

# Mathematics Curriculum Effects on Student Achievement in California

**Cory Koedel**

**Diyi Li**

*University of Missouri*

**Morgan S. Polikoff**

**Tenice Hardaway**

*University of Southern California*

**Stephani L. Wrabel**

*The RAND Corporation*

*We estimate relative achievement effects of the four most commonly adopted elementary mathematics textbooks in the fall of 2008 and fall of 2009 in California. Our findings indicate that one book, Houghton Mifflin’s California Math, is more effective than the other three, raising student achievement by 0.05 to 0.08 student-level standard deviations of the Grade 3 state standardized math test. We also estimate positive effects of California Math relative to the other textbooks in higher elementary grades. The differential effect of California Math is educationally meaningful, particularly given that it is a schoolwide effect and can be had at what is effectively zero marginal cost.*

**Keywords:** *curriculum, econometric analysis, economics of education, educational policy, mathematics education, quasiexperimental analysis*

SEVERAL recent experimental and quasiexperimental studies point toward differences in curriculum materials having educationally meaningful effects on student achievement (Agodini, Harris, Atkins-Burnett, Heaviside, & Novak, 2010; Bhatt & Koedel, 2012; Bhatt, Koedel, & Lehmann, 2013). Chingos and Whitehurst (2012) argue that relative to other potential educational interventions—and in particular, human resource interventions—making better-informed decisions about curriculum materials represents an easy, inexpensive, and quick way to raise student achievement. However, the extent to which educational administrators can improve achievement by selecting better curriculum materials is hampered by a general lack of information. Given the wide variety of materials from which decision makers can choose, and the wide variety of implementation contexts (e.g., high-/low-poverty schools, states with different curricular goals and assessments, etc.), the handful of available efficacy studies is far from sufficient to inform those charged with selecting curriculum materials on behalf of students.

We contribute to the sparse literature on the efficacy of curriculum materials by leveraging unique school-level data on textbook adoptions to estimate the relative effects on student achievement of four commonly used elementary

mathematics textbooks in California (we refer to curriculum materials as “curricula” and “textbooks” interchangeably throughout our study). In addition to adding to the literature that narrowly evaluates the effects of curriculum materials, our study also contributes to a larger literature on the effects of curriculum interventions broadly defined. For example, recent studies by Clotfelter, Ladd, and Vigdor (2015); Cortes, Goodman, and Nomi (2015); Domina, McEachin, Penner, and Penner (2015); and Dougherty, Goodman, Hill, Litke, and Page (2015) examine curricular interventions that intensify and/or modify the timing of exposure to mathematics course work. Jackson and Makarin (2016) study the effects of making “off-the-shelf” supplemental mathematics lessons available to teachers online. All of these studies identify large effects of curriculum-based interventions on student achievement. Echoing the sentiments of Chingos and Whitehurst (2012) and Bhatt and Koedel (2012), Jackson and Makarin conclude that their “highly scalable” intervention is more cost-effective than alternative policies aimed at improving teacher quality.

Textbook adoptions in California are reported by individual schools as a requirement of the 2004 *Eliezer Williams et al. v. State of California et al.* court ruling and resulting



legislation (<http://www.cde.ca.gov/eo/ce/wc/wmslawsuit.asp>). The plaintiff argued that low-income students do not have access to the same high-quality resources available to their higher-income peers. As a result of the *Williams* ruling, each school in the state is required to report on the presence of various educational resources, including textbooks. These data are kept in School Accountability Report Cards (SARCs) as PDF files available online from the California Department of Education (CDE). We manually collect textbook data from schools' SARCs and merge textbook adoptions with a longitudinal data file containing information about school achievement and characteristics. We use the merged file to perform a quasiexperimental evaluation of curriculum effects on Grade 3 state standardized assessments.

Our results indicate that one elementary mathematics textbook—*California Math*, published by Houghton Mifflin—outperformed the other three popular textbooks during the period we study. Specifically, we estimate that *California Math* increased student test scores by 0.05 to 0.08 student-level standard deviations on the Grade 3 test relative to the alternatives. We extend our analysis into Grades 4 and 5 and find that *California Math* increased math achievement in these grades as well, particularly in Grade 5. The effect of *California Math* is educationally meaningful, especially given the scope of the intervention and low cost of implementation. With regard to scope, curriculum effects apply on average across entire cohorts of students in schools. With regard to cost, as noted by Bhatt and Koedel (2012) and Chingos and Whitehurst (2012), the marginal cost of choosing one textbook over another is so small that it is effectively zero.<sup>1</sup>

The differential curriculum effects that we document in California are on the lower end of the range of estimates reported in similar recent studies, which have been between 0.08 and 0.17 student-level standard deviations (Agodini et al., 2010; Bhatt et al., 2013; Bhatt & Koedel, 2012). The fact that estimated differences in curriculum effects in California are smaller than in the handful of locales where other, similar evaluations have been conducted is interesting and worthy of further exploration. This could be due to differences in the curricula studied, evaluation contexts (including differences in the assessments used to gauge impact), or simply sampling variability. Ideally, our efficacy estimates could be compared to a much larger set of estimates for the same and different curricula, in similar and different evaluation contexts, to gain further insight into why the effect sizes vary. However, given that so few states collect textbook adoption data, and correspondingly there are so few studies of curricular efficacy, we can do little more than speculate as to the source of the differential results. Detailed investigations into the specific characteristics of curricula and the contexts in which they are used (e.g., different state standards) that drive differential curriculum effects are not currently possible because they are underidentified—that is,

there are too many potential explanations and too few achievement-based estimates of curricular efficacy. Our inability to contextualize our findings within a larger literature—which essentially does not exist—highlights the frustrating lack of information nationally about the effectiveness of different sets of curriculum materials.

## Background and Data

California has what is best described as a partially centralized curriculum adoption process. The important centralized feature is that the state initiates the process for a particular subject in a particular year by assembling a list of “state-approved” curriculum materials through an extensive process that is documented in state reports (e.g., CDE, 2009) and described in more detail below. The list goes out to districts, but it is advisory only. Districts can choose any curriculum materials they would like—on list or off—or they can choose not to adopt curriculum materials. Like other states with partially centralized adoption processes, districts in California adopt new curriculum materials in each subject on roughly the same schedule. In math, California districts have recently used a 6-year cycle (2008–2009 to 2014–2015), although again districts can choose when and whether to adopt. This cycle length is typical of other states. Districts are all prompted to move together by the state’s initiation of the adoption process, so the large majority of districts make adoption decisions in the years immediately following the state adoption.

We focus our analysis on elementary mathematics textbooks adopted in California schools in fall 2008 and fall 2009. Our curriculum materials data, which we collected from schools’ 2013 SARCs, include information on textbooks from this adoption cycle that were still in use in 2013 (only a small fraction of schools adopted a new textbook after fall 2009 and before the publication of the 2013 SARCs, which we drop; see Appendix A). The textbook adoption we study was intended for fall 2008, and the state-approved list was released in November of 2007, but based on data collected from individual schools’ SARCs, it appears that many schools and districts delayed the adoption 1 year. Thus, we refer to the adoption as occurring in 2009/2010 (for presentational convenience we refer to school years by the spring year throughout our study, e.g., 2009 for 2008–2009).

We merge information on schools’ curriculum adoptions with a longitudinal database containing school and district characteristics and achievement outcomes covering the school years 2003 to 2013, constructed based on publicly available data from the CDE. We supplement the CDE data with data from the U.S. Census on the median household income and education level in the local area for each school, linked at the zip code level. Achievement effects are estimated using school-average test scores on state standardized tests.<sup>2</sup> We focus most of the evaluation on Grade 3

achievement, which aligns our study with previous related work focusing on early primary grades (Agodini et al., 2010; Bhatt et al., 2013; Bhatt & Koedel, 2012). We also extend our analysis to examine curriculum effects on test scores in Grades 4 and 5.

Appendix Table A1 provides details about the construction of our analytic sample. There are several notable attrition points documented in the appendix. Most lost data are because of (a) incomplete information on the SARCs, (b) off-cycle curriculum adoptions, and (c) schools either using more than one textbook (treatment) in Grades 1 through 3 or reporting textbook usage in such a way that we cannot rule this out. With regard to the latter situation, although in principle these schools could be used to examine mixed-treatment effects, in practice there are too few observations for an effective analysis along these lines, so we simply drop them from the analytic sample.<sup>3</sup> Appendix A provides an extended discussion about the data.

After imposing the data restrictions, we are left with a sample of just over half of the schools in California. These schools clearly report which curriculum materials they use and use the same materials in Grades 1 through 3. Among them, 78% adopted one of these four textbooks: *enVision Math California*, published by Pearson Scott Foresman; *California Math*, published by Houghton Mifflin; *California Mathematics: Concepts, Skills, and Problem Solving*, published by McGraw Hill; and *California HSP Math*, published by Houghton Mifflin Harcourt. We focus our evaluation on these textbooks and the schools that adopted them. In total, this group initially included 2,281 California schools spread across 311 districts; however, after our analysis began, we also dropped data from the Los Angeles Unified School District (LAUSD) and Long Beach Unified School District (LBUSD). Both districts are much larger than all other districts in the state, which created comparability problems in our evaluation. After dropping LAUSD and LBUSD schools, our final analytic data set includes 1,878 California schools in 309 districts.<sup>4</sup>

Table 1 provides descriptive characteristics and sample sizes for all California schools in our initial universe and schools that were retained in our final analytic sample. We also report separate statistics for schools that adopted each of the four focal curricula. The initial universe of schools in column 1 includes all schools in the CDE data for which at least one Grade 3 test score is available during the years 2009 to 2013, school characteristics are available for either 2007 or 2008, and the highest grade is 8 or lower.<sup>5</sup> The table shows that schools in our analytic sample are negatively selected relative to all schools in the state but not substantially. Within our analytic sample, adopters of *California Math* are similar to, although somewhat more advantaged than, adopters of the other curricula. There is substantial distributional overlap in preadoption achievement and other school characteristics between *California*

*Math* adopters and the comparison schools, which facilitates our analysis as outlined below. This overlap is illustrated in Appendix Figure B1.

### *Focal Textbooks*

As noted above, the textbooks were adopted in either fall 2008 or fall 2009 (the state refers to these textbooks as being a part of the 2007 adoption cycle; see CDE, 2009). The adoption was to select books aligned with the state’s 1997 mathematics content standards and 2005 mathematics framework. The multistep adoption process, which is described in detail in the adoption report (CDE, 2009), included 14 content experts (university professors) and 141 instructional materials experts (K–12 educators) divided into 26 panels. The chosen books were required to meet criteria in five categories: mathematics content/alignment, program organization, assessment, universal access, and instructional planning and support. The final selections passed through a public comment period and were approved by the State Board of Education in winter 2007. The process is somewhat selective, as nearly a third of the curricula submitted by publishers were not approved.

In total, 10 textbooks for Grades K–3 were approved. We study four of these books, which we chose because they were the most popular. In addition to their popularity making these books the most policy-relevant ones to study, it also affords sufficient sample sizes to support our empirical evaluation. In Appendix B we briefly discuss the four studied books.

