

# Leveraging Complementarities between Teachers’ Content Knowledge and Pedagogical Supports: Experimental Evidence from India\*

Alejandro J. Ganimian<sup>†</sup>

Harvard University/  
New York University

Isaac M. Mbiti<sup>‡</sup>

University of Notre Dame

Abhilash Mishra<sup>§</sup>

Equitech Futures/  
University of Chicago

September 29, 2025

## Abstract

We present experimental evidence on a program in India that recruited college students (“fellows”) in math and science fields to teach in primary schools for a year. Fellows were younger, less educated, and less experienced than teachers, but they outperformed teachers by  $1.4\sigma$  on a test of content knowledge and pedagogy. They received a brief training, lesson scripts, and instructional coaches. During unannounced visits, fellows were no more likely to go to work or arrive early than control teachers. Yet, during announced observations, they fared  $0.73\sigma$  better on an index of positive instructional practices. After a year, their students scored  $0.34\sigma$  higher in math,  $0.22\sigma$  in science, and  $0.15\sigma$  in language than those taught by regular teachers. By the start of the next year, they still scored  $0.36\sigma$ ,  $0.14\sigma$ , and  $0.08\sigma$  higher in these subjects, respectively. They did not, however, differ on attitudes towards math and science, intelligence and math mindsets, or aspirations to pursue related careers. Evidence of compensatory behavior (more resources) in control classes suggests our estimates may understate true effects.

**JEL codes:** C93, I21, I22, I25

**Keywords:** student achievement, math and science, teaching effectiveness, professional development, coaching, lesson scripts, India.

---

\*We gratefully acknowledge the funding provided by the Abdul Latif Jameel Poverty Action Lab’s (J-PAL) Post-Primary Education initiative and the University of Chicago’s Kevin Xu Initiative on Science, Technology, and Global Development for this study. We thank Priya Pethe and Rahul Panda for making this study possible. We also thank Austin Dempewolff, Sharnic Djaker, Joshua Gilbert, Mustufa Patel, and especially Rashmi Menon, who provided excellent research assistance. We thank Felipe Barrera-Osorio, Andrew Ho, Dan Koretz, and seminar participants at NYU, UVA, Harvard, LSE, AEF, and SREE for useful comments. We registered this study with the AEA Trial Registry (RCT ID: AEARCTR-0002386). Research protocols were approved by the Institutional Review Boards at UVA and the Institute for Financial Management and Research. All views expressed are our own and not of the institutions with which we are affiliated.

<sup>†</sup>Visiting Associate Professor of Education, Harvard Graduate School of Education. Associate Professor of Applied Psychology and Economics, Steinhardt School of Culture, Education, and Human Development, New York University. [alejandroganimian@gse.harvard.edu](mailto:alejandroganimian@gse.harvard.edu).

<sup>‡</sup>Professor of Poverty and Education, Keough School of Global Affairs, University of Notre Dame. [imbiti@nd.edu](mailto:imbiti@nd.edu).

<sup>§</sup>Equitech Futures and University of Chicago. [abhilash@equitechfutures.com](mailto:abhilash@equitechfutures.com).

# 1 Introduction

There is mounting evidence on the individual and social returns to math and science skills. Across very different contexts, individuals with higher mathematical and scientific literacy complete more years of education, are more likely to get a job, and earn more once employed. For example, in the United States, each standard deviation (SD) in math test scores at the end of secondary school is associated with 12% higher earnings (Mulligan, 1999; Lazear, 2003).<sup>1</sup> In India, those who study math and science in high school complete more years of schooling and have 22% higher earnings than those in either business or humanities (Jain et al., 2022).<sup>2</sup> Further, a country’s scores on international assessments of math and science predict its level of innovation and economic growth (Bianchi and Giorcelli, 2020; Hanushek and Kimko, 2000; Barro, 2001; Woessmann, 2003; Jamison et al., 2007; Hanushek and Woessmann, 2008).

Teachers play a key role in the production of the math and science skills of their students. In Peru, each SD on a test of teachers’ content knowledge increased student achievement by about 0.1 SDs (Metzler and Woessmann, 2012).<sup>3</sup> Yet, their capacity on these subjects is low. Secondary-school students hoping to become teachers score below the national mean in math on global tests, and a fourth to a full SD below aspiring engineers (Bruns and Luque, 2014).<sup>4</sup> The traditional approach to this challenge, professional development, has often been ineffective (Popova et al., 2022), possibly because it tries to improve pedagogy without addressing gaps in teacher content knowledge. If the two are complements in the education-production function, focusing only on one of them will not be enough to raise student learning (Mbiti et al., 2019).

We evaluated a novel intervention in India that seeks to improve teacher capacity in math and science by simultaneously alleviating constraints on both content knowledge and pedagogy. The Science Education Initiative (SEI) recruits college students who are majoring in science, technology, engineering, and math (STEM) fields to teach for a year in public primary schools. These “fellows” are selected through a competitive process (involving exams, interviews, and demonstration lessons), they receive pre- and in-service training, scripts for every lesson, and periodic feedback from instructional coaches. They teach, on average, for six hours per week, replacing part of the lesson time allocated to math and science otherwise taught by teachers.<sup>5</sup> They receive a stipend of USD 4 per class, which is about 8% of the median teacher salary.

---

<sup>1</sup>See also Bishop (1989); O’Neill (1990); Grogger and Eide (1995); Murnane et al. (1995); Neal and Johnson (1996); Mulligan (1999); Altonji and Pierret (2001); Murnane et al. (2001) or Hanushek (2002) for a review.

<sup>2</sup>See also Boissiere et al. (1985); Alderman et al. (1996); Glewwe (1996); Angrist and Lavy (1997); Jolliffe (1998); Molina et al. (2020); Behrman et al. (2008) or Hanushek and Woessmann (2008) for a review.

<sup>3</sup>This finding is consistent with correlational evidence across 31 countries (Hanushek et al., 2019)

<sup>4</sup>Gaps in teachers’ content knowledge may be even more pronounced in developing settings. For example, a study in seven Sub-Saharan African countries found two thirds of grade 4 teachers could not solve an algebra problem intended for their students (Bold et al., 2017). See also Schmidt et al. (2007); Tatto et al. (2012).

<sup>5</sup>Study schools did not have timetables, so we cannot calculate the share of instructional time that shifted from teachers to fellows. If we assume a four-hour school day, six hours/week represents 30% of class time.

