

The Relative Effectiveness of Teachers and Learning Software: Evidence from a Field Experiment in El Salvador*

Konstantin Büchel^a, Martina Jakob^b,
Christoph Kühnhanss^b, Daniel Steffen^a, Aymo Brunetti^a

^aDepartment of Economics, University of Bern

^bInstitute of Sociology, University of Bern

November 3, 2020

Abstract

This study provides novel evidence on the relative effectiveness of computer-assisted learning (CAL) software and traditional teaching. Based on a randomized controlled trial in Salvadoran primary schools, we evaluate three interventions that aim to improve learning outcomes in mathematics: (i) teacher-led classes, (ii) CAL classes monitored by a technical supervisor, and (iii) CAL classes instructed by a teacher. We find that CAL lessons lead to larger improvements in students' mathematics skills than traditional teacher-centered classes. In addition, teachers add little to the effectiveness of learning software. Our results highlight the value of CAL approaches in an environment with poorly qualified teachers.

JEL classification: C93, I21, J24, O15.

Keywords: computer-assisted learning, productivity in education, teacher content knowledge.

*We are grateful to David Burgherr, Malin Frey, and Amélie Speiser who provided excellent research assistance and to Philippe Sasdi for coordinating data collection in Swiss primary schools. The project further benefited from invaluable feedback by Michael Gerfin, Ben Jann, Florian Keller, Ulf Liebe, Blaise Melly, Urs Moser, Adina Rom, Mauricio Romero, Erik Snowberg, and the participants at the NADEL Workshop (ETHZ), Brown Bag Seminar (Department of Economics, U. Bern), CRED Seminar (U. Bern), SEVAL Meeting (hosted by the SDC in Bern), Rational Choice Sociology Conference in Venice, the SSES Annual Congress in Geneva, the EEA Annual Meeting 2020, and the Annual DENS Meeting in St. Gallen. This study would not have been possible without the *Impact Award 2017* prize money awarded by the SDC and NADEL (ETHZ), and funding by the chair of Ben Jann (Institute of Sociology, University of Bern) as well as the IMG Stiftung. Martina Jakob and Christoph Kühnhanss disclose that they serve on a voluntary basis as president and vice-president of *Consciente – Unterstützungsverein für El Salvador (Schweiz)*. We received IRB approval from the Faculty of Business, Economics and Social Sciences at the University of Bern. A randomized controlled trials registry entry is available at: <https://www.socialscisceregistry.org/trials/2789>.

Contact Details (°Corresponding Author):

°Büchel: Univ. of Bern, Dept. of Economics, Schanzenneckstr. 1, CH-3001 Bern, konstantin.buechel@vwi.unibe.ch

Jakob: Univ. of Bern, Inst. of Sociology, Fabrikstr. 8, CH-3012 Bern, martina.jakob@soz.unibe.ch

Kühnhanss: Univ. of Bern, Inst. of Sociology, Fabrikstr. 8, CH-3012 Bern, christoph.kuehnhanss@soz.unibe.ch

Steffen: Univ. of Bern, Dept. of Economics, Schanzenneckstr. 1, CH-3001 Bern, daniel.steffen@vwi.unibe.ch

Brunetti: Univ. of Bern, Dept. of Economics, Schanzenneckstr. 1, CH-3001 Bern, aymo.brunetti@vwi.unibe.ch

1 Introduction

While net primary school enrollment rates in low-income countries climbed from 56% in 2000 to 81% in 2019, learning outcomes have failed to keep pace. Less than 15% of primary school children in low-income countries pass minimum proficiency thresholds in reading and math, compared to about 95% of pupils in high-income countries (World Bank, 2018, p. 8). Public schooling systems in developing countries face multiple challenges that curb their productivity. These include a mismatch between national curricula and student abilities (Pritchett and Beatty, 2015), large and heterogeneous classes (Mbiti, 2016; Glewwe and Muralidharan, 2016), and low levels of effort among poorly qualified teachers (Chaudhury et al., 2006; Bold et al., 2017a). A much-noticed approach to overcome these barriers is the use computer-assisted learning software (e.g. The Economist, 2017). Computer-assisted learning (CAL) has several advantages over traditional teaching methods, as it allows for self-paced learning that is tailored to the abilities of the student, provides instant feedback and is less sensitive to the motivation and skills of teachers. Previous studies on the impact of technology-based teaching methods on learning outcomes are encouraging. CAL interventions are usually found to improve students' test scores and seem to be particularly beneficial if the software is used to personalize instructions (for a review see Escueta et al., 2020).¹

Yet, most studies evaluate CAL lessons that were offered as a supplement to regular classes, meaning that beneficiaries experienced a considerable expansion of school time compared to the untreated students in the control group. Thus, it remains unclear whether learning gains are

¹Experimental studies on CAL interventions in low- and middle-income countries include Banerjee et al. (2007, math in Indian primary schools), Carrillo, Onofa and Ponce (2011, language and math in Ecuadorian primary schools), Yang et al. (2013, language and math in Chinese primary schools), Mo et al. (2015, math in Chinese primary schools), Lai et al. (2015, language and math in Chinese primary schools), and Muralidharan, Singh and Ganimian (2019, language and math with Indian secondary school pupils). They consistently report positive intent-to-treat estimates on learning outcomes that range between 0.1σ and 0.4σ .

actually attributable to the use of the software or if additional lessons conducted by a teacher might have produced similar or even better results.² In addition, there is little evidence on whether CAL is a substitute for teachers or if it is a complement to them. Finally, previous research has mostly evaluated specifically customized software which is only available in a limited number of languages. As a result, many policy-makers with an interest in implementing CAL cannot draw on software that is readily available and has been successfully evaluated.

Based on a randomized controlled trial, this paper examines the relative effectiveness of primary school math teachers and a freely available CAL software that features content in more than 30 languages. To disentangle the effects of additional teaching and the use of a learning software, the experimental design features three different treatments: The first treatment comprises additional math lessons instructed by a teacher (henceforth labeled as TEACHER). The second and third treatments are additional math lessons based on CAL software; one group of classes is monitored by technical supervisors (CAL + SUPERVISOR), while the other group is taught by teachers (CAL + TEACHER). Teachers had to be officially certified to teach math in primary schools, whereas supervisors were laymen instructed to provide no content-related help to students. CAL lessons were taught using an offline application of the “Khan Academy” platform, and the three treatment arms were implemented by the Swiss-Salvadoran NGO *Consciente*.

We conducted the experiment between February and October 2018 with a sample of 198 primary school classes spanning grades 3 to 6 in the rural district of Morazán, El Salvador. 29 out of 57

²To our knowledge, the only study that evaluates the effectiveness of CAL lessons as a substitute to regular teaching in the development context was conducted by Linden (2008) in India. While attending *additional* CAL lessons raised math scores of second and third graders, CAL had a negative impact when it *substituted* regular classes. As the author points out, the study sample only covers NGO-run schools with well trained staff and innovative teaching methods. While it is unclear whether these findings translate to the challenging contexts of public education systems in developing countries, they still raise doubts about the inherent benefits of technology-based instruction.

eligible schools were randomly selected for program participation. The 158 classes from these 29 schools were then randomly assigned to either Treatment 1 (i.e. TEACHER, 40 classes), Treatment 2 (i.e. CAL + SUPERVISOR, 39 classes), Treatment 3 (i.e. CAL + TEACHER, 39 classes) or a program school control group (40 classes). In the 28 non-program schools, a random sample of 40 classes was drawn resulting in a “pure” control group that is not subject to potential treatment externalities.

Our analysis establishes four key findings. *First*, the additional CAL classes had a considerable impact on students’ math skills. Being assigned to additional CAL lessons increased their math scores by 0.21σ (p-value<0.01) when overseen by a supervisor and by 0.24σ (p-value<0.01) when instructed by teachers. These intent-to-treat estimates, which reflect a program attendance rate of 59%, correspond to the average increase in math abilities over 0.6 school years. Using the treatment assignment as instrumental variable for attendance, we estimate that participating in all 46 additional CAL lessons (each lasting 90 minutes) translates to average learning gains of 0.38σ (p-value<0.01) and 0.40σ (p-value<0.01), respectively. This is equivalent to the average increase in math abilities during 1.1 school years.

Second, additional CAL lessons were more productive than the additional math lessons instructed by a teacher. The intent-to-treat effect of being assigned to additional teacher-led classes without CAL was 0.15σ (p-value=0.01). Hence, students assigned to CAL + TEACHER outperformed students assigned to TEACHER by 0.09σ (p-value=0.10). The CAL treatment overseen by technical supervisors (CAL + SUPERVISOR) was also more successful in raising student learning than traditional teaching, even though this difference is not statistically significant (p-value=0.24). The advantage of CAL lessons relative to teacher-centered lessons was most pronounced in the domain of number sense and elementary arithmetic, and less so with respect to geometry, measurement and data. Focusing on number sense and elementary arithmetic, the difference between the CAL and non-CAL treatments increases to 0.11σ (p-value=0.06) for CAL instructed by teachers and to 0.09σ

(p-value=0.12) for the CAL monitored by supervisors.

Third, we present multifaceted evidence that points to a low productivity of teachers. The difference in learning gains between program school control classes and those classes receiving additional teacher-centered math lessons was close to zero and statistically insignificant (p-value=0.78). Similarly, teachers did not provide much “value added” to the learning software: the estimated impact for CAL lessons instructed by teachers is slightly higher than for CAL lessons conducted by supervisors but the difference is negligible and statistically insignificant (p-value=0.65). Moreover, the productivity of teachers dropped as the complexity of concepts increased: The impact of additional math lessons instructed by teachers decreased with both the grade level and the baseline achievement of their students, while the effect of the CAL-based lessons was largely insensitive to students’ grades and initial ability levels.

To gain a better understanding of the mechanisms behind these patterns, we conducted a comprehensive teacher math assessment covering the primary school curriculum of El Salvador. This assessment documents inadequate content knowledge among the teachers hired by the NGO. The math content knowledge of the contract teachers is positively correlated with student learning gains in both traditional (coef.=0.08, p-value=0.28) and CAL-based math lessons (coef.=0.09, p-value=0.14), whereas the math score of technical supervisors is, as one would expect, virtually orthogonal to students’ learning gains (coef.<0.01, p-value=0.94). The obtained point estimates for teachers almost perfectly correspond to international evidence on the impact of teacher content knowledge on learning outcomes (for a review see Table 5 in Brunetti et al., 2020) and hence corroborate the hypothesis that the inadequate subject mastery of the contract teachers impaired the impact of additional math lessons.

Importantly, regular math teachers in local primary schools are even less proficient in their subject than the contract teachers hired by the NGO. While the median contract teacher correctly

answered 66% of 50 questions on an assessment covering the curriculum of grade levels two to six, the median regular teacher correctly solved only 44% of the same items. Potential productivity gains resulting from an introduction of CAL to regular classes thus are likely to be larger than suggested by our estimates.

Fourth, we document substantial treatment externalities. At endline, students in program school control classes outperformed students in pure control classes by 0.14σ (p-value=0.02), although they were only indirectly exposed to the three treatments. In particular, we find evidence for spillovers from the two CAL treatments. While we cannot comprehensively pin down the mechanisms at work, suggestive evidence points toward social learning. At the same time, the data rejects hypotheses operating via direct exposure of students in control classes to the treatments (i.e. non-compliance) or testable hypotheses on behavioral adjustments in response to the experimental design.

This study makes several contributions to the literature on educational interventions in developing countries. *First*, it improves our understanding of how CAL performs relative to alternative teaching models. To our knowledge, this is the first well-identified study assessing the value-added of CAL in a public school setting of a developing country. As opposed to Linden (2008), who documents a negative value-added of CAL in NGO-administered schools in India, our findings suggest that CAL has the potential to outperform traditional teacher-led instruction, especially if teachers are poorly qualified. According to our estimates, it would take a teacher at the 91th percentile of the local teacher ability distribution to achieve the same learning gains as observed in CAL lessons overseen by a supervisor. This corresponds to a teacher in the 75th ability percentile among the contract teachers hired for this experiment or to 88% correct answers in the administered teacher assessment. While CAL has been praised in terms of its individualized and interactive pedagogy (e.g. Banerjee et al., 2007; Muralidharan, Singh and Ganimian, 2019), these numbers highlight that it may also be a promising approach to mitigate the adverse effects of teachers' inadequate content

knowledge and pedagogical knowledge, as it has been recently documented for several developing countries (e.g. Bold et al., 2017a).

Second, we present the first experimental test of the complementarities between teachers and learning software. In our setting, teachers play a marginal role in the success of technology-based instruction, with CAL lessons being almost equally effective when conducted by a supervisor rather than an officially certified teacher. Thus, teachers and learning software are likely substitutes and not complements, at least in a public schooling system staffed by poorly qualified teachers. Only few experimental studies aspire to distinguish between complementary and substitutable inputs entering the educational production function; notable exceptions are recent papers by Mbiti et al. (2019) on the complementarity between school grants and teacher incentives in Tanzanian primary schools, and by Attanasio et al. (2014) on the complementarity between psychosocial stimulation programs and nutritional supplements in early childhood development.

Third, we contribute to the broader literature on treatment externalities (e.g. Miguel and Kremer, 2004; Baird et al., 2015). By including control classes from treatment schools as well as spatially separated pure control classes from non-treatment schools into our experimental design, this study provides a reasonable identification of potential externalities. Our findings underscore the importance of appreciating the possibility of externalities in the design of experimental evaluation studies, even when such effects appear unlikely at first sight. Moreover, the presence of positive treatment externalities provide a strong rationale in favor of scaling the evaluated program.

Finally, this study adds to the accumulated evidence on the effectiveness of CAL by evaluating a widely available off-the-shelf software. In contrast to software tested in previous evaluations, Khan Academy is free of charge and features extensive math content in more than 30 languages.³

³The full version is available in 16 languages including Spanish, and a subset of content is available in about 20 languages. Another off-the-shelf learning software that has been successfully evaluated is Mindspark (see Muralidharan, Singh and Ganimian, 2019), which operates in English and Hindi for math and language training.

Since the employed software is arguably one of the most important features of a CAL intervention, our findings bear direct policy relevance for educational decision-makers around the globe that are looking for a readily available learning software suitable in non-English speaking countries.

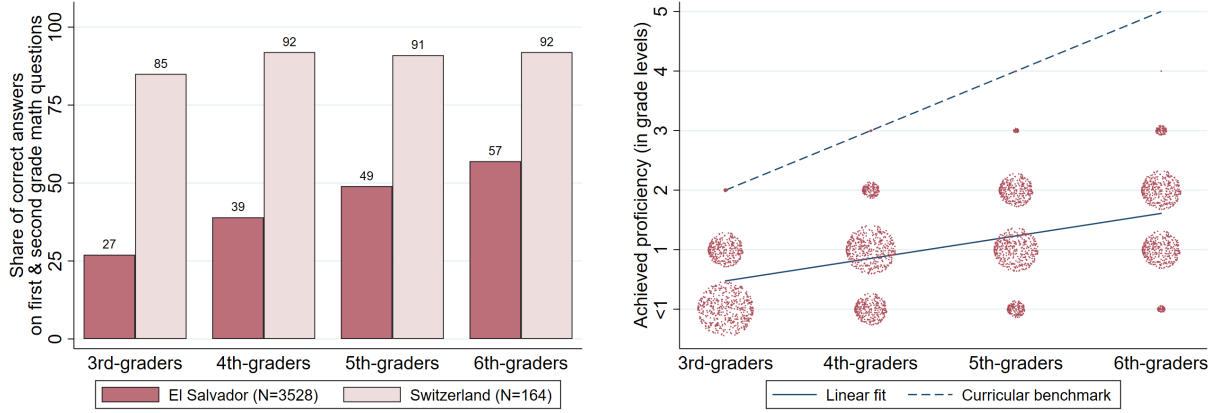
2 Context and Intervention

El Salvador is a lower middle-income country in Central America. The country’s net primary enrollment rates are estimated at 80%, which is 7 percentage points below the average of lower middle-income countries. While most children attend primary school, enrollment declines to 67% at the secondary level and to 28% in tertiary education.⁴

The department of Morazán is a poor and rural region in the northeast of the country with roughly 200,000 inhabitants. An average person in Morazán lives on 3.80 USD per day and, according to national definitions, almost 50% of the households face multifaceted poverty. While Morazán ranks second-last among all Salvadorian departments in terms of adult literacy, its secondary school students came forth in the 2018 “PAES” national examinations (DIGESTYC, 2018; MINED, 2018).