## **Empirical Strategy**

### *Methodological Overview*

We estimate the achievement effects of *California Math* relative to a composite alternative of the three other focal curricula using three related empirical strategies: (a) kernel matching, (b) common-support-restricted ordinary least squares (restricted OLS), and (c) “remnant”-based residualized matching. Despite the fact that adoption decisions commonly occur at the district level, we use schools as units of analysis in our study. There are several reasons for this. One is that although many districts are “uniform adopters” (also see Bhatt et al., 2013; Bhatt & Koedel, 2012)—that is, where all schools in the district adopt the same curriculum materials—some are not. Specifically, 16.5% of districts contain schools that adopt different curriculum materials in the elementary grades we study. Using schools as units of analysis allows us to include schools in non-uniform-adopting districts in a straightforward manner. Additional benefits of using schools instead of districts as the units of analysis in a similar evaluation are discussed by Bhatt and Koedel (2012).<sup>6</sup> Although we note these benefits of a school-level evaluation, we also cluster our standard errors at the district

TABLE 1

*Descriptive Statistics for California Schools, Our Full Analytic Sample, and by Textbook Adoption*

Variable	Within the analytic sample, by textbook						
	All schools	All schools without LAUSD or LBUSD	Analytic sample	<i>enVision Math California</i>	<i>California Math</i>	<i>California Mathematics: Concepts, Skills, and Problem Solving</i>	<i>California HSP Math</i>
<b>School outcomes</b>							
Preadoption Grade 3 math score	0.02	0.03	-0.03	-0.08	0.06	-0.09	-0.02
Preadoption Grade 3 ELA score	0.01	0.05	-0.03	-0.09	0.07	-0.10	-0.02
<b>School characteristics</b>							
% Female	48.7	48.7	48.7	48.5	48.9	48.8	48.7
% Economically disadvantaged	56.6	54.2	56.9	56.2	56.0	59.1	58.0
% English learner	29.3	28.3	29.5	30.2	28.0	29.9	30.3
% White	31.4	33.6	29.5	30.1	29.9	29.2	26.4
% Black	7.8	7.1	7.3	8.0	6.3	8.6	4.8
% Asian	8.4	8.8	7.5	7.6	7.2	7.5	8.4
% Hispanic	47.9	45.8	50.5	48.6	51.7	49.7	56.3
% Other	4.6	4.7	5.2	5.7	4.9	5.0	4.1
Enrollment	385.7	378.1	410.5	399.9	429.5	405.7	399.0
2008 adopter			50.2	49.0	53.7	53.5	36.2
<b>School-area characteristics (census)</b>							
Median household income (log)	11.0	11.0	10.8	10.7	10.9	10.8	10.9
Share low education	17.8	17.2	19.5	17.6	19.3	22.6	23.9
Share missing census data	3.1	3.3	1.8	2.7	1.2	0.8	0.0
<b>District outcomes</b>							
Preadoption Grade 3 math score	0.01	0.01	-0.02	0.03	0.00	-0.05	-0.05
Preadoption Grade 3 ELA score	0.02	0.02	-0.09	-0.03	-0.12	-0.09	-0.12
<b>District characteristics</b>							
Enrollment	5138.0	4438.9	5690.4	6404.0	6075.5	5279.0	4339.9
<i>n</i> (Schools)	5,494	4,931	1,878	710	602	389	177
<i>n</i> (Districts)	825	823	309	107	92	69	48

*Note.* The “all schools” sample is the universe of schools reported in Appendix Table A1. It includes schools in the California Department of Education data with characteristics from either 2007 or 2008, at least one Grade 3 test score from 2009 to 2013, and where the highest grade is 8 or lower. The descriptive statistics for the analytic sample in column 3 are a weighted average of the textbook-by-textbook statistics reported in columns 3 through 6. Note that some districts have a uniformly adopting school of more than one textbook; thus the sum of the district counts in the last four columns is greater than 309. LAUSD = Los Angeles Unified School District; LBUSD = Long Beach Unified School District; ELA = English language arts.

level throughout the analysis to reflect data dependence within districts across schools, including along the dimension of curriculum adoptions.

In the remainder of this section, we describe our methods within the context of our evaluation of Grade 3 test scores. The methods carry over directly when we extend our analysis to study test scores in Grades 4 and 5, as we discuss briefly when we present those results below.

*Matching.* Our matching estimators follow Bhatt and Koedel (2012) and draw on the larger matching literature to identify the approach best suited to our data (Caliendo & Kopeinig, 2008; Frölich, 2004). The key to identification is the conditional independence assumption (CIA), which requires potential outcomes to be independent of curriculum choice

conditional on observable information. Denoting potential outcomes by  $\{Y_0, Y_1\}$ , curriculum treatments by  $D \in \{0, 1\}$ , and  $X$  as a vector of (pretreatment) observable school, district, and local-area characteristics, the CIA can be written as

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

Conditional independence will not be satisfied if there is unobserved information influencing both treatments and outcomes. For example, if districts have access to information that is unobserved to the researcher,  $Z$ , such that  $P(D = 1 \mid X, Z) \neq P(D = 1 \mid X)$ , and the additional information in  $Z$  influences outcomes, matching estimates will be biased. We discuss the plausibility of the CIA in our context and provide evidence consistent with it being satisfied below.

We match schools using propensity scores (Rosenbaum & Rubin, 1983; Lechner, 2002). The propensity score model predicts whether each school adopted *California Math* as a function of a variety of school, district, and local-area characteristics. Specified as a probit, our propensity score model is as follows:

$$T_{sd} = \mathbf{X}_s \beta_1 + \mathbf{X}_d \beta_2 + \varepsilon_{sd}. \quad (2)$$

In Equation (2),  $T_{sd}$  is an indicator variable equal to 1 if school  $s$  in district  $d$  adopted *California Math* and 0 if it adopted one of the other focal curricula.  $\mathbf{X}_s$  and  $\mathbf{X}_d$  are vectors of school and district covariates, respectively, that include the variables listed in Table 1. For schools,  $\mathbf{X}_s$  includes preadoption student achievement in math and reading; the share of students by race, gender, language fluency, and socioeconomic disadvantage; school enrollment (cubic); and whether the school adopted new materials in 2008 or 2009. The vector  $\mathbf{X}_s$  also includes the log of median household income and the share of individuals over age 25 without a high school degree in the local area. These data are taken from the 2013 American Community Survey 5-year average (from the U.S. Census) and merged to schools by zip code.<sup>7</sup> The vector  $\mathbf{X}_d$  includes district-level preadoption achievement in math and reading, and enrollment (cubic).<sup>8</sup>

With the estimated propensity scores from Equation (2) in hand, we estimate the average treatment effect (ATE) of adopting *California Math*. Defining *California Math* as curriculum  $j$  and the composite alternative as curriculum  $m$ , where  $Y_j$  and  $Y_m$  are standardized test score outcomes for adopters of  $j$  and  $m$ , respectively, we estimate  $ATE_{j,m} \equiv E(Y_j - Y_m \mid D \in \{j, m\})$ . We use kernel matching estimators (with the Epanechnikov kernel), which construct the match for each “treated” school using a weighted average of “control” schools, and vice versa. The formula for our estimate of  $ATE_{j,m}$  is

$$\hat{\theta}_{j,m} = \frac{1}{N^S} \left[ \sum_{j \in N_j \cap S_p} \{Y_j - \sum_{m \in I_{0j} \cap S_p} W(j,m) Y_m\} - \sum_{m \in N_m \cap S_p} \{Y_m - \sum_{j \in I_{0m} \cap S_p} W(m,j) Y_j\} \right]. \quad (3)$$

In (3),  $N^S$  is the number of schools using  $j$  or  $m$  on the common support,  $S_p$ .  $I_{0j}$  indicates the schools that chose  $m$  in the neighborhood of observation  $j$ , and  $I_{0m}$  indicates the schools that chose  $j$  in the neighborhood of observation  $m$ . Neighborhoods are defined by a fixed-bandwidth parameter obtained via conventional cross-validation (as in Bhatt & Koedel, 2012).  $W(j, m)$  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 compute separate ATE estimates by year based on the distance from the adoption year using the formula in Equation (3). All of our standard errors are estimated via

bootstrapping using 250 replications and clustered at the district level (i.e., with district resampling). We omit a more detailed discussion of the matching estimators for brevity, but more information can be found in Caliendo and Kopeinig (2008); Heckman, Ichimura, and Todd (1997); and Mueser, Troske, and Gorislavsky (2007).

*Restricted OLS.* We also use restricted OLS models to estimate curriculum effects for schools on the common support of propensity scores. We use the same school and district characteristics taken from preadoption data in the OLS models as we use to match schools, allowing the coefficients to change over time as follows:

$$Y_{sdt} = \mathbf{X}_s \pi_{1t} + \mathbf{X}_d \pi_{2t} + T_{sd} \theta + u_{sdt}. \quad (4)$$

In Equation (4),  $Y_{sdt}$  is a Grade 3 math test score for school  $s$  in district  $d$  in year  $t$ ,  $\mathbf{X}_s$  and  $\mathbf{X}_d$  are the vectors of preadoption school/district characteristics that we use for matching as described above (these variables do not change over time),  $T_{sd}$  is an indicator equal to 1 if the school adopted *California Math*, and  $u_{sdt}$  is the error term. The coefficient vectors  $\pi_{1t}$  and  $\pi_{2t}$  allow the preadoption school and district characteristics to differentially predict achievement over time.

The OLS estimates are very similar to the matching estimates and rely on the same assumption of conditional independence for identification. The benefit of the OLS models is that they improve statistical precision by imposing a parametric form—linearity—on the outcome model. The cost is that if the linearity assumption is not justified, it could introduce bias (Black & Smith, 2004). In our application, where California schools and districts are diverse and we have small samples (at least by the standards of matching analyses, especially given the district clustering), the efficiency benefit of imposing linearity is substantial. This will become clear when we present our findings below. With regard to the potential for the linearity assumption to introduce bias into our estimates, we show results from falsification tests that provide no indication of bias in our OLS estimates.

*Remnant-Based Residualized Matching.* Our third approach uses remnant-based residualization as another way to improve statistical power. It blends aspects of the restricted-OLS and matching strategies. The fundamental idea, taken from Sales, Hansen, and Rowan (2014), is to pull in data from outside of the evaluation—that is, “remnant data”—to regression-adjust outcomes prior to matching. Sales et al. suggest several potential uses of remnant-based residualization; in our application, the appeal is that the procedure can remove noise from the outcome data before matching occurs, thereby improving the precision of our estimates. Our evaluation is particularly well suited for remnant-based residualization because we have access to substantial data from outside of the evaluation, for example, from schools

in California that use a textbook outside of the four focal curricula. We describe our method for using the remnant data in detail in Appendix C.

### *Conditional Independence*

All three approaches outlined above rely on the assumption of conditional independence for identification (the restricted-OLS and remnant-residualized matching methods further impose a functional form assumption on the outcome model to improve statistical power, utilizing either in-sample or out-of-sample data). Although conditional independence cannot be tested for directly, in this section we provide a brief intuitive case for why it may be a plausible assumption in our application. We also discuss falsification tests that we use to look for evidence of violations to the CIA.

One aspect of our evaluation that makes the CIA more plausible is that curriculum materials are adopted on behalf of large groups of students and teachers rather than being the product of individual choice. When individuals choose whether to seek treatment, characteristics that are difficult to observe, such as motivation, may influence treatment and outcomes. However, in the case of school- and district-level choices and conditional on the rich covariates to which we have access—preadoption test scores, demographic and socioeconomic status measures, and so on (see Table 1)—selection on unobservables that are uncorrelated with observables is more difficult. For example, consider two school districts that are similar demographically and located in zip codes with similar socioeconomic conditions. It is harder to imagine that there are substantial differences in group-average unobservable characteristics, like motivation or innate ability, across these districts that are not already accounted for by the group-level observed measures, certainly relative to the case of a treatment influenced by individual choice.

One could also argue that school- and district-level differences in teacher quality, which are not directly accounted for in our study, might lead to a violation of the CIA if teacher quality helps to determine curriculum adoptions. Research is quite clear that teacher quality affects student achievement (e.g., see Koedel, Mihaly, & Rockoff, 2015). However, many of the same arguments from the preceding paragraph apply—specifically, it would need to be the case that there are systematic differences in teacher quality across schools and districts after conditioning on the rich set of characteristics of schools and their local areas used in our study. It has also been documented in research that most of the variation in teacher quality occurs within schools (Aaronson, Barrow, & Sanders, 2007; Koedel & Betts, 2011), not across schools, let alone districts. The limited cross-school variation in teacher quality leaves less scope for systematic differences in the quality of teachers to lead to significant violations of the CIA.

As noted by Bhatt and Koedel (2012), perhaps the biggest conceptual threat to the CIA in curriculum evaluations is that some decision makers simply make better choices than others. For example, effective leaders might choose a more effective textbook and/or set up a more effective curriculum adoption process and also make other decisions that improve student outcomes. This would violate the CIA because “decision maker quality” is not observed in our data. Bhatt and Koedel discuss why this problem might be minor in practice. A key argument is that the curriculum adoption process is complex, and based on available documentation, it does not appear that any single decision maker has undue influence (Zeringue, Spencer, Mark, & Schwinden, 2010). This sentiment is supported by interviews of district administrators conducted by our research team.<sup>9</sup> However, we are not aware of any empirical evidence that can rule out the importance of district administrators in influencing curriculum adoptions, either directly or indirectly by establishing the adoption process. Therefore, although what little we know about the curriculum adoption process makes it seem less likely that differences in administrator quality will cause significant bias, such bias is arguably the biggest conceptual threat to the CIA in our application.

We aim to provide some formal evidence on the plausibility of the CIA by providing two different types of falsification estimates. The falsification estimates cannot be used to confirm the satisfaction of the CIA because, as noted above, it is not possible to confirm with certainty that the CIA is upheld. However, the falsification estimates can be used to look for evidence consistent with the CIA being violated. We describe the falsification tests in detail below, but the basic idea is to look for curriculum effects in situations where (a) we should not expect any effects at all or (b) we should expect small effects at most. If, for example, we estimate nonzero “curriculum effects” in situations where we know the effects should be zero, this would be a strong indication that the CIA is violated. Estimates from all of our falsification tests are as expected and provide no indication that the primary results are biased by unobserved selection. We elaborate on our falsification tests and their interpretation when we show the results below.