We conducted a randomized evaluation of this program in the city of Pune, which is in the state of Maharashtra. We randomly assigned 48 schools to a fellow in either grade 5 or 6.<sup>6</sup>

Fellows had less schooling than teachers, no teaching degrees, and had never taught before. Yet, they outscored teachers by 1.4 SDs ( $p < 0.01$ ) on a test of math and science knowledge, optimal approaches to teach concepts and procedures, and understanding of student errors. They also scored 0.73 SDs higher ( $p < 0.01$ ) on an index of positive practices during classroom observations (e.g., allowing students to ask questions or praising and encouraging students).

Consistent with these improvements in instruction, the intervention had moderate-to-large positive impacts on student achievement. Treatment students outperformed their control peers by 0.34 SDs in math and by 0.22 SDs in science ( $p < 0.01$  for both).<sup>7</sup> During the endline, we found that the organization had used items from the baseline test on instructional materials. Effects on student learning, however, cannot be solely explained by coaching on such items. Treatment students outperformed their control peers on items that were introduced at endline (and thus, were not susceptible to coaching).<sup>8</sup> They also fared better on a contemporaneous test administered to a random sub-sample, even if it was less well aligned with the curriculum.<sup>9</sup> And perhaps most compellingly, they still scored 0.36 SDs better in math ( $p < 0.01$ ) and 0.14 SDs in science ( $p < 0.05$ ) at the start of the next school year, two months after the program.<sup>10</sup> Both the endline and follow-up gains in math and science achievement were broad-based, with the treatment distributions stochastically dominating the control distributions in both rounds.

In fact, although fellows only taught math and science, the intervention had spillover effects on language test scores of 0.15 SDs at endline ( $p < 0.01$ ) and 0.08 SDs at follow-up ( $p < 0.05$ ). These effects show impacts in math and science did not come at the expense of other subjects. They also provide further evidence against item coaching explaining effects on achievement (SEI did not cover language, so items for this test could not have been subject to coaching).

We can rule out two key alternative mechanisms of impact other than improved instruction. First, it seems unlikely that the intervention raised test scores by increasing instructional time. Fellows were no more likely than teachers to be present during unannounced visits to schools (partly because attendance among control teachers was already high at an average of 84%). They were more likely than teachers to be found in the classroom (rather than elsewhere at school) conditional on being present (56% v. 4.8%,  $p < 0.01$ ), but during announced classroom observations, they spent a similar share of time teaching (75% v. 78%, non-significant [n.s.]). Treatment students were no more likely than their control peers to attend school, regardless of

---

<sup>6</sup>The 26 grade 5 and 25 grade 6 classes that received the intervention made up the treatment group and the rest served as controls. This randomization strategy, pioneered by Banerjee et al. (2007), ensures that all schools have an incentive to participate in all data-collection rounds. We discuss it in detail in section 2.3.

<sup>7</sup>These impacts are equivalent to gains of 6.7 and 3.8 percentage points in math and science, respectively.

<sup>8</sup>The impacts on these items were 0.04 SDs in math ( $p < 0.01$ ) and 0.02 SDs in science ( $p < 0.01$ ).

<sup>9</sup>The effects on this test were 0.09 SDs in math (non-significant) and 0.39 SDs in science ( $p < 0.01$ ).

<sup>10</sup>Once again, we find impacts on items first introduced in this round of data collection (and thus, not subject to coaching): 0.05 SDs in math ( $p < 0.01$ ) and 0.02 SDs in science ( $p < 0.05$ ).

how we measure student attendance (school records, self reports, and unannounced checks). They were no more prone to seek, or spend more time on, private tuition in math or science (almost a third of control students received tuition in math and over a fourth did in science).

Second, we find no evidence that test-score effects are due to increased student motivation. The program had null effects on indices of students' attitudes towards math and science (0.04 SDs, n.s.), views about the malleability of intelligence or math and science skills (-0.05 SDs, n.s.), and desire to take math and science courses or pursue related careers (0.04 SDs, n.s.).

Our study design allows us to estimate the combined effect of a bundle of four potentially complementary inputs: recruiting fellows, providing them with pre- and in-service training, scripting their instruction, and shifting some math and science lessons from teachers to fellows. We believe this is the parameter of policy interest, as it would make little sense to place college students inside classrooms and expect them to teach well without any training or supports. We also suspect that it is in part this training and supports that lead many college majors in math and science to apply to the program (instead of trying to enter teaching on their own), so removing them could have knock-on effects on the composition of the pool of applicants, rendering an evaluation of the change in the composition of instructors by itself challenging.<sup>11</sup>

We wanted to better understand, however, how scripts complemented fellows' knowledge, both because scripts were arguably the most important intensive support fellows received<sup>12</sup> and because they provided precisely what the fellows lacked: clear guidance on how to teach. We find some evidence for the complementarity between fellow knowledge and script pedagogy. First, fellows used the scripts, which indicates that they are useful: 82% of fellows wrote on the blackboard, 75% asked students questions, and 52% used materials as indicated in scripts during classroom observations. Conversely, only about a third of fellows changed, added to, or rephrased parts of the scripts, and less than a fifth actually excluded parts of the scripts. Second, fellows with higher scores on the test described above were more likely to use scripts, which suggests that scripts were complements (rather than substitutes) to fellows' expertise.<sup>13</sup> These results indicate scripts played an important role in the intervention. Consistent with these findings, we cannot reject the null hypothesis that impacts were the same across sites.<sup>14</sup>

Our first and most important contribution is to experimental evidence on how to improve the quality of instruction in math and science in low- and middle-income countries (LMICs).

---

<sup>11</sup>Further, several prior studies have already evaluated some of these components on their own, including: pre- and in-service teacher training (Yoshikawa et al., 2015; Loyalka et al., 2019; Albornoz et al., 2020), coaching (Cilliers et al., 2019, 2022), and lesson scripting (Piper et al., 2014; Gray-Lobe et al., 2022).

<sup>12</sup>Pre-service training only lasted five weeks, in-service training occurred once a week, and the frequency of coaching was variable. By contrast, fellows had a script for every lesson that they taught during the fellowship.

<sup>13</sup>Fellows who scored above the median were 4.2 pp. less likely than their below-median scoring peers to write in the blackboard and 9.9 pp. less likely to use materials as indicated in the script. Yet, they were 12 more likely to ask students questions and 23 pp. more likely to assign activities following the script. Further, they were 14.1 pp., 23 pp., and 5.1 less likely to change, add to, or exclude parts of the scripts.

<sup>14</sup>We show these results graphically and analytically, following von Hippel and Bellows (2018).