Our math assessments with 3,528 third to sixth graders conducted in February 2018 reveal that primary school children in Morazán barely grasp the most elementary concepts in math. Figure 1a shows that the share of correct answers to first and second grade questions increases from 27% among third graders to 57% among sixth graders, who by then should have attended more than 1,000 math lessons. To put these numbers into context, we conducted the same test with 164 pupils in Switzerland, who answered on average between 85% and 92% of the items correctly. Even the worst performing Swiss third grader outperformed the median sixth grader in Morazán. Similarly flat learning curves among primary school children have also been reported for other low-

⁴Enrollment statistics according to the *World Development Indicators* provided online by the World Bank, see <https://data.worldbank.org/indicator> (last access: 26.10.2019)



(a) Share of correct answers on 1st/2nd grade math questions among Salvadoran and Swiss pupils. (b) Assessed grade level in math among third to sixth graders in Morazán early in their school year.

Figure 1: Math learning outcomes in Morazán (Panels a & b) and Switzerland (Panel a).

Note: Panel (b) illustrates the achieved proficiency in math (measured in grade levels) among third to sixth graders in Morazán at the beginning of their school year. A student, each represented by a dot, needs to score at least 50% correct answers on grade specific items in order to reach the next proficiency level. Since the test was administered in the first weeks of their school year, a third grader answered first and second grade items and therefore may be assigned to grade level 2, 1 or <1 depending on her performance. The size of the bubbles are proportional to the number of students they represent. Further explanations are provided in appendix A.1. *Source:* Baseline data, February 2018.

and middle-income regions across the globe, including countries in Western Africa, Eastern Africa, Central America, and South East Asia (e.g. Beatty et al., 2018; PAL, 2020).

Several challenges that plague Morazán’s schooling system can account for its low productivity. For instance, our monitoring data from school visits reveal high rates of teacher absenteeism suggesting that, on average, 25% of regular lessons are canceled. Low teacher motivation mixes with outdated pedagogical techniques that essentially follow the logic of “copy, memorize, and reproduce”. And despite relatively small class sizes – the pupil-teacher ratio averages 28-to-1 at the national level and 19-to-1 in our sample – heterogeneous student performance and an overambitious curriculum make it difficult to teach at an appropriate level. As Figure 1b shows, third to sixth graders lag considerably behind the official curriculum and this gap widens as children move up to higher grade levels. Moreover, performance heterogeneity within classes is considerable. In the majority of classes, students’ math ability spans three grades or more (for further explanations see appendix A.1). In general, the public schooling system in El Salvador faces challenges similar to

those reported for other low- and middle income countries.⁵

The Salvadoran Ministry of Education has recently put considerable effort into addressing learning deficiencies in public schools. While primary schooling has been typically confined to either morning or afternoon lessons throughout El Salvador, a recent policy aims to extend school time over a full day and to complement traditional teaching with innovative learning approaches (MINED, 2013). The government hopes that longer schooldays will not only boost learning outcomes, but also shield children from the influence of criminal gangs. Within the scope of this countrywide program, the Ministry of Education seeks to cooperate with NGOs in order to collectively promote an open and flexible curriculum. While all schools received official instructions to expand their school days, most of them have not put the policy into practice due to a lack of resources to pay for further teaching staff.

Intervention. In this context, we evaluate the impact of an educational initiative on math abilities of primary school children of grades 3 to 6. The program features three intervention arms that offer two additional lessons of 90 minutes per week and almost double the beneficiaries' number of math lessons during the program phase. The first intervention arm comprises additional math lessons instructed by a teacher without using software. The second and third intervention arms are

⁵The pupil-teacher ratio in middle-income countries averages 24-to-1, while it climbs to 40-to-1 in low income countries (UNESCO, 2019); in some contexts, such as rural India, it can even reach 90-to-1 (Mbiti, 2016). Besides the large class size, students' abilities and preparation levels are often very heterogeneous, which is also the case in our data. For example, Muralidharan, Singh and Ganimian (2019) report for their sample of 116 Indian middle schools that students' ability in the median classroom spans four grades in both math and language, while we obtain three grade levels for primary schools. Moreover, Pritchett and Beatty (2015) show that the pace of learning is very slow in developing countries and that there is a mismatch between curriculum and student abilities. This is consistent with what we observe in Figure 1b. Finally, low teacher motivation is a well-known issue: Chaudhury et al. (2006) find that 19% of teachers in developing countries are absent during unannounced visits, while our monitoring data suggests that 25% of classes in Morazán's primary school are canceled.

additional math lessons based on computer-assisted learning software; one group of classes is taught by teachers, while the other group is monitored by supervisors.

The *CAL-lessons* in the second and third intervention arm were based on an offline application of the learning platform *Khan Academy*, which is known as *KA Lite*. This freely available software provides a wide range of instructional videos and exercises for every difficulty level. While the learning tool is not directly adaptive, it allows teachers to track the progress of each student and assign appropriate contents based on prior performance. To tailor instruction to students' learning levels, a set of working plans covering different content units was prepared. Based on a placement test, children received a plan that was viewed as adequate for their respective level and they could then proceed to subsequent plans at their own pace. Since one computer was available per student, each child could follow its individual learning path. Typically, students started with materials from lower grades and then slowly progressed towards contents corresponding to their actual grade level.

A similar methodology was used for the first intervention arm that features more traditional *math lessons instructed by a teacher*. According to their initial math skills, children were arranged in two different table groups where they worked on plans tailored to their ability. Teachers were instructed to explain important concepts, correct students' work at home and promote children (or entire table groups) to subsequent plans when appropriate. While this strategy only allows for a crude approximation of teaching to each child's ability level, it represents a degree of individualization that can realistically be achieved without the help of technology.

To pay credit to the *social component of learning*, all treatments combined individualized learning with educational games. For this purpose, a manual containing animation, concentration and math games was developed. The manual compiles simple techniques to promote students' collective learning as well as their motivation to participate in class. Games were usually played at the beginning or at the end of each session. While supervisors were instructed to use animation and

concentration games, teachers were additionally introduced to a series of math games.

The contracted *teachers* were required to be officially certified to instruct grades 3 to 6 in math. That is, they all possessed a university degree and had either completed a teacher education, or another study program combined with a one-year pedagogical course. Teachers were selected based on a brief math assessment and a job interview. They were employed on short-term contracts and earned 300 USD per month for assuming four classes.⁶ For lessons that were canceled, teachers received no remuneration. To optimize the comparability of treatments, all teachers were assigned an equal number of CAL and non-CAL lessons. Before and during the intervention, teachers were trained to operate the learning software and they reviewed mathematical concepts as well as central pedagogical strategies including the use of educational games. Teaching was tightly monitored by our partner NGO through monthly feedback meetings at the NGO’s headquarters and unannounced classroom visits during the implementation phase.

The *supervisors* received only technical training and were paid substantially less than teachers, that is 180 USD for taking care of four classes. They were required to have minimal IT skills and some experience in dealing with children, but no contracted supervisor possessed a degree in education or teaching credentials. During the intervention, supervisors were instructed to restrain from providing any content-specific help. Like teachers, supervisors were employed on short-term contracts and were paid conditional on the number of classes they conducted.

3 Research Design

This study is built around an RCT to identify the causal impact of the three interventions arms. It started in February 2018 with a baseline assessment and a survey covering all control and program classes. The additional math classes began in April 2018 and were implemented until the end of the

⁶This corresponds to 8×90 minutes of teaching per week, or – including preparatory work – to a 60% job.

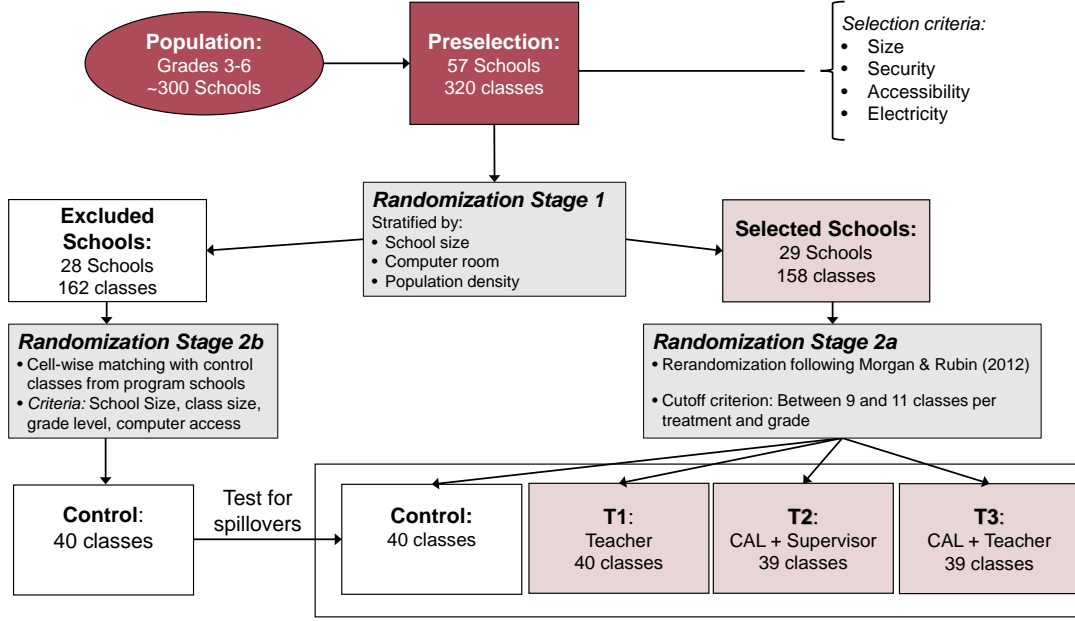


Figure 2: Sampling and randomization scheme.

school year in fall 2018; the school year in El Salvador starts in mid-January and ends in November. The endline tests took place in October 2018, six months after the start of the intervention. Again, all program and control classes took part in the endline tests.

3.1 Sampling and Randomization

Our sampling and randomization scheme has three layers, as exemplified in Figure 2. Starting point are all 302 primary schools in Morazán. In coordination with the NGO and the regional Ministry of Education, we defined the following eligibility criteria for a preselection of primary schools:

- **School size**, *eliminates 221 schools*: A school was considered too small, if it had integrated classes (across grades) or gaps in its grade structure (i.e. not at least one class per grade). This guarantees that every eligible school has at least four different classes in grades 3 to 6, and therefore can participate with at least (i) one CAL+TEACHER, (ii) one CAL+SUPERVISOR, (iii) one TEACHER, and (iv) one control class;

- ***Security***, *eliminates 14 of the remaining 81 schools*: Based on an assessment by the local staff and the regional Ministry of Education, schools located in areas dominated by criminal gangs were excluded due to security concerns;
- ***Accessibility***, *eliminates 7 of the remaining 67*: Schools where access by car is difficult were discarded. To inform this decision we relied on Google-Maps driving times and a validation by local staff and the regional Ministry of Education;
- ***Electricity***, *eliminates 3 of the remaining 60 schools*: Schools without a (close-by) power supply did not qualify for the program.

After this pre-selection, 57 schools with a total of 320 eligible classes and about 6,400 students remained in the sample. In *randomization stage 1*, 29 of the 57 schools were randomly chosen to participate in the program. To improve balance, the assignment was stratified by school size, local population density and students' access to a computer room.

In *randomization stage 2a*, we randomly assigned the 158 classes in the 29 selected program schools to the control group or one of the three intervention arms. Following Morgan and Rubin (2012) we re-run the randomization routine until the interventions were balanced across schools and grades. This mechanism assigned 39 classes to CAL+TEACHER, 39 classes to CAL+SUPERVISOR, 40 classes to TEACHER, and 40 classes to the control group. We account for the re-randomization procedure when comparing estimates within program schools by computing randomization inference test statistics based on 2,000 random draws subject to the identical cut-off criterion. Our choice to run 2,000 draws is guided by Young (2019, p. 572), who finds no appreciable change in rejection rates beyond this threshold. To implement the randomization tests we rely on Stata's RITEST-package developed by Hess (2017).

As prominently discussed in Miguel and Kremer (2004), interventions can have spillover effects

on non-participating students from the same school or area. A unique feature of our design allows us to estimate the size of such treatment externalities. For this purpose, in *randomization stage 2b*, 40 additional control classes from non-treatment schools were included in the study. These additional “pure” control classes are spatially separated from the intervention, and thus not affected by the NGO’s work. The *pure control classes* were randomly selected from the 28 control schools by matching them cell-wise to the distribution of control classes from program schools, accounting for school size, grade level, class size and students’ access to computers.

This procedure yields five different groups of primary school classes, namely the 39 or 40 classes assigned to each of the three treatment groups, 40 control classes from the 29 program schools, and 40 pure control classes from the 28 control schools.

3.2 Data

In the course of the evaluation, four types of data were gathered: *(i)* Math learning outcomes of students were assessed before and after the intervention, *(ii)* socio-demographic statistics stem from a survey that children answered prior to the baseline math assessment, *(iii)* administrative data on schools was collected between October 2017 and February 2018, and *(iv)* monitoring data was recorded during unannounced school visits throughout the program phase. Table 1 shows summary statistics for the main variables collected before the start of the program as well as absence rates at the endline and baseline assessment. In particular, it displays means and standard errors for the different variables by treatment status, and tests whether the mean is equal across the five groups.

While the treatment and control groups do not differ significantly on any observable dimension at baseline, Table 1 shows a sizeable increase in the absence rates between the baseline and the endline assessment. Before both rounds of data collection, we updated comprehensive class lists of registered pupils. This revealed that large numbers of children either migrated out of Morazán or

discontinued their education. We achieved an attendance of about 95% registered pupils in both rounds, but since classes shrank during the school year, the overall attrition at endline almost hits the 10% mark. Importantly, Table 1 does not point toward systematic differences in attrition rates by treatment status.⁷ Moreover, compliance with the experimental protocol was very good in the sense that only 38 out of 3197 students (i.e. 1.2%) within our estimation sample switched between different classes, grades or schools.

3.2.1 Math Learning Outcomes

The math assessments include 60 items covering the primary school curriculum of El Salvador. The weighting of questions across the three main topics (a) number sense & elementary arithmetic ($\sim 65\%$), (b) geometry & measurement ($\sim 30\%$), and (c) data & statistics ($\sim 5\%$) was closely aligned with the national curriculum. Moreover, we verified the appropriateness of each question through a careful mapping to the national curriculum and a feedback loop involving the regional Ministry of Education and local education experts. The math problems presented to the children mostly required a written answer (as opposed to a multiple choice format) and were inspired by El Salvador’s official textbooks as well as various international sources of student assessments. Section B in the appendix explains the design of our assessments step by step.

In the appendix, we further present detailed statistics on the distribution of student test scores and the difficulty of the items. Top or bottom coding is neither an issue with respect to students nor the selected items: Table B.2 shows that virtually all items (except one for fifth graders in the endline assessment) were at least once answered correctly or incorrectly. Likewise, Table B.1 documents that only about 0.5% of test-takers scored zero points, while nobody achieved the maximum score.

⁷We examine this more closely in Table A.1 in the appendix, confirming that the treatment status is not significantly correlated with presence at the endline test.