## **Results**

### *Six Pairwise Comparisons*

We compare *California Math* to a composite alternative of the three other focal curricula. To arrive at this final research design, we began by examining all six possible pairwise comparisons across the four curricula. After performing the six comparisons, it became clear that it would be difficult to obtain meaningful insight from them individually. Two issues arose: (a) covariate-by-covariate balance is mediocre in some of the pairwise comparisons with little scope for improvement given our small sample sizes (at least relative

to typical matching applications), and (b) statistical power is limited. The statistical-power issue is more problematic than we had anticipated because our point estimates suggest curriculum differences in California that are smaller than in previous, similar evaluations. Moreover, because of the diversity of curriculum materials adopted in California—which means that there are fewer districts adopting any single book—and the district-level clustering structure necessitated in the evaluation, our effective sample sizes in the pairwise comparisons are no larger in California than in previous studies in smaller states.

Despite their limitations, collectively, the pairwise comparisons point toward *California Math* being more effective than the other three textbooks. Moreover, they also point toward the other curricula having similar effects. It is these preliminary results that motivate our comparison of *California Math* to the composite alternative. Rebuilding the evaluation in this way is advantageous because it allows us to identify better matches for *California Math* schools and to perform a better-powered study of the effectiveness of *California Math* relative to the other three books. For interested readers, Appendix Table B1 presents disaggregated matching estimates for the six pairwise comparisons that led us to restructure our study to focus on *California Math*.<sup>10</sup>

#### *Comparison of California Math to the Composite Alternative*

*The Propensity Score.* The propensity score model in Equation (2) explains roughly 12% of the variance in curriculum adoptions between *California Math* and the composite alternative. The limited scope for observed selection into curriculum materials implied by this *R*-squared value is notable given that our covariates are strong predictors of achievement.<sup>11</sup> For interested readers, Appendix Table B2 reports results from the estimation of Equation (2) for our evaluation. The only statistically significant predictors of a *California Math* adoption are the linear, squared, and cubed district enrollment variables. Thus, collectively, the covariates do not predict adoptions of *California Math* well. The lack of predictive power in the selection model is consistent with qualitative accounts of the complexity of the curriculum adoption process and the lack of clear objectives and information to make decisions (Jobrack, 2011; Zeringue et al., 2010).

*Covariate Balance.* Table 2 presents information on covariate balance after matching. We report results for each year of the data panel separately, including both pre- and postadoption years. Subsequent tables follow a similar reporting format.

The results for each school are centered around the year of the curriculum adoption. In the case of a fall 2008 adoption, Year 1 indicates the 2008–2009 school year (the 1st

year the new book was used), Year 2 indicates the 2009–2010 school year, and so on; for a fall 2009 adoption, Year 1 indicates the 2009–2010 school year, and so on. We use data from 2 years preceding the adoption to match schools, as described above, so we do not perform any direct analysis in these years. Thus, the first preadoption year shown in Table 2 and subsequent tables is Year P3—3 years prior to adoption. For schools that first used the new books during the 2008–2009 school year, Year P3 is the 2005–2006 school year, Year P4 is 2004–2005, and so on.

Although we split out the data by year in Table 2, the years are strongly dependent. The practical implication is that balancing evidence from a 2nd year of data provides very little new information relative to what can be inferred from 1 year of data. Put differently, because the sample of schools is largely unchanged over time (except for changes due to school openings and closings and, small schools, data reporting issues) and the treatment designation does not change over time (i.e., adoptions are static), covariate balance should change very little from one year to the next. Nonetheless, for completeness, we show balancing results in Table 2 for each year.

As suggested by Smith and Todd (2005), we present results from several balancing tests. The results are presented tersely in Table 2 and expanded on in Appendices B and C. The first row of the table reports the number of unbalanced covariates using simple covariate-by-covariate *t* tests among the matched sample. These tests are the simplest balancing metrics we provide. There are 22 covariates in total, and none are unbalanced at the 5% level. In row 2, we report the average absolute standardized difference across all covariates (Rosenbaum & Rubin, 1985), which is small in all years, on the order of just 3% to 4%. Rows 3 and 4 show corroborating results from alternative, regression-based balancing tests proposed by Smith and Todd. These tests identify only marginally more covariates as unbalanced than would be expected by chance (one to three covariates); moreover, the average *p* values across all covariates from the Smith and Todd tests are essentially what would be expected from a balanced comparison (i.e.,  $\approx 0.50$ ).

Overall, we conclude that our comparison of *California Math* to the composite alternative is well balanced along the observable dimensions of our data. Again, we provide more detailed information about the balancing tests in Appendices B and C (Appendix B: supplementary balancing results, including covariate-by-covariate details; Appendix C: additional methodological information).

*Estimated Curriculum Effects for Grade 3.* Table 3 shows results for our comparison between *California Math* and the composite alternative for cohorts of students exposed to 1 to 3 years of the curriculum materials in Grade 3. The Year 1 results compare students who used the textbooks for Grade 3 only (and used previously adopted materials in Grades 1 and

TABLE 2

*Balancing Results for the Primary Comparison*

Variable	Year P6	Year P5	Year P4	Year P3	Year 1	Year 2	Year 3	Year 4
Treatment: <i>California Math</i>								
Control: Composite alternative								
No. unbalanced covariates, matched <i>t</i> tests (5%)	0	0	0	0	0	0	0	0
Mean absolute standardized difference of covariates (%)	2.8	3.0	3.0	3.6	3.5	3.4	3.1	3.6
No. unbalanced covariates, Smith-Todd regression tests (5%)	1	2	2	2	2	3	2	2
Average <i>p</i> value, Smith-Todd regression tests	0.50	0.50	0.51	0.48	0.48	0.46	0.50	0.49
No. of districts/schools ( <i>California Math</i> )	89/560	88/567	90/575	90/588	92/597	89/588	91/595	90/590
No. of districts/schools (composite alternative)	210/1,063	213/1,085	212/1,106	215/1,124	213/1,143	214/1,145	216/1,146	213/1,144

*Note.* There are 22 covariates included in the balancing tests. The sample size fluctuates year to year due to school openings and closings, and data reporting issues for small schools. Year 1 denotes the 1st year the new curriculum was adopted (e.g., the 2008–2009 school year for textbooks adopted in fall 2008), Year 2 denotes the 2nd year, and so on. Similarly, year P3 denotes the school year 3 years prior to the new curriculum being adopted (e.g., the 2005–2006 school year for textbooks adopted in fall 2008), Year P4 denotes the year 4 years prior, and so on. Note that there is a 2-year gap between Year P3 and Year 1. We use data from the two gap years to match schools as described in the text.

TABLE 3

*Effects of California Math on Grade 3 Mathematics Achievement for Exposed Cohorts Relative to the Composite Alternative, by Year After the Initial Adoption*

Variable	Year 1	Year 2	Year 3	Year 4
Treatment: <i>California Math</i>				
Control: Composite alternative				
Treatment effect: Kernel matching	0.063 (0.054)	0.083 (0.051)	0.061 (0.059)	0.070 (0.059)
Treatment effect: Restricted ordinary least squares	0.050 (0.019)**	0.064 (0.023)**	0.049 (0.023)**	0.058 (0.023)**
Treatment effect: Remnant-residualized matching	0.050 (0.020)**	0.065 (0.024)**	0.052 (0.024)**	0.060 (0.026)**
No. of districts/schools ( <i>California Math</i> )	92/597	89/588	91/595	90/590
No. of districts/schools (composite alternative)	213/1,143	214/1,145	216/1,146	213/1,144

*Note.* Standard errors are estimated by bootstrapping using 250 repetitions and clustered at the district level. Year 1 denotes the 1st year the new curriculum was adopted (e.g., the 2008–2009 school year for textbooks adopted in fall 2008), Year 2 denotes the 2nd year, etc. All estimates are converted from school-level standard deviation units to student-level standard deviation units by multiplying them by a factor of 0.45, which is the ratio of standard deviations of the school-average test score distribution to the student-level test score distribution averaged across our data panel, as reported in the text. This transformation has no bearing on the results qualitatively or quantitatively; the rescaling is performed only to improve comparability of our findings to those in other studies that report effect sizes in student-level standard deviation units.

\* $p \leq .10$ . \*\* $p \leq .05$ .

2), the Year 2 results show results for students who used the books in Grades 2 and 3, and the Year 3 and Year 4 results are for students who used the books in all three grades leading up to the Grade 3 test.

The point estimates from all three estimation strategies are similar and indicate effect sizes on the order of 0.05 to 0.08 student-level standard deviations of achievement.<sup>12</sup> However, the standard errors for the matching estimates are much larger than for the OLS or remnant-residualized estimates. The standard errors decrease using the latter two methods because the linear regression model removes

variation in outcomes attributable to observed covariates. The cost of the improved precision is that the linear specification may be wrong, which is a reason that the parametrically less-restrictive matching estimators are preferred conceptually. That said, our falsification tests below suggest that our use of the linear functional form to improve precision does not result in biased estimates.

It is somewhat surprising that the treatment effect estimates do not become more pronounced over time in Table 3. One might expect cohorts of students who are exposed to the curricula for all three years during Grades 1 through 3

(students in Year 3 or Year 4 in the table) to show larger test score differences than cohorts who are exposed for just 1 to 2 years (students in Year 1 or Year 2), but no such pattern emerges. There are a number of potential explanations. One possibility is that there is a dosage effect, but it is small enough that we lack the statistical power to detect it. Given the sizes of our standard errors, even in the OLS and remnant-residualized models, moderate dosage effects cannot be ruled out. Another possibility is that the most recent textbook is the dominant treatment. Given that even the Year 1 students used the new textbooks in Grade 3, which is the most recent year of instruction leading up to the Grade 3 test, it is possible that increased dosage in earlier grades is not important enough to show up in contemporary achievement results (if it matters at all). This explanation is consistent with numerous other studies showing fade-out in educational interventions (e.g., Chetty, Friedman, & Rockoff, 2014; Currie & Thomas, 2000; Deming, 2009; Krueger & Whitmore, 2001). Another explanation is that curriculum materials quality is not stable from grade to grade. We are not aware of any research that directly informs this hypothesis, but Bhatt et al. (2013) show that math curriculum effects can vary by subtopic, documenting at least one dimension of effect non-uniformity. Our analysis of Grades 4 and 5, shown below, is also suggestive of some grade-to-grade variability in the efficacy of *California Math*.

### Falsification Tests

In this section, we present results from falsification models designed to detect violations to our key identifying assumption, conditional independence. We estimate two types of falsification models. The first is a time-inconsistent model where we estimate curriculum effects on student achievement for cohorts of students who predate the adoption of the curriculum materials we study: specifically, students in the cohorts from 3 years preadoption (Year P3) to 6 six years preadoption (Year P6). If our matching and regression-adjusted models are resulting in truly balanced comparisons (on observed and *unobserved* dimensions), we would not expect to see achievement differences between cohorts of students in matched schools prior to the curriculum adoptions of interest. The second type of falsification model estimates math curriculum effects on contemporaneous achievement in English language arts (ELA). In these models, we cannot rule out nonzero curriculum effects because math curricula may have spillover effects, but we would anticipate smaller cross-subject effects.

One issue with the falsification models is that we do not know which curriculum materials were used by schools prior to the curriculum adoption we study. No such longitudinal data on curriculum materials adoptions exist, which points to an underlying problem with the state of curriculum data and research. We rely on lagged school- and district-level test scores to capture the impacts of previous curriculum

materials on achievement (and all other educational inputs that we do not observe, for that matter), despite our inability to directly observe these materials. Whether this strategy is sufficient is an empirical question, which our falsification tests are designed to inform. If lagged test scores (and our other controls) are not sufficient to control for previous curriculum effects, and if curriculum adoptions are correlated across cycles within schools (which seems likely, but again data are limited), serial correlation in adoptions would be expected to manifest itself in the falsification tests in the form of nonzero preadoption “curriculum effects.”

Table 4 shows the first set of falsification results from the time-inconsistent models of math achievement. Across all three estimation strategies and in all preadoption years, the false “effects” of *California Math* relative to the composite alternative are substantively small and far from statistical significance. This is as expected if our methods are sufficient to generate balanced comparisons. Table 5 shows the complementary falsification results using ELA achievement as the dependent variable. As is the case in Table 4, all of our estimates in Table 5 are small and statistically insignificant.<sup>13</sup> Figure 1 visually illustrates our treatment effect and falsification estimates side by side. The bars with asterisks are for estimates that are statistically distinguishable from zero at the 5% level.

