

# Effectiveness of Cognitive Tutor Algebra I at Scale

**John F. Pane**  
**Beth Ann Griffin**  
*RAND Corporation*

**Daniel F. McCaffrey**  
*Educational Testing Service*

**Rita Karam**  
*RAND Corporation*

*This article examines the effectiveness of a technology-based algebra curriculum in a wide variety of middle schools and high schools in seven states. Participating schools were matched into similar pairs and randomly assigned to either continue with the current algebra curriculum for 2 years or to adopt Cognitive Tutor Algebra I (CTAI), which uses a personalized, mastery-learning, blended-learning approach. Schools assigned to implement CTAI did so under conditions similar to schools that independently adopt it. Analysis of posttest outcomes on an algebra proficiency exam finds no effects in the first year of implementation, but finds evidence in support of positive effects in the second year. The estimated effect is statistically significant for high schools but not for middle schools; in both cases, the magnitude is sufficient to improve the median student's performance by approximately eight percentile points.*

**Keywords:** *algebra, effectiveness, education technology, blended learning, randomized controlled trial*

## Introduction

COGNITIVE Tutor Algebra I (CTAI), published by Carnegie Learning, Inc., is a first-year algebra course designed for students at a variety of ability and grade levels. The curriculum includes traditional textbook and workbook materials along with innovative automated tutoring software that provides self-paced individualized instruction and attempts to bring students to mastery of a topic before progressing to more advanced topics. On the basis of prior evidence of this curriculum's efficacy in some isolated

contexts, we conducted a large-scale randomized controlled trial (RCT) evaluation to estimate its effectiveness when implemented in a wide variety of natural school settings, in conditions similar to those of schools that independently adopt the curriculum.

This article reports CTAI's effects on student mathematics achievement and student confidence and attitudes about mathematics. As one of the few rigorous large-scale evaluations to date of interventions that use a personalized, mastery-learning, blended-learning (partly

online and partly classroom-based) approach, the study contributes to the evidence base on this increasingly popular approach to incorporating technology into instruction.

## Background and Context

Mathematics proficiency rates of students in the United States continue to be a concern for educators and policy makers. Although scores and proficiency rates on the National Assessment of Educational Progress have been on an upward trend since 1990, in 2011 only 35% of eighth-grade students performed at a level of proficient or higher (National Center for Education Statistics, 2011). Similarly, even though eighth-grade U.S. students have improved in international comparisons between 1995 and 2011, they continue to lag well behind the top-scoring countries (Mullis, Martin, Foy, & Alka, 2012). Moreover, low percentages of high school graduates demonstrate preparedness for higher education (American College Testing [ACT], 2012).

Educators are undertaking a variety of efforts to address these concerns, and many of these efforts are focused on algebra because it is considered a gateway course to higher level mathematics and science courses. At the same time, widespread availability of computers and network connections inside and outside of schools has drawn greater attention to technology-based course materials. There is some empirical evidence suggesting that technology-based curricula can help personalize students' learning experiences and facilitate the development of mathematical skills (Koedinger, Corbett, Ritter, & Shapiro, 2000; Ritter, Kulikowich, Lei, McGuire, & Morgan, 2007; Schacter, 1999; Wenglinsky, 1998). However, a meta-analysis conducted by the U.S. Department of Education (2010) concluded that evidence of the effects of online learning was very weak, with few rigorous controlled studies that enabled computing effect sizes. Nonetheless, that meta-analysis estimated that interventions combining online and face-to-face instruction in a blended-learning approach produced larger positive effects than either online or face-to-face instruction alone. This meta-analysis appears to have been influential in spurring widespread development and adoption of blended-learning approaches in recent years.

Stakeholders see a number of potential advantages to online or blended-learning courses (U.S. Department of Education, 2012). Such courses can provide access to high-quality instruction outside of normal school times and places. They can afford more efficient use of teaching resources by providing teachers with frequent and detailed information about the progress and struggles of each student and enabling them to provide focused attention to some students while others work online. Perhaps most importantly, they are seen as having the potential to improve student achievement by providing more engaging and personalized instruction and more immediate feedback. Moreover, many educators have called for curriculum and teaching methods that emphasize active learning, build on prior knowledge, and apply mathematical concepts to real-world problems (National Council of Teachers of Mathematics [NCTM], 2000).

In summary, the field is in the early stages of confirming whether these approaches produce positive effects. Where rigorous positive evidence does exist, it primarily comes from efficacy trials, which evaluate the intervention under optimal conditions in a limited context. Large-scale tests of effectiveness in diverse real-world school contexts, without any extraordinary effort to optimize implementation, like this study, have been rare.

## The Intervention

CTAI is a technology-based mathematics curriculum designed to promote student understanding of algebraic concepts and principles, to develop students' problem-solving skills, and to enable them to master higher order mathematical concepts (Ritter, Anderson, Koedinger, & Corbett, 2007). It is part of a broader set of curricula covering a number of secondary mathematics courses. In addition to textbook materials, each course includes an automated computer-based Cognitive Tutor (Anderson, Corbett, Koedinger, & Pelletier, 1995) that provides individualized instruction to address students' specific needs. The individualization is built into the software and is facilitated by detailed computational models of student thinking in a domain. Through the tutor, students work on

challenging problems that reflect real-world situations and provide opportunities for students to progress from concrete to abstract thinking. The company recommends that students spend 2 days per week of their class time using the computer-based, individualized one-on-one tutorial provided by the software while the teacher works with individual students as needed, and 3 days on classroom activities that are student-centered and involve group work and problem solving, guided by the teacher and the textbook but not using the software.

The CTAI software utilizes multiple representations, including diagrams, equations, and text, and concepts are often contextualized in real-world problem scenarios. As delivered, it is aligned to NCTM (2000) standards, and schools or districts can customize it to align with state or local standards. The software is available for students to use during class time and at other times during the school day, and outside of school from public libraries or homes, although this study did not collect data on how much this outside-of-class use occurred. Because students' progress in the software is self-paced and progression depends on mastery of the material, their work with the tutor might not be synchronized with material being covered in the classroom.

Classroom lessons address topics such as solving linear equations and systems of linear equations, mathematical modeling with linear and quadratic expressions, problem solving using proportional reasoning, and analyzing data and making predictions. During the lessons, students complete worksheets and other activities in which they record answers to questions posed in the problem scenarios and are encouraged to engage in a variety of problem-solving strategies such as breaking an unfamiliar problem down into simpler problems.

As part of implementation, teachers receive 4 days of training. During a 3-day session prior to the start of the school year, teachers are introduced to the curriculum materials, tutor software, and teacher tools, and given suggestions for how to implement the curriculum and make connections between the software and classroom instruction, to apply student-centered, standards-based instructional strategies including effective questioning strategies, and to learn

how to use data from the software to inform their instruction. The fourth training day occurs during the school year, when professional development staff observe classrooms, offer recommendations, and help teachers address any problems they are having with implementation. Teachers also receive a set of training materials, an implementation guide, and a book of resources and assessments.

