for the Implementation of Propensity Score Matching," IZA Discussion Paper No. 1588.

Chubb, John and Tom Loveless. 2002. *Bridging the Achievement Gap*, Brookings Institution Press, Washington, D.C.

Fan, Jianqing. 1993. "Local Linear Regression Smoothers and Their Minimax Efficiencies," *The Annals of Statistics*, 21, 196-216.

Finn, Chester. 2004. "The Mad, Mad World of Textbook Adoption" The Thomas B Fordham Foundation and Institute Report. Washington, D.C.

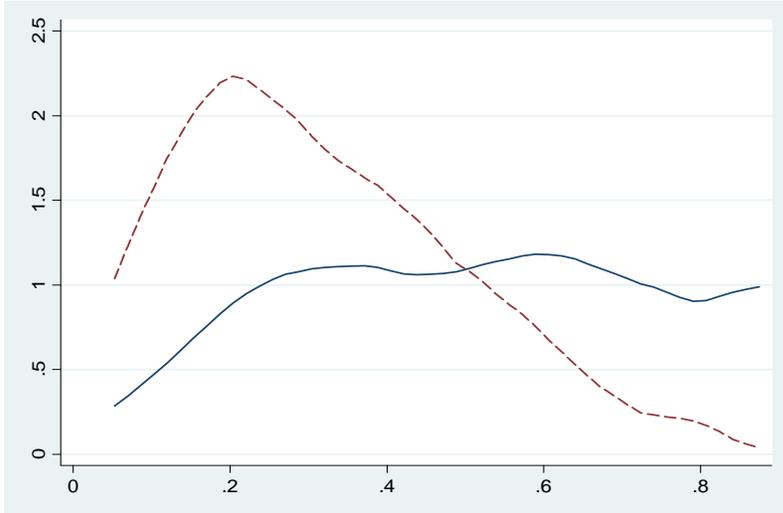
Frölich, Markus. 2004. "Finite-Sample Properties of Propensity-Score Matching and Weighting Estimators," *The Review of Economics and Statistics* 86 (1), 77-90.

- Fryer, Roland and Steven Levitt. 2006. "The Black-White Test Score Gap Through Third Grade," *American Law and Economics Review* 8 (2), 249-281.
- Ham, John and Xianghong Li and Patricia Reagan. 2006. "Propensity Score Matching, a Distance-Based Measure of Migration, and the Wage Growth of Young Men," Federal Reserve Bank, Staff Report No. 212.
- Heckman, James J. and Jeffrey A Smith. 1995. "Assessing the case for social experiments," *Journal of Economic Perspectives* 9 (2), 85-110.
- Heckman, James and Hidehiko Ichimura, and Petra Todd. 1997. "Matching as an Econometric Evaluation Estimator: Evidence from Evaluating a Job-Training Programme," *Review of Economic Studies* 64 (4), 261-294.
- Heckman, James and Hidehiko Ichimura and Petra Todd. 1998. "Matching As An Econometric Evaluation Estimator," *Review of Economic Studies* 65 (2), 261-294.
- Imbens, Guido. 2004. "Nonparametric Estimation of Average Treatment Effects Under Exogeneity: A Review," *The Review of Economics and Statistics* 86 (1), 4-29.
- Jacob, Brian and Lars Lefgren and David Sims. 2008. "The Persistence of Teacher-Induced Learning Gains," NBER Working Paper No. 14065.
- Koedel, Cory and Julian R. Betts. 2010. "Value-Added to What? How a Ceiling in the Testing Instrument Influences Value-Added Estimation," *Education Finance and Policy* 5 (1), 54-81.
- Kullback, Solomon and Richard Leibler. 1951. "On Information and Sufficiency," *Annals of Mathematical Statistics* 22 (1), 79-86.
- Lechner, Michael. 2002. "Program Heterogeneity and Propensity Score Matching: An Application to the Evaluation of Active Labor Market Policies," *The Review of Economics and Statistics* 84 (2), 205-220.
- Lechner, Michael. 2004. "Sequential Matching Estimation of Dynamic Causal Models," IZA Discussion Paper No. 1042.
- Lechner, Michael and Ruth Miquel (2010). "Identification of the Effects of Dynamic Treatments by Sequential Conditional Independence Assumptions," *Empirical Economics* 39, 111-137.
- Li, Qi and Jeff Racine. 2007. *Nonparametric Econometrics: Theory and Practices*, Princeton University Press, Princeton N.J.
- Ludwig, Jens and Douglas Miller. 2007. "Does Head Start Improve Children's Life Chances? Evidence from a Regression Discontinuity Design," *Quarterly Journal of Economics* 122 (1) 159-208.