Prior efforts have largely taken gaps in teacher capacity as given and focused on mitigating their effects on students by outsourcing some instructional responsibilities to expert teachers (e.g., through pre-recorded or live audios or videos, Naslund-Hadley et al., 2014; Fabregas, 2019; Navarro-Sola, 2021; Beg et al., 2022; de Barros, 2022) and devolving others to students (e.g., inquiry-based pedagogy, Beuermann et al., 2013; Bando et al., 2019, or peer-to-peer learning, Wachanga and Mwangi 2004; Ajaja and Eravwoke 2010; Berlinski and Busso 2017). Our study shows that it is possible to raise the capacity of the teaching labor force by recruiting individuals with subject-matter expertise and providing them with pedagogical supports. It is also among the first in this literature to include secondary-school grades and to explore the spillover effects of improvements in math and science instruction on language achievement.

Our second main contribution is to global evidence on the effectiveness of lesson scripts. Previous studies have focused on evaluating their impact as a *substitute* for teacher capacity: in settings where most teachers have low levels of education, scripts that tell them what to do at each stage of a lesson can ensure a minimum “floor” of instructional quality (Piper et al., 2018; Albornoz et al., 2020; Romero et al., 2020; Gray-Lobe et al., 2022; Eble et al., 2021).<sup>15</sup> Our study illustrates how scripts can also act as *complements*: for teachers with subject-matter expertise and low pedagogical training and experience, they provide much-needed scaffolding. It is also one of the first to document the extent to which instructors adhere to the different components of scripts and to examine which types of teachers are more likely to use them.

Lastly, our study offers a rare comprehensive account of teacher capacity in an LMIC. Internationally comparable data have traditionally tracked teachers’ education and experience (OECD, 2022; UNESCO, 2022; World Bank, 2023), but such metrics do a poor job predicting teaching effectiveness (Rockoff, 2004; Rivkin et al., 2005; Kane et al., 2008; Kane and Staiger, 2008; Rockoff et al., 2011; Kane and Staiger, 2012; Kane et al., 2013; Araujo et al., 2016). Studies in LMICs have drawn attention to the high teacher absence rates in these settings (Kremer et al., 2005; Chaudhury et al., 2006; Muralidharan et al., 2017), and more recently, to how lesson time is allocated (Abadzi, 2009; Bruns and Luque, 2014; Stallings et al., 2014). Yet, we still have a relatively narrow understanding of teachers’ competence in these contexts.<sup>16</sup> We draw on surveys to describe teachers’ background and credentials, assessments to measure their knowledge, school visits to track their attendance and punctuality, and class observations to document their lesson time allocation and the prevalence of instructional practices.<sup>17</sup>

---

<sup>15</sup>In fact, structured pedagogy more broadly was recently identified as one of the most cost- effective interventions to improve learning outcomes in LMICs (Akyaamong et al., 2023).

<sup>16</sup>International surveys of teachers (e.g., the International Evaluation Association’s Teacher Education and Development Study in Mathematics [TEDS-M] or the Organization for Economic Cooperation and Development’s Teaching and Learning International Survey [TALIS]) offer rich descriptions of teachers’ work, but only a handful of LMICs have ever participated in them (Tatto et al., 2012; OECD, 2019).

<sup>17</sup>For similar work, see Bhattacharjea et al. (2011); Bold et al. (2017); World Bank (2017).

## 2 Experiment

### 2.1 Context

We conducted our study in Pune, the second-largest city in the Indian state of Maharashtra (after its capital city, Mumbai). According to the latest census of India (conducted in 2011), there are 9.4 million people in Pune (NIEPA, 2017), rendering its population size comparable to that of countries such as Belarus, Honduras, and the United Arab Emirates (UN, 2019). Its school system is run by the Pune School Board (*Shikshan Mandal* or PSB) under the Pune Municipal Corporation (PMC). On the latest school year with available data, there were 3,473 primary-only schools (grades 1-5) and 1,941 additional primary schools with upper primary (grades 6-8), with 834,354 students enrolled in both types of primary schools (NIEPA, 2017).<sup>18</sup>

Nearly all primary-school aged children in Pune are enrolled in school: the gross enrollment rates are 110% in primary and 109% in upper primary.<sup>19</sup> Marathi is the language of instruction in 76% of primary-only schools and 55% of those with upper primary; English-medium schools account for 23% and 41% of these types of schools, respectively.<sup>20</sup> About 85% of primary-only schools and 63% of primary schools with upper primary are “government” (i.e., public) schools. There are, on average, 24 students per teacher in primary-only schools and 32 in primary schools with upper primary, and these numbers closely track class sizes (NIEPA, 2017).

The public sector employs most primary-school teachers: 61% of teachers in primary-only and 47% of those in schools with upper primary work in government schools. The vast majority of teachers in these types of schools are “regular” (i.e., tenured): 93% and 84%, respectively; the rest are hired on a renewable contract basis (NIEPA, 2017).

Most primary-school children in Pune lack basic math and science skills. According to a nationally representative student assessment administered by the central government, only 30% of fifth-graders in Pune could use arithmetic for daily situations, 38% could identify equivalent fractions, and 34% could estimate a volume. Results for science were equally discouraging: 23% could identify linkages between terrain, climate, and resources; 29% could group objects, materials, and activities according to properties such as shape, color, and sound; and 45% could estimate spatial quantities in simple standard units (NCERT, 2018).

---

<sup>18</sup>For reference, if Pune were a school system in the U.S., it would rank between the top two largest districts in number of students: New York City and Los Angeles Unified (NCES, 2018).

<sup>19</sup>The gross enrollment rate indicates the number of children enrolled at a given education level *irrespective of age*, divided by the number of children *of age* to attend this level and multiplied by 100. Gross enrollment rates often exceed 100% because the denominator includes both younger and older children. The net enrollment rate (the number of children enrolled at a given level who are of age to attend such level, divided by the total number of children of age for that level) is only reported for upper primary and it is 91% (NIEPA, 2017).

<sup>20</sup>According to anecdotal evidence, a non-trivial share of instruction in these schools is also in Marathi.

## 2.2 Sample

We selected 48 public and *Vidya Niketan* primary schools for this study.<sup>21</sup> We started with all 286 PSB-run primary schools. We excluded 118 schools away from the city center (because their location would have limited the capacity of the non-profit running the intervention to monitor its implementation), 30 Urdu-medium schools (because most fellows did not speak Urdu),<sup>22</sup> 59 schools where the PSB or other non-profits were conducting other programs (because we wanted to estimate of the effects of the intervention on its own), 20 schools with low enrollment (to minimize sampling error), and nine schools that already participated in the intervention (because we wanted to estimate the effects of the first year of the intervention). Our data-analytic sample includes 46 of the 48 sampled schools. Shortly after baseline, we had to drop two schools that could not be matched to any fellows based on their preferences.