To summarize, the theory of action underlying the CTAI curriculum incorporates constructivist principles and research on learning and memory in a blended-learning approach where technology and teachers work together to advance student learning. This theory considers student learning to be an active process of creating meaning from a variety of experiences, including engagement in content, inquiry, problem solving, explanation, and collaboration with peers. Classroom time is designed to provide students with such experiences. The technology also provides opportunities for active engagement by offering problems that are well matched to each individual student's level of mastery, providing immediate feedback, and advancing each student at an appropriate pace. Meanwhile, the software frees teachers to provide focused support and instruction to individuals. It also collects data on the whole class's understanding and difficulties, which can help guide teachers' classroom instruction. These approaches hold promise to be more effective than more traditional approaches (e.g., Hiebert et al., 2005; Pashler et al., 2007; Stecker, Fuchs, & Fuchs, 2005). This curriculum design is hypothesized to lead to greater student engagement, increased exposure to diverse applications of mathematics, and greater opportunity to practice problem-solving skills, which, in turn, are hypothesized to lead to increased student learning.

### **Prior Research on the Curriculum**

The most rigorous evidence of the efficacy of CTAI comes from an RCT of the curriculum in the Moore Independent School District in Moore, Oklahoma (Morgan & Ritter, 2002). The study randomly assigned ninth-grade students to classrooms using either CTAI or the district's existing traditional Algebra I curriculum, and controlled for possible teacher effects

by having some teachers teach both types of classes. The students using CTAI scored significantly higher on the Educational Testing Service's Algebra End-of-Course Assessment, received higher grades, and demonstrated more positive attitudes toward math than their counterparts using the traditional curriculum. The study measured a positive achievement effect of .23 standard deviation units.<sup>1</sup>

Two other experiments testing the efficacy of CTAI found negative, though not significant, treatment effects for the curriculum. These were an experiment testing CTAI among students in Grades 8 to 12 in five schools in Hawaii<sup>2</sup> (Cabalo & Vu, 2007) and a 2-year multisite experiment sponsored by the U.S. Department of Education that evaluated CTAI in Grades 8 to 12 along with nine other reading and mathematics software products (Campuzano, Dynarski, Agodini, & Rall, 2009). The latter study did find a significant positive effect for experienced teachers using algebra software in the study's second year, although the result was for a composite of products that included CTAI. In addition, an experiment testing the efficacy of the related Cognitive Tutor Geometry curriculum in Grades 9 to 12 found a significant negative effect on student achievement (Pane, McCaffrey, Slaughter, Steele, & Ikemoto, 2010). The U.S. Department of Education's What Works Clearinghouse (WWC) reviewed these and other studies and concluded that CTAI and related curricula have "mixed effects on mathematics achievement for high school students" based on a medium to large evidence base (WWC, 2013).

### Design of the Effectiveness Trial

This study sought to implement the CTAI curriculum under conditions similar to those that exist when schools independently decide to obtain the curriculum from Carnegie Learning. Allowing the study schools to implement the curriculum authentically influenced our choice of research design and required us to consider several challenging pragmatic issues to retain experimental control.

The first major challenge was how to design the study so that it minimized the disruption to normal school operations such as the assignment

of students and teachers to curriculum and classes, and the ability of teachers to interact freely with colleagues in their school. To meet this challenge, we chose to randomize assignment at the school level, as this would not require any modifications to the activities within schools other than the implementation of the CTAI curriculum. Moreover, randomization of schools makes it easier to prevent crossover between the treatment and control conditions.

A second consideration was the necessity for the study to administer an algebra posttest, because assessment programs vary across states and across districts or schools within states. Combining the results from many different assessments would have been problematic; moreover, not all high school students are tested in mathematics every year and the amount of algebra content on existing assessments might be small, making it more difficult to detect intervention effects. However, administering an algebra assessment to all students in the school, regardless of grade level or the level of mathematics attained, would have seemed illogical and would have been disruptive, imposing a recruiting challenge because few schools would have been willing to participate in the study. Schoolwide testing would also have been cost-infeasible given the available funding for this study. Thus, we determined that the study would administer algebra proficiency posttests to only study participants, ruling out some analytic approaches that require measuring outcomes for all students in the school.

Third, and perhaps most challenging, was defining and holding constant the study population for the duration of the experiment, given that students who enroll in algebra come from a range of grade levels and schools typically do not have firm rules about algebra course taking. Generally, the majority of students take algebra in 9th grade, more advanced students take algebra in 8th grade or earlier, and lower achieving students might enroll in 10th grade or later. Schools also sometimes spread the curriculum over 2 years for lower-performing students. The definitions of these groups such as lower achieving, mainstream, or advanced/honors/gifted are not precisely specified and may drift over time, and the placement of individuals into these groups is often partly based on the discretion of

families, teachers, principals, or guidance counselors. Moreover, at the time of the study, some states and districts were establishing new policies to encourage more students to enroll in algebra in 8th grade, and some of the schools participating in the study were subject to these policies. For these reasons, authentic implementation would not enable the research team to dictate exactly which students would take algebra or in what grade, and thus the study design avoided doing so, allowing schools to assign students to algebra classes according to their normal routines.

A fourth consideration was the need to allow teachers time to prepare to teach the new curriculum under a typical timeline. Meeting this goal required that randomization of schools would occur well before the start of the school year so treatment teachers could receive curriculum training, and schools could install the curriculum software along with any hardware necessary to support it. Moreover, the study intended to conduct the experiment for 2 years in each school to be able to capture any improvement in implementation in the second year. This meant that second-year classes would begin more than a year after randomization. Thus, because of natural student mobility and the necessity for students to meet prerequisites, it was not considered authentic to require school officials to define, prior to randomization, the precise set of students to take algebra during the 2 years of the study. Instead, we sought to allow schools to retain discretion regarding when to enroll students in algebra and, as a consequence, the resulting population of students taking algebra at any particular time. Such discretion, left unchecked, would enable schools to change algebra enrollment patterns in response to their randomly assigned experimental condition, that is, after the treatment is assigned. For instance, treatment schools might view the adoption of the new algebra curriculum as an opportunity to change algebra enrollment patterns. Such enrollment changes in response to assigned experimental condition subvert the control that the researcher is trying to establish through randomization. Randomization is intended to ensure that the preexisting characteristics of the students in the treatment and control groups are unrelated to experimental conditions, but the

authentic discretion held by schools in this study could lead to systematic differences between groups.

Ultimately, we settled on an approach where, prior to randomization, schools specified *schema* for selecting the students who would participate for both years. The schema identified a set of criteria schools would apply to determine the students who would participate each year. Schools specified this by answering a series of questions on a form we provided to them as part of their enrollment into the study. The form asked whether the school would include all algebra classes in the study, or if not, which classes would be included or excluded according to teacher or variant of the class (i.e., gifted, honors, remedial, etc.) or other criteria. The form also inquired about the ability levels of the students in the designated classes (low, average, or high performing—check all that apply); and how many classes and teachers were estimated to participate. By answering these questions, schools specified an exact set of rules for selecting the types of students who would participate even though they did not specify the specific individuals. We monitored, and to the extent possible enforced, adherence to the schema throughout the study. We judged that this method achieved an appropriate balance between granting schools authentic discretion over algebra enrollment and providing sufficient control over experimental assignment to treatment, while avoiding serious feasibility issues. Moreover, we could include the schema as one of the characteristics used in blocking schools for randomization, helping to ensure that similar types of students would be participating in the treatment and control schools.

## Method

The project conducted two parallel experiments, one in middle schools and one in high schools. The study was designed and powered to examine these groups separately because the populations of students taking algebra in middle schools (Grade 8 or earlier) is generally higher achieving than the population of students taking algebra in high schools (Grades 9–12) and the curriculum might have different effects with these student populations or their schools.