Table 1: Balance at baseline and absence rates during assessments

	(1)	(2)	(3)	(4)	(5)	(6)
Panel A:	T1: Math	T2: CAL	T3: CAL	Within school	Pure control	
Math scores (N=3528)	w. teacher	w. supervisor	w. teacher	controls	classes	p-value
%-share correct answers	30.33 (1.80)	33.47 (1.90)	31.97 (2.07)	32.60 (1.32)	30.80 (2.00)	0.45
Std. IRT math score	0.01 (0.14)	0.18 (0.14)	0.08 (0.16)	0.08 (0.10)	0.00 (0.15)	0.72
Panel B: Sociodemographics (N=3528)						
Female student	0.50 (0.03)	0.52 (0.04)	0.55 (0.04)	0.51 (0.03)	0.49 (0.04)	0.43
Student age	-0.09 (0.08)	-0.01 (0.09)	0.02 (0.09)	-0.03 (0.06)	-0.03 (0.09)	0.70
Household size	5.56 (0.13)	5.61 (0.12)	5.57 (0.12)	5.55 (0.08)	5.50 (0.12)	0.92
Household assets index	0.55 (0.02)	0.55 (0.02)	0.54 (0.02)	0.56 (0.02)	0.56 (0.02)	0.88
Panel C: Class room variables and absence rates during assessments (N=198)						
Class size	18.40 (1.37)	19.33 (1.35)	18.69 (1.37)	18.13 (0.96)	18.32 (1.54)	0.92
Female teacher	0.80 (0.10)	0.77 (0.10)	0.77 (0.10)	0.73 (0.07)	0.55 (0.11)	0.14
Absence rate at baseline (%)	3.88 (1.33)	3.15 (1.16)	5.39 (1.74)	4.39 (0.95)	3.38 (1.15)	0.59
Absence rate at endline (%)	9.09 (2.09)	9.72 (2.04)	10.50 (2.18)	9.99 (1.63)	8.10 (2.00)	0.72
Panel D: School variables (N=49)				Treatment schools	Pure control schools	p-value
# classes grade 3–6				5.48 (0.43)	6.25 (0.76)	0.32
Computer lab				0.79 (0.08)	0.75 (0.13)	0.73
Local population density				0.18 (0.01)	0.19 (0.02)	0.63

Notes: This table presents the mean and standard error of the mean (in parenthesis) for several characteristics of students (Panels A & B), class rooms (Panel C), and schools (Panel D), across treatment groups. The student sample consists of all students tested by the research team during the baseline survey in February 2018. Column 6 shows the p-value from testing whether the mean is equal across all treatment groups. IRT-scores are standardized such that $\mu = 0$ and $\sigma = 1$ for the pure control group. The household asset index measures what share of the following assets a household owns: Books, electricity, television, washmachine, computer, internet and car. Local population density is the municipality's population density measured in 1000 inhabitants per km². Standard errors are clustered at the class level in Panels A & B, and at the school level in Panel C. * p<0.10, ** p<0.05, *** p<0.01.

In general, the assessments seem to nicely capture the different performance levels, with the scores being roughly normally distributed around a median of 30% (3rd graders) to 40% (6th graders) correct answers (see Figure B.2).

A particularly nice feature of our math assessments is that they allow us to project all outcomes on a common ability scale by drawing on techniques from psychology, that is *Item Response Theory (IRT)* (e.g. de Ayala, 2009). This implies that we can directly compare children across grades and express their learning gains between the baseline and the endline assessment in terms of how many additional school years would be required to reproduce the same effect. The conversion of our estimates into program effects measured in terms of additional school years is explained in the appendix B.

3.2.2 Socio-Demographic Survey

The socio-demographic survey was distributed 15 minutes before the baseline math assessment began. It asked students about their age, gender, household composition, household assets and parental education. Since literacy can be an issue, questions were illustrated with pictures and the enumerators helped children to understand and answer them correctly.

3.2.3 Administrative Data on Schools

In the run-up to the study, we collected various administrative data on Morazán’s school. While the government gathers thematically vast information on the school environment through a paper-and-pencil survey administered to school principals, the data turned out to be of rather poor quality. To obtain usable information on the class structure, enumerators had to call each school at the beginning of the school year, because the planning data from official sources was too unreliable. Moreover, the paper-and-pencil surveys left many missing values, so that we had to discard most items due to an insufficient coverage. We therefore decided to use a minimal set of school level

variables, which were either comprehensively available, relatively cheap to supplement, or essential for the study. These include the number of grade 3 to grade 6 classes (i.e. school size), information on the presence of gangs (i.e. security at school), accessibility measures based on Google-Map estimates and validated by local staff, power supply according to the administrative survey and validated via phone calls, student access to computer labs according to the administrative survey and validated via phone calls, and local population density from the National Bureau of Statistics.

3.2.4 Monitoring Data

From May to September 2018, NGO staff made on average five unannounced school visits (about 1000 visits in total) to collect monitoring data. These visits covered regular lessons in *pure control schools* as well as both regular and treatment lessons in *program schools*. The enumerators collected data on teacher attendance, student attendance, computer usage, and the implementation of the additional math lessons in the afternoon.

4 Results

4.1 The Overall Program Effects

We begin by estimating *intent to treat* (ITT) effects of being assigned to one of the three programs (i.e. $\beta_{T1}, \beta_{T2}, \beta_{T3}$) or the program school control classes (i.e. β_{CX}) using

$$Y_{icsk}^{EL} = \alpha + \beta_{T1}T1_{csk} + \beta_{T2}T2_{csk} + \beta_{T3}T3_{csk} + \beta_{CX}CX_{csk} + \delta Y_{icsk}^{BL} + X'_{icsk}\gamma + V'_{csk}\lambda + \phi_k + \epsilon_{1icsk}, \quad (1)$$

where Y_{icsk}^{EL} is the endline math score of student i in class c , school s and stratum k ; math scores are either measured as percentage of correct answers or as the IRT-score normalized to $\mu=0$ and $\sigma=1$ based on the baseline score of the pure control group. The binary treatment indicators are

defined as follows: $T1$ equals one for those assigned to extra math lessons conducted by a teacher, $T2$ equals one for those assigned to extra CAL lessons overseen by a supervisor, $T3$ equals one for those assigned to extra CAL lessons instructed by a teacher, and CX equals one for those assigned to program school control classes that are potentially subject to externalities. Our control variables include Y_{icsk}^{BL} which stands for the baseline math score, X_{icsk} representing a set of student-level control variables (i.e. age normalized by the average age at the same grade level, gender, household size and household assets), and V_{csk} comprising a set of classroom-level variables (i.e. indicator for grade level, class size and teacher gender). Finally, ϕ_k stands for k strata fixed effects and ϵ_{1icsk} represents the error term.

The upper panel of Table 2 displays the program effect relative to pure control classes (i.e. $\hat{\beta}_{T1}$, $\hat{\beta}_{T2}$, $\hat{\beta}_{T3}$ and $\hat{\beta}_{CX}$) and the lower panel of Table 2 presents estimates for the pairwise differences between the three treatment groups in program schools. The lower panel reports p-values obtained from a randomization inference test statistic based on 2,000 random draws subject to the identical cut-off criterion as used in our re-randomization scheme (see section 3). In the upper panel, however, p-values are based on traditional clustered standard errors, since the assignment to program schools and pure control schools did not involve re-randomization.⁸

Students who were assigned to one of the treatments perform significantly better in the end-line assessment than students assigned to the pure control classes. Compared to the pure control

⁸Moreover, we cannot properly apply randomization inference to the upper panel due to missing information on ability levels of non-selected classes from pure control schools. As we show in appendix A.3, randomization inference in the upper panel is based on draws that include on average 37% missing data points. Consequently, p-values obtained from these randomization tests increase by a factor of about 5 to 10 compared to p-values from traditional inference with clustered standard errors. While this is clearly too conservative, our main conclusion are not altered when we apply randomization inference to the upper panel (see Table A.3). The only notable difference is that program externalities, captured by β_{CX} , turn insignificant with p-values around 0.13. When we apply traditional inference to the lower panel, as shown in Table A.2, changes in p-values are very small and do not show a clear pattern.

Table 2: ITT-Estimates on the effects of the different interventions on children's math scores

	Percent Correct		IRT-Scores	
	(1)	(2)	(5)	(6)
T1: Lessons with Teacher	2.904*** (0.005)	2.643** (0.012)	0.165*** (0.006)	0.152** (0.013)
T2: CAL-Lessons with Supervisor	4.095*** (0.000)	3.869*** (0.000)	0.226*** (0.000)	0.214*** (0.000)
T3: CAL-Lessons with Teacher	4.554*** (0.000)	4.328*** (0.000)	0.250*** (0.000)	0.238*** (0.000)
CX: Control Classes for Externalities	2.595** (0.011)	2.407** (0.017)	0.147** (0.013)	0.137** (0.020)
$\beta_{T4} := \beta_{T2} - \beta_{T1} = 0$	1.191	1.226	0.061	0.063
p-value ($\beta_{T4}=0$)	(0.214)	(0.194)	(0.268)	(0.244)
$\beta_{T5} := \beta_{T3} - \beta_{T1} = 0$	1.650*	1.686*	0.084	0.086
p-value ($\beta_{T5}=0$)	(0.069)	(0.059)	(0.117)	(0.102)
$\beta_{T6} := \beta_{T3} - \beta_{T2} = 0$	0.459	0.460	0.024	0.023
p-value ($\beta_{T6}=0$)	(0.618)	(0.615)	(0.650)	(0.653)
Adjusted R ²	0.66	0.67	0.69	0.70
Observations	3197	3197	3197	3197
Individual & Classroom Controls	No	Yes	No	Yes
Baseline Score	Yes	Yes	Yes	Yes
Stratum & Grade FE	Yes	Yes	Yes	Yes

Notes: In the upper panel (coef. $\beta_{T1} - \beta_{CX}$), p-values are based on traditional clustered standard errors. In the lower panel (coef. $\beta_{T4} - \beta_{T6}$), p-values are based on a two-sided randomization inference test statistic that the placebo coefficients are larger than the actual; randomization inference is based on 2000 random draws.

* p<0.10, ** p<0.05, *** p<0.01.

students, participants assigned to extra classes with math teachers (i.e. $T1$) perform 2.6 percentage points or 0.15σ better, students assigned to CAL classes with supervisors (i.e. $T2$) perform about 3.9 percentage points or 0.22σ better, and students assigned to CAL classes with a teacher (i.e. $T3$) perform 4.3 percentage points or 0.24σ better. Remarkably, students in control classes within program schools (i.e. CX) also score 2.4 percentage points or 0.14σ higher than students in pure control classes. As we discuss in section 5.1, our analysis points towards spillovers from CAL-lessons to program school control classes, while we find no evidence for direct exposure of control units (i.e. non-compliance) or behavioral changes at the level of students, regular teachers

or school administrations.

Finally, we test whether the observed gaps in the endline performance of students assigned to one of the three treatments (defined as β_{T4} , β_{T5} , and β_{T6}) are statistically different from zero. While we find that the two CAL treatments outperform additional math classes, only the difference between additional math classes and CAL classes conducted by a teacher is statistically significant at the 10%-level: students assigned to CAL+TEACHER outperform students assigned to TEACHER by 1.7 percentage points or 0.085σ with p-values ranging from 0.059 to 0.117.

On the one hand, this is novel evidence that CAL delivers sizable learning gains in a Latin American context using off-the-shelf learning software: Expressing the estimates in terms of school years suggests that the effect of the CAL interventions is equivalent to the average student’s progress in 0.6 to 0.7 school years (see appendix B for details on this conversion). On the other hand, traditional math classes conducted by teachers are relatively ineffective compared to additional math lessons with CAL-software: In comparison to the program school control classes, boosting the supply of conventional math lessons by roughly 80% delivered no measurable impact. Importantly, the performance difference between CAL classes taught by teachers and additional teacher-centered math classes is statistically (marginally) significant. We interpret this as suggestive evidence that the learning gains reported in a series of CAL-evaluations can – at least partially – be attributed to the learning software and not necessarily to the increase in the number of math lessons.

4.2 Heterogeneity Analysis

We now examine effect heterogeneity along several dimensions. We first decompose program effects by subtopics, before we explore effect heterogeneity along baseline ability, grade level and class size.

4.2.1 Program Effects by Subtopic

In this subsection, we explore the impact of the three interventions on learning outcomes by topics. In accordance with the official curriculum, 65% of the items cover number sense and arithmetic (NSEA), 30% of the items cover geometry and measurement (GEOM), and 5% of the items cover data, probability and statistics (DSP). In particular, we re-estimate equation (1) but calculate separate math scores based on (i) NSEA-questions and (ii) GEOM- as well as DSP-questions.

The ITT-effects on students' NSEA skills are shown in Table 3. We find that both CAL treatments had a more pronounced effect on the NSEA score than on the overall math ability. Students who were assigned to CAL classes with supervisors score 4.6 percentage points or 0.24σ higher in NSEA questions than students assigned to pure control classes; this is an increase of 10% to 20% compared to the overall impact reported in Table 2. The NSEA math score of students assigned to CAL classes with teachers is 4.9 percentage points or 0.26σ higher than the score of students assigned to pure control classes; again this effect is 10% to 15% larger compared to estimates based on all questions. Since the impact on the NSEA math score remains about the same for students receiving additional math classes instructed by teachers, the gap between CAL and conventional teaching widens.

When we compare the learning gains attributed to CAL with the gains attributed to the additional math classes without software the differences range from 1.7 to 2.1 percentage points or from 0.092σ to 0.115σ . The corresponding p-values lie between 0.046 and 0.055 for the CAL classes with teachers and between 0.093 and 0.129 for CAL classes with supervisors. Hence, when focusing on NSEA questions, the overall pattern remains qualitatively similar to the estimations including all subject domains, but the gap between the two CAL treatments and additional math classes in the traditional sense (i.e. without the use of software) becomes more pronounced.

Table 4 shows the results that are based on GEOM- and DSP-items. Focusing on these topics

Table 3: ITT-Estimates on the effects of the interventions on children's *NSEA*-scores

	Percent Correct		IRT-Scores	
	(1)	(2)	(5)	(6)
T1: Lessons with Teacher	3.174*** (0.002)	2.849*** (0.006)	0.166*** (0.006)	0.146** (0.013)
T2: CAL-Lessons with Supervisor	4.907*** (0.000)	4.581*** (0.000)	0.258*** (0.000)	0.238*** (0.000)
T3: CAL-Lessons with Teacher	5.225*** (0.000)	4.895*** (0.000)	0.279*** (0.000)	0.259*** (0.000)
CX: Control Classes for Externalities	2.711*** (0.008)	2.463** (0.012)	0.145** (0.013)	0.130** (0.020)
$\beta_{T4} := \beta_{T2} - \beta_{T1} = 0$	1.733	1.732*	0.092	0.091
p-value ($\beta_{T4}=0$)	(0.103)	(0.093)	(0.129)	(0.115)
$\beta_{T5} := \beta_{T3} - \beta_{T1} = 0$	2.051**	2.047**	0.113*	0.112*
p-value ($\beta_{T5}=0$)	(0.046)	(0.047)	(0.051)	(0.055)
$\beta_{T6} := \beta_{T3} - \beta_{T2} = 0$	0.318	0.315	0.021	0.021
p-value ($\beta_6=0$)	(0.750)	(0.752)	(0.706)	(0.714)
Adjusted R ²	0.62	0.63	0.65	0.65
Observations	3197	3197	3197	3197
Individual & Classroom Controls	No	Yes	No	Yes
Baseline Score	Yes	Yes	Yes	Yes
Stratum & Grade FE	Yes	Yes	Yes	Yes

Notes: In the upper panel (coef. $\beta_{T1} - \beta_{CX}$), p-values are based on traditional clustered standard errors. In the lower panel (coef. $\beta_{T4} - \beta_{T6}$), p-values are based on a two-sided randomization inference test statistic that the placebo coefficients are larger than the actual; randomization inference is based on 2000 random draws.

* p<0.10, ** p<0.05, *** p<0.01.

reduces the impact of both CAL treatments. The effects compared to pure control classes remain significant but they decrease considerably in magnitude. The results show, for instance, that additional CAL lessons conducted by a teacher increase the NSEA-score by about 5 percentage points, while the increase in the combined GEOM- and DSP-score is only 3.5 percentage points. Since this drop is less pronounced for those classes receiving additional math lessons instructed by a teacher, the within treatment school comparisons yield insignificant effects.

These results show that computer-assisted learning software can be a valuable substitute to traditional teaching, but its impact seems to be sensitive to the concepts that are taught. While we

Table 4: ITT-Estimates on the effects of the interventions on children’s *GEOM* & *DSP*-scores

	Percent Correct		IRT-Scores	
	(1)	(2)	(5)	(6)
T1: Lessons with Teacher	2.433*	2.132*	0.155**	0.140*
	(0.055)	(0.093)	(0.035)	(0.057)
T2: CAL-Lessons with Supervisor	3.207***	3.014**	0.196***	0.187***
	(0.009)	(0.014)	(0.006)	(0.009)
T3: CAL-Lessons with Teacher	3.646***	3.472***	0.201***	0.193**
	(0.006)	(0.008)	(0.008)	(0.010)
CX: Control Classes for Externalities	2.773**	2.561**	0.159**	0.149**
	(0.032)	(0.048)	(0.036)	(0.050)
$\beta_{T4} := \beta_{T2} - \beta_{T1} = 0$	0.775	0.882	0.041	0.047
p-value ($\beta_{T4}=0$)	(0.498)	(0.432)	(0.543)	(0.464)
$\beta_{T5} := \beta_{T3} - \beta_{T1} = 0$	1.213	1.340	0.046	0.053
p-value ($\beta_{T5}=0$)	(0.279)	(0.221)	(0.481)	(0.412)
$\beta_{T6} := \beta_{T3} - \beta_{T2} = 0$	0.438	0.458	0.005	0.006
p-value ($\beta_6=0$)	(0.692)	(0.669)	(0.934)	(0.926)
Adjusted R ²	0.46	0.47	0.49	0.50
Observations	3197	3197	3197	3197
Individual & Classroom Controls	No	Yes	No	Yes
Baseline Score	Yes	Yes	Yes	Yes
Stratum & Grade FE	Yes	Yes	Yes	Yes

Notes: In the upper panel (coef. $\beta_{T1} - \beta_{CX}$), p-values are based on traditional clustered standard errors. In the lower panel (coef. $\beta_{T4} - \beta_{T6}$), p-values are based on a two-sided randomization inference test statistic that the placebo coefficients are larger than the actual; randomization inference is based on 2000 random draws.