## 2.3 Randomization

We randomly assigned the 48 sampled schools to receive the intervention in grades 5 or 6. This process resulted in 26 grade 5 and 25 grade 6 treatment classrooms and 26 grade 5 and 25 grade 6 control classrooms, such that all schools had at least one classroom with a fellow.<sup>23</sup> This randomization strategy, pioneered by Banerjee et al. (2007), seeks to minimize the risk of differential attrition (i.e., schools without any intervention dropping out before the endline).

We also randomly assigned fellows to schools, conditional on their preference set. First, we grouped fellows based on their preferred neighborhood, school shift (morning or afternoon), and medium of instruction (English or Marathi). Then, we ran 17 separate lotteries—one per preference set (e.g., one lottery for neighborhood 1, morning shift, English medium schools).

Table 1 presents summary statistics on students and compares the characteristics and achievement of students between experimental groups. The mean control-group student was 11 years old, which is expected given that the sample is split between grades 5 and 6. Most control students (69%) speak Marathi at home, some (17%) speak Hindi, and few (1%) speak English. Less than two-thirds of them have mothers who completed primary school and more than three-fourths have fathers who reached this level. Nearly all of them (90%) have a TV, but fewer have Internet (43%), a desk (28%), a computer (20%), or their own room (17%),

---

<sup>21</sup>These schools are publicly funded and managed like regular public schools, but they have more autonomy over school-management decisions.

<sup>22</sup>Note, however, that Urdu-medium schools account for about 1% of primary-only schools and 3% of schools with upper primary in Pune (NIEPA, 2017).

<sup>23</sup>Two schools had two grade 5 and two grade 6 classrooms and two other schools had *either* two grade 5 *or* two grade 6 classrooms. In both cases, we assigned all classrooms within the same grade to the same experimental group to prevent contamination, which is particularly likely to occur within the same grade (e.g., a grade 5 fellow in a treatment classroom sharing materials with the grade 5 teacher in a control classroom). All other schools had one treatment and one control classroom (either grade 5 or grade 6).

suggesting that the schools where SEI places its fellows serve relatively low-income families. Yet, two in three of these students attend tuition in math and one in four does so in science.<sup>24</sup>

We find no systematic differences in the characteristics or achievement of students across experimental groups. By chance, treatment students had lower achievement in math ( $p < 0.1$ ), so we estimate student-level effects accounting for baseline test scores.<sup>25</sup> Yet, if this difference were to affect our estimates, it should lead us to underestimate the effect of the intervention.

## 2.4 Attrition

We see no difference in attrition across control and treatment groups at endline or follow-up (Table A.1). We also see no differences in the composition of the sample across experimental groups along students’ age, sex, or baseline test score, with all interactions between treatment assignment and these characteristics statistically being insignificant in a model for attrition. Further, we see no evidence of differential attrition between control and treatment groups on observed baseline characteristics according to a joint test of significance across all interactions.

## 2.5 Intervention

The program we studied is the Science Education Initiative (SEI) undergraduate fellowship.<sup>26</sup> SEI is a Pune-based non-profit organization dedicated to improving math and science learning and this fellowship was its flagship program. Since 2014, it has placed 200 fellows in 110 classes. It recruits, trains, and supports college majors (“fellows”) in science, technology, engineering, and math to teach math and science in schools serving disadvantaged students for one year.

The intervention has four main components. The first one is changing the composition of math and science instructors. Fellows are selected through a four-stage process, including: (a) a test of content knowledge based on the math and science curricula for grades 5 to 10;<sup>27</sup> (b) a brief (5- to 7-minute) demonstration lesson on a topic chosen by each applicant; (c) an interview to assess applicants’ scientific aptitude, leadership skills, and motivation to teach; and (d) a longer (15- to 20-minute) demonstration lesson on a topic chosen by SEI.<sup>28</sup>

Fellows differ from teachers in background, education, experience, and knowledge (Table 2). They are less likely to be female (63% v. 76%) and nearly half their age (21 v. 42 years old). Only one in three fellows has a bachelor’s degree (most are still completing it), compared to

---

<sup>24</sup>Private tuition is common in urban India, even among low-income families (Berry and Mukherjee, 2019).

<sup>25</sup>All impact estimates without accounting for baseline are identical in sign and statistical significance and very similar in magnitude. They are available upon request.

<sup>26</sup>SEI also has a fellowship for college graduates, which we did not evaluate in this study.

<sup>27</sup>These curricula are jointly determined by the National Council of Educational Research and Training and the Board of Education of Maharashtra.

<sup>28</sup>There are no language requirements, but applicants proficient in Marathi are prioritized, given that most primary schools in Pune are Marathi-medium schools.



eight in ten teachers. By the end of the study, the average fellow had completed their first year of teaching, whereas the average teacher had accrued 18 years of teaching experience, five of which had been at their current school, and two of which focused on math or science. Yet, despite their lower education and experience, fellows outperformed teachers by 14 pp. on a test of content knowledge, instructional practices, and student misconceptions.<sup>29</sup>

The second component of the intervention is providing pre- and in-service teacher training. Before they start teaching, fellows must complete a three-week course on pedagogy, classroom management, and the math and science curriculum. Once they begin the fellowship, they are also enrolled in a bachelor’s in education paid for by SEI at a local teacher-training college.

The third component is scripting instruction. Fellows are given lesson scripts for every lesson, which specify the topics to be taught on each day, the materials to be used, the words to be written in the blackboard, the activities to complete, and the questions to ask students.

Most fellows wrote on the board and asked students questions as specified on the script, but only half used the suggested materials and activities, and a fifth to a third made changes (adding to, rephrasing, or excluding parts of the lesson script; see Table A.3 in appendix A).<sup>30</sup> Interestingly, scripts served as a complement to subject-matter expertise: fellows who scored above the median in the test of content knowledge, instructional practices, and understanding of students’ misconceptions mentioned above were *more* likely to follow the scripts.

The fourth component is substituting part of the lesson time otherwise assigned to teachers. Schools in Pune do not have timetables that specify the number of minutes each lesson is supposed to last or even the number of lessons assigned to each subject, so we cannot calculate the share of lesson time in math and science for which fellows substituted regular teachers. Yet, we know fellows were expected to teach three times per week, and that on each occasion, they taught both math and science for 120 minutes (i.e., 24 hours per month) for one year. As compensation, fellows received a stipend of INR 3,000 (USD 47) per month.<sup>31</sup>

### 3 Data

We conducted four data-collection rounds: a baseline at the start of the school year (to check the comparability of experimental groups, increase the precision of our impact estimates, and test for heterogeneous effects); a midline during the year (to identify potential impact mechanisms); an endline at the end of the year (to estimate the impact of the intervention); and a follow-up at the start of the following school year (to check whether effects persisted). We list the dates, instruments administered, and participation rates per round in Table A.2.

---

<sup>29</sup>We describe this test in section 3. Fellows also vary less in their scores (see Figure A.1 in appendix A).