### *Study Setting*

The study was conducted in 73 high schools and 74 middle schools in 51 school districts in seven states. Participating schools include urban, suburban, and rural public schools, and some Catholic Diocese parochial schools. The sites include city districts in Texas, Connecticut, New Jersey, and Alabama, suburban districts near Detroit, Michigan, generally rural districts in Kentucky, and districts throughout Louisiana. Each school participated for 2 years. All sites participated in both the middle school and high school arms of the study except Alabama (middle school only).

### *Study Population*

We define the study population as all students present at the time of the pretest as well as any additional students who entered the study later and remained for the posttest. Nearly 18,700 students in Grades 9 through 12 participated in the high school study, with 89% of the participants in 9th grade. Nearly 6,800 students in Grades 6 through 8 participated in the middle school study, with more than 99% of them in 8th grade.

### *Research Design*

The study used a pair-matched cluster randomized design to assign schools to study condition. Schools within each state were matched into pairs on a number of criteria, including school-level variables and the achievement profile of students targeted for participation, as specified by the schema that schools prepared as part of enrollment in the study. The full set of variables used in matching is shown in Table 1. To select the pairs for matching, we first calculated the distance between each pair of schools in a site, where distance equaled the sum of the absolute values of the standardized differences between schools on the matching variables. The algorithm then selected the set of pairs that minimized the average distance between schools in the pairs under the constraint that each school was in only one of the selected pairs. Schools were randomized in the spring prior to their first year of implementation.

Schools randomized to the treatment group implemented the CTAI curriculum and those

assigned to the control group continued to use their existing Algebra I curriculum. Those were published by Prentice Hall, Glencoe, and McDougal Littell. Assignments to treatment or control groups continued for 2 academic years in each school.

### *Data Collection*

The study-administered pretests and posttests from the CTB/McGraw-Hill Acuity series. About 3 weeks after the start of the algebra course, the study administered the Algebra Readiness Exam as a pretest. This is a 40-item multiple-choice assessment designed for students who have completed Grades 6 through 11 to assess their preparation in the skills necessary for successful performance in algebra. At the end of the course, the study administered the Algebra Proficiency Exam as a posttest. This is a 32-item multiple-choice assessment designed to measure mastery of Algebra I content knowledge at the end of the course. The exams were scored using a three-parameter item response theory (IRT) model, and posttest scores were standardized within the population analyzed to have a mean of 0 and standard deviation of 1, enabling regression coefficients to be read as standardized effect sizes.

We collected additional administrative data from district or state sources. These consisted of sociodemographic information, including race/ethnicity; gender; socioeconomic status, as indicated by eligibility for the federal free or reduced-price meal program (FRL); whether the student was an English-language learner (ELL); and whether the student was in special education or gifted programs. The administrative data also included state test scores from the 2 years prior to enrollment in the study for each student.

At the end of the algebra course, the study also administered a 17-item student survey that measured student opinions about mathematics, computers, the algebra course they just completed, and their future schooling plans. The survey was derived from Fennema and Sherman (1976) and a similar survey that RAND previously developed and administered in an efficacy study of Cognitive Tutor Geometry (Pane et al., 2010). From this survey, we derived scales on mathematics confidence ( $\alpha = .84$ ),<sup>3</sup> utility of

TABLE 1

*Variables Used for Matching and School-Level Balance After Randomization*

	High school study		Middle school study	
	Control group	Treatment group	Control group	Treatment group
Number of schools	35	33	37	37
School size	852	825	836	780
Race/ethnicity				
American Indian/Alaskan Native	0.00	0.00	0.00	0.00
Asian/Pacific Islander	0.02	0.01	0.03	0.02
Black non-Hispanic	0.34	0.41*	0.30	0.30
Hispanic	0.12	0.13	0.24	0.27
White non-Hispanic	0.52	0.45*	0.44	0.41
Student demographics				
Eligible for free or reduced-price lunch	0.38	0.44	0.63	0.69
Classified as English-learner	0.09	0.09	0.10	0.10
Proficiency rates on state assessments				
Reading 2004	0.52	0.50	0.70	0.66*
Reading 2005	0.55	0.51	0.75	0.72*
Reading 2006	0.54	0.55	0.75	0.75
Mathematics 2004	0.47	0.42	0.51	0.52
Mathematics 2005	0.47	0.43	0.54	0.50*
Mathematics 2006	0.44	0.42	0.58	0.57
Intended study population as defined by schema				
All classes	0.55	0.61	0.68	0.65
Identified teachers only	0.42	0.39	0.39	0.42
Includes low-achieving students	0.61	0.55	0.03	0.00
Includes average-achieving students	0.52	0.48	0.16	0.19
Includes high-achieving students	0.33	0.35	0.52	0.48
Number of classes projected to participate	7.04	6.56	2.30	2.03
Number of students projected to participate	161	153	56	48

*Note.* Five participating high schools are not represented in this table due to lack of school-level variables. These schools were matched by hand into a pair in Kentucky and a triplet in Michigan based on the shared characteristic of being alternative schools.

\* $p < .05$  for difference between treatment and control groups, calculated with a permutation test (see text; all  $p > .01$ ).

mathematics for the future ( $\alpha = .77$ ), technology confidence and enjoyment ( $\alpha = .73$ ), and opinion about the course ( $\alpha = .86$ ). Two remaining items, not part of scales, asked opinions of the utility of computers in learning math and future schooling plans.

### *Statistical Analyses*

To assess success of the random assignment blocked by matched schools, we examined pre-treatment group balance at the school level using a permutation test (Efron & Tibshirani, 1993). The permutation test is an approach for calculating the probability of obtaining the observed differences between the intervention and control schools by chance that does not rely on model assumptions

like a Wald test of whether the intervention and control schools are significantly different from each other. To calculate the  $p$  values in a permutation test, we randomly reassigned the treatment assignment indicators to schools following the randomization design (i.e., after schools were matched into pairs), to simulate the group differences with alternative realizations of the randomization. We repeated this process 10,000 times, and for each resampled data set, we calculated the group differences (mean differences for continuous measures and McNemar's test statistic for binary measures). The  $p$  value is the proportion of times the difference from the resampled data equals or exceeds the observed difference between the two groups from the randomization that actually occurred in our study.

Later, after the sample was defined, we assessed pretreatment group balance at the student level, using the same hierarchical model as described below (Model 1) for measuring posttest outcomes, by substituting the IRT pretest score for the outcome and determining if the treatment indicator was significantly associated with the pretest. Because this analysis revealed pretreatment differences between the treatment and control groups, described below in the “Results” section, we also fit covariate-adjusted models and models using prognostic scores (Hansen, 2008) that attempt to restore group balance, described below.

Including error-prone covariates in a regression model to control for systematic group differences without a correction for measurement error will generally result in biased estimates of all model parameters including the treatment indicator (Greene, 2003; Lockwood & McCaffrey, in press). To avoid this potential bias given the apparent purposive selection of students for CTAI classes, we used regression calibration (Carroll, Ruppert, Stefanski, & Crainiceanu, 2006) to correct for measurement error in the pretests. We implemented regression calibration by replacing each error-prone test score with a random variable drawn from the conditional distribution of the corresponding error-free test score given the error-prone test score and all other model covariates. The required conditional distribution was constructed assuming a linear relationship between the error-free scores and the other model covariates and used the conditional standard error of measure of each test score. To support standard error estimation, we imputed 20 values of the error-free test score for each student. An additional benefit of this approach was that we were able to impute missing pretest values for approximately 18% of the high school sample and 7% of the middle school sample who were absent for the study-administered pretest.