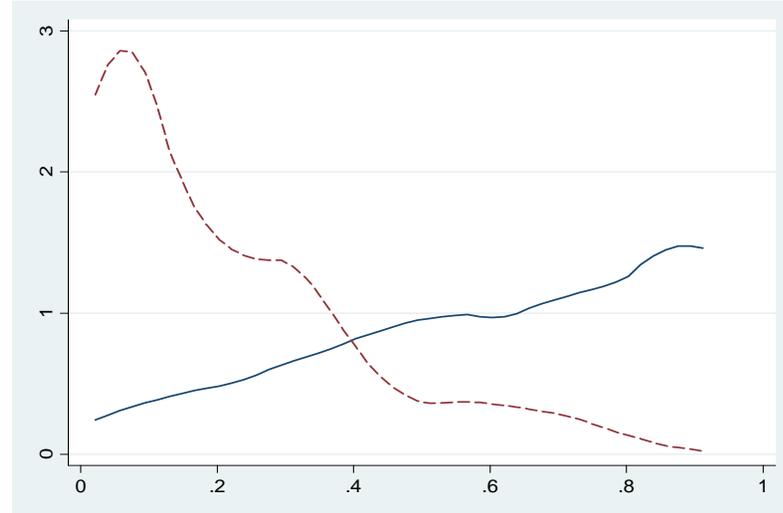
- Manski, Charles. 1996. "Learning About Treatment Effects from Experiments with Random Assignment of Treatments," *The Journal of Human Resources* 31 (4), 709-733.
- Millimet, Daniel L. and Rusty Tchernis. 2009. "On the Specification of Propensity Scores, with Applications to the Analysis of Trade Policies," *Journal of Business & Economic Statistics* 27 (3), 397-415.
- Mueser, Peter R. and Kenneth R. Troske and Alexey Gorislavsky. 2007. "Using State Administrative Data to Measure Program Performance," *The Review of Economics and Statistics* 89 (4), 761-83.
- National Research Board. 2004. *On Evaluating Curricular Effectiveness: Judging the quality of K-12 Mathematics Evaluations*, The National Academies Press, Washington DC.
- Resendez, Miriam and Mariam Azin. 2007. "The Relationship Between Using Saxon Elementary and Middle-School Math and Student Performance on California Statewide Assessments," Planning Research Evaluation Services.
- Rosenbaum, Paul R. and Donald B. Rubin. 1983. "The Central Role of the Propensity Score in Observational Studies for Causal Effects," *Biometrika* 70 (1), 41-55.
- Rosenbaum, Paul R. and Donald B. Rubin. 1985. "The Bias due to Incomplete Matching," *Biometrika* 41 (1), 103-116.
- Slavin, Robert E. and Cynthia Lake. 2008. "Effective Programs in Elementary Mathematics: A Best-Evidence Synthesis," *Review of Educational Research*, 78 (3), 427-515.
- Smith, Jeffrey and Petra Todd. 2005. "Rejoinder," *Journal of Econometrics* 125 (2), 365-375.
- Tulley, Michael. 1989. "The Pros and Cons of State-Level Textbook Adoption," *Publishing Research Quarterly*, 5 (2).
- U.S. Department of Health and Human Services, Administration for Children and Families. 2010. *Head Start Impact Study*. Final Report. Washington, DC.
- What Works Clearinghouse. 2007. Topic Report: Elementary School Math. Available at: http://ies.ed.gov/ncee/wwc/reports/elementary_math/topic
- Zhao, Zhong. 2004. "Using Matching to Estimate Treatment Effects: Data Requirements, Matching Metrics, and Monte Carlo Evidence" *The Review of Economics and Statistics* 86 (1) 91-107.

Figure 1. Probability Density Functions for Estimated Propensity Scores for Treatment and Control Units on the Common Support in Each Comparison Using 2001 Data (Solid Lines are Treatment Densities, Dashed Lines are Control Densities).