<sup>30</sup>We describe the announced classroom observations on which these figures are based in section 3.

<sup>31</sup>At the time of the study, the median teacher salary in Pune was about INR 50,000 (USD 601) per month.

### 3.1 Instructor surveys and assessments

We administered a survey and an assessment to regular teachers and fellows to describe how the intervention shifted the composition of instructors for some of the math and science lessons. We administered both instruments at endline, so their results may partly reflect the impact of the program on instructors (e.g., lesson scripts may have affected fellows’ pedagogical skills).

We decided which domains to assess based on previous tests for educators (e.g., Hill et al., 2005; Bhattacharjea et al., 2011; Gitomer et al., 2014; Bold et al., 2017; World Bank, 2017). We assessed instructors’ knowledge of the math and science curriculum for grades 5 and 6, the instructional practices they considered optimal to teach concepts and processes, and their understanding of the misconceptions that may explain mistakes that students frequently make. The test had 36 multiple-choice items from both domestic and international instruments.<sup>32</sup>

### 3.2 Student assessments

We administered assessments of math and science to measure the impact of the intervention. We also administered assessments of language to explore potential spillovers from the program. We assessed all three subjects on all data-collection rounds (baseline, endline, and follow-up). We decided which topics and skills to assess in each subject, and the share of items for each topic and skill, based on the frameworks for two international assessments (IEA, 2015, 2017). We tested three topics in math (numbers, geometry and measurement, data) and science (life, physical, and earth sciences) and three skills in both (knowledge, application, and reasoning). In language, we tested vocabulary, grammar, and reading, and whether students could retrieve explicit information, make inferences, and interpret and integrate ideas from written texts. Each test had 30 multiple-choice items, with repeated items across rounds for linking results. We used items from domestic and international assessments and prior impact evaluations from a wide range of difficulty levels to prevent students answering none or all questions correctly.<sup>33</sup>

During the endline, we discovered that some of the items in from our baseline assessments had been used in SEI’s instructional materials, even if we had instructed them not to do so.<sup>34</sup> Specifically, we identified 15 items in the program’s practice tests, but it is possible (and we had no way of verifying whether) more items were featured in other instructional materials. To address this issue in a timely manner, we administered an “audit” test of math and science

---

<sup>32</sup>The items on content knowledge were sampled from the student assessments. Those on instructional practices presented objectives for hypothetical lessons and asked respondents to choose their preferred approach to pursue those goals. Those on student misconceptions presented mistakes that students made and asked respondents to identify the most likely underlying reason for students’ misunderstanding.

<sup>33</sup>Figures A.2 shows the distributions of proportion-correct scores on these assessments. As the figure shows, we see little evidence of floor or ceiling effects.

<sup>34</sup>Due to time constraints at the start of our study, we asked fellows to administer our baseline assessments (we administered all other rounds). We asked them not to review the assessments or take any copies home.

from a contemporaneous evaluation (Gray-Lobe et al., 2022) to a random subset of students.<sup>35</sup> In section 5, we present impacts separately on the items included in the baseline (and could have been compromised) and the ones that were not (and could not have been compromised). We also show effects on language, which fellows did not teach (and had no way of coaching). As we discuss, we see evidence of both item coaching and of learning beyond such coaching.

### 3.3 Student surveys

We administered a short student survey at baseline focusing on background characteristics (e.g., sex and socio-economic status) to test for heterogeneous effects, and a longer one at endline measuring constructs that may be affected by the intervention (e.g., attitudes towards math and science, intelligence and math mindsets, and educational and career aspirations).

### 3.4 Unannounced school visits

We conducted unannounced visits to school during the school year to estimate the impact of the intervention on instructor attendance and punctuality by comparing SEI fellows to PMC teachers in control and treatment classes. We did not announce these visits to minimize the chances that instructors would attend school, or would do so earlier than usual, because of us. We also collected administrative data and counted the number of students in the classroom to estimate the impact of the intervention on student attendance and punctuality.<sup>36</sup>

### 3.5 Announced classroom observations

We conducted announced classroom observations during the school year to estimate the impact of the intervention on instructor lesson-time allocation by comparing SEI fellows to PMC teachers in control and treatment classes.<sup>37</sup> We announced our observations because we were interested in how teachers used lesson time when they attended.<sup>38</sup> We also collected data on whether fellows and control teachers engaged in certain practices during the lesson.

---

<sup>35</sup>These assessments were less closely aligned with the math and science curriculum for grades 5 and 6 in Pune, but we use them to confirm impacts on achievement cannot be solely explained by item coaching.

<sup>36</sup>We supplemented these measures with survey-based measures from endline.

<sup>37</sup>We adapted a classroom-observation protocol that has been widely used in LMICs, including India (Stallings, 1977; Bruns and Luque, 2014; Sankar and Linden, 2014; Stallings et al., 2014; World Bank, 2017).

<sup>38</sup>As stated above, schools in Pune do not have timetables. However, based on conversations with principals, we observed PMC teachers for 30 minutes per subject, which seemed to be the modal duration of a lesson. We observed fellows for their entire 120-minute slot, during which they taught both math and science.

## 4 Empirical strategy

We estimate the intent-to-treat effect of the offer of the intervention by fitting the model:

$$Y_{igs}^t = \alpha_{r(g)} + Y_{igs}^{t=0}\gamma + T_g'\beta + \epsilon_{igs} \quad (1)$$

where  $Y_{igs}^t$  is the outcome of interest for student  $i$  in grade  $g$  and school  $s$  at endline ( $t = 1$ ) or follow-up ( $t = 2$ );  $Y_{igs}^{t=0}$  is a measure of that outcome at baseline (when available);  $r(g)$  is the randomization stratum of grade  $g$  and  $\alpha_{r(g)}$  is the corresponding stratum fixed effect;  $T_g$  is an indicator variable for random assignment to the intervention; and  $\epsilon_{igs}$  is an error term. The parameter of interest is  $\beta$ , which captures the causal effect of the offer of the intervention. We estimate equation (1) by ordinary least-squares regression. We use cluster-robust standard errors to account for within-school correlations across students in outcomes.

We also fit variations of this model in which outcomes are measured at the instructor level (to estimate the impact of the intervention on instructors) and models that interact the intervention indicators with student and teacher covariates (to test for heterogeneous effects).

## 5 Results

### 5.1 Student achievement

After one year, the intervention had moderate-to-large effects on math and science test scores. As Table 3 shows, the endline scores of treatment students were 0.247 SDs higher in math and 0.207 SDs higher in science than those of their control peers ( $p < 0.01$ ). If we account for baseline scores, estimates are slightly higher (0.340 and 0.216 SDs, respectively,  $p < 0.01$ ).<sup>39</sup> We observe improvements in all three content domains (numbers, geometry and measurement, and data display) and two of the three cognitive domains (knowing and applying) assessed in math, and all content domains (life, earth, and physical science) and cognitive domains in science (knowing, applying, and reasoning).