### Extensions

#### *Results for Grades 4 and 5*

Figures 2 and 3 show results from the full replication of our methods applied to Grade 4 and Grade 5 test scores, respectively. We follow analogous procedures as in the Grade 3 analysis to produce the results.<sup>14</sup> Like with the sample we constructed for our analysis of Grade 3 scores, our Grade 4 and Grade 5 samples are well balanced between *California Math* and composite-alternative schools. This is not surprising because the samples are essentially the same.<sup>15</sup> We do not report the balancing details for the Grade 4 and Grade 5 samples for brevity, but they are available upon request. The falsification results for the higher grades are shown in Figures 2 and 3; as in Figure 1, they provide no indication of selection bias in our primary estimates.

Taken together, our findings in Figures 2 and 3 are broadly consistent with the interpretation that *California Math* outperformed the composite of the other three focal curricula in California. Specifically, the estimates in Grade 4 are all nominally positive and sometimes statistically significant, and the estimates in Grade 5 are larger than in Grade 3 and statistically significant in all posttreatment years, at least using restricted OLS and remnant-based residualization. There is no evidence to suggest a negative relative effect of *California Math* in any grade or year.

The effect sizes we estimate are largest in Grade 5, but it is somewhat puzzling that they are smallest in Grade 4. The up-and-down pattern of estimates holds even for cohorts

TABLE 4

*Falsification Results: California Math “Effects” on Grade 3 Mathematics Achievement for Cohorts of Students in Years Prior to the 2008/2009 Adoption Cycle*

Variable	Year P3	Year P4	Year P5	Year P6
Treatment: <i>California Math</i>				
Control: Composite alternative				
Treatment effect: Kernel matching	0.010 (0.053)	0.023 (0.051)	0.017 (0.054)	0.017 (0.056)
Treatment effect: Restricted ordinary least squares	0.002 (0.014)	0.015 (0.015)	0.006 (0.015)	0.010 (0.018)
Treatment effect: Remnant-residualized matching	0.001 (0.014)	0.018 (0.018)	0.014 (0.020)	0.011 (0.024)
No. of districts/schools ( <i>California Math</i> )	89/560	88/567	90/575	90/588
No. of districts/schools (composite alternative)	210/1,063	213/1,085	212/1,106	215/1,124

*Note.* Standard errors are estimated by bootstrapping using 250 repetitions and clustered at the district level. Year P3 denotes the school year 3 years prior to the new curriculum being adopted (e.g., the 2005–2006 school year for textbooks adopted in fall 2008), Year P4 denotes the year 4 years prior, and so on. Data from the 2 years preceding the adoption are used to match schools and thus not analyzed directly. All estimates are converted from school-level standard deviation units to student-level standard deviation units by multiplying them by a factor of 0.45, which is the ratio of standard deviations of the school-average test score distribution to the student-level test score distribution in math averaged across our data panel, as reported in the text. This transformation has no bearing on the results qualitatively or quantitatively; the rescaling is performed only to improve comparability of our findings to those in other studies that report effect sizes in student-level standard deviation units.

TABLE 5

*Falsification Results: California Math “Effects” on Grade 3 English Language Arts (ELA) Achievement for Exposed and Unexposed Cohorts*

Variable	Year P6	Year P5	Year P4	Year P3	Year 1	Year 2	Year 3	Year 4
Treatment: <i>California Math</i>								
Control: Composite alternative								
Treatment effect: Kernel matching	0.002 (0.064)	0.016 (0.060)	0.015 (0.058)	0.016 (0.057)	0.027 (0.061)	0.043 (0.056)	0.014 (0.066)	0.020 (0.064)
Treatment effect: Restricted ordinary least squares	–0.000 (0.016)	0.013 (0.016)	0.008 (0.014)	0.008 (0.013)	0.012 (0.017)	0.019 (0.021)	–0.003 (0.020)	0.004 (0.022)
Treatment effect: Remnant-residualized matching	–0.005 (0.021)	0.012 (0.018)	0.006 (0.015)	0.004 (0.015)	0.012 (0.017)	0.020 (0.022)	0.001 (0.020)	0.006 (0.026)
No. of districts/schools ( <i>California Math</i> )	89/560	88/567	90/575	90/588	92/597	89/588	91/595	90/590
No. of districts/schools (composite alternative)	210/1,063	213/1,085	212/1,106	215/1,124	213/1,143	214/1,145	216/1,146	213/1,143

*Note.* Standard errors are estimated by bootstrapping using 250 repetitions and clustered at the district level. Year P3 denotes the school year 3 years prior to the new curriculum being adopted (e.g., the 2005–2006 school year for textbooks adopted in fall 2008), Year P4 denotes the year 4 years prior, and so on. Year 1 denotes the 1st year the new curriculum was adopted (e.g., the 2008–2009 school year for textbooks adopted in fall 2008), Year 2 denotes the 2nd year, and so on. All estimates are converted from school-level standard deviation units to student-level standard deviation units by multiplying them by a factor of 0.47, which is the ratio of standard deviations of the school-average test score distribution to the student-level test score distribution in ELA averaged across our data panel. This transformation has no bearing on the results qualitatively or quantitatively; the rescaling is performed only to improve comparability of our findings to those in other studies that report effect sizes in student-level standard deviation units.

who were exposed to *California Math* in all three grades, which indicates something other than a linearly progressing dosage effect.<sup>16</sup> This is consistent with the results in Table 3, which also show no evidence of dosage effects for cohorts with differential exposure to *California Math* in the early primary grades.

Unfortunately, because the literature on curricular efficacy is so thin, there is little prior evidence on which we can

draw to gain inference about dosage effects. In similar previous studies, there is suggestive evidence of increased effect sizes for greater dosages in the early primary grades, but no study finds a statistically significant effect of longer exposure to a curriculum that is more effective on average (Agodini et al., 2010; Bhatt et al., 2013; Bhatt & Koedel, 2012). Our study, which provides the longest range of curricular-efficacy estimates in the literature to date (up to 4

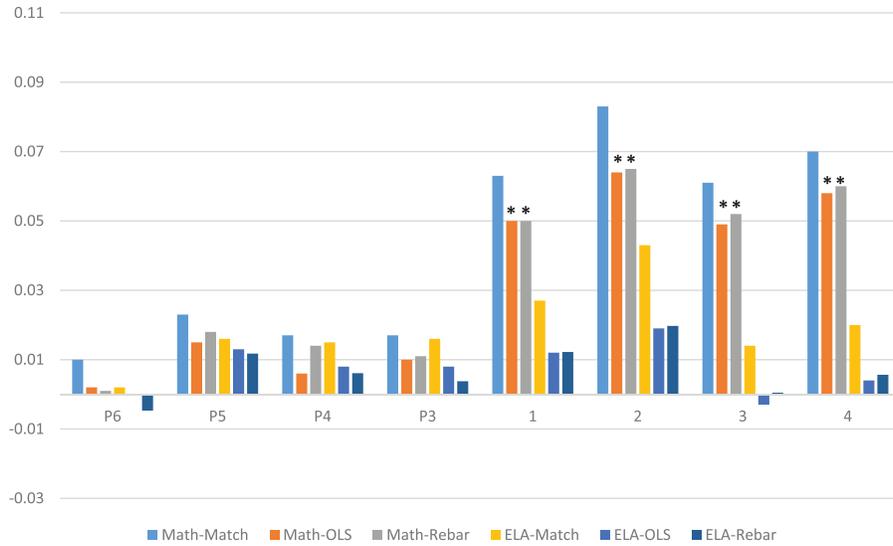


FIGURE 1. *Effects of California Math relative to the composite alternative on Grade 3 test scores, over time and using different estimators.*  
*Note.* Each bar shows an estimate reported in the preceding tables. All estimates are converted to student-level standard deviation units. Bars with asterisks are for estimates that are statistically distinguishable from zero at the 5% level. Years P6 to P3 are pretreatment years; Years 1 to 4 are posttreatment years.

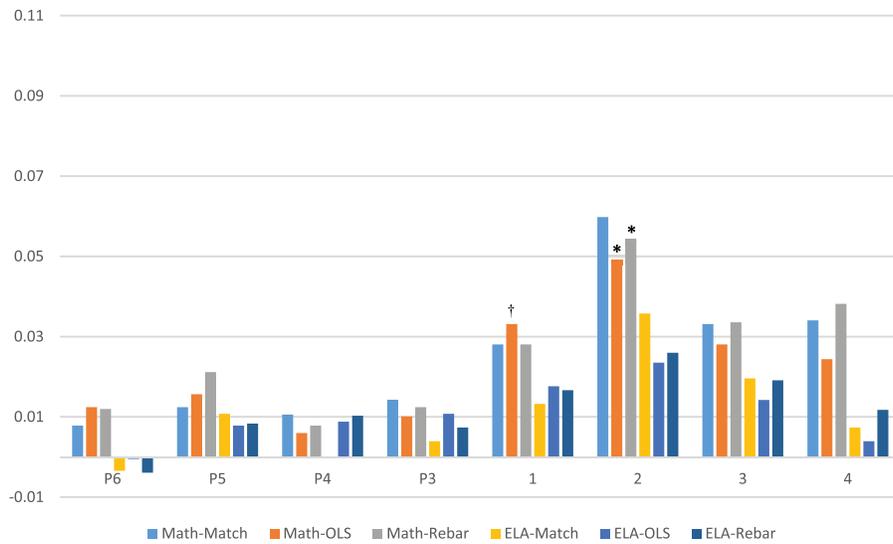


FIGURE 2. *Effects of California Math relative to the composite alternative on Grade 4 test scores, over time and using different estimators.*  
*Note.* All estimates are converted to student-level standard deviation units. Bars with asterisks (\*) are for estimates that are statistically distinguishable from zero at the 5% level; † indicates statistical significance at the 10% level. Years P6 to P3 are pretreatment years; Years 1 to 4 are posttreatment years. They Year 2, Grade 4 cohort in the posttreatment period corresponds to the Year 1, Grade 3 cohort; the Year 3, Grade 4 cohort corresponds to the Year 2, Grade 3 cohort; and so on.

consecutive years of use for the cohorts we follow the longest), can be characterized similarly. On the one hand, suggestive evidence of positive dosage effects across four different studies is more compelling than suggestive evidence from any single study, but on the other hand, it is interesting that evidence of dosage effects is not stronger. The dip in our estimates in Grade 4 for *California Math* suggests a potential mechanism worthy of additional exploration: the presence of grade-to-grade variability in the relative efficacy of curriculum materials.<sup>16</sup> More research is needed to

understand why dosage effects are not stronger than they appear in the handful of available studies, which has implications for understanding the scope for the adoption of more-effective curriculum materials to raise student achievement.

### *Effect Heterogeneity*

Finally, we also briefly consider the potential for curriculum effects to be heterogeneous across different student

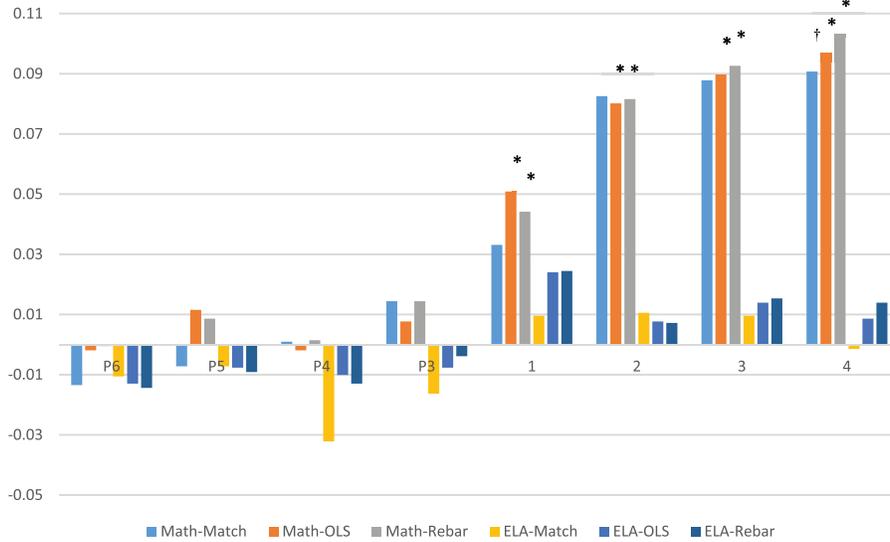


FIGURE 3. *Effects of California Math relative to the composite alternative on Grade 5 test scores, over time and using different estimators.*  
*Note.* All estimates are converted to student-level standard deviation units. Bars with asterisks (\*) are for estimates that are statistically distinguishable from zero at the 5% level; † indicates statistical significance at the 10% level. Years P6 to P3 are pretreatment years; Years 1 to 4 are posttreatment years. They Year 3, Grade 5 cohort in the posttreatment period corresponds to the Year 1, Grade 3 cohort; the Year 4, Grade 5 cohort corresponds to the Year 2, Grade 3 cohort.

subpopulations. We focus on students who differ by whether they are identified as socioeconomically disadvantaged by the CDE.<sup>17</sup> The CDE reports school-level test scores separately for students who differ by socioeconomic-disadvantage status (among other subgroups), which facilitates the heterogeneity analysis. Note that the sample of schools used for the heterogeneity analysis differs from the main sample because some schools do not have enough students in one or both subgroups (socioeconomically disadvantaged and non-socioeconomically disadvantaged) to report subgroup-specific test scores.