Treatment: Silver-Burdett Ginn (B) Control: Saxon (A)



Treatment: Scott-Foresman (C) Control: Saxon (A)



Treatment: Scott-Foresman (C) Control: Silver-Burdett Ginn (B)

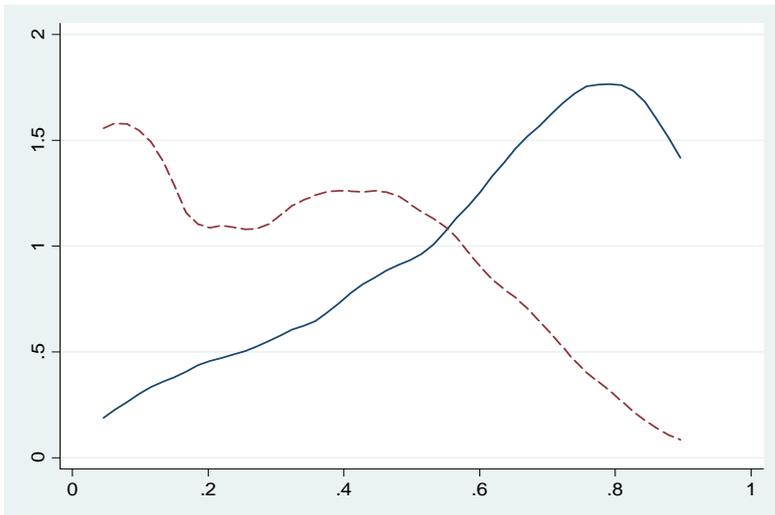


Table 1. Average Characteristics of Schools and Districts, by Adopted Curriculum (1997 values)

	Sample Average	Saxon (A)	Silver (B)	Scott (C)
<u>School-Level Outcomes</u>				
Attendance Rate	96.2	96.3 ^a	96.1 ^a	96.3
Grade-3 Math Test Score	496.6	496.5	494.2 ^c	499.7 ^c
Grade-3 Language Test Score	496.7	496.1	495.8	498.7
<u>School-Level Characteristics</u>				
<i>Percent Free Lunch</i>	27.4	24.7 ^{a,b}	28.5 ^a	30.5 ^b
<i>Percent Reduced Lunch</i>	6.7	7.1 ^a	6.3 ^a	6.6
<i>Percent Not Fluent in English</i>	1.2	0.7 ^a	1.7 ^a	1.2
<i>Percent Language Minority</i>	2.6	1.8 ^a	3.9 ^a	2.6
<i>Percent White</i>	91.3	95.4 ^{a,b}	88.0 ^a	88.4 ^b
<i>Percent Black</i>	5.6	2.3 ^{a,b}	7.2 ^{a,c}	9.2 ^{b,c}
<i>Percent Asian</i>	0.7	0.4 ^{a,b}	0.9 ^a	1.1 ^b
<i>Percent Hispanic</i>	2.2	1.8 ^{a,b}	3.7 ^{a,c}	1.1 ^{b,c}
<i>Percent American Indian</i>	0.2	0.1	0.2	0.2
<i>Enrollment (logs)</i>	5.95	5.92	5.97	5.96
N (Schools)	716	311	221	184
<u>District-Level Outcomes</u>				
Attendance Rate	95.8	95.7 ^b	95.8	96.1 ^b
Grade-3 Math Test Score	498.1	495.8 ^b	498.1 ^{a,c}	506.9 ^b
Grade-3 Language Test Score	498.9	496.5 ^{a,b}	500.6 ^a	505.6 ^b
<u>District-Level Characteristics</u>				
<i>Enrollment (logs)</i>	7.72	7.60 ^{a,b}	7.8 ^{a,c}	8.2 ^{b,c}
<i>Total Per-Pupil Revenue (logs)</i>	8.83	8.81 ^b	8.84	8.87 ^b
<i>Local Per-Pupil Revenue (logs)</i>	7.24	7.18 ^b	7.24 ^c	7.47 ^{b,c}
<u>Census Information (District Level)</u>				
Median Household Income (logs)	10.81	10.8 ^{a,b}	10.8 ^{a,c}	10.9 ^{b,c}
Share of Population with Low Education	18.2	18.8 ^b	19.2 ^c	14.3 ^{b,c}
N (Districts)	213	124	56	33

^a Indicates statistically significant difference at the 10% level between Saxon and Silver-Burdett Ginn adopters.

^b Indicates statistically significant difference at the 10% level between Saxon and Scott-Foresman adopters.

^c Indicates statistically significant difference at the 10% level between Silver-Burdett Ginn and Scott-Foresman adopters.

Note: The propensity-score specification also uses italicized information from 1998 – differences in means for these years are not reported for brevity.

Table 2. Balancing details for the 32 covariates included in the multinomial probit specification.