* p<0.10, ** p<0.05, *** p<0.01.

obtain a consistently positive value-added of CAL + TEACHER and CAL + SUPERVISOR relative to TEACHER, the measured differences seem primarily driven by the pronounced improvements in the domains of number sense and elementary arithmetic. The CAL interventions were less successful in shifting abilities to solve questions on geometry, measurement, data and statistics: The difference in point estimates (see $\hat{\beta}_{T4}$ and $\hat{\beta}_{T5}$) decrease by about 20%, and the p-values clearly exceed the 0.1-threshold for statistical significance.

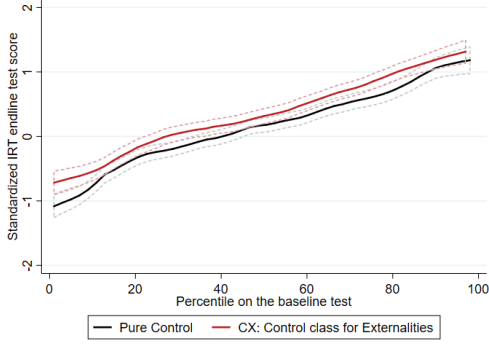
Overall, this sub-analysis also points toward a strikingly low productivity of teachers: Regardless of the domain, classes receiving additional math lessons conducted by teachers do not perform better

than control classes subject to externalities. However, the pronounced differences across domains also suggest that CAL may not be well-equipped to substitute for all aspects of the complex task a teacher is expected to perform. While it may be relatively easy to automate the correction of errors in simple arithmetic exercises, evaluating students’ progress and providing helpful feedback on tasks that require creativity or connected thinking may be much harder for a computer. Moreover, CAL may face a difficult job in connecting instructed concepts to real world experiences. While a teacher can, for example, distribute rulers to make students measure different objects in the classroom, pure CAL instruction is limited to what can be achieved with a two-dimensional screen. This suggests that a blended learning approach, where CAL is combined with active teachers who focus their engagement on tangible tasks may be a promising way to go.

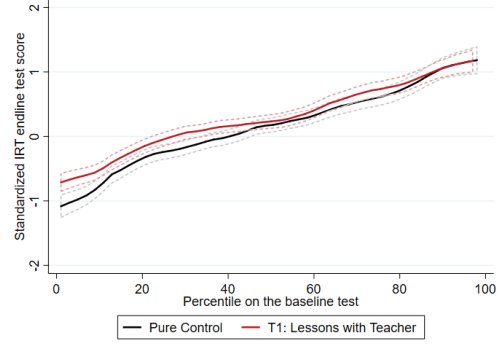
4.2.2 Effect Heterogeneity by Baseline Ability, Grade Level and Class Size

We continue the heterogeneity analysis by discussing Figure 3, which plots kernel-weighted locally-smoothed means of the endline test score at each percentile of the baseline test score by treatment status. Figure 3a shows that endline tests scores in the control group for spillovers are slightly higher than those in the pure control group at all percentiles of the baseline test score, but the 95% confidence bands mostly overlap. Comparing pure control classes to the TEACHER classes in Figure 3b shows that the latter outperform the former at low percentiles of the baseline score, while there is no difference at higher percentiles. Both CAL intervention groups, as illustrated in Figures 3c and 3d, achieve considerably higher endline scores than pure control classes across all percentiles in the baseline achievement, although the gap seems to narrow at higher percentiles in the CAL + TEACHER group.

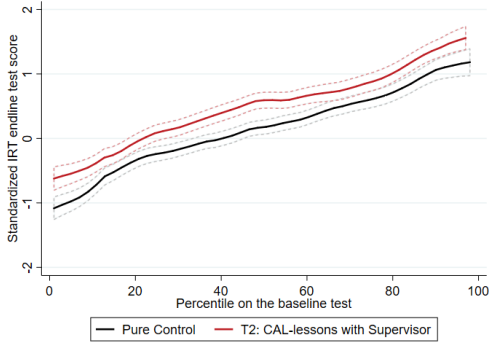
In a next step, we examine the functional relation between treatment effects and baseline achievement more thoroughly. We further investigate whether the reported effects vary by grade level or



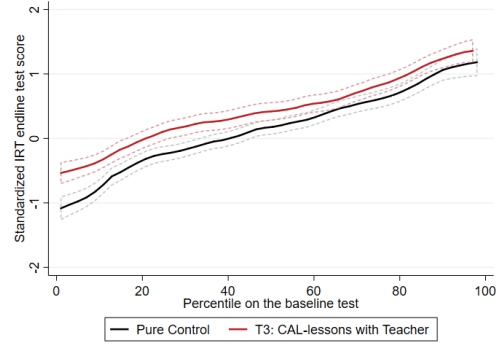
(a) Spillover vs. pure control classes



(b) TEACHER vs. pure control classes



(c) CAL + SUPERVISOR vs. pure control classes



(d) CAL + TEACHER vs. pure control classes

Figure 3: Endline test scores by treatment status and baseline percentiles.

Note: The figures present kernel-weighted local mean smoothed plots relating endline test scores to percentiles in the baseline achievement by treatment status alongside 95% confidence bands.

by class size. To do so, we estimate

$$\begin{aligned}
Y_{icsk}^{EL} = & \alpha + \beta_{T1}T1_{csk} + \beta_{T2}T2_{csk} + \beta_{T3}T3_{csk} + \beta_{CX}CX_{csk} \\
& + \theta_1(T1_{csk} \times Var_{icsk}) + \theta_2(T2_{csk} \times Var_{icsk}) \\
& + \theta_3(T3_{csk} \times Var_{icsk}) + \theta_{CX}(CX_{csk} \times Var_{icsk}) \\
& + \delta Y_{icsk}^{BL} + X'_{icsk}\gamma + V'_{csk}\lambda + \phi_k + \epsilon_{2icsk}
\end{aligned} \tag{2}$$

where $(T_{csk} \times Var_{icsk})$ is the interaction of the treatment dummy with the variable of interest (i.e. baseline math score, grade level and class size). Except for the four interaction terms, equation (2) is identical to our benchmark estimation equation, i.e. equation (1).

In terms of baseline math ability, the regression analysis confirms our visual analysis of Figure 3.

Table 5: Effect heterogeneity along baseline ability, grade level and class size.

<i>Treatment indicators interacted with:</i>	Baseline Math Score	Grade Level	Class Size (log)
<i>Dependent variable: Std. IRT-Score</i>	(1)	(2)	(3)
T1: Lessons with Teacher \times Var.	−0.105*** (0.004)	−0.140*** (0.000)	−0.437*** (0.004)
T2: CAL-Lessons with Supervisor \times Var.	−0.014 (0.741)	−0.052 (0.250)	−0.109 (0.434)
T3: CAL-Lessons with Teacher \times Var.	−0.038 (0.284)	−0.058 (0.181)	−0.270* (0.052)
CX: Classes exposed to Externalities \times Var.	−0.004 (0.913)	−0.023 (0.675)	−0.118 (0.482)
Adjusted R ²	0.70	0.70	0.70
Observations	3197	3197	3197
Treatment Indicators	Yes	Yes	Yes
Baseline Score	Yes	Yes	Yes
Individual & Classroom Controls	Yes	Yes	Yes
Stratum & Grade Level FE	Yes	Yes	Yes

Notes: All interaction variables are demeaned, so that the main effects of the treatment indicators remain unaltered. p-values are based on class-level clustered standard errors and are shown in parentheses.

* p<0.10, ** p<0.05, *** p<0.01.

Regarding the effect of additional math classes instructed by teachers, the effect size and baseline achievement are indeed negatively correlated (see column 1 in Table 5). This suggests that teachers were more effective in improving the performance of children with low math ability than those children who performed well in the baseline assessment. The regression also yields negative signs for the interaction between the baseline math score and T2 (i.e. CAL + SUPERVISOR) as well as T3 (i.e. CAL + TEACHER), but the p-values do not reach the 10%-threshold. Hence, the benefit of attending CAL-based lessons was independent of initial ability levels, while the effectiveness of teachers without software was particularly low among well-performing students.

A similar pattern emerges when we study effect heterogeneity by grade level of the participating students (see column 2 in Table 5). The effects of the CAL treatments do not significantly vary with the grade level of students, but we find that additional math lessons taught by a teacher are less effective in higher grades. This corroborates the finding that without the help of learning software, teachers in Morazán seem to be least effective when explaining more complex concepts.

Finally, we find that large class sizes reduce the effectiveness of teachers (see column 3 in Table 5), no matter whether they use CAL software or not. This pattern does not emerge for CAL classes overseen by supervisors, which seems plausible since supervisors were directed to refrain from explaining math contents but solely provided technical assistance. Comparing the point estimates of the interaction terms of the two treatments conducted by teachers, we find that the effect of traditional classes ($\hat{\theta}_1=0.436$, p-value=0.005) is more sensitive to class size than the effect of CAL-lessons instructed by teachers ($\hat{\theta}_3=0.270$, p-value=0.052). Overall, this confirms the notion that computer-based learning can mitigate the problems related to large class sizes (e.g. Banerjee and Duflo, 2011; Muralidharan, Singh and Ganimian, 2019).

4.3 Program Attendance and IV-Estimates

Our benchmark analysis focuses on ITT-estimates that do not account for the actual attendance rate of students in the additional math lessons. In this section, we present data on the overall compliance, examine the correlation between individual attendance and endline scores, and finally discuss instrumental variable estimates for the impact of the three interventions assuming full attendance.

Figure 4 plots the distribution in attendance rates across all eligible students. With an average attendance rate of 59%, participation of students was a weak spot of the program. Attendance rates slightly varied across the three treatments, although the differences are statistically insignificant: Additional CAL classes instructed by teachers achieved the highest participation (60%), followed by additional classes instructed by teachers (59%) and CAL classes conducted by a supervisor (57%).

The individual attendance rate of students is strongly correlated with their performance in the endline math assessment, as one would expect considering that the programs successfully increased math learning outcomes.

Figure 5 plots the residual endline IRT-score (net of all control variables including baseline

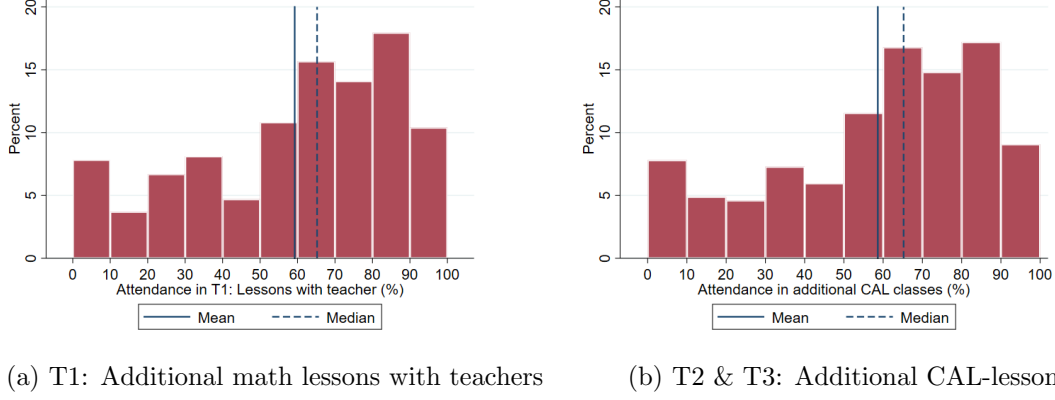


Figure 4: Attendance of students in additional math lessons.

scores) on the y-axis, and the attendance rates of the students on the x-axis. We aggregated the individual data points into 15 bins in order to improve readability, and plot the correlation by treatment type. Figure 5a covers the students that were assigned to additional math classes taught by teachers, while Figure 5b illustrates the correlation between attendance and residual endline scores for the two CAL interventions.⁹

We next appraise the question, how much children would have learned had they fully participated in the additional math lessons they were offered. To do so, we estimate an IV-model, with the first-stage estimation being specified as

$$Att_{icsk}^{T=t} = \alpha + \pi_1 T1_{csk} + \pi_2 T2_{csk} + \pi_3 T3_{csk} + \delta Y_{icsk}^{BL} + X'_{icsk} \gamma + V'_{csk} \lambda + \phi_k + \epsilon_{3icsk} \quad for \quad t \in [1, 2, 3] \quad (3)$$

where $Att_{icsk}^{T=t}$ is student's i attendance rate in treatment t and takes values between 0 and 1.

All other variables are defined as in the benchmark estimation equation, i.e. equation (1). In the

⁹Regressing endline IRT scores on attendance rates (continuous between 0 and 1), baseline scores, individual and classroom controls yields the following correlations between attendance and performance: $\hat{\gamma}_{T1}=0.40$ (t -value=5.0); $\hat{\gamma}_{T2}=0.56$ (t -value=4.2); $\hat{\gamma}_{T3}=0.55$ (t -value=3.6). Including a quadratic term we get: $\hat{\gamma}_{T1}^1=-0.53$ (t -value=-1.9), $\hat{\gamma}_{T1}^2=0.89$ (t -value=3.1); $\hat{\gamma}_{T2}^1=0.56$ (t -value=1.1), $\hat{\gamma}_{T2}^2=0.01$ (t -value=0.0); $\hat{\gamma}_{T3}^1=-0.41$ (t -value=-1.0), $\hat{\gamma}_{T3}^2=0.94$ (t -value=2.1).

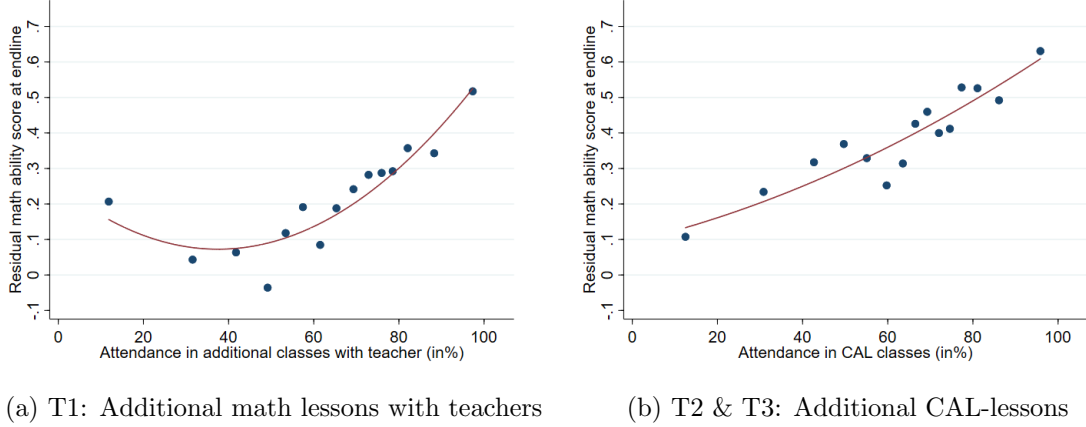


Figure 5: Residual endline test scores and attendance in additional math lessons.

Note: The figures present the partial correlation between individual attendance rates and residual endline test scores after controlling for baseline scores, individual and classroom characteristics. To ease readability, we aggregated individual data points into 15 bins.

second stage, we replace the binary treatment indicators with the predicted attendance rates from

stage 1, i.e. $\widehat{Att}_{icsk}^{T=t}$, and estimate

$$Y_{icsk}^{EL} = \alpha + \beta_1 \widehat{Att}_{icsk}^{T=1} + \beta_2 \widehat{Att}_{icsk}^{T=2} + \beta_3 \widehat{Att}_{icsk}^{T=3} + \delta Y_{icsk}^{BL} + X'_{icsk} \gamma + V'_{csk} \lambda + \phi_k + \epsilon_{4icsk}. \quad (4)$$

In order to interpret $\hat{\beta}_1$, $\hat{\beta}_2$, and $\hat{\beta}_3$ as the treatment effects of attending all 46 additional math lessons, we have to impose two (restrictive) properties that go beyond the standard monotonicity and independence assumptions (see Angrist and Pischke, 2008; Muralidharan, Singh and Ganimian, 2019). *First*, the treatment effect needs to be homogenous across students. *Second*, the functional form between attendance and math score gains has to be linear.