Table 6 shows results from models that follow the same estimation procedures as outlined above but run separately using school-average test scores for the student subgroups as the outcomes of interest. The table combines falsification estimates and treatment effect estimates—that is, it shows all at once what we report in Tables 3 and 4 for the primary analysis. As in the primary analysis, the falsification tests are generally as expected, although some are nominally larger than their analogs in the main tables. In terms of the main results, the treatment-effect estimates for Years 1 to 4 in Table 6 suggest that the positive relative effect of *California Math* is driven disproportionately by its effect on achievement for socioeconomically disadvantaged students. The effect sizes for these students are typically larger than for their nondisadvantaged counterparts and are more consistent in terms of size and statistical significance (e.g., the estimates in Years 3 and 4 for nondisadvantaged students, although positive and with standard errors that cannot rule out consistent effects over time, imply declining effects relative to Years 1 and 2).

Importantly, the results in Table 6 do not suggest a trade-off between groups associated with adopting *California Math* in the sense that neither group is made worse off. Noting that many of the estimates between advantaged and disadvantaged students are close (and some even overlap), especially when one accounts for our standard errors, we do not interpret the suggestive evidence of heterogeneous effects in Table 6 too strongly but identify this as an area of potentially interesting future research.<sup>18</sup>

## Conclusion

We use unique school-level data on curriculum adoptions in California to estimate the achievement effects of *California Math* relative to a composite alternative consisting of *enVision Math California*; *California Mathematics: Concepts, Skills, and Problem Solving*; and *California HSP Math*. We find that *California Math* outperformed the composite alternative curriculum. The differential effect in Grade 3 is on the order of 0.05 to 0.08 student-level standard deviations of the state standardized assessment in mathematics, which would move the 50th percentile school to roughly the 54th to 56th percentile in the school-level distribution of average test scores.<sup>19</sup> The Grade 5 estimates are suggestively larger, but although the Grade 4 estimates are always nominally positive, they are smaller and mostly insignificant. A potential explanation for the grade-by-grade variability in our estimates that merits attention in future research is that the effects of curriculum materials vary across grades.

Our estimates imply that *California Math* has an economically and educationally meaningful positive effect on

TABLE 6

*Heterogeneous Effects of California Math on Grade 3 Mathematics Achievement Relative to the Composite Alternative, for Exposed and Unexposed Cohorts by Student Socioeconomic Advantage Status*

Variable	Year P6	Year P5	Year P4	Year P3	Year 1	Year 2	Year 3	Year 4
Treatment: <i>California Math</i>								
Control: Composite alternative								
Socioeconomically nondisadvantaged test scores								
Treatment effect: Kernel matching	0.077 (0.047)	0.032 (0.048)	0.051 (0.050)	0.032 (0.049)	0.104 (0.048)**	0.084 (0.043)**	0.046 (0.043)	0.035 (0.055)
Treatment effect: Restricted ordinary least squares	0.022 (0.022)	-0.004 (0.022)	0.025 (0.022)	0.009 (0.022)	0.045 (0.027)*	0.054 (0.028)*	0.032 (0.032)	0.022 (0.030)
Treatment effect: Remnant-residualized matching	0.023 (0.030)	0.002 (0.024)	0.017 (0.022)	0.000 (0.024)	0.055 (0.027)**	0.062 (0.028)**	0.045 (0.035)	0.027 (0.039)
Socioeconomically disadvantaged test scores								
Treatment effect: Kernel matching	0.024 (0.034)	0.023 (0.027)	0.014 (0.032)	0.010 (0.032)	0.075 (0.032)**	0.087 (0.032)**	0.078 (0.034)**	0.109 (0.038)**
Treatment effect: Restricted ordinary least squares	0.002 (0.020)	0.005 (0.019)	0.007 (0.017)	0.001 (0.018)	0.059 (0.023)**	0.068 (0.025)**	0.062 (0.025)**	0.083 (0.026)**
Treatment effect: Remnant-residualized matching	0.015 (0.028)	0.020 (0.029)	0.006 (0.020)	0.000 (0.022)	0.063 (0.027)**	0.069 (0.028)**	0.067 (0.031)**	0.089 (0.030)**
Sample sizes								
Socioeconomically nondisadvantaged models								
No. of districts/schools ( <i>California Math</i> )	72/452	71/453	70/442	71/454	65/406	69/397	72/407	68/380
No. of districts/schools (composite alternative)	181/942	177/915	181/915	181/914	173/885	173/838	169/800	165/755
Socioeconomically disadvantaged models								
No. of districts/schools ( <i>California Math</i> )	80/473	83/478	81/495	83/511	83/532	82/532	86/540	83/548
No. of districts/schools (composite alternative)	183/883	186/929	187/950	192/972	189/1,010	192/1,009	194/1,022	193/1,015

*Note.* This table replicates the results in Tables 3 and 4 but using separate school-level achievement measures for students identified as either socioeconomically disadvantaged or non-socioeconomically disadvantaged by the California Department of Education. The school sample sizes vary between the groups because not all schools have enough students in each group for reporting. Standard errors are estimated by bootstrapping using 250 repetitions and clustered at the district level. Year P3 denotes the school year 3 years prior to the new curriculum being adopted (e.g., the 2005–2006 school year for textbooks adopted in fall 2008), Year P4 denotes the year 4 years prior, and so on. Year 1 denotes the 1st year the new curriculum was adopted (e.g., the 2008–2009 school year for textbooks adopted in fall 2008), Year 2 denotes the 2nd year, and so on. All estimates are converted from school-level standard deviation units to student-level standard deviation units by multiplying them by a factor of 0.45, which is the ratio of standard deviations of the school-average test score distribution to the student-level test score distribution averaged across our data panel. This transformation has no bearing on the results qualitatively or quantitatively; the rescaling is performed only to improve comparability of our findings to those in other studies that report effect sizes in student-level standard deviation units. \* $p \leq .10$ . \*\* $p \leq .05$ .

student achievement relative to the popular alternatives we consider. The effect is particularly notable given that (a) it is a schoolwide effect and thus applies, on average, to each student in a treated school and (b) the marginal cost of choosing one curriculum over another is so small as to be effectively zero (Bhatt & Koedel, 2012; Chingos & Whitehurst, 2012). Note that an alternative intervention targeted at 10% of the student population would need to have an effect 10 times as large as the *California Math* effect to generate as large an increase in student achievement overall (ignoring spillovers). Of course, providing empirical evidence of differential curriculum effects on

student achievement is just one step toward giving districts comprehensive information to inform textbook adoptions. The weight given to these results (and results from future, related studies) in the adoption process will depend on how decision makers value achievement as measured by state assessments. The more weight given to achievement, the more appealing *California Math* becomes relative to the alternative textbook options we consider.

We can only speculate as to why we find smaller differential curriculum effects in California than in previous studies.<sup>20</sup> Candidate explanations include that the curriculum materials in California are more similar to each other than the

curriculum materials that have been evaluated previously, the context in California is such that curriculum effects are smaller (e.g., curricular objectives, assessments, etc.), or simply sampling variance. Our ability to gain inference into the mechanisms underlying the differential curriculum effects that we estimate here, and related estimates elsewhere, is limited by the lack of a larger literature within which our findings can be contextualized. In analytic terms, empirical analyses that would aim to link specific textbook characteristics and/or contextual evaluation factors (such as the assessment used) to efficacy estimates are currently hampered by underidentification—there are too many potential explanatory factors and too few efficacy estimates.

Ours is one of only a small handful of rigorous studies to test for impacts of textbooks on student achievement in mathematics. Moreover, we are not aware of any similar studies in a subject outside of mathematics. By this point, we believe our methods are well established enough that it would be straightforward to apply them in other contexts if textbook data were available. By replicating this study across states and within states over time, we could begin to gather enough impact data to explore variation in curricular impact estimates as a function of features thought to matter (e.g., textbook content, alignment to standards, approach to teaching the subject, etc.). However, currently there is not enough efficacy information to support such investigations, and in the meantime, studies like ours contribute evidence for specific sets of materials and can be used to inform contemporary curriculum adoption decisions, even if the features that make some curricula outperform others remain unidentified.

We conclude by reiterating the calls made by Bhatt and Koedel (2012) and Chingos and Whitehurst (2012) for improved efforts to collect data on curriculum materials. Curriculum materials are a substantial input into educational production, and data consistently point toward high curriculum materials usage by students and teachers in the Common Core era (Opfer, Kaufman, & Thompson, 2016; Perry et al., 2015). However, it remains the case that in nearly all states, which curriculum materials are being used by which schools is not tracked. Even in California, where reporting on curriculum materials is the law, we found that information provided by a significant fraction of schools does not actually identify the curriculum materials being used, which suggests little oversight of the data. This much is for certain: With no data, we are committed to leaving educational decision makers to adopt curricula without efficacy evidence.

## Appendix A

### *Data Appendix*

Appendix Table A1 documents attrition from our data set beginning with a universe of California elementary schools in the California Department of Education data with characteristics from either 2007 or 2008, at least one Grade 3 test score from 2009 to 2013, and where the highest grade is 8 or lower. In the text of this appendix, we also briefly elaborate on the three key attrition points.

First, although California provides a SARC template for schools, which some follow, the quality of information about curriculum materials reported on the SARCs varies greatly. Curriculum materials information was either not reported

TABLE A1  
*Construction of the Analytic Sample*

Variable	Schools	% of total	Districts	% of total
Initial universe	5,494		825	
Reason for data loss				
No record in textbook file	−339	6.2	−32	3.9
Indeterminate textbook information	−804	14.6	−134	16.2
Adoption year other than 2008 or 2009	−876	15.9	−119	14.4
Non-uniform adopter (or uncertain), Grades 1–3	−481	8.8	−54	6.6
Grade-span conflict between CDE and SARC data	−33	0.6	−17	2.1
Missing school/district outcome data	−48	0.9	−19	2.3
Missing district/school covariate data	0	0	0	0
Did not use one of the four focal curricula	−632	11.5	−139	16.8
Initial analytic sample	2,281	41.5	311	37.7
Drop LAUSD and LBUSD	−403	7.3	−2	0.2
Final analytic sample	1,878	34.2	309	37.5

*Note.* The initial universe includes all schools in the CDE data with characteristics from either 2007 or 2008, at least one grade-3 test score from 2009–2013, and where the highest graded is 8 or lower. CDE = California Department of Education; SARC = School Accountability Report Card; LAUSD = Los Angeles Unified School District; LBUSD = Long Beach Unified School District.

(perhaps because no book was used in some cases), or reported in such a way that the actual textbook used is indeterminate, for 20.8% of elementary schools in the state. As an example of an indeterminate report, a district might list only a publisher's name for a publisher that produced multiple state-approved textbooks (e.g., list "Houghton Mifflin," which published both *Harcourt California HSP Math* and *California Math*). In such a case, if no other information is provided, the actual textbook cannot be determined. We drop all schools from the sample that report no textbook information or indeterminate information.

A second notable reason schools were removed from the analytic sample is that they report a curriculum adoption year other than 2008 or 2009 on the 2013 SARC. Appendix Table A1 shows that this applies to approximately 15.9% of schools. Schools may have delayed adoptions beyond 2009 for a variety of reasons, including budgetary issues or a lack of need. As an example of the latter, a school may have adopted off cycle in a recent year prior to 2009/2010 and thus may not have needed to adopt new materials on the standard timeline.

A third significant source of attrition from our data set, conditional on schools adopting textbooks in 2008 or 2009 and reporting identifiable materials, is that we drop approximately 8% of schools that either (a) explicitly indicate using more than one textbook in Grades 1 to 3 or (b) indicate using more than one textbook in the school and where the SARC was ambiguous about which curriculum materials were used in which grades. The reason for this restriction is that we focus primarily on estimating achievement effects on Grade 3 mathematics tests. Schools that use more than one textbook in Grades 1 to 3 have mixed treatments. As noted in the text, although in principle these schools could be used to examine mixed-treatment effects, in practice there are too few observations for an effective analysis along these lines, so we simply drop them from the analytic sample.

## Appendix B

### *Supplementary Materials and Results*

*Focal Curricula.* In this section, we briefly describe the four books, drawing on available data from the What Works Clearinghouse (WWC), the state adoption report, and available web materials. All of the textbooks we study are the California editions of their respective book series. Because some of the information available online describes the national or Common Core versions of these series, we cannot always be confident that it applies to the California versions we study. We are hampered in our descriptions by the fact that there is little or no publicly available information about the differences between state-specific and national versions of textbooks.