	1992	1993	1994	1995	<i>1996</i>	<i>1999</i>	2000	2001	2002	2003	2004	2005	2006	2007	2008
<u>Silver (B) to Saxon (A)</u>															
# of unbalanced covariates (p-values below 0.05/0.10)	1/4	0/4	0/3	0/2	<i>0/2</i>	<i>0/2</i>	0/0	0/0	0/2	0/1	0/0	0/0	0/2	1/2	1/3
Average p-value from balancing tests, all covariates	0.55	0.55	0.55	0.55	<i>0.55</i>	<i>0.55</i>	0.56	0.56	0.56	0.56	0.57	0.58	0.58	0.57	0.53
Mean Standardized Diff	3.4	2.9	3.8	3.9	<i>3.3</i>	<i>3.5</i>	3.6	3.3	3.5	3.3	3.0	3.7	3.9	4.2	4.7
<u>Scott (C) to Saxon (A)</u>															
# of unbalanced covariates (p-values below 0.05/0.10)	2/4	4/6	3/6	4/6	<i>3/5</i>	<i>3/6</i>	3/5	3/5	3/6	5/5	3/5	4/5	5/5	5/5	3/4
Average p-value from balancing tests, all covariates	0.48	0.49	0.49	0.48	<i>0.50</i>	<i>0.49</i>	0.48	0.48	0.49	0.44	0.45	0.46	0.47	0.47	0.46
Mean Standardized Diff	8.5	5.9	6.1	6.2	<i>6.1</i>	<i>6.0</i>	6.6	6.3	6.0	7.2	7.6	7.4	7.6	7.6	8.2
<u>Scott (C) to Silver (B)</u>															
# of unbalanced covariates (p-values below 0.05/0.10)	2/5	2/5	2/5	2/5	<i>1/3</i>	<i>0/4</i>	0/4	1/4	1/4	0/4	0/4	0/4	1/4	3/5	2/4
Average p-value from balancing tests, all covariates	0.48	0.47	0.44	0.46	<i>0.51</i>	<i>0.50</i>	0.50	0.49	0.50	0.52	0.54	0.54	0.50	0.51	0.54
Mean Standardized Diff	9.6	10.2	8.8	9.3	<i>9.3</i>	<i>9.2</i>	9.5	9.7	9.8	10.1	10.6	10.6	10.8	10.6	10.8

Note: Columns in italics are for years that are contiguous to the years from which the matching criteria are drawn. Results reported using the samples of treatments and controls that are on the common support in each year for the kernel-matching estimators. The numbers of covariates that fail the balancing tests at the 5 percent level are a subset of those that fail at the 10 percent level.

Table 3. Kullback-Leibler (KL) Information Criteria by Curriculum Comparison.

<u>Comparison</u>	<u>KL Information Criterion</u>
B and A	0.63
C and A	1.58
C and B	0.91

Note: Based on 2001 sample of schools.

Table 4. Estimates of Math Curricular Effectiveness on Math Test Scores for Partially and Fully-Exposed Cohorts. All Comparisons.

	1999	2000	2001	2002	2003	2004	2005	2006
<u>Treatment: B Control: A</u>								
OLS	0.124 (0.105)	0.162 (0.101)	0.354 (0.095)**	0.356 (0.087)**	0.374 (0.099)**	0.268 (0.131)*	0.293 (0.104)**	0.250 (0.110)*
Kernel Matching	0.144 (0.139)	0.191 (0.145)	0.396 (0.125)**	0.400 (0.102)**	0.401 (0.116)**	0.279 (0.135)*	0.318 (0.130)**	0.253 (0.132)†
LLR Matching	0.154 (0.184)	0.173 (0.153)	0.397 (0.117)**	0.398 (0.122)**	0.398 (0.126)**	0.273 (0.147)†	0.308 (0.138)*	0.259 (0.134)†
<u>Treatment: C Control: A</u>								
OLS	0.130 (0.120)	-0.013 (0.134)	0.187 (0.104)†	0.261 (0.096)**	0.208 (0.110)†	0.014 (0.119)	0.109 (0.104)	0.183 (0.119)
Kernel Matching	0.117 (0.169)	0.010 (0.184)	0.215 (0.158)	0.270 (0.122)*	0.272 (0.124)*	-0.042 (0.118)	0.113 (0.187)	0.150 (0.187)
LLR Matching	0.128 (0.248)	0.135 (0.269)	0.169 (0.220)	0.295 (0.156)†	0.301 (0.199)	0.032 (0.224)	0.085 (0.243)	0.141 (0.354)
<u>Treatment: C Control: B</u>								
OLS	0.008 (0.100)	-0.160 (0.123)	-0.100 (0.117)	-0.186 (0.129)	-0.285 (0.166)†	-0.271 (0.162)†	-0.181 (0.129)	-0.083 (0.139)
Kernel Matching	-0.088 (0.255)	-0.237 (0.274)	-0.165 (0.230)	-0.164 (0.183)	-0.331 (0.193)†	-0.275 (0.204)	-0.208 (0.239)	-0.148 (0.249)
LLR Matching	-0.072 (0.657)	-0.230 (0.531)	-0.149 (0.652)	-0.122 (0.898)	-0.302 (0.358)	-0.236 (0.219)	-0.163 (0.484)	-0.163 (0.798)
N(A)	309	307	307	305	300	294	286	287
N(B)	220	219	219	213	213	212	210	207
N(C)	184	182	182	181	176	174	169	163