Our data suggest that these two additional assumptions may be violated and that the IV-estimates are potentially *downward* biased. *Effect homogeneity* seems questionable, since the impacts of the interventions are homogenous (in case of both CAL treatments) or decreasing (in case of the TEACHER treatment) in initial ability, even though attendance rates are positively correlated with baseline scores. Attending an additional math lesson thus had a stronger effect on low ability than high ability students. Hence, the IV-estimates might understate the true effect under

Table 6: IV-Estimates: Program effects with full participation

	Percent Correct		Std. IRT-Scores	
	(1)	(2)	(3)	(4)
T1: Lessons with Teacher	5.066*** (0.002)	4.739*** (0.004)	0.286*** (0.002)	0.269*** (0.005)
T2: CAL-lessons with Supervisor	7.104*** (0.000)	6.859*** (0.000)	0.390*** (0.000)	0.378*** (0.000)
T3: CAL-lessons with Teacher	7.517*** (0.000)	7.236*** (0.000)	0.411*** (0.000)	0.396*** (0.000)
Kleibergen-Paap F-statistic	214.45	193.47	213.78	192.89
Adjusted R ²	0.65	0.66	0.69	0.69
Observations	2570	2570	2570	2570
Baseline Score	Yes	Yes	Yes	Yes
Individual & Classroom Controls	No	Yes	No	Yes
Stratum & Grade Level FE	Yes	Yes	Yes	Yes

Notes: p-values are based on class-level clustered standard errors and are shown in parentheses.

* p<0.10, ** p<0.05, *** p<0.01.

full participation. Moreover, the functional form between attendance and ability gains appears to be (slightly) convex rather than linear, suggesting that children experienced increasing returns to attending the additional lessons. Again this would lead to a downward bias in the IV-estimates.

Table 6 presents the IV-estimates, which can be interpreted as the (potentially downward biased) treatment effects of attending all 46 additional math lessons. Attending the full CAL program during the intervention period leads to an increase in the endline score of about 7 percentage points or 0.38σ to 0.41σ , which is comparable in magnitude to effects of technology-aided instruction found in India, where Muralidharan, Singh and Ganimian (2019) report average learning gains in math of 0.6 standard deviations for 90 days attendance at CAL learning centers.

An increase in math ability of 0.4σ is about equivalent to the average student's progress in 1.1 school years. This translation of the average treatment effects under full compliance into school-year equivalents should be read with caution, however: On the one hand, school year equivalents do not only represent what children learn in their regular math classes at school but also reflect age-based cognitive development, learning at home or spillovers from other subjects. On the other

hand, our monitoring data suggest that about 25% of regular lessons are canceled due to teacher absenteeism and that children miss another 10%. Hence, compliance in regular classes is far from perfect, which complicates statements concerning the relative effectiveness of the additional CAL-based lessons compared to regular math lessons based on these estimates.

5 Discussion

5.1 Treatment Externalities

Our research design allows us to quantify spillovers on non-treated classes in program schools. As discussed in section 4.1, we find positive and significant externalities: Students assigned to control classes in program schools scored about 0.14σ higher in the endline assessment than students assigned to pure control classes. This effect is comparable in magnitude to the treatment effect for additional math lessons instructed by teachers. While we do not have rigorous experimental evidence to pin down the mechanisms with certainty, the data we collected from different sources allows for a discussion of what may (or may not) explain these externalities. In the following we distinguish between three broad explanations: *(i)* direct exposure of students in control classes, *(ii)* behavioral adjustments to the experimental design, and *(iii)* social learning among peers.

Direct Exposure. We begin with examining the hypothesis that control students in program schools may have been directly exposed to one of the treatments, either by (illicitly) participating in the additional math lessons, by targeted migration and class changes, or by using CAL-software in regular lessons or at home.

To prevent direct exposure of control students to the treatments, the implementing NGO instructed contract teachers and supervisors to confine access to children that were registered as participants. Our monitoring data shows high compliance with this directive, as unauthorized

participation was only recorded during 6 out of about 750 unannounced visits in treatment classes.

Likewise, we aimed to eliminate any incentives to change classes or schools and therefore barred students that changed into treatment classes during the school year from attending the additional math lessons. Only 38 students in our estimation sample changed classes or schools during the program and excluding these students from the estimation models leaves the results unchanged.

Control students in program classes may also have been exposed to the learning software in regular classes or at home. Again, our data suggests otherwise: The enumerators recorded computer usage in only 5 out of about 1,000 regular class visits. Similarly, computer usage at home is an unlikely candidate to account for treatment externalities: According to our socio-demographic survey, only 576 students (about 18%) live in a household that owns a computer with internet access and this asset class is not correlated with learning outcomes in the endline assessment.

Behavioral Adjustments to the Experimental Design. There are several ways in which students or school staff might have inadvertently adjusted their behavior to the experiment (see appendix A.4 for a more detailed discussion). *First*, the presence of the NGO could have *incentivized* schools to make a good impression, for instance to be allowed to keep the IT equipment after the intervention or to be considered for future collaborations. Our data does not support this claim, however, as neither teacher nor student attendance is higher in program schools (see Table A.4) and the number of installed computers is uncorrelated with student performance (see Table 7).

Second, the difference between control classes within and outside program schools may be driven by *John Henry effects*, a bias induced from reactive behavior of the control group to overcome the *disadvantage* of not being treated (e.g. Glennerster and Takavarasha, 2013). If such behavior arises within program schools, but not in spatially separated pure control schools, it could account for the observed externalities. This mechanism has similar implications, but is distinguishable from those discussed in the previous paragraph. While the previous paragraph explores the possibility

of a general boost in student or teacher motivation across all groups in treatment schools, the John Henry effect would only operate for the control group. As shown in columns (3) and (6) of Table A.4, the data reject this hypothesis.

Third, the experiment might have induced behavioral changes in response to being *observed*, so called *Hawthorne effects* (e.g. Levitt and List, 2011). If being part of an experiment was more salient to subjects in program schools, they might have worked harder as a response to being observed, producing the pattern we find in our data. This seems unlikely for several reasons. Most importantly, the monitoring process was structured homogeneously, meaning that enumerators visited program and control schools with the same frequency and followed the same observational procedure. Moreover, only few studies provide evidence for the presence of Hawthorne effects in the context of educational interventions, even though the topic received considerable attention (e.g. Adair, 1984; Adair, Sharpe and Huynh, 1989; Krueger, 1999).

Finally, divergent behavioral responses by treatment status might only have occurred during the math assessment. A large body of literature shows that *test-taking motivation* can have profound effects on low-stakes test results (e.g. Silm, Pedaste and Täht, 2020). Since the implementation of our experiment did not hint at any personal or institutional rewards for participants, it seems unlikely that the treatment status systematically influenced students’ test-taking effort. One may further hypothesize that students in program schools put more effort into the tests because they perceived it as more “purposeful” when other classes of the same grade also participated in the examinations. We test this claim by interacting the control classes with a binary indicator equaling one for classes in schools that have other classes of the same grade level that took the test (almost exclusively satisfied in program schools), but do not find a significant correlation. Finally, even if motivational effects were present, one would also expect them to have influenced performance during the baseline assessment, which would cancel out any potential bias operating via this channel.

Table 7: Externality channel: Proxies for social learning and in-kind incentives

<i>Dependent variable: Std. IRT Score</i>	Treatment Intensity		Installed NGO computers	
<i>CX-indicator interacted with:</i>	All Treatments	CAL	Per Student	Total
	(1)	(2)	(3)	(4)
CX: Control Classes for Externalities	0.146** (0.019)	0.135** (0.023)	0.142** (0.020)	0.146** (0.037)
CX: Control Classes for Externalities \times Var.	0.010 (0.290)	0.015*** (0.001)	0.031 (0.950)	0.001 (0.865)
Adjusted R ²	0.73	0.74	0.73	0.73
Observations	1279	1279	1279	1279
Individual & Classroom Controls	Yes	Yes	Yes	Yes
Baseline Score	Yes	Yes	Yes	Yes
Stratum & Grade FE	Yes	Yes	Yes	Yes

Notes: Treatment intensity defined as share of treated students in a school. p-values are based on school-level clustered standard errors and are shown in parentheses. * p<0.10, ** p<0.05, *** p<0.01.

Social Learning among Peers. The treatment externalities may also stem from *peer effects*, as participants could have shared their knowledge with schoolmates from other classes. Results in columns (1) and (2) of Table 7 suggest that this may have been the case: What explains part of the performance differential between program school control classes and pure control classes is the share of children that participated in the CAL treatments. One explanation is that the learning gains produced by CAL were passed on by the participants to their peers from non-treated classes. Another explanation for this pattern would be that hosting many CAL classes went along with a more generous furnishing of computer-labs by the NGO, which might have incentivized school staff to make a good impression with the NGO so that they could keep the equipment even after the NGO-run program expired. As discussed above, columns (3) and (4) in Table 7 show no relevant correlation between the number of NGO computers installed in a school and the endline performance of students. Hence, the interpretation that CAL beneficiaries passed on their learning gains to their peers seems more plausible than behavioral adjustments in prospect of being donated equipment. This finding is consistent with a literature of peer-effects that documents how the performance of each student affects achievements of her classmates (see Sacerdote, 2011).

Summarizing Remarks. Although we cannot comprehensively pin down the channels through which the observed externalities operate, *social learning among peers* is the mechanism that can be reconciled best with the data at hand. In contrast, we are confident to rule out *direct exposure* of control units to the evaluated treatments, as our data documents excellent compliance with the experimental protocol. *Behavioral adjustments to the experimental design* may unfold in many ways, which makes it difficult to track them exhaustively. We tested several potential channels of this category, but the data consistently rejects the underlying hypotheses.¹⁰

5.2 Cost-Effectiveness

Since all three interventions were assessed within the same context and framework, we can directly compare their cost-effectiveness. The bulk of expenditures comes from salaries to teachers and supervisors (65% for TEACHER, 41% for CAL + SUPERVISOR, and 51% for CAL + TEACHER). The two computer treatments additionally entail costs for acquiring the IT equipment, shipping

¹⁰The discussed channels imply three competing interpretations of our results: *First*, the program may have unfolded treatment externalities in the narrower sense due to social learning among peers. This is supported by our data as we indeed find a correlation between the number of students attending CAL lessons and the performance of their peers in control classes. *Second*, the observed pattern may not result from actual spillovers, but from a biased estimate for the program school control classes caused by John Henry effects. John Henry effects may operate via test effort on the assessment day (not testable with our data) or throughout the implementation phase (no supporting evidence in our data), but – in either case – they would *not* affect the interpretation of the three treatment estimates. A *third* possibility is that the pure control classes do not constitute a sound counterfactual because students from program and control schools differed systematically in their test taking effort. In this the case, the performance differential between the two control groups would have to be subtracted from the impact estimates for each treatment, roughly halving the impact the two CAL treatments (T2: p-values=0.10-0.13; T3: p-values=0.04-0.06) and virtually eliminating the impact of the teacher treatment (T1: p-values=0.73-0.79). While we do not find any supporting evidence for this claim, we do not possess the data to rule it out with certainty. Since such behavior might occur in a large number of educational field experiments, incorporating measures of test taking effort in future experiments may yield important insights about this potential methodological artifact.

it to El Salvador and maintaining it. Since our partner NGO acquired most computers as in-kind donations, the factual IT-related costs incurred by the NGO (about 18 USD per computer) provide a poor guidance for educational policy-makers aiming to implement CAL interventions at scale. To make the cost-effectiveness calculations more insightful for a generic setting, we assume costs of 200 USD per work station and an average of five years of usage time.

Based on these assumptions for the costs of the computer hardware, the cost accounting of our partner NGO, and the guidelines developed by Dhaliwal et al. (2014), we estimate the cost per child to be 44 USD for TEACHER, 43 USD for CAL + SUPERVISOR, and 56 USD for CAL + TEACHER. Assuming a linear dose-response-relationship, TEACHER can thus be expected to yield a 0.35σ increase in test scores per 100 USD, while investing the same amount of money in CAL lessons would produce 0.49σ and 0.43σ , respectively. This implies that even when the computers have to be acquired at a considerable price, the two CAL interventions outperform additional teacher-led classes in terms of cost-effectiveness. Moreover, hiring lower-paid supervisors rather than officially certified teachers to conduct the CAL classes might be slightly more cost-effective, as supervisor were paid only about 60% of a teacher’s wage. These conclusions should be interpreted with caution: Not only is precision impaired by the statistical uncertainty of our estimates, but the relative cost-effectiveness also depends on contextual factors such as local wages.

5.3 The Role of Teacher Ability

Multifaceted evidence derived in our analysis points to a relatively low productivity of teachers. *First*, the difference in learning gains between program school control classes and classes receiving additional teacher-centered math lessons is close to zero and statistically insignificant (p-values around 0.7). Similarly, teachers do not seem to add much to the effect of computer-assisted learning lessons: The estimated impact for CAL lessons instructed by teachers is only marginally and

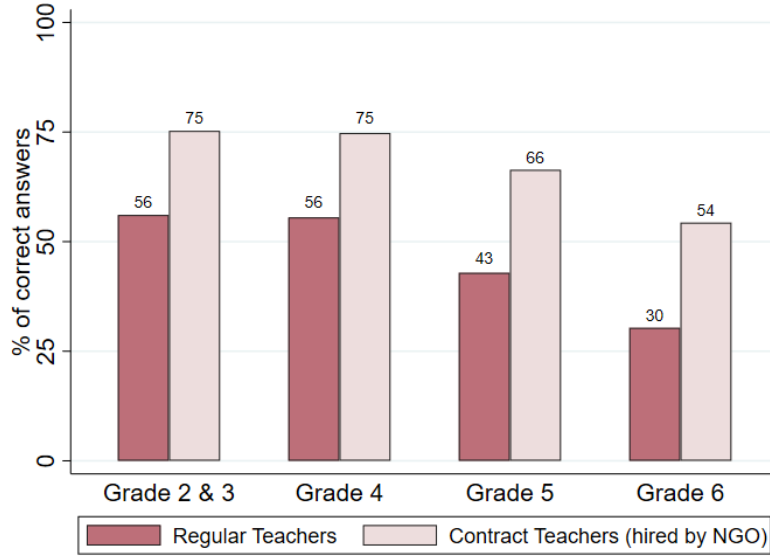


Figure 6: Math proficiency among regular teachers and teachers hired for additional math lessons. *Note:* The graph shows the share of correct answers on questions covering the official math curriculum of grades 2 & 3, grade 4, grade 5, and grade 6. This data was collected after the endline assessment for students in late 2018 and early 2019. The sample includes all program teachers as well as a representative sample of regular primary school teachers teaching math in grades 3 to 6 in the department of Morazán. *Source:* Brunetti et al. (2020).

insignificantly higher than that of CAL lessons conducted by supervisors (p-values around 0.6). *Second*, the heterogeneity analysis shows that the productivity of teachers declines as the complexity of concepts increases: The impact of the additional math lessons instructed by a teacher is decreasing in both the grade level as well as the baseline achievement of their students. *Third*, both CAL interventions outperform the additional math lessons instructed by teachers: The point estimates of the CAL interventions are consistently larger, and their impact neither decreases with student baseline performance nor grade level. Hence, it appears that in our setting, learning software is more productive in teaching basic math than officially certified teachers, especially as the complexity of the content increases.

In order to analyze the root cause of the low productivity of teachers, we asked the instructors hired by the NGO to participate in a 90 minutes math assessment covering the primary school curriculum of grades 2 to grade 6. Moreover, we administered the same assessment to a representative sample of regular math teachers of grade 3 to grade 6 classes allowing us to learn how

the contract teachers compare to the regular teaching staff (see Brunetti et al., 2020, for further details). Figure 6 illustrates the main insights from this assessment: Primary school math teachers in the department of Morazán do not master large parts of the content they are supposed to teach. The contract teachers hired by the NGO answered on average only 75% of the second and third grade questions correctly and this share declines to 54% for the sixth grade questions. Hence, even for the simplest questions, the average contract teacher does not meet the minimum proficiency of 80% correct answers as advocated by the World Bank (see Bold et al., 2017a; World Bank, 2018).