Pearson Scott Foresman's *enVision Math California* is an early edition of the *enVision* series that is still marketed and sold by Pearson as Common Core and Texas editions. According to the WWC, *enVision* aims to help students develop an understanding of mathematics concepts through problem-based instruction, small-group interaction, and visual learning, with a focus on reasoning and modeling. Each lesson is intended to include small-group problem solving. The book's lead author, Randall Charles, was a coauthor of the National Council of Teachers of Mathematics' Focal Points, widely considered a reform-oriented mathematics document. Despite its seemingly reform-oriented description, analyses of other editions of *enVision* (the Common Core and Florida Grade 4 versions) found them to be typical in terms of their cognitive-demand coverage and far below the level of cognitive demand emphasized in the standards (Polikoff, 2015). The California state adoption report indicates that this curriculum met all five evaluative criteria.

We have far less information about the other three textbooks. The California state adoption report indicates that all three meet the five evaluative criteria (California Department of Education, 2009). Houghton Mifflin's *California Math* and Harcourt's *California HSP Math* are both updated versions of textbooks previously adopted by the state in the 2001 adoption, whereas McGraw Hill's *California Mathematics* was not adopted previously by the state. Other than this, we were unable to find information about the Houghton Mifflin and Harcourt books. McGraw Hill's *California Mathematics* has an evaluation report (Papa & Brown, 2007) that describes the book as including both conceptual understanding and guided practice and argues that it aligns with what is known about effective mathematics instruction. McGraw Hill does not appear to have published any books in this series since 2009. In the conclusion of the paper, we discuss the challenge of characterizing these textbooks, and correspondingly, in interpreting our results based on student achievement in terms of their content and form.

*Pairwise Comparisons.* Appendix Table B1 summarizes initial results from the six pairwise comparisons. The first three comparisons involve what becomes the focal curriculum in our analysis: *California Math*. *California Math* is the treatment curriculum in the first comparison, and the control curriculum in the other two (we use the convention of defining the most-adopted book as the "control" curriculum in each pairwise comparison). Notice that we obtain fairly large point estimates in all three comparisons involving *California Math*, and all three comparisons suggest that *California Math* is more effective. For the comparisons involving the other curricula, our point estimates are consistently small and do not suggest differential effects.

TABLE B1

*Balance and Estimation Results for the Six Initial Pairwise Comparisons During Treatment Years*

Variable	Estimated treatment effects and balancing results by year after adoption			
	Year 1	Year 2	Year 3	Year 4
Comparison 1				
Treatment: <i>California Math</i>				
Control: <i>enVision Math</i>				
Treatment effect (kernel matching)	0.048	0.059	0.041	0.054
	(0.063)	(0.066)	(0.058)	(0.061)
No. unbalanced covariates, matched <i>t</i> tests (5%)	2	3	2	2
Mean standardized difference of covariates	6.0	5.8	6.1	6.1
No. unbalanced covariates, Smith-Todd (5%)	3	4	4	3
Average <i>p</i> value, Smith-Todd	0.41	0.28	0.43	0.40
Comparison 2				
Treatment: <i>California Mathematics: Concepts, Skills, and Problem Solving</i>				
Control: <i>California Math</i>				
Treatment effect (kernel matching)	-0.087	-0.152	-0.110	-0.091
	(0.072)	(0.077)**	(0.072)	(0.076)
No. unbalanced covariates, matched <i>t</i> tests (5%)	0	0	0	0
Mean standardized difference of covariates	4.8	5.5	4.3	4.1
No. unbalanced covariates, Smith-Todd (5%)	11	5	3	3
Average <i>p</i> value, Smith-Todd	0.21	0.30	0.31	0.31
Comparison 3				
Treatment: <i>California HSP Math</i>				
Control: <i>California Math</i>				
Treatment effect (kernel matching)	-0.063	-0.065	-0.039	-0.059
	(0.063)	(0.057)	(0.072)	(0.076)
No. unbalanced covariates, matched <i>t</i> tests (5%)	0	0	0	0
Mean standardized difference of covariates	6.2	6.4	5.9	5.6
No. unbalanced covariates, Smith-Todd (5%)	5	5	5	4
Average <i>p</i> value, Smith-Todd	0.28	0.27	0.28	0.29
Comparison 4				
Treatment: <i>California Mathematics: Concepts, Skills, and Problem Solving</i>				
Control: <i>enVision Math</i>				
Treatment effect (kernel matching)	0.010	-0.017	-0.003	0.005
	(0.066)	(0.066)	(0.058)	(0.065)
No. unbalanced covariates, matched <i>t</i> tests (5%)	0	0	0	0
Mean standardized difference of covariates	5.6	5.2	5.4	5.2
No. unbalanced covariates, Smith-Todd (5%)	3	3	3	3
Average <i>p</i> value, Smith-Todd	0.52	0.56	0.56	0.57
Comparison 5				
Treatment: <i>California HSP Math</i>				
Control: <i>enVision Math</i>				
Treatment effect (kernel matching)	0.065	0.004	-0.009	0.013
	(0.090)	(0.078)	(0.106)	(0.104)
No. unbalanced covariates, matched <i>t</i> tests (5%)	0	0	0	0
Mean standardized difference of covariates	6.4	6.7	6.5	6.3
No. unbalanced covariates, Smith-Todd (5%)	4	4	4	4
Average <i>p</i> value, Smith-Todd	0.39	0.39	0.39	0.39
Comparison 6				
Treatment: <i>California HSP Math</i>				
Control: <i>California Mathematics: Concepts, Skills, and Problem Solving</i>				
Treatment effect (kernel matching)	0.016	0.021	0.058	0.028
	(0.091)	(0.079)	(0.081)	(0.083)
No. unbalanced covariates, matched <i>t</i> tests (5%)	0	0	0	0
Mean standardized difference of covariates	4.1	4.2	3.9	4.2
No. unbalanced covariates, Smith-Todd (5%)	6	6	6	6
Average <i>p</i> value, Smith-Todd	0.31	0.34	0.33	0.34
No. of districts/schools				
<i>enVision Math California</i>	106/706	106/707	107/707	105/706
<i>California Math</i>	92/602	89/593	91/600	90/599
<i>California Mathematics: Concepts, Skills and Problem Solving</i>	67/387	69/389	69/389	69/389
<i>California HSP Math</i>	48/177	47/176	48/177	47/176

*Note.* The balancing tests report results based on the same 22 matching covariates used in each pairwise comparison. Standard errors for matching estimators are estimated by bootstrapping using 250 repetitions and clustered at the district level. Year 1 denotes the 1st year the new curriculum was adopted (e.g., the 2008–2009 school year for textbooks adopted in fall 2008), Year 2 denotes the 2nd year, and so on. All estimates are converted from school-level standard deviation units to student-level standard deviation units by multiplying them by a factor of 0.45, which is the ratio of standard deviations of the school-average test score distribution to the student-level test score distribution in math averaged across our data panel, as reported in the text. This transformation has no bearing on the results qualitatively or quantitatively; the rescaling is performed only to improve comparability of our findings to those in other studies that report effect sizes in student-level standard deviation units.

\* $p \leq .10$ . \*\* $p \leq .05$ .

Like for the primary comparison in the text, we report balancing information in several ways for each pairwise comparison in Appendix Table B1. As is clear from the table, a limitation of most of the pairwise comparisons is that the balancing results, although not indicative of egregious imbalance, are also not particularly compelling. Covariate balance using the matched *t* tests generally looks good, but the mean standardized difference for several of the pairwise comparisons is large, and certainly much larger than in the comparison between *California Math* and the composite alternative. In all pairwise comparisons, the Smith and Todd (2005) regression tests indicate imbalance in one form or another (i.e., either too many unbalanced covariates and/or average *p* values that are too low).

As noted in the text, our small sample sizes in the pairwise comparisons (relative to sample sizes more typical of matching analyses in other contexts) limit our ability to improve covariate balance separately for each comparison. Thus, based on these initial results, and the suggestion that *California Math* is more effective than the other three textbooks (which all appear to be similarly effective), we focus our main evaluation on comparing *California Math* to a composite of the other three popular curricula. Reducing the dimensionality of the comparison in this way yields a more effective matching procedure, which can be seen by comparing the balance statistics shown in Appendix Table B1 for the pairwise comparisons to the analogous numbers for the composite comparison in the main text (Table 2). The falsification tests shown in the main text offer additional evidence consistent with our final evaluation of *California Math* being balanced.

*Matching Details for the Primary Comparison and Overlap of Propensity Scores.* Appendix Tables B2 and B3 report details about the matching procedure for the primary comparison between *California Math* and the composite alternative. First, Table B2 shows the output from the initial selection model from which the propensity scores are generated to give a sense of which covariates predict the adoption of *California Math*. The only statistically significant covariates are the three terms for district enrollment (linear, quadratic, cubic).

Second, Table B3 shows covariate-by-covariate balancing results to complement the aggregate reporting in Table 2. For brevity, we show covariate-by-covariate balance using the Year 1 sample of schools and districts only (recall from the text that the balancing results fluctuate mildly from year to year because of sample changes due to building openings and closings and data reporting issues for small schools).

Figure B1 shows the distributional overlap in propensity scores between *California Math* (treatments) and other focal-curricula adopters (controls). The propensity scores are summary measures of school and district characteristics, weighted by their predictive influence over the adoption of

TABLE B2  
*Probit Coefficients From the Propensity Score Model Predicting the Adoption of California Math Instead of the Composite Alternative*

Variable	Coefficient
Data quality indicator	-1.028 (0.715)
Census data missing indicator	-1.777 (4.580)
Fall 2008 adoption	0.305 (0.222)
School average math score (standardized)	-0.049 (0.059)
School average ELA score (standardized)	0.189 (0.126)
District average math score (standardized)	-0.078 (0.434)
District average ELA score (standardized)	0.338 (0.410)
Share female	0.438 (1.148)
Share socioeconomically disadvantaged	0.633 (0.673)
Share African American	-1.328 (1.296)
Share Asian	-0.729 (0.807)
Share White	-0.738 (0.781)
Share Other	-1.539 (1.721)
Share English learner	-0.886 (0.879)
School enrollment (1,000s)	4.091 (2.640)
School enrollment squared (1,000s)	-0.00615 (0.00473)
School enrollment cubed (1,000s)	0.00000302 (0.00000251)
District enrollment (1,000s)	-0.167 (0.089)*
District enrollment squared (1,000s)	0.0000172 (0.00000710)**
District enrollment cubed (1,000s)	0.00000000382 (0.00000000150)**
Share low education (U.S. Census)	-0.004 (0.011)
Median household income (U.S. Census)	-0.153 (0.398)
Constant	2.263 (4.450)
Pseudo <i>R</i> -squared	0.1221
<i>N</i> (total)	1,878

*Note.* The data quality indicator is set to 1 if the sum of student subgroups does not equal total enrollment as reported by the California Department of Education. This was not an issue for most schools, and even when it was, inequalities were small. Ex post, this variable has no bearing on our findings, and all of our results are robust to excluding it. The omitted student categories are the share male, socioeconomically disadvantaged, Hispanic, and non-English learner. ELA = English language arts.

\**p* ≤ .10. \*\**p* ≤ .05.

TABLE B3

*Covariate-by-Covariate Balancing Details for the Comparison between California Math and the Composite Alternative, Year 1 Sample*

Variable	Matched <i>t</i> test, significant difference	Standardized difference	Smith-Todd test, significant difference	Smith-Todd <i>p</i> value
Data quality indicator	No	12.2	No	.85
Census data missing indicator	No	2.6	No	.27
Fall 2008 adoption	No	2.7	No	.69
School average math score	No	1.8	No	.50
School average ELA score	No	4.3	No	.77
District average math score	No	1.3	No	.68
District average ELA score	No	5.5	No	.92
Share female	No	1.1	Yes	.02
Share socioeconomically disadvantaged	No	-4.4	No	.97
Share African American	No	-0.3	No	.53
Share Asian	No	-2.1	No	.29
Share White	No	7.4	No	.91
Share Other	No	-1.9	No	.45
Share English learner	No	-5.9	No	.73
School enrollment	No	2.6	No	.24
School enrollment squared	No	1.5	No	.15
School enrollment cubed	No	0.7	No	.13
District enrollment	No	3.1	Yes	.01
District enrollment squared	No	3.4	No	.10
District enrollment cubed	No	4.3	No	.14
Share low education (Census)	No	-5.5	No	.79
Median household income (Census)	No	-2.3	No	.36

*Note.* This table provides full details for the balancing results shown in Table 2 for Year 1. Detailed balancing results for other years are substantively similar. The average absolute standardized difference reported in Table 2 is the average of the absolute values of the standardized differences reported in this table.