These effects are not entirely explained by fellows coaching students on specific items.<sup>40</sup> As Table 4 shows, effects on items that first appeared at baseline are larger (7.1 pp. in math and 5.4 pp. in science, respectively,  $p < 0.01$ ) than those on those that we introduced at endline (2.8 pp. in math and 1.8 pp. in science,  $p < 0.01$  and  $p < 0.05$ , respectively), and for

---

<sup>39</sup>These effects are equivalent to 6.7 pp. in percent-correct scores in math and 3.8 pp. in science (Table A.8). We obtain very similar results if we scale the endline test scores using a two-parameter logistic IRT model to account for differences in students' latent ability and item characteristics (Table A.9).

<sup>40</sup>As we discuss in section 3.2, right before the endline, we noticed that the nonprofit running the program had included some of our baseline items in its materials.

most specifications, we can reject the null that effects on these two sets of items are equal.<sup>41</sup> Yet, the impacts on the latter are statistically significant, suggesting actual learning gains.

A parallel test that we administered at endline offers further evidence that coaching does not fully account for the impacts.<sup>42</sup> As panel B of Table 3 shows, treatment students also outperformed control students in this “audit” test by 0.046 SDs in math and 0.367 SDs in science. Only the effects on science are statistically significant ( $p < 0.01$ ), but this may be because the test may have overlapped more with the material taught in that subject.

The fact that treatment students still fared better than their control peers at the start of the following school year confirms that the intervention impacted learning beyond coaching. As panel C of Table 3 shows, treatment students scored 0.283 SDs higher than control students in math and 0.120 SDs in science or 0.356 SDs and 0.142 SDs accounting for baseline scores (in both specifications,  $p < 0.01$  and  $p < 0.05$ , respectively).

The intervention also had spillover effects on language. As Table 3 indicates, the treatment group outperformed the control group by 0.132 SDs or 0.151 SDs once we account for baseline scores ( $p < 0.05$  and  $p < 0.01$ , respectively), even if fellows did not teach this subject.<sup>43</sup> We see effects in all three content domains (vocabulary, grammar, and reading) and two of the three cognitive domains (retrieving information and making inferences) assessed in language. In fact, treatment students still outperformed their control counterparts in language at the start of the following school year by 0.059 SDs, but the effect is only statistically significant if we account for baseline (0.081 SDs,  $p < 0.05$ ).

The intervention improved test scores across all levels of the achievement distribution. Quantile treatment effect plots show that the treatment distributions first-order stochastically dominate the control distributions on the endline, audit, and follow-up tests, suggesting that the intervention led to broad-based gains (Figure A.3). Non-parametric estimates of average treatment effects at each percentile of the baseline scores also show large positive impacts for all three rounds of data collection across the full range of baseline achievement (Figure A.4). Consistent with these results, we find no heterogeneous effects by baseline scores (Table A.10).

We find some evidence that the intervention was more effective for students with higher socio-economic status. In our estimation of endline effects, the interaction term between the treatment and the first principal component of a principal-component analysis of home assets (SES index) is positive and statistically significant for all subjects, and the p-value of the sum of the main and interaction effects of this index is statistically significant for math and science (Table A.11, panel A, cols. 1, 3, and 5). Once we account for baseline performance, however,

---

<sup>41</sup>Remarkably, the differences between effects on repeated and non-repeated items in math and science remain stable at followup, suggesting the effects of coaching persisted (Table 5).

<sup>42</sup>As we mention in section 3.2, this test drew on an entirely different item bank and was used in a concurrent impact evaluation (Gray-Lobe et al., 2022).

<sup>43</sup>This effect translates into 3.5 pp. (Table A.8). IRT-scaled impacts are very close (Table A.9). Effects on language, however, are no longer statistically significant at follow-up.

the p-value of the sum is no longer statistically significant for any subject (cols. 2, 4, and 6). Further, we do not see any evidence of heterogeneity by SES in the endline audit (panel B). The interaction between the treatment and SES index is positive and statistically significant for the follow-up (panel C, cols. 1, 3, and 5), but this interaction and the p-value of the sum only remain statistically significant for science after we control for baseline (cols. 2, 4, and 6).

We do not find evidence of heterogeneous effects by other student or teacher characteristics. The interaction between treatment and an indicator for female students is positive for all subjects at endline, but it is statistically significant only for math once we account for baseline (Table A.12, panel A) and statistically insignificant in the other two rounds (panels B and C). The interaction between treatment and an indicator for students from scheduled castes or tribes is negative but statistically insignificant for nearly all subjects and rounds (Table A.13). Lastly, students do not benefit more from the intervention when they have an instructor of the same sex (Table A.14) or one who scores higher on the written assessment (Table A.15).<sup>44</sup>

## 6 Mechanisms

### 6.1 Instructor pedagogical practices

The intervention seems to have impacted test scores mainly by improving the pedagogical practices instructors used—specifically, by increasing the frequency of positive practices. As panel B of Table 6 shows, control teachers used closed- and open-ended questions (63%), asked students to explain their answers (47%), corrected wrong answers (72%), allowed students to ask questions (39%), provided individual help (68%), assigned homework (48%), and praised or encouraged students (76%). Fellows were even more prone to pursue these practices. In fact, they scored 0.73 SDs higher on a composite index of all positive practices ( $p < 0.01$ ).

The introduction of fellows did not affect the frequency of negative practices, which were already rare. As panel A shows, only 19% of control teachers taught from the same spot and 5% or fewer remained sitting down, used their phone, got upset at incorrect answers, or was aggressive towards students. Fellows were 13 pp. less likely than control teachers to teach from the same spot ( $p < 0.05$ ), but they were also 12 pp. more likely to get upset at incorrect answers ( $p < 0.1$ ), and they did not differ from control teachers on a composite index of all negative practices, or reduce treatment teachers’ engagement in these practices.

---

<sup>44</sup>In fact, if we estimate the effect of the program only among high-scoring instructors (i.e., PMC teachers and SEI fellows), the point estimates (Table A.16) resemble the average effects of the intervention (Table 3), suggesting that effects are not explained solely by a selection effect.