These insights raise the question whether the teachers hired for the intervention have a particularly low proficiency in math – which could explain why they are not part of the publicly employed teaching staff. Figure 6 suggests otherwise: Regular teachers performed considerably worse than the contract teachers, as they achieved on average only 56% correct answers on second and third grade questions and 30% on items pertaining to the sixth grade curriculum.¹¹

Inadequate content knowledge of teachers likely puts a binding constraint on their productivity. To test this claim we conduct two additional analysis using the data on the math ability of contracted instructors as well as regular math teachers:¹² *First*, we re-estimate equation (2) and interact the three treatment dummies with the instructors’ standardized math ability score. The math content knowledge of teachers is correlated with student learning gains in both traditional ($\hat{\theta}_1=0.08$,

¹¹Note that the implementing NGO administered a very short math assessment in the hiring process in order to eliminate the least qualified candidates. Moreover, the hired teachers participated in several workshop to prepare them for the teaching assignment. Since the assessment reported in Figure 6 was conducted *after* the intervention finished, it is likely that the NGO’s selection process and the additional training for the contract teachers partly explains the pronounced differences in content knowledge between the regular teachers and the contract teachers.

¹²This additional analysis comes with some caveats: The experimental protocol did not take teacher ability into consideration, which is why teachers were assessed in the aftermath of the field experiment. Moreover, the number of different instructors was not optimized with respect to statistical power. The implementing NGO hired 23 teachers and 15 instructors to conduct the additional math lessons; all contract teachers instructed both traditional lessons (T1) and CAL-based courses (T3) so that an average contract teacher was responsible for two class per treatment arm.

p-value=0.28) and CAL-based math lessons ($\hat{\theta}_3=0.09$, p-value=0.14), whereas the math score of supervisors is virtually orthogonal to learning gains ($\hat{\theta}_2 < 0.01$, p-value=0.94). Since supervisors did not provide math related explanations, it makes sense that their math ability does not moderate the impact of CAL-based lessons. *Second*, we correlate the standardized math ability of regular teachers with the learning gains of their students between the baseline and the endline assessment. Depending on the model specification, the point estimates vary between 0.09 and 0.12 and are significant at the 0.05 level or higher (see Table 4 in Brunetti et al., 2020). Although these are purely observational estimates, they are not only very similar to the point estimates obtained for the impact of teacher ability in the program classes reported above, but also to quasi-experimental evidence established in studies from various countries: The benchmark estimates for the annual impact of one standard deviation in additional teacher content knowledge on standardized learning outcomes of children are 0.09 for math in Peruvian primary schools (Metzler and Woessmann, 2012), 0.09 for math in Pakistani primary schools (Bau and Das, 2020), and 0.07 for math and language in Eastern African primary schools (Bold et al., 2019). Overall, these consistently positive point estimates for teacher content knowledge corroborate the hypothesis that the contract teachers' poor subject mastery impaired the impact of the evaluated teacher-centered math lessons. Hence, inadequate content knowledge is a plausible factor that helps to explain the low productivity of teachers reported in this study.

In view of drawing general conclusion for the effectiveness of additional math lessons instructed by regular teachers, the results reported in Figure 6 are particularly grim. The relatively low impacts found for the additional math lessons instructed by contract teachers may be too optimistic when aiming for a scale-up with regular teachers, who have on average a lower math proficiency than the contract teachers hired by the implementing NGO. According to our most reliable estimates, it would take a teacher with 88% correct answers on the administered assessment to teach additional

math classes so that the attending students improve on average as much as students attending CAL lessons overseen by a supervisor. This score corresponds to the 75th math ability percentile among the hired contract teachers and the 91st math ability percentile among regular primary school teachers in Morazán.¹³

These results highlight how learning software can compensate for the poor content knowledge of teaching staff. Earlier contributions on the value of computer-assisted learning emphasized its advantages in terms of mitigating issues of large class sizes and the challenges of “teaching at the right level” (e.g. Banerjee and Duflo, 2011; Muralidharan, Singh and Ganimian, 2019). While our heterogeneity analysis corroborates this line of reasoning, this section showed that CAL can help to remedy shortcomings related to low teacher ability.

6 Conclusion

Computer-assisted learning (CAL) is widely perceived as a promising approach to address the low quality of teaching in developing countries. While encouraging, previous research is inconclusive regarding the value of technology-based instruction relative to traditional teaching and has little to say on the complementarities between teachers and learning software. The evidence presented in this paper suggests that CAL can not only produce substantial learning gains, but may also outperform traditional instruction. In our setting, this relative advantage seems to be driven by a mismatch between teacher qualification and the complexity of the concepts they have to teach: Under traditional teaching models, children are unlikely to learn what their teachers fail to understand, while CAL allows them to make progress beyond their teachers’ content knowledge. Overall,

¹³For this back-of-the-envelope calculation, we use 0.063σ for the impact difference between teacher-based additional lessons and CAL lessons monitored by supervisors (see Table 2), 0.09 for the conversion factor of standardized teacher content knowledge on standardized student ability, and data on content knowledge of teachers presented in Figure 6.

our findings point to an alarmingly low productivity of teachers. Not only is the effect of additional teacher-led instruction comparatively low (and might be partly if not completely attributable to treatment externalities), but poorly qualified teachers also do little to improve the productivity of CAL lessons. In light of the fact that they do not master a substantial share of the contents they are required to teach, these results are hardly surprising.

Promoting the targeted use of computers may therefore be an attractive option for governments and NGOs operating in settings with low teacher quality. When teachers are struggling with the concepts they have to teach, learning software can be an important remedy allowing them to improve the quality of their teaching. Another approach would be to invest in the skills of teachers, for instance by offering professional development programs: Teachers may not make much of a difference when they do not master what their students are supposed to learn, but vast empirical evidence from developed countries suggests that they can matter a great deal when they are well prepared and adequately qualified (Rockoff, 2004; Chetty, Friedman and Rockoff, 2014). Hence, gaining a better understanding of how teachers' preparedness, and particularly their content knowledge, can be improved is likely to yield large social returns. Since hardly any rigorous evidence on this aspect is available (Muralidharan, 2017; Bold et al., 2017a), we teamed up with the same implementing partner to examine whether computer-assisted learning software can help to advance the content knowledge of teachers and therewith their productivity in the classroom (see Brunetti et al., 2019).

References

- Adair, John. 1984. "The Hawthorne Effect: A Reconsideration of the Methodological Artifact." *Journal of Applied Psychology* 69(2):334–345.
- Adair, John, Donald Sharpe and Cam-Loi Huynh. 1989. "Hawthorne Control Procedures in Educational Experiments: A Reconsideration of Their Use and Effectiveness." *Review of Educational Research* 59(2):215–228.
- Angrist, Joshua and Jörn-Steffen Pischke. 2008. *Mostly Harmless Econometrics. An Empiricist's Companion*. Princeton: Princeton University Press.
- Attanasio, Orazio, Camila Fernández, Emla Fitzsimons, Sally Grantham-McGregor, Costas Meghir and Marta Rubio-Codina. 2014. "Using the Infrastructure of a Conditional Cash Transfer Program to Deliver a Scalable Integrated Early Child Development Program in Colombia: Cluster Randomized Controlled Trial." *BMJ* 349:1–12.
- Baird, Sarah, Aislinn Bohren, Craig McIntosh and Berk Özler. 2015. "Designing Experiments to Measure Spillover Effects." PIER Working Paper No. 15-021. URL: www.ssrn.com/ (last access: 10.03.2020).
- Banerjee, Abhijit and Esther Duflo. 2011. *Poor Economics*. London: Penguin Books.
- Banerjee, Abhijit, Shawn Cole, Esther Duflo and Leigh Linden. 2007. "Remedying Education: Evidence from Two Randomized Experiments in India." *The Quarterly Journal of Economics* 122(3):1235–1264.
- Bau, Natalie and Jishnu Das. 2020. "Teacher Value-Added in a Low-Income Country." *American Economic Journal: Economic Policy* 12(1):62–96.
- Beatty, Amanda, Emilie Berkhout, Luhur Bima, Thomas Coen, Menno Pradhan and Daniel Suryadarma. 2018. "Indonesia Got Schooled: 15 Years of Rising Enrolment and Flat Learning Profiles." RISE Working Paper No. 18/026. URL: <https://riseprogramme.org> (last access: 07.10.2020).
- Bold, Tessa, Deon Filmer, Ezequiel Molina and Jakob Svensson. 2019. "The Lost Human Capital: Teacher Knowledge and Student Achievement in Africa." World Bank Policy Research Working Paper No. 8849. URL: www.worldbank.org/ (last access: 10.03.2020).
- Bold, Tessa, Deon Filmer, Gayle Martin, Ezequiel Molina, Brian Stacy, Christophe Rockmore, Jakob Svensson and Waly Wane. 2017a. "Enrollment without Learning: Teacher Effort, Knowledge and Skill in Primary Schools in Africa." *Journal of Economic Perspectives* 31(4):185–204.
- Brunetti, Aymo, Konstantin Büchel, Martina Jakob, Ben Jann, Christoph Kühnhanss and Daniel Steffen. 2020. "Teacher Content Knowledge in Developing Countries: Evidence from a Math Assessment in El Salvador." Working Paper No. 2005, Department of Economics, University of Bern. URL: www.vwi.unibe.ch (last access: 10.03.2020).

- Brunetti, Aymo, Konstantin Büchel, Martina Jakob, Ben Jann, Christoph Kühnhanss and Daniel Steffen. 2019. “How Effective Are Computer-Based Teacher Training Programs? Evidence from a Randomized Controlled Trial in El Salvador.” *AEA RCT Registry*: <https://doi.org/10.1257/rct.4092-2.0> (last access: 05.03.2020).
- Carrillo, Paul, Mercedes Onofa and Juan Ponce. 2011. “Information Technology and Student Achievement: Evidence from a Randomized Experiment in Ecuador.” IDB Working Paper Series No. 223. URL: www.econstor.eu (last access: 10.03.2020).
- Chaudhury, Nazmul, Jeffrey Hammer, Michael Kremer, Karthik Muralidharan and Halsey Rogers. 2006. “Missing in Action: Teacher and Health Worker Absence in Developing Countries.” *Journal of Economic Perspectives* 20(1):91–116.
- Chetty, Raj, John Friedman and Jonah Rockoff. 2014. “Measuring the Impacts of Teachers II: Teacher Value-Added and Student Outcomes in Adulthood.” *American Economic Review* 104(9):2633–2679.
- de Ayala, R.J. 2009. *The Theory and Practice of Item Response Theory*. New York: Guilford Press.
- Dhaliwal, Iqbal, Esther Duflo, Rachel Glennerster and Caitlin Tulloch. 2014. Comparative Cost-effectiveness Analysis to Inform Policy in Developing Countries: a General Framework with Applications for Education. In *Education Policy in Developing Countries*, ed. Paul Glewwe. Chicago and London: University of Chicago Press pp. 285–338.
- DIGESTYC, Direccion General de Estadistica y Censos El Salvador. 2018. “Encuesta de Hogares de Direccion General de Estadistica y Censos 2017 (EHPM).” Online available, URL: www.digestyc.gob.sv (last access: 25.07.2018).
- Escueta, Maya, Andre Nickow, Philip Oreopoulos and Vincent Quant. 2020. “Upgrading Education with Technology: Insights from Experimental Research.” *Journal of Economic Literature* forthcoming.
- Glennerster, Rachel and Kudzai Takavarasha. 2013. *Running Randomized Evaluations: A Practical Guide*. Princeton: Princeton University Press.
- Glewwe, Paul and Karthik Muralidharan. 2016. Improving Education Outcomes in Developing Countries: Evidence, Knowledge Gaps and Policy Implications. In *Handbook of the Economics of Education*, ed. Eric Hanushek, Stephen Machin and Ludger Woessmann. Amsterdam: Elsevier pp. 653–743.
- Hess, Simon. 2017. “Randomization Inference with Stata: A Guide and Software.” *The Stata Journal* 17(3):630–651.
- Krueger, Alan. 1999. “Experimental Estimates of Education Production Functions.” *Quarterly Journal of Economics* 114(2):497–532.
- Lai, Fang, Renfu Luo, Lixiu Zhang, Xinzhe Huang and Scott Rozelle. 2015. “Does Computer-Assisted Learning Improve Learning Outcomes? Evidence from a Randomized Experiment in Migrant Schools in Beijing.” *Economics of Education Review* 47(1):34–48.

- Levitt, Steven and John List. 2011. “Was There Really a Hawthorne Effect at the Hawthorne Plant? An Analysis of the Original Illumination Experiments.” *American Economic Journal: Applied Economics* 3(1):224–238.
- Levitt, Steven, John List, Susanne Neckermann and Sally Sadoff. 2016. “The Behavioralist Goes to School: Leveraging Behavioral Economics to Improve Educational Performance.” *American Economic Journal: Economic Policy* 8(4):183–219.
- Linden, Leigh. 2008. “Complement or Substitute? The Effect of Technology on Student Achievement in India.” infoDev Working Paper No. 17. URL: www.worldbank.org/ (last access: 10.03.2020).
- Mbiti, Isaac. 2016. “The Need of Accountability in Education in Developing Countries.” *Journal of Economic Perspectives* 30(3):109–132.
- Mbiti, Isaac, Karthik Muralidharan, Mauricio Romero, Youdi Schipper, Constantine Manda and Rakesh Rajani. 2019. “Inputs, Incentives, and Complementarities in Education: Experimental Evidence from Tanzania.” *The Quarterly Journal of Economics* 134(3):1627–1673.
- Metzler, Johannes and Ludger Woessmann. 2012. “The Impact of Teacher Subject Knowledge on Student Achievement: Evidence from Within-Teacher Within-Student Variation.” *Journal of Development Economics* 99(2):486–496.
- Miguel, Edward and Michael Kremer. 2004. “Worms: Identifying Impacts on Education and Health in the Presence of Treatment Externalities.” *Econometrica* 72(1):159–217.
- MINED, Ministerio de la Educacion de El Salvador. 2013. “Elementos para el Desarrollo del Modelo Pedagogico del Sistema Educativo Nacional – Escuela Inclusiva de Tiempo Pleno.” Online available, URL: <https://www.mined.gob.sv/jdownloads/Institucional/modelopedagogico.pdf> (last access: 14.01.2018).
- MINED, Ministerio de la Educacion de El Salvador. 2018. “Informe de resultados: PAES 2018.” Online available, URL: <https://www.mined.gob.sv> (last access: 17.06.2020).
- Mo, Di, Linxiu Zhang, Jiafu Wang, Weiming Huang, Yao Shi, Matthew Boswell and Scott Rozelle. 2015. “Persistence of Learning Gains from Computer Assisted Learning: Experimental Evidence from China.” *Journal of Computer Assisted Learning* 31:562–581.
- Morgan, Kari and Donald Rubin. 2012. “Rerandomization to Improve Covariate Balance in Experiments.” *The Annals of Statistics* 40(2):1263–1282.
- Muralidharan, Karthik. 2017. Field Experiments in Education in Developing Countries. In *Handbook of Economic Field Experiments*. Amsterdam: Elsevier pp. 323–385.
- Muralidharan, Karthik, Abhijeet Singh and Alejandro Ganimian. 2019. “Disrupting Education? Experimental Evidence on Technology-Aided Instruction in India.” *American Economic Review* 109(4):1426–1460.
- Pritchett, Land and Amanda Beatty. 2015. “Slow Down, You’re Going Too Fast: Matching Curricula to Student Skill Levels.” *International Journal of Educational Development* 40:276–288.