We also used as covariates students’ achievement test scores from the 2 years prior to their participation in the study. Because tests differ across states, our models include interaction terms between the prior test scores and state indicators, allowing for different relationships between prior achievement and the posttest by state. Within states, prior tests additionally

differ across grades, and students have different numbers of prior tests due to whether the state administers tests in a particular grade, student absences and transfers, and students who were retained in or skipped grades. Consequently, our models allow the coefficients on prior tests to depend not only on state but also on the grade and year in which a prior test was completed and on the number of prior tests available for the student. We did not use the prior test scores available for students with uncommon test taking patterns; instead these students were treated as having missing prior test scores. Overall, about 92% of the available scores were used in our models.

To enable models to control for any potential imbalance remaining after covariate adjustment, we estimated students’ expected posttest outcomes using prognostic scores (Hansen, 2008). Prognostic scores, which are defined as the predicted value of an outcome conditional on individuals being in the control condition, have been shown to be a useful tool for handling covariate imbalance in cases when multivariate adjustment is not sufficient (Arbogast & Ray, 2011; Hansen, 2008; Miettinen, 1976). Prognostic scores collapse the covariates of a study into a single measure that summarizes the covariates’ association with potential responses (here, potential posttest scores) had each study participant been in the control condition. They are particularly useful in settings like our study where information about the outcome’s relationship with covariates in the control condition is readily available. To estimate prognostic scores for each student in our study, we fit linear regression models to posttest scores in the control group sample, controlling for imputed pretest scores, prior state test scores, and student sociodemographic measures (race/ethnicity, gender, ELL, FRL, special education or gifted status, and grade level) along with the appropriate missingness indicators. These models were then used to calculate predicted posttest values for all students in both the treatment and control groups. Twenty prognostic scores were calculated for each student, 1 corresponding to each of the imputed pretest scores.

To estimate the impact of the treatment on student mathematics achievement and student confidence and attitudes about mathematics, we



compared the performance of the experimental (CTAI) and control (standard algebra) groups on the posttest scores and survey items. Specifically, we fit the following hierarchical linear models (Raudenbush & Bryk, 2002) for student posttest scores. For each model, let  $y_{ijk1}$  denote the score on the outcome for the  $k$ th student in classroom  $j = 1$  to  $J_i$  for school  $i = 1$  to  $I$  where  $J_i$  denotes the number of classrooms in school  $i$  and  $I$  denotes the total number of schools in the analysis. Let  $y_{ijk0}$  denote a student's score on one imputed version of the IRT pretest, centered to have mean 0. At the first level, the student level, we model the student's score at the end of the course using four different specifications. In Model 1, we only include the overall classroom mean ( $\mu_{ij}$ ) and a student-specific residual error term ( $\varepsilon_{ijk}$ ) at Level 1:

$$y_{ijk1} = \mu_{ij} + \varepsilon_{ijk}, \quad (\text{Model 1})$$

where the  $\varepsilon_{ijk}$  are independent  $N(0, \sigma^2)$  random variables. The mean and the slope parameters are classroom specific and specified via the classroom level model:  $\mu_{ij} = \gamma_{0i} + \eta_{ij0}$  with the term  $\eta_{ij0}$  representing random classroom effects that are assumed to be normal with mean zero and an unknown variance. The term  $\gamma_{0i}$  is specified in a school-level model:

$$\gamma_{0i} = \omega_{00} + \omega_{01}T_i + \omega_{02[i]} + \zeta_{i0} \quad (\text{school-level model})$$

where  $T_i$  indicates the school's treatment assignment (0 for traditional curriculum and 1 for the CTAI curriculum),  $\omega_{02[i]}$  denote the fixed effects for matched pairs corresponding to the pair for school  $i$ , and  $\zeta_{i0}$  denotes normally distributed random school error terms each with mean zero and variance  $\tau^2$ . The random school and classroom effects allow for the inherent clustering of outcomes at these levels of the hierarchical sample. The models were fit and parameter testing was doing using the lme command in R.

The effect of the CTAI curriculum on student achievement was tested by testing the null hypothesis that  $\omega_{01} = 0$  with two-tailed test and a .05 level of significance. Because the pretest scores are centered to have mean zero,  $\omega_{01}$  estimates the average effect of the intervention on students in the study.

Additional models use the same structure but extend the Level 1 model to increasingly control for variables that potentially confound the effect of the intervention and to increase precision of the estimated treatment effects. In Model 2, students' imputed pretest scores ( $y_{ijk0}$ ) are included as covariates, such that

$$y_{ijk1} = \mu_{ij} + \beta y_{ijk0} + \varepsilon_{ijk}. \quad (\text{Model 2})$$

Standard errors for multiple imputation results are calculated as described in Rubin (1987), with adjustments to degrees of freedom and significance tests based on Barnard and Rubin (1999).

In Model 3, we add more student covariates, including prior state test scores and sociodemographic measures along with appropriate missingness indicators. Namely, we have,

$$y_{ijk1} = \mu_{ij} + \beta y_{ijk0} + \alpha'x_{ijk0} + \varepsilon_{ijk}, \quad (\text{Model 3})$$

where  $x_{ijk0}$  is the vector of the student variables and  $\alpha$  is the vector of corresponding regression coefficients.

Finally, Model 4 additionally includes four dummy indicators of prognostic score quintiles within each matched pair in the Level 1 model.

$$y_{ijk1} = \mu_{ij} + \beta y_{ijk0} + \alpha'x_{ijk0} + \sum_{jm} \gamma_{jm} p_{mijk} + \varepsilon_{ijk}, \quad (\text{Model 4})$$

where  $p_{mijk}$  denotes the indicator for whether the student fell into the  $m$ th quintile ( $m = 1, \dots, 4$ ) and  $\gamma_{mi}$  denotes the regression coefficient for each pair by quintile dummy.

Finally, to examine treatment effects by student ability level, we supplemented Model 4 to include interaction terms between the prognostic score quintiles and the treatment indicator. We used quintiles of the prognostic score in our last two models rather than the continuous values to allow for nonlinearities between prognostic scores and the outcome and to facilitate the examination of nonlinear interaction effects between the prognostic scores and treatment.

Our analysis plan specified that we would analyze results separately by cohort (the first or second year of implementation in the school) to allow for the possibility that implementation might be better the second year due to teachers gaining experience with the curriculum.