Notes: Bolded columns are for the fully-exposed cohorts. Matching estimators impose the common support restriction. Standard errors in parentheses are clustered at the district level for all estimates, and bootstrapped using 250 repetitions for the matching estimators. N(A) refers to the number of schools in our sample that use curriculum A (Saxon), and similarly for N(B) and N(C). Estimates are provided in terms of the school-level distribution of test scores. Dividing the estimates by three will roughly return estimates in the more common student-level metric (see the discussion in the text and Appendix Table A.2 for more details).

** Denotes statistical significance at the 1 percent level or better

* Denotes statistical significance at the 5 percent level or better

† Denotes statistical significance at the 10 percent level or better

Table 5. Average 2004 Curriculum Adoptions in Math by District for the Four Most Common Curricula from the 2004 Adoption Cycle.

		1998 Uniform Math Adoptions – Grades 1 Through 3				
		All	Saxon (A)	Silver-Burdett Ginn (B)	Scott-Foresman (C)	Other
<u>2004 Math Adoptions</u>						
Grade 1						
	Saxon	0.48	0.76	0.25	0.12	0.21
	Harcourt	0.19	0.07	0.32	0.35	0.24
	Houghton Mifflin	0.10	0.06	0.11	0.21	0.15
	Scott-Foresman	0.09	0.07	0.07	0.15	0.18
Grade 2						
	Saxon	0.48	0.77	0.25	0.09	0.24
	Harcourt	0.19	0.08	0.32	0.35	0.21
	Houghton Mifflin	0.10	0.06	0.11	0.21	0.15
	Scott-Foresman	0.09	0.05	0.07	0.18	0.18
Grade 3						
	Saxon	0.48	0.76	0.23	0.09	0.24
	Harcourt	0.18	0.08	0.32	0.35	0.21
	Houghton Mifflin	0.12	0.07	0.14	0.21	0.15
	Scott-Foresman	0.09	0.06	0.05	0.21	0.15
Grade 4						
	Saxon	0.47	0.73	0.21	0.12	0.24
	Harcourt	0.18	0.09	0.30	0.35	0.21
	Houghton Mifflin	0.12	0.09	0.12	0.18	0.15
	Scott-Foresman	0.11	0.07	0.11	0.21	0.15
Grade 5						
	Saxon	0.47	0.74	0.21	0.18	0.22
	Harcourt	0.18	0.08	0.30	0.32	0.22
	Houghton Mifflin	0.10	0.07	0.11	0.18	0.16
	Scott-Foresman	0.11	0.08	0.11	0.18	0.16
N		286	128	57	34	33

Notes: N indicates the number districts where we observe a 2004 math-curriculum adoption and at least one grade-3 math test score between 1998 and 2008. The “other” category includes all districts that did not adopt any of the “big three” curricula in any grade during the 1998 adoption cycle. Districts that adopted at least one of the big-three curricula *non-uniformly* during the 1998 adoption cycle are included only in the “all” category.

Table 6. Falsification Estimates of Math Curricular Effectiveness, Estimated Using Math Test Scores for Grade-3 Cohorts Who Were Never Exposed to the Curricula of Interest. All Comparisons.