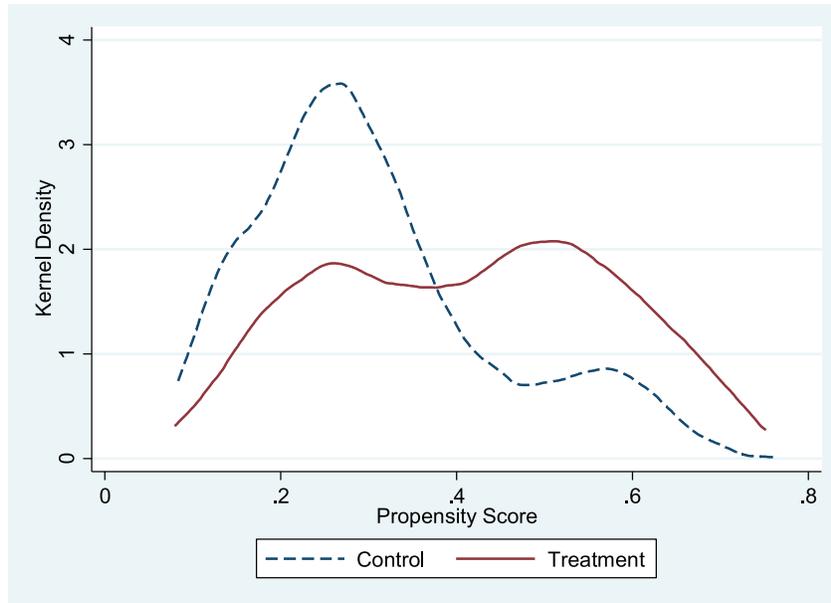


FIGURE B1. *Kernel densities of estimated propensity scores for treatment (California Math) and control (composite alternative) schools on the common support, Grade 3 math.*

*California Math*. In any program evaluation where treatment is predicted at least to some degree by observable characteristics, treatment units will have higher propensity scores on average than controls, as in the case in Figure B1. However, the figure shows considerable overlap in the distributions of propensity scores for treatment and control schools, which is conducive to our matching evaluation.

## Appendix C

### Technical Appendix

*Remnant-Based Residualization.* The “remnant” sample we use for remnant-based residualization includes data from all schools in California that adopted a new curriculum in fall 2008 or fall 2009 uniformly but chose a curriculum other than one of the four primary textbooks (there are 632 such schools per Appendix Table A1). Thus, these schools are outside of our evaluation sample. Following Sales, Hansen, and Rowan (2014), we start by estimating the following linear regression model using the remnant data:

$$Y_{sdt} = \mathbf{X}_s \boldsymbol{\alpha}_{1t} + \mathbf{X}_d \boldsymbol{\alpha}_{2t} + \eta_{sdt}. \quad (\text{C1})$$

In Equation (C1),  $Y_{sdt}$  is a Grade 3 math test score for school  $s$  in district  $d$  in year  $t$ , and  $\mathbf{X}_s$  and  $\mathbf{X}_d$  are defined as above.<sup>22</sup> After estimating Equation (C1), we store the coefficient estimates  $\hat{\boldsymbol{\alpha}}_{1t}$  and  $\hat{\boldsymbol{\alpha}}_{2t}$  and construct the following residualized test score outcome for each school *in our analytic sample* in each year:

$$Q_{sdt} = Y_{sdt} - (\mathbf{X}_s \hat{\boldsymbol{\alpha}}_{1t} + \mathbf{X}_d \hat{\boldsymbol{\alpha}}_{2t}). \quad (\text{C2})$$

In Equation (C2),  $Y_{sdt}$  is the Grade 3 test score for school  $s$  in district  $d$  in year  $t$  for a school that adopted one of the four primary curricula.  $\mathbf{X}_s$  and  $\mathbf{X}_d$  continue to be defined as above.  $\hat{\boldsymbol{\alpha}}_{1t}$  and  $\hat{\boldsymbol{\alpha}}_{2t}$  are out-of-sample parameter estimates based on the remnant data that link the preadoption school and district characteristics to test score outcomes by year. Intuitively, Equation (C2) can be described as specifying a set of general relationships between school/district characteristics and test score outcomes in California as defined by  $\hat{\boldsymbol{\alpha}}_{1t}$  and  $\hat{\boldsymbol{\alpha}}_{2t}$ , and netting the influence of these characteristics out of the outcome data.

We implement the matching procedure as described by Equation (3) in the main text using the residualized outcomes,  $Q_{sdt}$ , in place of the actual outcomes. This procedure is very similar to restricted ordinary least squares (OLS), with the added benefit that the adjustment parameters  $\hat{\boldsymbol{\alpha}}_{1t}$  and  $\hat{\boldsymbol{\alpha}}_{2t}$  are estimated entirely out of sample. Using an out-of-sample “training set” for the outcome model has several conceptual benefits over using in-sample data (as was the case with OLS) as described by Sales et al. (2014). In our application, it addresses the concern that bias could be introduced by the OLS models if the covariate coefficients are

disproportionately influenced by schools in the control condition, which dominate our sample. This in turn would result in asymmetric overfitting of the outcome model, potentially causing bias.<sup>23</sup>

A concern with remnant-based residualization is that the relationships between school/district characteristics and test scores may be different in the analytic sample and the remnant sample. Although in such a scenario the adjustment parameters  $\hat{\boldsymbol{\alpha}}_{1t}$  and  $\hat{\boldsymbol{\alpha}}_{2t}$  will be less useful, Sales et al. (2014) show that the procedure still improves inference, albeit by less. In practice, if  $\hat{\boldsymbol{\alpha}}_{1t}$  and  $\hat{\boldsymbol{\alpha}}_{2t}$  measure a relatively constant set of relationships between characteristics and outcomes in California schools within years, remnant-based residualization and restricted OLS should return similar results. This is the case in our application as shown in the main text.

*Balancing Test Details.* This section elaborates on the balancing results shown in Table 2. The first row of Table 2 reports the number of unbalanced covariates using simple covariate-by-covariate  $t$  tests among the matched sample. A covariate where the difference between treatment and control values is significant at the 5% level is reported as unbalanced. We use 22 covariates in total to match schools, and none are individually unbalanced at the 5% level within the matched sample based on the  $t$  tests. This indicates that the unconditional differences in school characteristics shown in Table 1 disappear completely in the matched comparisons.<sup>24</sup>

In row 2, we report the average absolute standardized difference across all covariates. Following Rosenbaum and Rubin (1985), the formula for the absolute standardized difference for covariate  $X_k$  is given by

$$SDIFF(X_k) = \frac{\left| \frac{1}{N^S} \left[ \sum_{j \in N_t \cap S_p} \{X_{kj} - \sum_{m \in I_{t,j} \cap S_p} W(j,m) X_{km}\} - \sum_{m \in N_m \cap S_p} \{X_{km} - \sum_{j \in I_{t,m} \cap S_p} W(m,j) X_{kj}\} \right] \right|}{\sqrt{\frac{Var(X_{kj}) + Var(X_{km})}{2}}} * 100. \quad (\text{C3})$$

The numerator in Equation (C3) is analogous to the formula for our matching estimators in Equation (3) where we replace  $Y$  with  $X_k$  and take the absolute value (note the denominator is calculated using the full sample). The absolute average standardized difference is complementary to the covariate-by-covariate  $t$  tests reported in the first row of the table. Beyond measuring purely statistical differences as with the  $t$  tests, the absolute average standardized difference provides an indication of the magnitude of potential imbalance.

A weakness of reporting on standardized differences is that there is no clear rule by which to judge the results. Rosenbaum and Rubin (1985) suggest that a value of 20 is large, although recent studies have applied more stringent criteria (e.g., Sianesi, 2004). The average absolute standardized differences that we report in Table 2 are quite small

compared to similar estimates reported in other studies, on the order of just 3% to 4% across the pre- and postadoption years of our data panel. This corroborates the result from the  $t$  tests that the covariates are well balanced between *California Math* adopters and other schools. In Appendix Table B3, we report standardized differences on a covariate-by-covariate basis for interested readers.

Rows 3 and 4 of Table 2 show results from alternative, regression-based balancing tests proposed by Smith and Todd (2005). Like with the standardized difference measure, we perform the regression test for each covariate in each year and aggregate the results. Specifically, we estimate the following regression on a covariate-by-covariate basis:

$$X_{ik} = \beta_0 + \beta_1 p_i + \beta_2 p_i^2 + \beta_3 p_i^3 + \beta_4 p_i^4 + \beta_5 D_i + \beta_6 D_i p_i + \beta_7 D_i p_i^2 + \beta_8 D_i p_i^3 + \beta_9 D_i p_i^4 + \xi_{ik}. \quad (\text{C4})$$

In Equation (C4),  $X_{ik}$  represents a covariate from the propensity score specification for school  $i$ ,  $p_i$  is the estimated propensity score, and  $D_i$  is an indicator variable equal to 1 if the school adopted *California Math* and 0 otherwise. The test for balance is for whether the coefficients  $\beta_5$  to  $\beta_9$  are jointly equal to 0—that is, whether treatment predicts the  $X$ s conditional on a quartic of the propensity score.<sup>25</sup>

We report the number of unbalanced covariates at the 5% level and the average  $p$  value from the joint test of significance for  $\beta_5$  to  $\beta_9$  across the 22 covariates in each year. Although we see marginally more unbalanced covariates than would be expected by chance using the Smith-Todd tests (two to three per year), the implied level of imbalance is small. Moreover, the average  $p$  values from the regression tests are consistently around 0.50 across the covariates in each year, which is as expected in a balanced comparison.<sup>26</sup>

### Acknowledgment

This study is based on work supported by the National Science Foundation under Grant No. 1445654 and the Smith Richardson Foundation. Any opinions, findings, and conclusions or recommendations expressed in this study are those of the author(s) and do not necessarily reflect the views of the funders.

### Notes

1. We do not know the list price of the textbooks we study, but research indicates that most textbooks are approximately the same unit cost. The elementary mathematics books in Boser, Chingos, and Straus (2015) cost an average of \$34 per pupil, or approximately 0.32% of per-pupil spending (the true per-pupil expenditure is even lower because textbooks are used for multiple years).

2. Access to student-level test scores would offer little additional value for our evaluation because the curriculum adoption data are at the school level. It is also unlikely that student-level data on test scores and curriculum exposure (we are not aware of the latter existing anywhere in the United States), even if available, would meaningfully improve inference from our evaluation, because very few schools report using more than one set of curriculum materials

in the same grade (these schools are a small subsample of “non-uniform” adopters reported in Appendix Table A1). This implies limited treatment variability within schools that could be exploited with student-level data.

3. When we extend our analysis to Grades 4 and 5, we also extend the restriction of constant materials usage to Grades 4 and 5. Most schools that used constant materials in Grades 1 to 3 used the same materials in Grades 4 and 5, but there is a small amount of sample attrition owing to this issue in the later grades.

4. Districtwide enrollment in Los Angeles Unified School District (LAUSD) is at least an order of magnitude higher than in any other individual district in the sample (the next largest district in California, San Diego, is roughly one fifth the size of LAUSD but is dropped from the sample because it adopted off cycle); enrollment in Long Beach Unified School District is 50% larger than the next largest district. The adoption outcomes of these districts disproportionately affect the matching model we use to construct observationally equivalent comparisons (the model is shown below). Modeling adjustments can be made (at a cost) to reduce the disproportionate effect of these districts, but the obvious noncomparability problem remains, and for this reason, we exclude schools from these districts from the evaluation.

5. Prior to merging in the curriculum data, these are the minimal conditions for inclusion into our analysis.

6. Briefly, these benefits include (a) there are more schools than districts, which allows us to construct better unit-level matches, and (b) performing our analysis at the school level allows us to directly control for school- and district-specific factors that may influence adoption decisions, whereas it is not clear how one would control for disaggregated school characteristics and their potential role in adoption decisions if the match were performed at the district level.

7. We also include a binary variable to indicate California Department of Education (CDE) data quality for individual schools and an indicator for missing census data. The CDE data quality indicator is equal to 1 if the enrollment counts by subgroup (e.g., by race, gender, etc.) do not exactly match total reported enrollment for schools. For most schools, the subgroup enrollments sum to total enrollment, and this variable is of no practical consequence in our analysis (i.e., if we omit the variable entirely, our results are unchanged).

8. The non-test score school and district covariates are averaged over the 2 years immediately prior to the adoption of the new materials, and the test score covariates are from 2 years before the adoption. We follow Bhatt and Koedel (2012) in not using test score information from the year immediately before the new books were adopted because this information would not have been available to decision makers at the time of the decision per the above discussion. That said, none of our findings are substantively affected if we include lagged test score information from the year just before adoption into the selection models. Moreover, we have considered the robustness of our findings to expanding the set of controls at the district level, including the use of more of the school-level analogs (e.g., district racial composition shares) and district per-pupil spending. Including these additional covariates in our models leads to negligible changes in our estimates.