- PAL, People's Action for Learning Network. 2020. "International Common Assessment of Numeracy." Online available, URL: <https://palnetwork.org/ican/> (last access: 14.10.2020).
- Rios, Joseph. 2020. "Improving Test-Taking Effort in Low-Stakes Group-Based Educational Testing: A Meta-Analysis of Interventions." *Applied Measurement in Education*, forthcoming.
- Rockoff, Jonah. 2004. "The Impact of Individual Teachers on Student Achievement: Evidence from Panel Data." *American Economic Review Papers and Proceedings* 94(2):247–252.
- Sacerdote, Bruce. 2011. Peer Effects in Education: How Might They Work, How Big Are They and How Much Do We Know Thus Far? In *Handbook of the Economics of Education*, ed. Eric Hanushek, Stephen Machin and Ludger Woessmann. Amsterdam: Elsevier pp. 249–277.
- Silm, Gerli, Margus Pedaste and Karin Täht. 2020. "The Relationship between Performance and Test-Taking Effort when Measured with Self-Report or Time-Based Instruments: A Meta-Analytic Review." *Educational Research Review* 31:1–22.
- UNESCO, United Nations Educational, Scientific and Cultural Organization. 2019. "UNESCO Institute for Statistics Database." Online available, URL: <http://data.uis.unesco.org/> (last access: 04.12.2019).
- The Economist. 2017. "Technology is Transforming What Happens When a Child Goes to School." published in *Briefing* section of the print edition under headline "Machine Learning", July 22nd 2017.
- World Bank. 2018. *World Development Report 2018: Learning to Realize Education's Promise*. Washington D.C.: World Bank.
- Yang, Yihua, Linxiu Zhang, Junxia Zeng, Xiaopeng Pang, Fang Lai and Scott Rozelle. 2013. "Computers and the Academic Performance of Elementary School-Aged Girls in China's Poor Communities." *Computers & Education* 60(1):335–346.
- Young, Alwyn. 2019. "Channeling Fisher: Randomization Tests and the Statistical Insignificance of Seemingly Significant Experimental Results." *The Quarterly Journal of Economics* 134(2):557–598.

A Appendix: Additional Analyses

A.1 Learning Gap and Grade Level Heterogeneity in our Sample

In order to examine the learning gap and grade level heterogeneity in our sample of primary school pupils, we follow the approach by Muralidharan, Singh and Ganimian (2019) and convert the pupils' performance in the baseline assessment into a proficiency measure expressed in grade levels. As point of origin, we calculate for each participant her share of correct answers by item grade level. The score that a child obtains in our discrete proficiency measure is determined by those grade specific set of items, where the child scores at least 50% correct answers. To be assigned to a certain grade level, a participant needs to reach the 50%-threshold that corresponds with said grade level and all preceding grades. For example, a fourth grader that scored 80% on first grade items, 55% on second grade items and 40% on third grade items would be assigned to a second grade proficiency level. Participants answering less than 50% of first grade items correctly, are assigned to grade level <1.

Based on the previously specified measure, which is plotted in Figure 1b, we obtain a performance gap of two grades between the best and worst student in the *median* class of our sample. By construction the *mean* in the within-class performance range is lowest in third grade classes (about 1.3, i.e. the math abilities of students' within the same class cover on average 2.3 grades) and highest in sixth grade classes (about 2.4). A simple regression analysis also confirms that within-class variation is substantial, as classroom fixed effects only account for about 25% of the total variation at a certain grade level.

A.2 Attrition

In Table A.1 we examine whether the attrition at endline is correlated with the treatment status. To do so, we present results based on Linear Probability Models in columns (1) to (3), and on Logit Models in columns (4) to (6). The results unequivocal suggest, that the probability to miss the endline test did not depend on the treatment status.

Table A.1: Differences in attrition across treatments

<i>Dependent var.: Attrition at endline</i>	OLS			Logit		
	(1)	(2)	(3)	(4)	(5)	(6)
T1: Lessons with Teacher	0.018 (0.302)		0.017 (0.318)	0.224 (0.298)		0.224 (0.293)
T2: CAL-Lessons with Supervisor	0.021 (0.203)		0.026 (0.115)	0.263 (0.202)		0.330 (0.116)
T3: CAL-Lessons with Teacher	0.023 (0.226)		0.025 (0.190)	0.280 (0.215)		0.315 (0.173)
CX: Control Classes for Externalities	0.019 (0.307)		0.022 (0.237)	0.236 (0.298)		0.282 (0.228)
Baseline math score		-0.002*** (0.000)	-0.002*** (0.000)		-0.024*** (0.000)	-0.024*** (0.000)
Adjusted R ²	0.00	0.01	0.01	-	-	-
Pseudo R ²	-	-	-	0.00	0.02	0.02
Observations	3528	3528	3528	3528	3528	3528

Notes: p-values (in parentheses) are based on class-level clustered standard errors. * p<0.10, ** p<0.05, *** p<0.01.

A.3 Method of Inference and Robustness of our Results

As explained in section 4.1, we apply two methods of inference. When we assess the impact of the different treatments relative to the children in pure control classes, the reported p-values are based on class-level clustered standard errors. Inference on within program school comparisons between the different treatments (including control classes subject to externalities), however, are based on a randomization inference test statistic with 2,000 random draws subject to the identical cut-off criterion as used in our re-randomization scheme.

This mixed estimation approach directly follows from our two-step randomization design (see Figure 2). Randomization inference is indispensable when comparing experimental groups within program schools since the underlying assignment process involved re-randomization. Conversely, selection of program schools and pure control schools was not based on re-randomization, making the use of randomization inference less critical.

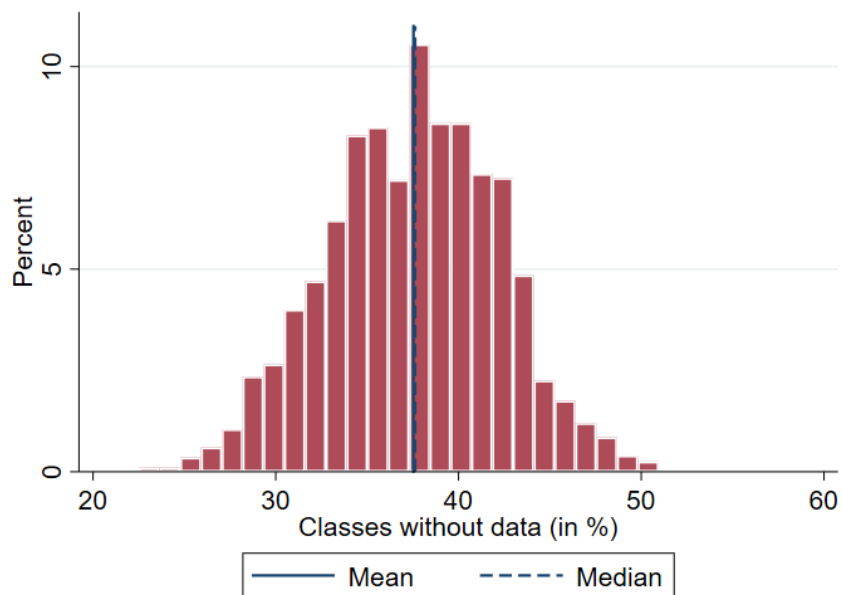


Figure A.1: Full re-randomization (incl. steps 1, 2a, and 2b) and the share of classes without data points ($N=2000$ draws).

Notes: This graph plots the distribution of the share of missing data points, when we conduct randomization inference by reiterating both stages of our randomization procedure. The large number of missing data points weakens the precision of our estimates, which explains why the p-values in the upper panel of Table A.2 increase by a factor of 5 to 10 compared to the p-values in Table 2.

While randomization inference is also preferable for assignment processes based on plain (or stratified) randomization (e.g. Young, 2019), its application is problematic in our case due to a particular feature of our study design: Out of the 162 eligible classes in pure control schools, we only collected data for a random sample of 40 classes. Implementing randomization inference for both stages of the randomization process thus comes with the downside that each draw will contain a considerable number of classes that did not participate in the assessments. As illustrated in Figure A.1, re-iterating the full randomization procedure yields an average of 37% of classes without data per draw. Even though missing data points in the replication procedure create an artificial loss of statistical power, we present the respective estimates as a conservative reference point.

Table A.2: ITT-Estimates on the effects of the different interventions on children's math scores with p-values based on clustered standard errors

	Percent Correct		IRT-Scores	
	(1)	(2)	(5)	(6)
T1: Lessons with Teachers	2.904*** (0.005)	2.643** (0.012)	0.165*** (0.006)	0.152** (0.013)
T2: CAL-Lessons with Supervisor	4.095*** (0.000)	3.869*** (0.000)	0.226*** (0.000)	0.214*** (0.000)
T3: CAL-Lessons with Teacher	4.554*** (0.000)	4.328*** (0.000)	0.250*** (0.000)	0.238*** (0.000)
CX: Control Classes for Externalities	2.595** (0.011)	2.407** (0.017)	0.147** (0.013)	0.137** (0.020)
$\beta_{T4} := \beta_{T2} - \beta_{T1} = 0$	1.191	1.226	0.061	0.063
p-value ($\beta_{T4}=0$)	(0.203)	(0.180)	(0.267)	(0.241)
$\beta_{T5} := \beta_{T3} - \beta_{T1} = 0$	1.650*	1.686*	0.084	0.086*
p-value ($\beta_{T5}=0$)	(0.080)	(0.063)	(0.115)	(0.093)
$\beta_{T6} := \beta_{T3} - \beta_{T2} = 0$	0.459	0.460	0.024	0.023
p-value ($\beta_{T6}=0$)	(0.606)	(0.599)	(0.637)	(0.636)
Adjusted R ²	0.66	0.67	0.69	0.70
Observations	3197	3197	3197	3197
Individual & Classroom Controls	No	Yes	No	Yes
Baseline Score	Yes	Yes	Yes	Yes
Stratum & Grade FE	Yes	Yes	Yes	Yes

Notes: p-values based on traditional clustered standard errors in parentheses.

* p<0.10, ** p<0.05, *** p<0.01.

To assess the robustness of our results with respect to the method of inference, we report three versions of our benchmark analysis: In Table 2, the upper panel p-values are based on class-level clustered standard errors, while we run randomization tests in the lower panel. Table A.2 replicates these results, but inference is consistently based on class-level clustered standard errors. Finally, Table A.3 presents the main results with p-values based on a full randomization tests.

Reassuringly, our main conclusion do not depend on the method of inference. When we apply traditional inference to the lower panel, as in Table A.3, changes in p-values are very small and do not show a clear pattern. And despite losing a lot of power when applying randomization inference to the upper panel, as in Table A.2, the only notable difference is, that the program externalities captured by β_{CX} turn insignificant with p-values around 0.13.

Table A.3: ITT-Estimates on the effects of the different interventions on children’s math scores with p-values based on randomization inference

	Percent Correct		IRT-Scores	
	(1)	(2)	(5)	(6)
T1: Lessons with Teacher	2.904*	2.643*	0.165*	0.152*
	(0.073)	(0.089)	(0.083)	(0.097)
T2: CAL-Lessons with Supervisor	4.095***	3.869**	0.226**	0.214**
	(0.009)	(0.013)	(0.015)	(0.018)
T3: CAL-Lessons with Teacher	4.554***	4.328***	0.250***	0.238**
	(0.006)	(0.006)	(0.007)	(0.011)
CX: Control Classes for Externalities	2.595	2.407	0.147	0.137
	(0.117)	(0.136)	(0.120)	(0.140)
Adjusted R ²	0.66	0.67	0.69	0.70
Observations	3197	3197	3197	3197
Individual & Classroom Controls	No	Yes	No	Yes
Baseline Score	Yes	Yes	Yes	Yes
Stratum & Grade FE	Yes	Yes	Yes	Yes

Notes: p-values based on a two-sided randomization inference test statistic that the placebo coefficients are larger than the actual are shown in parentheses. The p-values were computed based on 2000 random draws.

* p<0.10, ** p<0.05, *** p<0.01.

A.4 A Detailed Account on Behavioral Adjustments to the Experimental Design

This section presents more exhaustive considerations regarding behavioral adjustments to the experimental design than the main text. It follows the same structure, but discusses our four sets of explanations more extensively: *(i)* unintended incentives to improve performance at the school level, *(ii)* John Henry effects, *(iii)* Hawthorne effects, and *(iv)* systematic variation in test effort.

The presence of the NGO might have induced a general motivational boost in participating schools. This is, school staff might have been incentivized to make a good impression to be allowed to keep the IT equipment or to remain part of the project after the completion of its trial phase. We first examine this reasoning by using class cancellation rates and attendance rates as proxies for the effort by school staff/teachers, and then continue by testing whether a more generous supply of computer hardware raised performance in control classes; the absolute and relative number of supplied computers varied across program schools due to heterogeneity in both class sizes as well the number of classes assigned to CAL-treatments. Contrary to expectations, cancellation rates appear to be slightly higher in program schools than in control schools although the difference is not statistically significant (see columns 4 & 5 in Table A.4).¹⁴ Similarly, student attendance rates do not point towards intensified efforts in program schools, as the estimated differences in columns (1) and (2) of Table A.4 yield p-values larger than 0.8. Finally, we test whether a more

¹⁴The project could also have affected class cancellation rates directly, e.g. via space limitations induced through the additional lessons at the expense of regular classes. Furthermore, differences in cancellation rates may be an artifact of the data collection process. To keep expenses low, we randomly selected entire *schools* rather than *classes* to be visited on a given day. Thus, enumerators had to record data from all classes on grades 3–6 in program schools, but only up to two classes during control school visits. One could hypothesize that, in control schools, data collectors were more inclined to wait for the teacher to turn up, while, in program schools, they moved on to the next class.

Table A.4: Externality channel: Motivation proxied with class attendance and cancellations.

<i>Dependent variable:</i>	Student Attendance (%)			Class Cancellations (%)		
	(1)	(2)	(3)	(4)	(5)	(6)
Program Schools	−0.304 (0.891)	−0.287 (0.896)	−0.978 (0.648)	6.879 (0.238)	6.537 (0.264)	8.215 (0.140)
Adjusted R ²	0.07	0.06	0.00	0.08	0.08	0.07
Observations	198	198	80	198	198	80
Control Classes Only	No	No	Yes	No	No	Yes
Classroom Controls	No	Yes	Yes	No	Yes	Yes
Stratum & Grade FE	Yes	Yes	Yes	Yes	Yes	Yes

Notes: p-values are based on school-level clustered standard errors and are shown in parentheses.

* p<0.10, ** p<0.05, *** p<0.01.

generous furnishing of computer-labs by the NGO has pushed schools to better performances that is not necessarily reflected in attendance and cancellation rates. Consistent with the previous results, columns (3) and (4) in Table 7 show no relevant correlation between the number of NGO computers installed in a school and the endline performance of students in control classes.

The difference between control classes within and outside program schools may also be driven by so-called *John Henry effects*. These refer to biases stemming from reactive behavior of the control group. In our experiment, students in the control group may have worked harder to overcome their disadvantage of not receiving the additional math lessons. Similarly, teachers could have redirected resources and effort towards control classes to compensate them for their relative deprivation. For example, teachers may have given more weight to math relative to other subjects when attending control classes. One could hypothesize that such behavior only arises within (but not across) schools, as students from the same school are more likely to be a relevant reference group. Limiting our analysis to the control group, we can, once again, compare student and teacher motivation between different experimental groups. As shown in columns (3) and (6) of Table A.4, limiting the analysis to the control classes does not alter our conclusions: The difference in class cancellation rates between program school control classes and pure control classes is small and remains aloof from any conventional level of statistical significance. The same is true for students' attendance rates.

Another factor potentially driving the difference between control classes within and outside program schools might be behavioral changes in response to being *observed*. Such *Hawthorne effects* would explain the measured difference in performance across the two control groups, if control units in program schools responded more strongly to the fact that they were part of an experiment than children and teachers in pure control schools. This seems unlikely for four reasons: *First*, the experiment was explained to parents and teachers before schools and classes were randomized into treatment and control groups, so that all subjects shared the same awareness about the field experiment when its implementation started. *Second*, control students received a similar degree of attention by the experimenters, no matter whether they attended program or control schools; in the monitoring process, enumerators visited both set of schools with the same frequency and followed the same procedure. Even though enumerators collected data from a larger number of classes within program schools, the time they spent with each school's headmaster and within a given classroom did not vary systematically by the schools' program status. *Third*, the results on student and teacher attendance rates in Table A.4 do not support the hypothesis that subjects in program schools had a stronger tendency to please experimenters, or what Levitt and List (2011)

describe as “experimenter demand effects”. *Finally*, very few of a considerable number of studies confirm the presence of Hawthorne effects in the context of educational interventions (e.g. Adair, 1984; Adair, Sharpe and Huynh, 1989; Krueger, 1999).¹⁵

A possibly remaining methodological artifact might be systematic differences in the participants’ *test-taking motivation*. A large body of literature, recently reviewed in Silm, Pedaste and Täht (2020) and Rios (2020), shows that test-taking motivation can have profound effects on low-stakes test results. Hence, if motivation of test-takers varied between control and program schools, this may (partly) explain the measured differences in learning gains between the two control groups.