TABLE 2  
*High School Study Attrition and Group Balance*

	Treatment group		Control group		Group difference <sup>a</sup>	<i>p</i> value
	<i>n</i>	Pretest mean	<i>n</i>	Pretest mean		
Cohort 1						
Eligible sample	4,541	−.468	5,014	−.347	−.194	.033*
Attrition	1,330	−.673	1,723	−.485	−.155	.091
Final sample	3,211	−.365	3,291	−.258	−.144	.187
Attrition rate	29.3%		34.4%			
Cohort 2						
Eligible sample	3,990	−.390	5,146	−.359	−.099	.284
Attrition	1,058	−.591	1,135	−.619	−.067	.529
Final sample	2,932	−.302	4,011	−.270	−.111	.276
Attrition rate	26.5%		22.1%			
Both cohorts						
Eligible sample	8,531	−.432	10,160	−.353	−.139	.082
Attrition	2,388	−.637	2,858	−.538	−.106	.205
Final sample	6,143	−.335	7,302	−.265	−.116	.188
Attrition rate	28.0%		28.1%			

*Note.* Eligible sample is defined as students present at pretest or entering the study after pretest. Attrition is defined as the portion of the eligible sample that did not take the posttest.

<sup>a</sup>Standardized mean difference in pretest scores between treatment and control groups (negative indicates treatment scored lower than control) as calculated by a model that includes fixed effects for randomization pairs and random effects for classrooms within schools.

Additional analyses examined whether treatment was associated with a number of secondary outcomes measuring student attitudes and confidence toward mathematics and technology. Specifically, models like Model 3 were fit to each of the following survey scales: mathematics confidence, utility of mathematics for the future, technology confidence and enjoyment, utility of computers in learning math, opinion about the course, and future schooling plans.

Finally, additional analyses explored whether there was any evidence of interactions with treatment by site or, for Cohort 2 only, whether the teacher was new to the study or was in the study the previous year. The latter analysis explores whether implementation of CTAI might be more effective during the second-year teachers use it, by comparing with teachers in the control group who were also present both years.

## Results

Table 1 summarizes balance on school-level variables between the treatment and control groups. Variables examined include school size;

percentages of the student population by race/ethnicity, eligibility for free or reduced-price lunch, or classified as English-learners; school-wide proficiency rates on state mathematics and reading assessments for the three most recent years; and variables defined by the schema schools planned to use to select their intended study populations. Randomization achieved balance on most, but not all of these variables. In the high school study, the overall student bodies of schools assigned to treatment had a greater proportion of Black students and fewer White students; in addition, the percentage of students classified as proficient was lower the previous 3 years in mathematics, and 2 of the previous 3 years in reading, although the proficiency differences were not significant. In the middle school study, the overall student bodies of schools assigned to treatment had significantly fewer students classified as proficient in reading in 2004 and both reading and mathematics in 2005.

Table 2 summarizes information about student participants in the high school sample. The final sample included 13,445 students after approximately 28% attrition from both treatment and

TABLE 3  
*Middle School Study Attrition and Group Balance*

	Treatment group		Control group		Group difference <sup>a</sup>	<i>p</i> value
	<i>n</i>	Pretest mean	<i>n</i>	Pretest mean		
Cohort 1						
Eligible sample	1,681	.265	1,743	.654	−.296	.021*
Attrition	169	−.067	212	.318	−.295	.135
Final sample	1,512	.306	1,531	.706	−.312	.016*
Attrition rate		10.1%		12.2%		
Cohort 2						
Eligible sample	1,534	.400	1,828	.738	−.422	.003**
Attrition	170	0.015	295	.733	−.392	.012*
Final sample	1,364	.449	1,533	.739	−.347	.016*
Attrition rate		11.1%		16.1%		
Both cohorts						
Eligible sample	3,215	.331	3,571	.698	−.366	.003**
Attrition	339	−.026	507	.560	−.385	.007**
Final sample	2,876	.376	3,064	.722	−.328	.007**
Attrition rate		10.5%		14.2%		

*Note.* Attrition is defined as the portion of the eligible sample that did not take the posttest.  
<sup>a</sup>Standardized mean difference in pretest scores between treatment and control groups (negative indicates treatment scored lower than control) as calculated by a model that includes fixed effects for randomization pairs and random effects for classrooms within schools. Eligible sample is defined as students present at pretest or entering the study after pretest.  
 \*Significance at the  $p < .05$  level. \*\*Significance at the  $p < .01$  level.

control groups. In Cohort 1, the treatment group scored .14 standard deviation units lower than the control group on the pretest ( $p = .19$ ), and in Cohort 2, the treatment group scored .11 lower ( $p = .28$ ). Similarly, Table 3 summarizes information about student participants in the middle school sample, which included 5,940 students after attrition of approximately 11% in the treatment group and 14% in the control group. In Cohort 1 of this study, the treatment group scored .31 standard deviation units lower than the control group on the pretest ( $p = .02$ ), and in Cohort 2, the treatment group scored .35 lower ( $p = .02$ ). Similar group differences are also apparent on students' prior state test scores (not shown in tables).

Table 4 summarizes the results for the high school study. Models consistently estimated negative treatment effects for Cohort 1, ranging from .10 to .19 standard deviation units and not significant. In contrast, models for Cohort 2 consistently estimated positive treatment effects, ranging from .14 to .21 standard deviation units; results for Models 2 through 4 are all below the .05 level of significance.

Similarly, Table 5 summarizes the results for the middle school study. Here again, models

estimated treatment effects for Cohort 1 that are not significant. These estimates are near zero in all the models that adjust for pretest scores. For Cohort 2, the unadjusted estimate is negative, and the estimates became positive with covariate adjustment. Although these estimates for middle school Cohort 2 are not significant, the estimated treatment effects are similar in magnitude to those found in the high school study.

Figures 1 and 2 show the estimated treatment effects for each of the prognostic score quintiles from our regression models that included interaction terms between the quintile indicators and treatment. For Cohort 2 in both the middle schools and high schools, effects are stable across the prognostic score quintiles and the interactions between the quintiles and the treatment indicator are not significant (joint Wald test  $p$  values = .51 and 1.00, respectively). Conversely, there is highly significant evidence of moderation by the prognostic score quintiles in middle school Cohort 1 (joint Wald test  $p$  value  $< .001$ ; see Figure 2) that indicates that there were potentially moderately large positive treatment effects in the lowest quintile and small negative effects of treatment in the highest two

TABLE 4  
*High School Study Treatment Effect Estimates*

Model	Cohort 1				Cohort 2			
	Estimate	SE	<i>t</i> value	<i>p</i> value	Estimate	SE	<i>t</i> value	<i>p</i> value
1	-.19	.12	-1.68	.10	.14	.12	1.20	.24
2	-.12	.10	-1.20	.24	.19	.09	2.05	.05 <sup>a,*</sup>
3	-.10	.10	-0.97	.34	.22	.09	2.33	.03*
4	-.10	.10	-1.02	.31	.21	.10	2.23	.03*

*Note.* Model 1 estimates group differences without any covariates; Model 2 includes regression-calibrated pretest scores; Model 3 also includes additional student covariates; and Model 4 includes all covariates as well as prognostic score quintiles.

<sup>a</sup>Value is less than .05 before rounding.

\*Significance at the  $p < .05$  level.

TABLE 5  
*Middle School Study Treatment Effect Estimates*

Model	Cohort 1				Cohort 2			
	Estimate	SE	<i>t</i> value	<i>p</i> value	Estimate	SE	<i>t</i> value	<i>p</i> value
1	-.20	.15	-1.34	.19	-.07	.17	-0.41	.69
2	.00	.10	0.01	.99	.17	.13	1.30	.21
3	-.03	.11	-0.24	.81	.19	.13	1.44	.16
4	.01	.11	0.11	.91	.19	.14	1.38	.17

*Note.* Model 1 estimates group differences without any covariates; Model 2 includes regression-calibrated pretest scores; Model 3 also includes additional student covariates; and Model 4 includes all covariates as well as prognostic score quintiles.