## 7 Robustness checks

### 7.1 Instructor attendance, punctuality, and lesson-time allocation

The intervention did not raise students’ achievement simply by increasing the frequency with which they saw their instructors. Fellows were no more likely than control teachers to go to work or arrive on time—partly, because the latter were already doing so at fairly high rates. As panel A of Table 7 shows, 84% of control teachers were present during unannounced visits, compared to 72% of fellows, and this difference was not statistically significant. In fact, fellows were 30 pp. *less* likely to arrive on time than control teachers ( $p < 0.01$ ).<sup>45</sup>

The introduction of fellows did not increase the attendance or punctuality of teachers in treatment classes either. Treatment teachers were slightly more likely than their control counterparts to be present (88% v. 84%) and less likely to arrive on time (71% v. 74%) during the unannounced visits, but neither difference was statistically significant.

Conditional on being at school, however, fellows were far more likely to be in their class. As panel B shows, only 4.8% of control teachers were in their classroom, compared to 56% of fellows, a difference of almost 51 pp. ( $p < 0.01$ ). Fellows did not impact the likelihood of treatment teachers of being in their classroom, which was nearly identical (4.7%).

The intervention did not raise achievement merely by increasing the share of lesson time devoted to instruction either—largely, because control teachers were already teaching for most of their lessons. As panel A of Table 8 shows, the average control teacher devoted 78% of their lesson to instruction during announced classroom observations, compared to 75% for the average fellow, and the difference between these groups was not statistically significant.<sup>46</sup>

In fact, fellows allocated instructional time similarly to teachers in control classes. Control teachers spent most of their lesson time lecturing and explaining (34%), asking and answering questions (18%), and assigning students classwork (13%; Table A.4 in appendix A). The corresponding figures for fellows were nearly identical (34%, 15%, and 16%, respectively). In fact, fellows also resembled teachers in their use of class management and off task time (Tables A.5-A.6), suggesting that time allocation was not a primary mechanism of impact.

Treatment teachers, however, rarely engaged with the fellows while they were teaching. The typical teacher in this group spent 8% of a fellow’s lesson teaching, 21% managing the class, and 69% being off task—61 pp. more than the typical control peer ( $p < 0.01$ ).

---

<sup>45</sup>Unlike teachers, fellows were only expected to be at school when teaching. Thus, their lower attendance rate may simply reflect that they were expected to be at school less frequently than teachers.

<sup>46</sup>As discussed in section 2.5, while PMC teachers teach either math or science in 30-minute lessons, SEI fellows teach both math and science combined in 120-minute lessons. Therefore, while we compare these groups focusing on the proportion of lesson time devoted to each type of activity, we also report the number of minutes devoted to each category in panel B of Table 8.

## 7.2 Student attitudes, mindsets, and aspirations

The intervention did not improve students’ achievement by increasing their motivation to learn math and science. Roughly eight out of ten control students liked learning math and science and almost as many found them useful (Table A.17, panel A). It did not decrease students’ negative feelings towards these subjects either. A fourth of control students felt nervous about them, a third wished they did not have to study these topics, and a fifth gave up when the material is difficult, but fellows did not change this pattern. In fact, treatment and control score perform similarly on a composite index of these measures.

Fixed mindsets were common among control students and fellows did not dispel them. About half of control students believed intelligence and math skills cannot be developed, and four in ten thought boys are both more intelligent and skilled at math. The intervention had a negative but statistically insignificant effect on all of these beliefs (Table A.17, panel B).

Lastly, the intervention did not make students more likely to want to stay in school, study math and science, or pursue a STEM-related job. Over two thirds of control students planned to pursue post-secondary education, three fourths wanted to study a STEM subject in high school, and a third aspired to a STEM-related job. The intervention had statistically insignificant effects on these metrics (Table A.17, panel C).<sup>47</sup>

Consistent with these results, we find no evidence that impacts on achievement are driven by role-model effects. As we would expect any potential role-model effects to be stronger when the fellow “looks like” the student, we test for heterogeneity in student test-scores by student-instructor sex match in Table A.14. However, we do not find any statistically significant effects on the interaction between treatment and sex match.

## 7.3 Student attendance

The intervention did not impact students’ achievement by increasing their attendance to school. We measured student attendance in three different ways: by calculating the share of present students observed during unannounced visits, digitizing school attendance records, and asking students how often they were late to or missed school. None of these approaches suggest the intervention impacted student attendance (Table A.22).

## 7.4 Student demand for tuition

The intervention did not impact achievement by making students more likely to seek private tuition. Between a fifth and a third of students in our sample receive tuition in one subject.

---

<sup>47</sup>We did not find any evidence of heterogeneous effects on any of these sets of outcomes by students’ sex, caste, socio-economic status, or baseline achievement (Tables A.18-A.21).



The intervention did not impact their propensity to seek tuition in a target (math and science) or non-target subject (English and Urdu) or the amount of tuition received (Table A.23).

## 7.5 Miscellaneous

Given the lower performance of teachers relative to fellows in our assessments, and the positive relationship between teacher content knowledge and student achievement, it is possible that the results reflect the presence of more knowledgeable instructors (i.e., fellows) rather than any other facet of SEI’s program. To ameliorate this concern, we restrict our analysis to the set of teachers and fellows where the distribution of assessment scores overlap (see Figure A.1). This restriction essentially drops the lowest scoring PMC teachers (and their students) from the analysis. In Table A.16 we find that, across all specifications, the treatment effects are similar in magnitude to our full-sample estimates, suggesting that the results are not driven by differences in the teaching assessment.

One important feature of scripted instruction is that if it is implemented with high fidelity, it can reduce the heterogeneity in teaching quality across schools (Gray-Lobe et al., 2022). Additionally, if scripts are adhered to, an individual teacher’s content knowledge will be less predictive of student outcomes. We test these implications in Figure A.5 and Table A.15.

Following von Hippel and Bellows (2018), we construct “caterpillar plots” of the school-specific (ITT) treatment effects and conduct a joint test against the null hypothesis of homogeneity. To ameliorate concerns about multiple hypothesis testing, we pool science and math scores into a composite test and the von Hippel and Bellows (2018) procedure computes both standard and Bonferroni confidence intervals to account for the number of school-specific effects. For the endline and the follow-up we fail to reject the null hypothesis of homogeneity, but we reject it for the audit. Across the three sets of student tests, we always fail to reject the null for science, but we reject the null in the math audit and the math follow-up (see Figures A.5-A.7). Overall, these results provide some suggestive evidence on the importance of scripted instruction, especially in science.

We examine the heterogeneity in student outcomes by teacher and fellow content score. In our empirical specification we use the assessment score of fellows in treatment classrooms, and the teacher’s score in control classrooms to better reflect who the actual instructor was. The results are reported in Table A.15. For the endline, teacher and fellow scores are associated with higher student test scores, although this effect is imprecisely estimated in three (out of six) specifications. Further, we also find that treatment dampens this effect, but this dampening is imprecisely estimated in five (out of six) specifications. While the patterns observed in the endline data are broadly consistent with the scripted instruction prediction, we do not find similar patterns in the endline audit or the follow-up.