The psychometric literature coarsely distinguishes between the following criteria that impact test-taking motivation: (i) stakes involved in test performance (e.g. grades or hard incentives), (ii) behavior of proctors, (iii) design of the test (e.g. difficulty of items), and (iv) the purposefulness of test (e.g. motivational instructions). We can rule out three of the four factors as potential performance wedges between the two control groups: The stakes were uniformly low, the test design only varied by grade level but not treatment status, and all proctors not only followed the same instructions but also supervised test-taking across multiple schools regardless of their treatment status. What potentially remains are vague differences in (iv) the purposefulness of the test, where one can distinguish between “personal” (usually more important) and “institutional” (typically less relevant) conditions. Regarding “personal” conditions, students in program schools might have perceived competition as fiercer, since only in large program schools multiple classes of the *same* grade participated in the assessment. We test this claim by interacting the control classes for externalities with a binary indicator equaling one for classes in schools that have other classes of the *same* grade level that took the test, but do not find a significant correlation (p-value=0.73, results not shown).¹⁶ In a similar vein, John Henry effects may have only kicked in on the test day, so that control students in program schools made an extra effort to balance out their disadvantage of fewer math lessons; unfortunately, this is a hypothesis we cannot test. Under “institutional” conditions, we subsume motivational instructions by school staff who potentially encouraged children more strongly, if they thought such (last-minute) briefings help to cast a good light on their school. We did not observe such behavior, however, and it is questionable whether such briefings could sufficiently shift student motivation. Levitt et al. (2016) show that *personal* rewards delivered *immediately after* assessments significantly impact student effort, while *personal* rewards delivered with a *delay of one month* did not change performance. The implementation of our experiment did neither hint at *personal* nor *immediate* rewards for participants. Hence, it is unlikely that institutional conditions systematically influenced students’ test-taking effort.

¹⁵Adair (1984) reviews 13 educational studies (7 on children and 6 on college students), of which only four studies, all of them analyzing the behavior of college students, produce evidence for Hawthorne effects. Adair, Sharpe and Huynh (1989) expand the meta-analysis to 86 studies on educational interventions. They classify less than 10% of the reviewed studies as a “Hawthorne case” and obtain an average Hawthorne effect across all studies of 0.01 (confidence interval: -0.07 and +0.08). In an experimental study on the role of class size, Krueger (1999) demonstrates a negative impact of larger classes on learning progress, while finding no indication for a potential Hawthorne effect on teachers or students. Somewhat tellingly, Levitt and List (2011) conclude from a re-examination of the *original* data collected at the Hawthorne plants in the 1930s, that “the most important lesson to be learned from the original Hawthorne experiments is the power of a good story. The mythology surrounding the Hawthorne experiments arose largely absent careful data analysis, and has persisted for decades even in the face of strong evidence against it [...]”

¹⁶This basically replicates the analysis of Table 7 using the binary indicator for multiple classes of the *same* grade level as interaction variable. We also re-estimate equation (2) using this proxy for perceived in-school competition but consistently reject the underlying hypothesis.

B Appendix: Measuring and Converting Learning Outcomes

To measure math skills of third to sixth graders, we conducted two standardized math assessments during the school year 2018. Both assessments include 60 items and were designed as follows:

1. We summarized the Salvadoran math curriculum for grades 1–6 along the three topics (a.) number sense & arithmetic, (b.) geometry & measurement, and (c.) data & probability.
2. We then mapped test items from various sources on the Salvadoran curriculum. These sources are (a.) official text books of El Salvador, (b.) publicly available items from the STAR¹⁷ evaluations in California, (c.) publicly available items from the VERA¹⁸ evaluations in Germany, and (d.) exercises from the Swiss textbook MATHWELT.
3. We then gathered pilot data on 180 test items answered by 600 Salvadoran pupils in October 2017 and estimated the difficulty and discrimination parameters of test questions based on *Item Response Theory* (e.g. de Ayala, 2009).

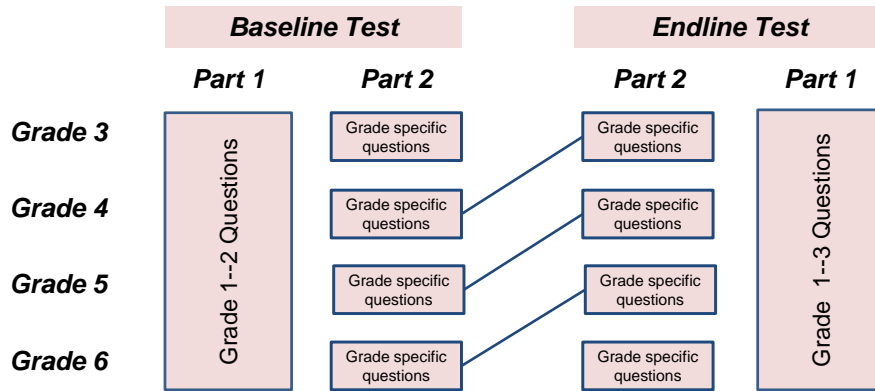


Figure B.1: Stylized illustration of the assessment design.
Note: Each part covers 30 items, adding up to 60 items per wave.

4. Finally, we designed paper and pencil maths tests using insights from step 3. The 60 items are selected such that they reflect the weighting in the official curriculum: 60–65% number sense & arithmetic, 30% geometry & measurement, 5–10% data & probability. Most items required a written answer, while the share of multiple choice questions varied between 10% and 15% depending on grade level. Figure B.1 illustrates how the math assessments at baseline and endline were structured and linked. Both assessments had two parts, with the first part being answered by all children independent of their grade. Moreover, the grade specific second part of 3rd/4th/5th graders in the endline assessment included many baseline questions of the 4th/5th/6th graders. This linking across grades and waves was essential to infer a commonly scaled ability score, i.e. the IRT scores.

¹⁷Further information on the Standardized Testing and Reporting (STAR) programme in California is available online: www.cde.ca.gov/re/pr/star.asp (last accessed: 14.01.2018).

¹⁸VERA is coordinated by the Institut für Qualitätsentwicklung im Bildungswesen (IQB), see www.iqb.hu-berlin.de/vera (last accessed: 14.01.2018).

Diagnostics. Table B.1 shows summary statistics on test items for each grade and wave of the assessment. In Table B.2 and Figure B.2, similar statistics are displayed for students’ percentage scores. As can be seen, our test is not subject to relevant floor or ceiling effects: Hardly any students could not answer a single question on a given assessment and not a single student scored all items correctly. Similarly, only one item was not solved by anyone and no question could be answered by all students. On average, students gave correct answers to about 25-43% of the questions in a test booklet (column 2 in Tables B.1 and B.2). Figure B.3a shows the corresponding IRT-based test information function for the entire assessment, i.e. for all grades and waves combined (see below for details on IRT). As can be seen, our test is very informative for students across all ability levels. However, the assessment is skewed towards high difficulty levels, meaning that it allows to differentiate very precisely among high-achieving, but less precisely among low-achieving students. Ideally, the precision (or “information”) of an assessment is highest around $\Theta = 0$ where most students are located (see Figure B.3b). This implies that, on average, students should be able to answer about 50% of the test items. This reflects our decision to construct the assessment based on the official Salvadoran curriculum in spite of the mismatch between the curriculum and students’ actual ability levels. Consequently, most of the included items could be answered by less than half of the students. While this curriculum-based approach allows for a more meaningful interpretation, it leads to a loss in terms of test information. Nevertheless, sufficient questions of differing difficulty levels are covered so that our item battery provides a reliable measurement instrument.

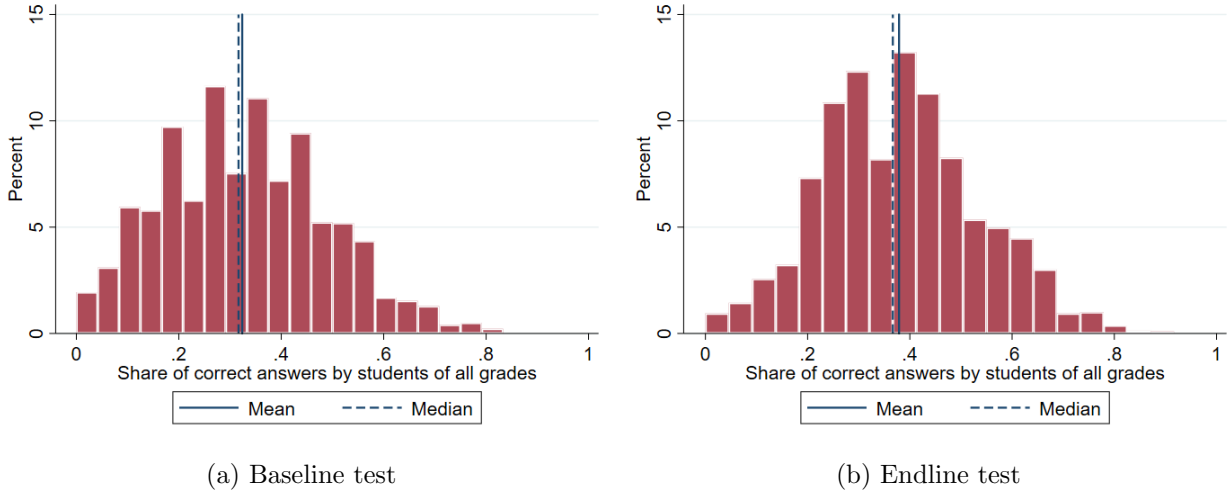


Figure B.2: Distribution of percentage scores across students

Calculating IRT-Scores. Our math assessments allows us to project all outcomes on a common ability scale by using Item Response Theory. Instead of summing up the correct answers to a total score taken to represent a person’s ability, Item Response Theory proposes a probabilistic estimation procedure. Ability is then viewed as a latent variable influencing the responses of each individual to each item through a probabilistic process: The higher a person’s ability and the lower the difficulty of a particular test item, the higher the probability of a correct answer. In the simplest form of the model, the probability that individual i succeeds on item j can be expressed as

$$Pr(success_{ij}|b_j, \theta_i) = \frac{\exp\{a(\theta_i - b_j)\}}{1 + \exp\{a(\theta_i - b_j)\}}$$

with θ_i denoting the ability of student i , and b_j representing the difficulty of item j .

Table B.1: Item diagnostic: The distribution of correct answers across items

a. Baseline	Share of correct answers across items (in %)					
	Minimum	Mean	Median	Maximum	Share 0% ^a	Share 100% ^b
3rd Graders	0.4	24.9	18.3	87.3	0.0	0.0
4th Graders	2.4	30.9	25.5	94.2	0.0	0.0
5th Graders	0.4	34.9	26.6	96.6	0.0	0.0
6th Graders	0.4	38.7	27.4	96.4	0.0	0.0
b. Endline	Minimum	Mean	Median	Maximum	Share 0% ^a	Share 100% ^b
3rd Graders	0.9	34.1	23.5	95.8	0.0	0.0
4th Graders	0.5	36.0	31.0	98.0	0.0	0.0
5th Graders	0.0	38.9	32.3	98.8	1.7	0.0
6th Graders	1.3	42.6	37.2	98.9	0.0	0.0

Notes: The share of correct answers bases on those students that participated in both assessments, and hence constitute the main estimation sample. *a. Share 0%:* This column displays the share of items with zero correct answers. *b. Share 100%:* This column displays the share of items that were answered correctly by all test-takers.

Table B.2: Item diagnostic: The distribution of percentage scores across students

a. Baseline	Percentage score across students (in %)					
	Minimum	Mean	Median	Maximum	Share 0% ^a	Share 100% ^b
3rd Graders	0.0	24.9	21.7	78.3	0.9	0.0
4th Graders	0.0	30.9	28.3	83.3	0.6	0.0
5th Graders	0.0	34.9	35.0	80.0	0.2	0.0
6th Graders	1.7	38.7	38.3	80.0	0.0	0.0
b. Endline	Minimum	Mean	Median	Maximum	Share 0% ^a	Share 100% ^b
3rd Graders	0.0	34.1	33.3	83.3	0.8	0.0
4th Graders	0.0	36.0	35.0	91.7	0.2	0.0
5th Graders	0.0	38.9	38.3	81.7	0.1	0.0
6th Graders	0.0	42.6	40.0	90.0	0.1	0.0

Notes: The distribution of percentage scores bases on those students that participated in both assessments, and hence constitute the main estimation sample. *a. Share 0%:* This column displays the share of students that answered zero questions correctly. *b. Share 100%:* This column displays the share of students that answered all questions correctly.

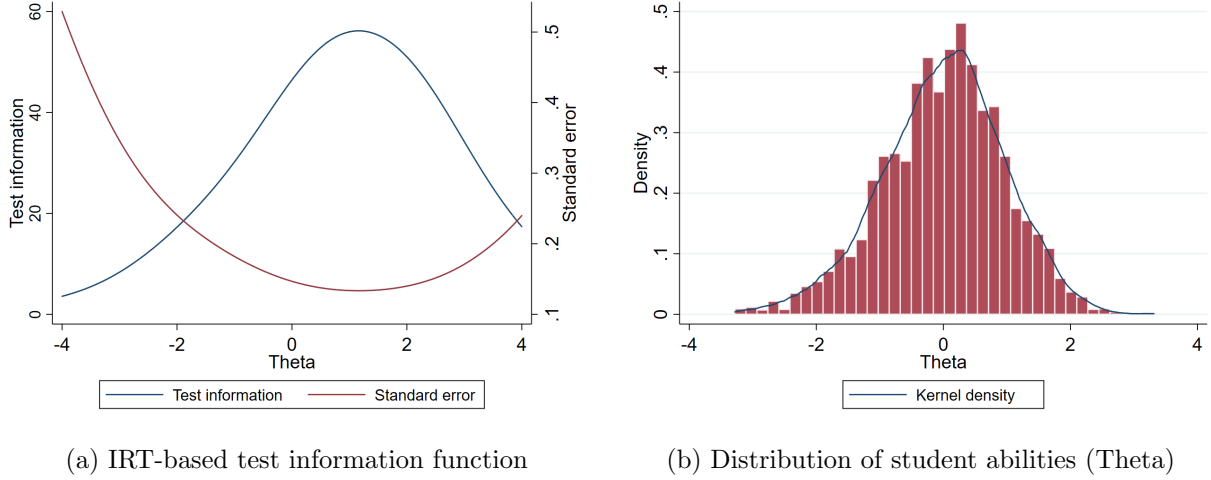


Figure B.3: Test information figure and distribution of students' abilities.

In this so-called *one-parameter model*, the probability that an individual correctly solves a particular item is thus a logistic function of the distance between the ability level of that individual and the difficulty of the item. Ability levels for each person and difficulties for all items can be computed through joint maximum likelihood estimation. IRT has many advantages over classical test theory. It tends to produce more reliable ability estimates, allows to link the scores of different individuals in different tests through overlapping items, and can help to better understand and improve the quality of a test (e.g. de Ayala, 2009).

As illustrated in Figure B.1 a selection of items overlap *(i)* between the baseline and endline assessments and *(ii)* across test booklets of different grades within an assessment wave. This allowed us to project the performance in the baseline and endline assessment onto a common scale through the estimation of an IRT one-parameter model. This procedure yields for every student i two ability estimates, namely one for the baseline assessment, i.e. θ_i^{BL} , and one for the endline assessment, i.e. θ_i^{EL} . The latter serves as outcome variable in the regression models that are labeled with “IRT-Scores”.

Converting IRT-Scores to School Year Equivalents. To allow for an intuitive interpretation, IRT scores can be represented as school year equivalents. For this purpose, we re-scale ability estimates based on between-grade ability differences. To obtain between-grade ability differences, we use the standardized endline IRT score from 651 pure control students (Y^{EL}) and regress it on the students' grade level (GL), i.e.

$$Y_{ic}^{EL} = \alpha + \pi GL_c + \epsilon_{ic}. \quad (\text{B.1})$$

We obtain $\hat{\pi}=0.36$ (p-value<0.01) for the average math ability progression between consecutive grade levels. This means that the average ability difference between third and fourth graders, fourth and fifth graders, and fifth and sixth graders in October 2018 equaled 0.36σ .

The estimated program effects can be divided by this average ability difference between adjacent grades and then be interpreted as proportion of the students' average progress during one school year. Note, however, that ability differences between grades do not only represent what children learn in their regular math classes at school but also reflect age-based cognitive development, learning at home or spillovers from other subjects (e.g. literacy or science).