	1992	1993	1994	1995	1996	2007	2008
<u>Treatment: B Control: A</u>							
Kernel Matching	-0.120 (0.112)	0.072 (0.135)	-0.019 (0.120)	0.079 (0.137)	0.094 (0.129)	0.091 (0.117)	0.192 (0.127)
<u>Treatment: C Control: A</u>							
Kernel Matching	-0.326 (0.162)*	-0.046 (0.174)	-0.011 (0.146)	-0.035 (0.186)	-0.045 (0.153)	-0.020 (0.157)	-0.050 (0.270)
<u>Treatment: C Control: B</u>							
Kernel Matching	-0.171 (0.274)	0.077 (0.277)	0.032 (0.237)	0.072 (0.294)	-0.066 (0.280)	-0.147 (0.202)	-0.235 (0.263)
N(A)	301	304	304	306	308	284	280
N(B)	209	210	213	216	220	205	201
N(C)	179	179	182	182	182	163	162

Notes: Matching estimators impose the common support restriction. Standard errors in parentheses are clustered at the district level and bootstrapped using 250 repetitions. See the notes for Table 4 for details on how to interpret the estimates.

** Denotes statistical significance at the 1 percent level or better

* Denotes statistical significance at the 5 percent level or better

† Denotes statistical significance at the 10 percent level or better

Table 7. Falsification Estimates of Math Curricular Effectiveness, Estimated Using Math Test Scores for Grade-6 Cohorts who were Never Exposed to the Curricula of Interest. Comparison of B and A only.

	1992	1993	1994	1995	1996	1999	2000	2001
<u>Treatment: B Control: A</u>								
Kernel Matching	-0.126 (0.155)	-0.290 (0.165)†	-0.055 (0.158)	-0.133 (0.139)	0.045 (0.142)	0.016 (0.177)	-0.190 (0.151)	-0.100 (0.130)
N(A)	205	208	213	213	218	212	205	204
N(B)	117	118	122	125	127	122	120	120

Notes: Matching estimators impose the common support restriction. Standard errors in parentheses are clustered at the district level and bootstrapped using 250 repetitions. See the notes for Table 4 for details on how to interpret the estimates.

** Denotes statistical significance at the 1 percent level or better

* Denotes statistical significance at the 5 percent level or better

† Denotes statistical significance at the 10 percent level or better

Table 8. Estimates of Math Curricular Effectiveness, Estimated Using Reading Test Scores for all Grade-3 Cohorts. All Comparisons.

	1992	1993	1994	1995	1996	1999	2000	2001	2002	2003	2004	2005	2006	2007	2008
<u>Treatment: B Control: A</u>															
Kernel Matching	-0.152 (0.110)	-0.036 (0.131)	-0.078 (0.130)	0.082 (0.129)	0.135 (0.141)	0.160 (0.141)	0.186 (0.146)	0.197 (0.146)	0.228 (0.123)†	0.150 (0.120)	0.084 (0.127)	0.027 (0.148)	0.044 (0.122)	-0.084 (0.118)	0.069 (0.117)
<u>Treatment: C Control: A</u>															
Kernel Matching	-0.200 (0.156)	-0.126 (0.175)	-0.107 (0.178)	-0.154 (0.206)	-0.161 (0.178)	0.023 (0.207)	0.048 (0.237)	0.036 (0.221)	-0.043 (0.159)	0.037 (0.177)	-0.030 (0.205)	-0.080 (0.205)	0.184 (0.208)	0.028 (0.205)	0.084 (0.212)
<u>Treatment: C Control: B</u>															
Kernel Matching	-0.023 (0.294)	0.149 (0.305)	0.118 (0.268)	0.009 (0.288)	-0.179 (0.290)	-0.172 (0.281)	-0.143 (0.321)	-0.166 (0.289)	-0.222 (0.245)	-0.125 (0.259)	-0.095 (0.213)	-0.020 (0.297)	0.113 (0.297)	0.065 (0.262)	-0.014 (0.303)
N(A)	301	304	304	306	308	309	307	307	305	300	294	286	287	284	280
N(B)	209	210	213	216	220	220	219	219	213	213	212	210	207	205	201
N(C)	179	179	182	182	182	184	182	182	181	176	174	169	163	163	162

Notes: Bolded columns are for the fully-exposed cohorts. Matching estimators impose the common support restriction. Standard errors in parentheses are clustered at the district level and bootstrapped using 250 repetitions. See the notes for Table 4 for details on how to interpret the estimates.

** Denotes statistical significance at the 1 percent level or better

* Denotes statistical significance at the 5 percent level or better

† Denotes statistical significance at the 10 percent level or better

Table 9. Persistence Effects. Estimated Curriculum Effects for Grade-6 Cohorts who were Partially or Fully Exposed. Comparison of B and A only.