9. We interviewed 21 district administrators from across California about the curriculum adoption processes in their districts. These interviews confirm the complexity of the adoption process and indicate that decisions are driven by committees made

up mostly of teachers. In none of the districts was there evidence of a strong decision maker.

10. Appendix Table B1 presents a lot of information tersely. It will be easier to interpret after reading the remainder of this section.

11. There are several ways to empirically verify this statement, but we must be careful to not contaminate the predictive power of our covariates with their predictive power over curriculum materials. As one straightforward data point, we use the remnant sample and estimate an achievement model during the 1st year of a new adoption using our matching covariates. The  $R$ -squared from this regression is 0.74.

12. The analysis is performed using school-level achievement measures. Effect sizes are converted into student-level standard deviation units, which are more commonly reported for other educational interventions in the literature, by multiplying them ex post by the ratio  $\sigma_s/\sigma_i$ , where  $\sigma_s$  is the standard deviation of the distribution of school-averaged math test scores and  $\sigma_i$  is the standard deviation of the distribution of student-level scores. We calculate  $\sigma_s$  using data from all reporting schools in California each year;  $\sigma_i$  is provided for all students by the CDE in annual reports. This conversion follows the procedure of Bhatt and Koedel (2012). The ratio  $\sigma_s/\sigma_i$  averaged across years in our data panel is 0.45 (the ratio varies very little from year to year).

13. For the remnant-residualized matching estimates in English language arts (ELA), we reestimate Equation (C1) using ELA scores from the remnant sample to obtain appropriate adjustment parameters analogously to the procedure for math scores described in Appendix C.

14. With obvious appropriate adjustments; for example, in the matching model for the Grade 4 analysis, we match schools on Grade 4 test scores.

15. Specifically, we lose 2.5% of the initial Grade 3 sample of schools in the Grade 4 analysis and another 1.5% when we move to the Grade 5 analysis. The small data loss is attributable to schools that did not have preadoption test scores in Grades 4 and/or 5 and schools that did not continue to uniformly adopt a focal curriculum past Grade 3. The schools dropped from the sample as we move to higher grades are much smaller than the typical California school.

16. There is some overlap in the samples between grades. For example, the Year 2, Grade 4 cohort is the same as the Year 1, Grade 3 cohort; the Year 3, Grade 4 cohort is the same as the Year 2, Grade 3 cohort; and the Year 4, Grade 4 cohort is the same as the Year 3, Grade 3 cohort. Similarly, there are two overlapping cohorts between the Grade 3 and Grade 5 results.

17. However, we caution against overinterpreting this one result, which may be unique to the particular curricula we evaluate or could be the product of sampling variability. Note that in some years, the Grade 4 estimates are substantially smaller than the Grade 5 and Grade 3 estimates, but in other years, they are quite close, especially given the sizes of our standard errors.

18. A student is identified as socioeconomically disadvantaged by the CDE if either (a) both of his or her parents do not have a high school diploma or (b) he or she is eligible for free/reduced-price lunch.

19. In future work (by us or others), we hope to learn more about the extent of curriculum effect heterogeneity. Given the primary purpose of our paper is to perform the overall comparison, we do not have the capacity to give this question the full attention it deserves. We simply apply our models to the subgroup scores. Although this seems sufficient at a quick pass, it would be useful

to examine the potential for curriculum effect heterogeneity in more detail, along with the modeling assumptions. For example, although our falsification estimates provide no obvious indications of problems in Table 6, for any two subgroups  $S1$  and  $S2$ , even if the condition  $Y_0, Y_1 \perp D | X$  is satisfied, it need not guarantee  $Y_0^{S1}, Y_1^{S1} \perp D | X$  and  $Y_0^{S2}, Y_1^{S2} \perp D | X$ . Plausible mechanisms that might lead to a violation of conditional independence in the subgroup analyses when there is no conditional independence assumption violation overall are not obvious but merit additional attention. We save this deeper dive for future research.

20. Per the conversions used in this paper as described in the table notes, a 0.05/0.08 student-level standard deviation move corresponds to a 0.11/0.18 school-level standard deviation move.

21. Note that in addition to our comparison centered on *California Math* yielding a smaller differential effect size, the suggestive results from our initial pairwise comparisons (Appendix Table B1) imply that the other three curricula are similarly effective.

22. The use of covariates from before the 2009/2010 adoptions is not particularly important given that none of these schools used any of the curricula of interest, but we follow the same timing convention as in other parts of our analysis for consistency. We obtain similar results if we estimate Equation (C2) using data from different years.

23. We are not aware of a specific example of this particular problem causing bias, but the possibility is implied in related findings by Hansen (2008), who shows that bias can be caused when observations in one condition (either treatment or control) disproportionately supply identifying variation for covariates.

24. The covariates are as listed in Table 1. As noted above, we also use cubics in school and district enrollment and include a variable to indicate CDE data quality for individual schools.

25. We cluster our standard errors in Equation (C4) at the district level to mimic the conditions of our primary analysis in the main text.

26. We report covariate-by-covariate balancing results for the primary comparison in Appendix Table B3.

## References

- Aaronson, D., Barrow, L., & Sander, W. (2007). Teachers and student achievement in the Chicago public high schools. *Journal of Labor Economics*, 25(1), 95–135.
- Agodini, R., Harris, B., Atkins-Burnett, S., Heaviside, S., & Novak, T. (2010). *Achievement effects of four early elementary school math curricula: Findings for first and second graders* (NCEE 2011-4001). Washington, DC: National Center for Education Evaluation and Regional Assistance, U.S. Department of Education, Institute of Education Sciences.
- Bhatt, R., & Koedel, C. (2012). Large-scale evaluations of curricular effectiveness: The case of elementary mathematics in Indiana. *Educational Evaluation and Policy Analysis*, 34, 391–412.
- Bhatt, R., Koedel, C., & Lehmann, D. (2013). Is curriculum quality uniform? Evidence from Florida. *Economics of Education Review*, 34(1), 107–121.
- Black, D., & Smith, J. (2004). How robust is the evidence on the effects of college quality? Evidence from matching. *Journal of Econometrics*, 121(2), 99–124.
- Boser, U., Chingos, M., & Straus, C. (2015). *The hidden value of curriculum reform: Do states and districts receive the most bang for their curricular buck?* Washington, DC: Center for American Progress.

- Caliendo, M., & Kopeinig, S. (2008). Some practical guidance for the implementation of propensity score matching. *Journal of Economic Surveys*, 22(1), 31–72.
- California Department of Education. (2009). *2007 mathematics primary adoption report*. Sacramento, CA: California Department of Education Press.
- Chetty, R., Friedman, J. N., & Rockoff, J. E. (2014). Measuring the impacts of teachers I: Evaluating bias in teacher value-added estimates. *American Economic Review*, 104(9), 2593–2632.
- Chingos, M. M., & Whitehurst, G. J. (2012). *Choosing blindly: Instructional materials, teacher effectiveness and the Common Core*. Policy report, Brown Center on Education Policy, Washington, DC.
- Clotfelter, C. T., Ladd, H. F., & Vigdor, J. L. (2015). The aftermath of accelerating algebra: Evidence from district policy initiatives. *Journal of Human Resources*, 50(1), 159–188.
- Cortes, K., Goodman, J., & Nomi, T. (2015). Intensive math instruction and educational attainment: Long-run impacts of double-dose algebra. *Journal of Human Resources*, 50(1), 108–158.
- Currie, J., & Thomas, D. (2000). School quality and the longer-term effects of Head Start. *Journal of Human Resources*, 35(4), 755–774.
- Deming, D. J. (2009). Early childhood intervention and life-cycle skill development: Evidence from Head Start. *American Economic Journal: Applied Economics*, 1(3), 111–134.
- Domina, T., McEachin, A., Penner, A., & Penner, E. (2015). Aiming high and falling short: California's eighth-grade algebra-for-all effort. *Educational Evaluation and Policy Analysis*, 37(3), 275–295.
- Dougherty, S., Goodman, J., Hill, D., Litke, E., & Page, L. (2015). Middle school math acceleration and equitable access to 8th grade algebra: Evidence from the Wake County public school system. *Educational Evaluation and Policy Analysis*, 37(1), 80S–101S.
- Frölich, M. (2004). Finite-sample properties of propensity-score matching and weighting estimators. *Review of Economics and Statistics*, 86(1), 77–90.
- Hansen, B. B. (2008). The prognostic analogue of the propensity score. *Biometrika*, 95(2), 481–488.
- Heckman, J., Ichimura, H., & Todd, P. (1997). Matching as an econometric evaluation estimator: Evidence from evaluating a job-training programme. *Review of Economic Studies*, 64(4), 261–294.
- Jackson, K., & Makarin, A. (2016). *Simplifying teaching: A field experiment with online “off-the-shelf” lessons* (NBER Working Paper No. 22398). Cambridge, MA: National Bureau of Economic Research.
- Jobrack, B. (2011). *Tyranny of the textbook: An insider exposes how educational materials undermine reforms*. Lanham, MD: Rowman & Littlefield.
- Koedel, C., & Betts, J. R. (2011). Does student sorting invalidate value-added models of teacher effectiveness? An extended analysis of the Rothstein critique. *Education Finance and Policy*, 6(1), 18–42.
- Koedel, C., Mihaly, K., & Rockoff, J.E. (2015). Value-added modeling: A review. *Economics of Education Review*, 47, 180–195.
- Krueger, A. B., & Whitmore, D. M. (2001). The effect of attending a small class in the early grades on college test taking and middle school test results: Evidence from Project STAR. *Economic Journal*, 111(468), 1–28.
- Lechner, M. (2002). Program heterogeneity and propensity score matching: An application to the evaluation of active labor market policies. *Review of Economics and Statistics*, 84(2), 205–220.
- Mueser, P. R., Troske, K. R., & Gorislavsky, A. (2007). Using state administrative data to measure program performance. *Review of Economics and Statistics*, 89(4), 761–783.
- Opfer, V. D., Kaufman, J. H., & Thompson, L. E. (2016). *Implementation of K–12 state standards for mathematics and English language arts and literacy*. Santa Monica, CA: RAND.
- Papa, R., & Brown, R. (2007). *The Research Base for California Mathematics: Concepts, Skills, and Problem Solving*. Columbus, OH: The McGraw-Hill Companies.
- Perry, R. R., Finkelstein, N. D., Seago, N., Heredia, A., Sobolew-Shubin, S., & Carroll, C. (2015). *Taking stock of Common Core math implementation: Supporting teachers to shift instruction*. San Francisco, CA: WestEd.
- Polikoff, M. S. (2015). How well aligned are textbooks to the Common Core standards in mathematics? *American Educational Research Journal*, 52(6), 1185–1211.
- Rosenbaum, P. R., & Rubin, D. B. (1983). The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1), 41–55.
- Rosenbaum, P. R., & Rubin, D. B. (1985). The bias due to incomplete matching. *Biometrika*, 41(1), 103–116.
- Sales, A., Hansen, B. B., & Rowan, B. (2014). *Rebar: Reinforcing a matching estimator with predictions from high-dimensional covariates*. Unpublished manuscript
- Sianesi, B. (2004). An evaluation of the Swedish system of active labor market programs in the 1990s. *Review of Economics and Statistics*, 86(1), 133–155.
- Smith, J., & Todd, P. (2005). Rejoinder. *Journal of Econometrics*, 125(2), 365–375.
- Zeringue, J. K., Spencer, D., Mark, J., & Schwinden, K. (2010, April). *Influences on mathematics textbook selection: What really matters?* Paper presented at the Research Pre-session of the National Council of Teachers of Mathematics, San Diego, CA. Retrieved from [http://mcc.edc.org/pdf/Final\\_Draft\\_Research\\_Pre-session\\_2010.pdf](http://mcc.edc.org/pdf/Final_Draft_Research_Pre-session_2010.pdf)

## Authors

CORY KOEDEL is an associate professor of economics and public policy at the University of Missouri. His research interests include teacher quality and compensation, curricular effectiveness, and the efficacy of higher education institutions.

DIYI LI is a PhD candidate in economics at the University of Missouri. His early research focuses on curriculum evaluation and higher education.

MORGAN S. POLIKOFF is an associate professor of education at the University of Southern California Rossier School of Education. He studies the design, implementation, and effects of standards, assessment, and accountability policies.

TENICE HARDAWAY is a graduate student at University of Southern California Rossier School of Education. Her research focuses on K–12 school policy and reform with particular emphasis on school choice reforms.

STEPHANI L. WRABEL is an associate policy researcher at the RAND Corporation. Her research is focused on school accountability policy, student mobility, and military-connected students and schools.