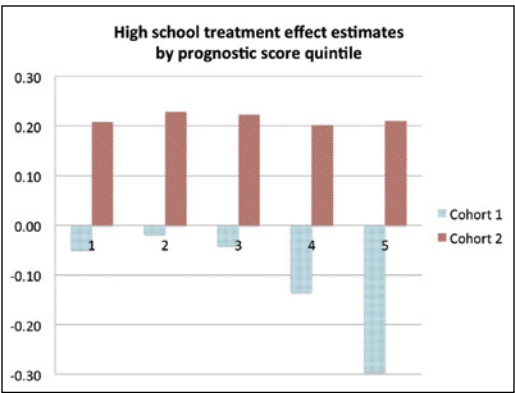


FIGURE 1. *Estimated treatment effects within the five prognostic score quintiles (1 = lowest and 5 = highest) for the high school study cohorts, calculated using a variant of Model 4 with interaction terms between the prognostic score quintiles and the treatment indicator.*

quintiles. Nonetheless, it is important to note that all of the quintile-specific treatment effect estimates for middle school Cohort 1 have confidence intervals that contain 0. Similar negative effects of treatment in the highest two quintiles

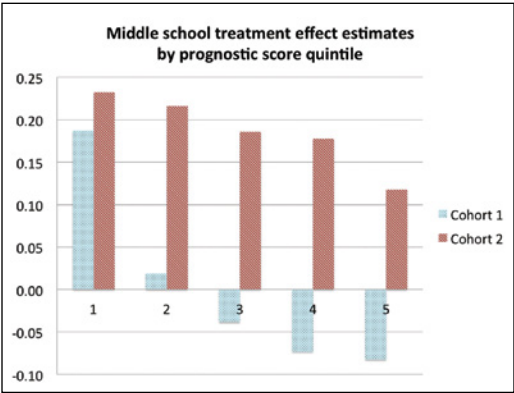


FIGURE 2. *Estimated treatment effects within the five prognostic score quintiles (1 = lowest and 5 = highest) for the middle school study cohorts, calculated using a variant of Model 4 with interaction terms between the prognostic score quintiles and the treatment indicator.*

were found for high school Cohort 1 students (see Figure 1); however, the joint Wald test for this cohort is not significant ( $p = .30$ ).

Analyses involving secondary outcomes on student attitudes and confidence in mathematics

and technology only revealed one significant relationship. Across all cohorts, students in the treatment condition reported significantly higher mean scores on the item that asked about the utility of computers in learning math.

Finally, we found no significant interactions between the treatment and treatment sites nor treatment and the indicator of whether a teacher was previously in the study.

## Discussion

It is necessary to be cautious in interpreting these results because mean student pretest scores were lower in the treatment group than the control group. The pretreatment differences, at  $-.14$  and  $-.11$  standard deviation units in the two cohorts in the high school study, are within the limit of  $-.25$  set by the WWC (U.S. Department of Education, 2011) for acceptable pretreatment differences. However, the differences are much larger in the middle school study. At  $-.31$  and  $-.35$  standard deviation units for the two cohorts, they exceed the WWC guideline and raise concern about the validity of the middle school experiment.

Our assessment of postrandomization balance on school-level variables found imbalance on some variables suggesting the treatment groups in both studies may have been disadvantaged through the luck of random assignment. Such differences are not unexpected when randomizing with such small sample sizes (at the school level); however, this cannot fully explain the magnitude of pretreatment differences on student pretests in the middle school study. Schools may have differentially exercised discretion over the population of students taking algebra and participating in the study depending on whether they were assigned to the treatment or control groups. As discussed above, the study was designed to allow for this discretion within limits; schools specified schema to describe the population of students in the study, but not the precise set of students. We monitored adherence to these schema and did not detect serious non-compliance. Nonetheless, schools may have worked within the constraints of the schema to shape the population of students participating in the study after becoming aware of treatment assignment. This may have occurred if, for

example, treatment schools believed the new curriculum would have positive effects, causing them to encourage more low-performing students to enroll in CTAI classes. Conversely, they may have believed the curriculum was better for lower-performing students and thus might have enrolled higher performing students in algebra classes that were not using CTAI and thus not part of the study. We explored these or other potential explanations with school officials and they uniformly expressed no awareness of such deliberate inclusion or exclusion of students.

For whatever reason, in both years CTAI students underperformed control group students on the pretest by a modest and nonsignificant amount in the high school study and a greater, significant amount in the middle school study. This could potentially raise concerns of bias in the impact estimates for both studies even though pretreatment group differences in the high school study are well within levels generally considered acceptable (e.g., U.S. Department of Education, 2011). Models 2 through 4 attempt to address the bias and improve precision through covariate and prognostic score adjustments.

The first thing to note is that for both studies and both years the treatment effect estimates from Model 1, which has no covariate adjustments, are substantially lower than for the models with adjustments, as would be expected if pretreatment differences were biasing the estimate. Second, in year 2 of the high school study, all three of the adjusted estimates are positive, of similar magnitude, and significant at the .05 level. Results from Models 2 through 4 for year 2 of the middle school study are also all positive and of similar magnitude to each other and to the effects for the year 2 high school cohort. However, the middle school results are not significant because of the smaller sample size<sup>4</sup>; enrollment in algebra is less common in middle school, consequently the middle school cohort of students is only about one third the size of the high school cohort.

Together the results provide strong evidence in support of a positive effect for CTAI in high schools in year 2: When we control for the selection of lower achieving students in the treatment group, apparent on the basis of their

prior achievement, we find positive effects that are robust to model specification. The replicated estimate in the underpowered middle school study is suggestive that the result may apply there as well. To help interpret the importance of these second-year effects, consider a student who would score at the 50th percentile of the posttest distribution if they were in the control group; an effect size of .20 is equivalent to having that student score at the 58th percentile if they were in the treatment group.

Another way to gauge the magnitude of an effect of .20 is to consider the amount of mathematics growth typically seen from one grade to the next on nationally normed assessments. Lipsey et al. (2012) report typical mathematics achievement gains of .32 for 8th grade, .22 for 9th grade, and .25 for 10th grade (most students in our two studies were in these grades). These data suggest that a treatment effect of .20 is nearly comparable in size with the amount of mathematics growth typically seen over a full year for 9th or 10th graders, or about two thirds of the amount typically seen for 8th graders. However, it worth noting that the typical gains cited here are from spring to spring (e.g., spring of 7th grade to spring of 8th grade for 8th graders) and thus might include the effects of any decline in achievement over the summer when school is in recess. Such a summer decline would not be captured in our effect estimates because we measured fall to spring gains.

Examination of prognostic score quintiles suggests that the positive effect for high school Cohort 2 was relatively uniform for students of all ability levels. Estimates for all five quintiles were in the range of .20 to .23. For middle score Cohort 2, the effect estimates are .23 for the lowest performing quintile and decrease for students of greater ability to .12 for the highest quintile. However, results from the two studies are not directly comparable because, as is evident in mean pretest scores in Tables 2 and 3, middle school students in the study are much higher achieving than their high school counterparts.

Finally, although the positive results for Cohort 2 varied from site to site, lack of statistical significance does not support any attempt to interpret this variation.