	2002	2003	2004	2005	2006	2007	2008
<u>Treatment: B Control: A</u>							
Kernel Matching	-0.064 (0.151)	0.141 (0.146)	0.156 (0.199)	0.077 (0.173)	0.007 (0.150)	-0.023 (0.169)	-0.016 (0.159)
N(A)	200	189	174	165	163	160	156
N(B)	118	115	105	101	97	94	93

Notes: Bolded columns are for the fully-exposed cohorts. Matching estimators impose the common support restriction. Standard errors in parentheses are clustered at the district level and bootstrapped using 250 repetitions. See the notes for Table 4 for details on how to interpret the estimates.

** Denotes statistical significance at the 1 percent level or better

* Denotes statistical significance at the 5 percent level or better

† Denotes statistical significance at the 10 percent level or better

Appendix A
Supplementary Tables

Appendix Table A.1. Data Sample Details.

	<u>Schools</u>	<u>% of Universe</u>	<u>Districts</u>	<u>% of Universe</u>
Universe*	1115		294	
<u>Missing Information:</u>				
District-reported curriculum adoption	3	0.3	3	1.0
District outcome variables (1997)	2	0.2	2	0.7
School outcome variables (1997)	23	2.2	1	0.3
District finance/enrollment data (1997, 1998)	2	0.2	1	0.3
School enrollment/demographic data (1997, 1998)	82	7.3	12	4.0
Did not use one of the primary curricula in grades one, two or three	211	18.9	38	12.9
Used only primary curricula, but did not uniformly adopt	76	6.8	24	8.2
<i>Final Sample</i>	<i>716</i>	<i>64.2</i>	<i>213</i>	<i>72.4</i>

* The universe consist of those schools and districts for which any information was reported in 1997, and at least one grade-3 math test score was reported for an exposed cohort (1999-2006).

Appendix Table A.2. Scaling Factors Used to Convert Estimation Metric from School-Level Distribution to Individual-Level Distribution for Grade-3 Math Scores.

Year	Standard Deviation of Distribution of School Scores	Standard Deviation of Distribution of Individual Scores	Approximate Scaling Factor
1992	2.8	N/A	N/A
1993	2.9	N/A	N/A
1994	2.8	N/A	N/A
1995	2.8	N/A	N/A
1996	1.9	N/A	N/A
1999	21.3	N/A	N/A
2000	20.5	61.0	0.34
2001	21.0	61.4	0.34
2002	19.9	59.7	0.33
2003	20.7	60.9	0.34
2004	22.5	63.1	0.36
2005	21.0	62.2	0.34
2006	20.0	64.3	0.31
2007	21.3	65.4	0.33
2008	22.5	63.7	0.35

Appendix B Bandwidth Selection

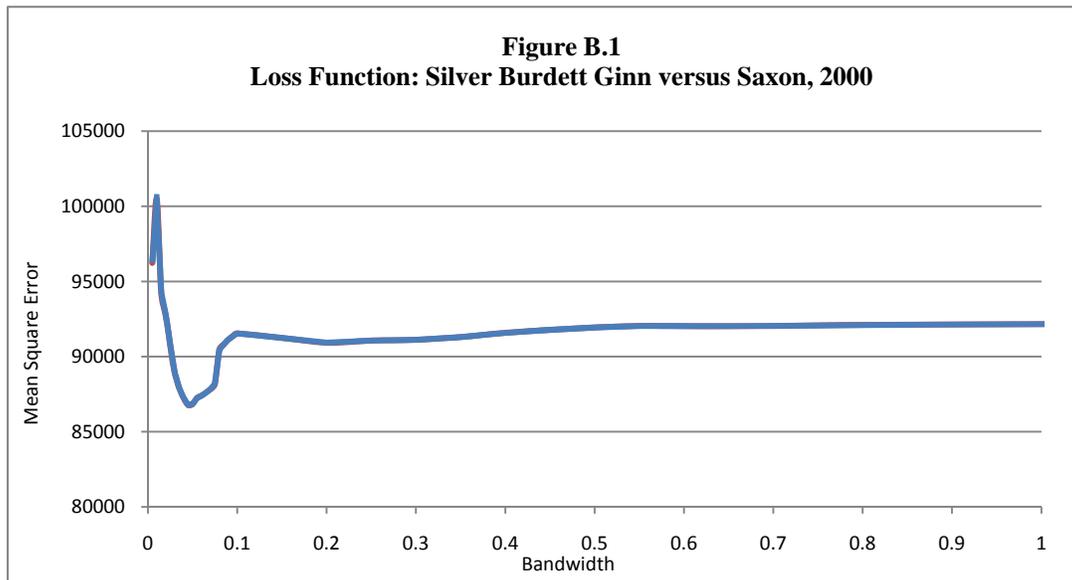
We use standard leave-one-out cross validation (C-V) to obtain fixed bandwidths for the kernel and LLR matching estimators. The grid search for kernel and LLR matching is over the range (0.005, 2.0). Using Frölich’s (2004) notation, the C-V approach selects the optimal bandwidth, h_{CV} , by solving the following minimization problem for control observations:³²

$$h_{CV} = \arg \min_{(h)} \sum_{q=1}^Q (Y_q - \hat{m}_{-q}(p_q))^2$$

where q indexes the sample of control units, Y is the outcome (test score) and $\hat{m}_{-q}(p_q)$ is the estimate of the mean outcome among the control observations, excluding observation q , conditional on the estimated propensity score for unit q .

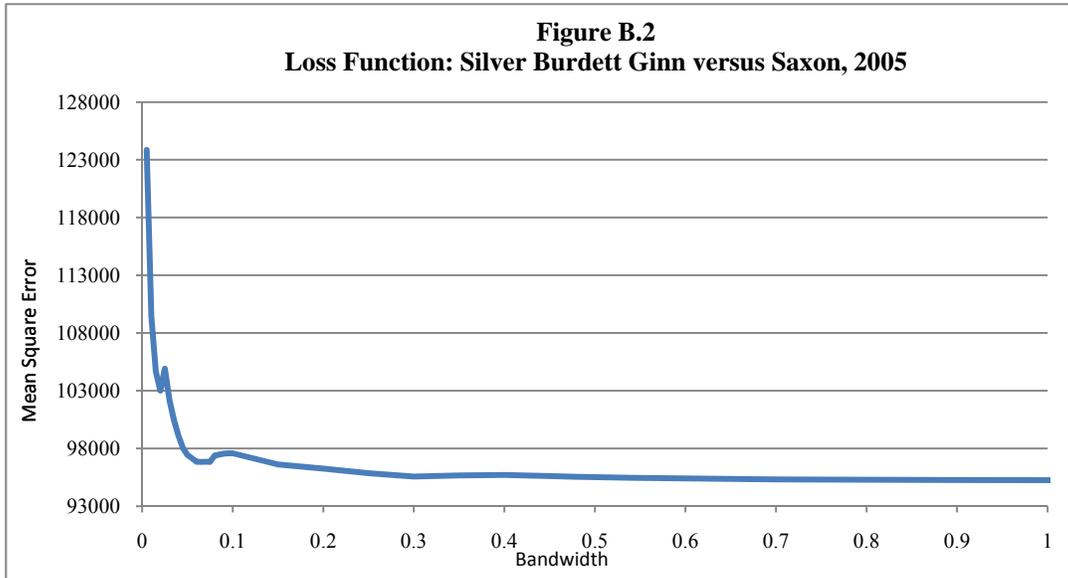
As has been reported in other contexts (see, for example, Ludwig and Miller, 2007), the loss function used to select the bandwidth is fairly flat in most of our comparisons. Therefore, we use a combination of conventional C-V and “visual inspection” to identify the appropriate bandwidth for each of our matching estimators.

First, Figure B.1 illustrates a case where cross-validation produces a clear bandwidth choice at the global minimum of the loss function, for our comparison between B and A in 2000 using the kernel matching estimator. In this case we use the bandwidth at the global minimum, 0.048.



³² In our case the definition of “treatment” and “control” is arbitrary and therefore, we could use either group. We use the largest group in each comparison as the control group.

Next, Figure B.2 illustrates a case where cross-validation suggests an optimal bandwidth at the edge of our grid search, for our comparison between B and A in 2005 using the kernel matching estimator. For this comparison we use a bandwidth of 0.062, which occurs just prior to the narrowly upward sloping portion of the curve.



We describe our bandwidth selection procedure for the comparison in Figure B.2 as a combination of cross-validation and visual inspection. Because the flat region of the curve has a mild negative slope, the mechanical application of the C-V procedure would produce a bandwidth at the edge of our grid search, 2.0. However, by visual inspection, we can see that there is very little difference in the loss function between the bandwidth determined mechanically by the C-V procedure and a much narrower bandwidth selected after the initial drop in the loss function. We ultimately use the narrower bandwidth in this and similar cases because the efficiency gains associated with the wider bandwidth will be minimal, and the narrower bandwidth should reduce bias in the estimates.