In contrast to the positive results for the second year of implementation, treatment effect

estimates are not significant the first year. The estimates are negative in the high school study and near zero in the middle school study. Examination of the effects by prognostic score quintiles suggest that the poor results in the first year may have been concentrated among higher performing students in both middle schools and high schools.

It is quite interesting that significant positive effects emerge in high school Cohort 2 after the negative (though not significant) results the prior year. One potential explanation is that teachers improved their implementation of CTAI after a year of experience using it. We explored this by dividing the Cohort 2 teachers into two groups by whether they were also in the study the prior year, and testing whether treatment effects were more positive among those present the prior year. This analysis produced mixed, nonsignificant results that thus do not lend support to this hypothesis. This question can be further informed by examining implementation data the study collected through teacher surveys and site visits.

Karam, Pane, Griffin, and Slaughter (2013) explores a variety of implementation questions, including the extent to which teachers report implementing the curriculum as specified by the CTAI developers, the relationship between implementation variables and outcomes, and how implementation changed over time. We identified the set of activities and practices recommended by the developer for high fidelity implementation, specifying not only whether particular activities should be present, but also the recommended amount of time or emphasis it should receive. As examples, the developer recommends that students spend 40% of instructional time using the CTAI software, and recommends that students spend less than 10% of instructional time taking notes. The article finds that although teachers generally implemented all components of the program, they did not emphasize the components in the same proportions as recommended, either underimplementing or overimplementing relative to the recommendation. As a result, no teacher received a perfect fidelity score, and on average fidelity scores were moderate. The fidelity scores did not change significantly between cohorts and analyses did not identify significant associations between fidelity scores and student outcomes.

With the exception of some components related to the CTAI software and materials, most implementation measures (e.g., the amount of time students spent taking notes) were also measured in control group classrooms. Among this subset of components, Karam et al. (2013) find that treatment and control group teachers reported contrasts in instructional practices that appear to have been induced by the curriculum. Relative to the control group, treatment group teachers reported less implementation of traditional practices such as lecturing with students taking notes and greater implementation of student-centered practices such as facilitating student work or assigning students to work in groups and give presentations. These differences are aligned with the recommendations of the developer for implementing CTAI. The article also finds that these contrasts in instructional practices were greater the first year than they were the second year, suggesting that treatment group teachers reverted somewhat back toward the more traditional practices over time. One interpretation of these results is that the reversion was an adaptation in response to poor results the first year, and that the adaptation was productive in that it coincided with positive treatment effect estimates for Cohort 2 students.

We found no meaningful effects of CTAI on secondary outcomes such as student attitudes and confidence in any of the cohorts.

Daugherty, Phillips, Pane, and Karam (2012) examines the costs of CTAI relative to the curricula in use in the control group schools in this study. While that analysis finds that CTAI is substantially more expensive, the cost must be weighed alongside the benefits reported herein. Educators may judge that the positive effects are large enough to warrant the additional cost.

## Conclusion

This large-scale effectiveness trial of CTAI finds a significant positive effect in high schools in the second year of implementation, relative to similar schools that continued to implement a variety of existing textbook-based algebra curricula. The effect size of approximately .20 is educationally meaningful—equivalent to moving an Algebra I student from the 50th to the

58th percentile. This positive result is important for educators and policy makers who are seeking interventions to improve Algebra I achievement, and is particularly notable because it was obtained in an effectiveness trial, where broad variety of schools implemented the curriculum without extraordinary support. The results may also be of broader potential interest because this curriculum uses technology to enable a personalized, blended-learning approach. As one of the first large-scale effectiveness trials of this type of intervention, the results help to inform researchers and practitioners whether this may be a productive way to employ technology to improve student achievement in mathematics or other subjects.

## Authors' Note

The opinions expressed are those of the authors and do not represent views of the Institute of Education Sciences or the U.S. Department of Education.

## Acknowledgments

The authors would like to express their gratitude to the many individuals who helped to make this study possible: Current and former RAND researchers Andrea Phillips, Abby Robyn, Nidhi Kalra, Regan Main, and J. R. Lockwood made extraordinary contributions, as did administrative assistants Crystal Baksis and Melanie Rote. Also playing important research or support roles were Scott Ashwood, Diane Bronowicz, Richard Bowman, Jaime Connors, Lindsay Daugherty, Maria Edelen, Amy Haas, Ann Haas, Laura Hamilton, Mark Hanson, Gina Ikemoto, Brian McInnis, Scott Naftel, Lawrence Painter, Louis Ramirez, Mary Ellen Slaughter, Anisah Waite, and Deborah Wesley. JoAnn Arcement, Michelle Auster, John Dilegghio, Barbara Dilegghio, Gayle Glusman, Kathy Hughes, Gary Kubina, and Diana Perez provided essential coordination and support at the seven research sites. We offer our special thanks to the state and district administrators who helped with recruiting, support, and data access, including David Akridge, Gary Asmus, Luellen Bledsoe, Rebecca Feola, Debbie Ferry, Michael Henderson, Kevin Hill, Monica Kendall, Gayle Kirwan, Brian McCarty, Karen Mohr, Ricardo Rosa, Marianne Srock, Liz Storey, and Kelly Trlica. We thank the teachers and principals in the participating schools, without whose participation this research would not have been possible. We also thank Carnegie Learning personnel who supported the study, including Sandy Bartle,

Steve Fancsali, Joseph Goins, Christy McGuire, Tristan Nixon, Steve Ritter, and Sean Sykes, as well as the company's field support and training staff. Finally, we thank Paco Martorell, who reviewed and provided helpful feedback on an earlier draft of this article, and the anonymous reviewers of our submission to this journal. The first author assumes full responsibility for any omissions from this acknowledgment, and offers his apologies and appreciation to those persons as well.

### Declaration of Conflicting Interests

The author(s) declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.

### Funding

The author(s) disclosed receipt of the following financial support for the research, authorship, and/or publication of this article: The research reported in this article was supported by the Institute of Education Sciences, U.S. Department of Education, through Grant R305A070185 to the RAND Corporation.

### Notes

1. Shown is the adjusted effect size calculated by the What Works Clearinghouse (2004).
2. The study also included some students enrolled in a remedial course at a community college.
3. We report Cronbach's alpha, a measure of internal consistency reliability.
4. The analysis for middle school Cohort 2 had power to detect an effect of approximately .26 or larger.

### References

- American College Testing. (2012). *The reality of college readiness: National*. Iowa City, IA: Author.
- Anderson, J. R., Corbett, A. T., Koedinger, K. R., & Pelletier, R. (1995). Cognitive tutors: Lessons learned. *Journal of the Learning Sciences*, 4, 167–207.
- Arbogast, P. G., & Ray, W. A. (2011). Performance of disease risk scores, propensity scores, and traditional multivariable outcome regression in the presence of multiple confounders. *American Journal of Epidemiology*, 174, 613–620.
- Barnard, J., & Rubin, D. B. (1999). Small-sample degrees of freedom with multiple imputation. *Biometrika*, 86, 948–955.
- Cabalo, J. V., & Vu, M.-T. (2007). *Comparative effectiveness of Carnegie Learning's Cognitive Tutor Algebra I Curriculum: A report of a randomized experiment in the Maui School District*. Palo Alto, CA: Empirical Education.
- Campuzano, L., Dynarski, M., Agodini, R., & Rall, K. (2009). *Effectiveness of reading and mathematics software products: Findings from two student cohorts (NCEE 2009-4041)*. Washington, DC: National Center for Educational Evaluation and Regional Assistance, Institute of Education Sciences, U.S. Department of Education.
- Carroll, R. J., Ruppert, D., Stefanski, L. A., & Crainiceanu, C. M. (2006). *Measurement error in nonlinear models: A modern perspective* (2nd ed.). London, England: Chapman & Hall.
- Daugherty, L., Phillips, A., Pane, J. F., & Karam, R. (2012). *Analysis of Costs in an Algebra I Curriculum Effectiveness Study*. RAND Corporation.
- Efron, B., & Tibshirani, R. J. (1993). *An introduction to the bootstrap*. Boca Raton, FL: Chapman & Hall.
- Fennema, E., & Sherman, J. A. (1976). Fennema-Sherman mathematics attitude scales: Instruments designed to measure attitudes toward the learning of mathematics by males and females. *JSAS Catalog of Selected Documents of Psychology*, 6(1), 31.
- Greene, W. H. (2003). *Econometric analysis* (5th ed.). Upper Saddle River, NJ: Prentice Hall.
- Hansen, B. B. (2008). The prognostic analogue of the propensity score. *Biometrika*, 95, 481–488.
- Hiebert, J., Stigler, J. W., Jacobs, J. K., Givvin, K. B., Garnier, H., Smith, M., . . . Gallimore, R. (2005). Mathematics teaching in the United States today (and tomorrow): Results from the TIMSS 1999 Video Study. *Educational Evaluation and Policy Analysis*, 27, 111–132.
- Karam, R., Pane, J. F., Griffin, B. A., & Slaughter, M. E. (2013). *Evaluating Cognitive Tutor Algebra I Curricula at Scale: Focus on implementation*. Manuscript submitted for publication.
- Koedinger, K. R., Corbett, A. T., Ritter, S., & Shapiro, L. J. (2000). *Carnegie Learning's Cognitive Tutor: Summary research results*. Pittsburgh, PA: Carnegie Learning.
- Lipsey, M. W., Puzio, K., Yun, C., Hebert, M. A., Steinka-Fry, K., Cole, M. W., . . . Busick, M. D. (2012). *Translating the statistical representation of the effects of education interventions into more readily interpretable forms* (NCSER 2013-3000). Washington, DC: National Center for Special Education Research, Institute of Education Sciences, U.S. Department of Education.



- Lockwood, J. R., & McCaffrey, D. F. (in press). Should nonlinear functions of test scores be used as covariates in a regression model? In R. Lissitz (Ed.), *Value-added modeling and growth modeling with particular application to teacher and school effectiveness*. Charlotte, NC: Information Age Publishing.
- Miettinen, O. S. (1976). Stratification by a multivariate confounder score. *American Journal of Epidemiology*, 104, 609–620.
- Morgan, P., & Ritter, S. (2002). *An experimental study of the effects of Cognitive Tutor Algebra I on student knowledge and attitude*. Pittsburgh, PA: Carnegie Learning.
- Mullis, I. V. S., Martin, M. O., Foy, P., & Alka, A. (2012). *TIMSS 2011 International Results in Mathematics*. Chestnut Hill, MA: TIMSS & PIRLS International Study Center, Lynch School of Education, Boston College.
- National Center for Education Statistics. (2011). *The Nation's Report Card: Mathematics 2011* (NCES 2012–458). Washington, DC: Institute of Education Sciences, U.S. Department of Education.
- National Council of Teachers of Mathematics. (2000). *Principles and Standards for School Mathematics*. Reston, VA: Author.
- Pane, J. F., McCaffrey, D. F., Slaughter, M. E., Steele, J. L., & Ikemoto, G. S. (2010). An experiment to evaluate the efficacy of Cognitive Tutor Geometry. *Journal of Research on Educational Effectiveness*, 3, 254–281.
- Pashler, H., Bain, P., Bottge, B., Graesser, A., Koedinger, K., McDaniel, M., & Metcalfe, J. (2007). *Organizing instruction and study to improve student learning* (NCER 2007-2004). Washington, DC: National Center for Education Research, Institute of Education Sciences, U.S. Department of Education.
- Raudenbush, S. W., & Bryk, A. S. (2002). *Hierarchical linear models*. Thousand Oaks, CA: SAGE.
- Ritter, S., Anderson, J. R., Koedinger, K. R., & Corbett, A. (2007). The Cognitive Tutor: Applied research in mathematics education. *Psychonomics Bulletin & Review*, 14, 249–255.
- Ritter, S., Kulikowich, J., Lei, P., McGuire, C. L., & Morgan, P. (2007). What evidence matters? A randomized field trial of Cognitive Tutor Algebra I. *Supporting Learning Flow Through Integrative Technologies*, 162, 13–20.
- Rubin, D. B. (1987). *Multiple imputation for nonresponse in surveys*. New York, NY: John Wiley & Sons.
- Schacter, J. (1999). *The impact of education technology on student achievement: What the most current research has to say*. Santa Monica, CA: Milken Exchange on Educational Technology, Milken Family Foundation.
- Stecker, P. M., Fuchs, L. S., & Fuchs, D. (2005). Using curriculum-based measurement to improve student achievement: Review of research. *Psychology in the Schools*, 42, 795–819.
- U.S. Department of Education. (2010). *Evaluation of evidence-based practices in online learning: A meta-analysis and review of online learning studies*. Washington, DC: Office of Planning, Evaluation, and Policy Development.
- U.S. Department of Education. (2011). *What Works Clearinghouse: Procedures and standards handbook* (Version 2.1). Washington, DC: Institute of Education Sciences.
- U.S. Department of Education. (2012). *Understanding the implications of online learning for educational productivity*. Washington, DC: Office of Educational Technology.
- Wenglinsky, H. (1998). *Does it compute? The relationship between educational technology and student achievement in mathematics*. Princeton, NJ: Educational Testing Service Policy Information Center.
- What Works Clearinghouse. (2004). *WWC Intervention Report: Cognitive Tutor®*. Washington, DC: U.S. Department of Education, Institute of Education Sciences.
- What Works Clearinghouse. (2013). *WWC Intervention Report: Carnegie Learning Curricula and Cognitive Tutor®*. Washington, DC: U.S. Department of Education, Institute of Education Sciences.

## Authors

JOHN F. PANE is a senior scientist at RAND. He uses experimental and rigorous quasi-experimental methods to study the implementation and effectiveness of innovations in education, particularly those involving technology.

BETH ANN GRIFFIN is a senior statistician at RAND. Her statistical research focuses on causal effects estimation when using observational data and innovative techniques for designing randomized studies. She works in a wide-range of applied areas, including education, public health, and criminal justice.

DANIEL F. MCCAFFREY is a principal research scientist at Educational Testing Service. He conducted this research while employed at RAND. His current research interests include value-added modeling and the measurement of teaching.

RITA KARAM is a policy researcher in RAND's Education unit. She has extensive experience in investigating educational reforms and programs, including math, science, and literacy programs. Her research focuses on measuring program implementation and examining implementation effects on school processes

and student learning. She holds a PhD in educational policy from the University of California, Riverside.

Manuscript received March 15, 2013

Revision received July 12, 2013

Accepted September 6, 2013