Although these results provide some suggestive evidence consistent with scripting, especially in science, they also potentially shed light on the implementation fidelity. In Table A.3 we use data from announced classroom observations of fellows and document their script adherence on multiple dimensions. Overall, we find moderate levels of adherence. For instance, three-quarters of fellows asked students the questions in the script, about half of fellows used materials from the script, and a third added or expanded to the script. The patterns in Figure A.5 and Table A.15 are consistent with these moderate (to low) levels of script adherence.

## 7.6 Compensation for control classes

Whenever resources (e.g., fellows) are (randomly) allocated to some classes and not others, principal may try to compensate the non-selected classes with other resources (e.g., materials). To explore this possibility, we leveraged the information on class materials collected during the announced observations to compare the availability of materials across experimental groups. Control classrooms tend to have more materials than treatment classrooms, and they are statistically significantly more likely to have textbooks for teachers (52 pp.) and students (46 pp.;  $p < 0.01$  in both cases) and math or science equipment (14 pp.,  $p < 0.05$ ; Table A.24). These results suggest that our estimates may understate the true impact of the intervention (i.e., the impact if all resources were equally distributed across control and treatment classes).

## 8 Conclusion

Developing math and science skills can help individuals improve their economic prospects. Yet, there is very little evidence on how to improve instructional quality in these subjects—especially, in contexts with low teacher capacity, which are more prevalent among LMICs.

We present experimental evidence on a novel approach to raise teacher capacity in math and science: hiring college students in these fields and providing them pedagogical support. We find that the introduction of these fellows results in moderate-to-large gains in student achievement in both of the subjects that they teach as well as in other subjects, and that those gains persist over time, influencing their students’ preparation level for the next school year. We also identify potential mechanisms of impact (e.g., increase in instructors’ likelihood to be in their class, increase in their use of positive pedagogical practices) and rule out confounders (e.g., increase in instructor or student attendance or in student demand for private tuition).

These findings shed new light on the importance of teachers’ subject-matter expertise. Until recently, its value had been largely inferred from correlations with students’ performance on standardized tests (Hill et al., 2005; Santibanez, 2006; Rockoff et al., 2011; Kane et al., 2013; Gitomer et al., 2014; Bold et al., 2017; Cruz-Aguayo et al., 2017; Hanushek et al., 2019). We only know one quasi-experimental study on this question (Metzler and Woessmann, 2012).

We provide experimental evidence that hiring instructors with high levels of content knowledge to replace teachers with lower levels for part of the school week raises student achievement.<sup>48</sup> The fellows in our study received considerable pedagogical support (e.g., pre-service training, lesson scripts, instructional coaches), so we cannot attribute their impact to their knowledge. Yet, if this intervention were to be taken to scale, fellows would likely receive such supports, so the parameter of policy interest is the estimated effect of combining fellows and supports.

Our results also add nuance to education-policy debates on the merits on lesson scripts. Opponents often argue that any gains from standardizing instruction must be weighted against the ensuing losses to teachers' professional autonomy (Valencia et al., 2006; Dresser, 2012). Yet, we found that the vast majority of fellows adhered to lesson scripts, suggesting that they may be a useful complement even for individuals with high levels of subject-matter expertise. In fact, the fellows who performed best in the assessment of content and pedagogical knowledge were less likely to omit, change, or add to the material in such scripts, implying that the fellows who were most likely to exercise discretion were precisely those least well equipped to do so. Future research should explore the extent to which this pattern is observed in other settings, and if so, how scripts can best encourage instructors to leverage their expertise appropriately.

Having demonstrated that this approach works, a logical next question is whether it could be implemented with comparable levels of fidelity and improve student achievement at scale. Of the reasons identified by List (2022) for why interventions that are successful in efficacy trials fail to sustain gains at scale, we are least concerned with the credibility of the evidence. Our randomized evaluation offers one of the most rigorous studies to date demonstrating that hiring individuals with subject-matter expertise can improve student achievement in a LMIC. The fact that we relied on a convenience sample, however, suggests that this intervention has more chances of scaling in urban, medium-sized, English-language public and charter schools. The motivation and capacity of the implementing organization seems much harder to replicate. SEI not only recruited, selected, and supported fellows; it also established partnership with the local governments to place them in schools and with teacher-training colleges for certification; and it created detailed lesson scripts and hired instructional coaches for pedagogical support.<sup>49</sup>

Yet, perhaps the most important challenge to take this intervention to scale is to identify enough college students in math and science who are willing and able to become teachers. According to data for the 2020-2021 year, there were 32.7 million undergraduate students in India, 16% of whom studied science and 12% of whom were studying engineering (MoE, 2022). In fact, science and engineering were the second and fourth most popular majors, respectively. By comparison, there were 247,236 students in teacher-training colleges, which means that

---

<sup>48</sup>While fellows were expected to teach alongside regular teachers, as we show in section 5, the latter were typically not engaged in instructional tasks while the former were teaching.

<sup>49</sup>Although this combination of know-how may be rare on a global scale, some of India's most successful education organizations already engage in similar tasks (e.g., Pratham, Educational Initiatives, Central Square Foundation, and their partner institutions), so such talent may be easier to find within the country.

there were 21 science majors and 16 engineer majors per aspiring teacher across the country, suggesting that there is not a shortage of college students who can teach math and science.<sup>50</sup> Female enrollment in science and engineering is relatively high (52% and 29%, respectively), so getting more undergraduate students in these fields to enter teaching would not necessarily alter the sex breakdown in the profession. The question is whether they would *want* to teach. We do not have data on the career preferences of science and engineering majors across India. The only comparable program (Teach for India) has recruited 4,000 fellows since 2008,<sup>51</sup> but unlike the SEI fellowship, this program does not intend to keep its graduates in the profession. These figures, however, suggest that there is scope to increase the number of undergraduates in math and science who could enter teaching and raise the achievement of Indian students.

---

<sup>50</sup>This is true even if we include students in education, who account for 5.3% of undergraduates nationwide.

<sup>51</sup>See Teach for India's website: <https://www.teachforindia.org/fellowship>.

# References

- Abadzi, H. (2009). Instructional time loss in developing countries: Concepts, measurement, and implications. *The World Bank Research Observer* 24(2), 267–290.
- Ajaja, O. P. and O. U. Eravwoke (2010). Effects of cooperative learning strategy on junior secondary school students achievement in integrated science. *The Electronic Journal for Research in Science & Mathematics Education* 14(1).
- Akyeampong, T., T. Andrabi, A. Banerjee, R. Banerji, S. Dynarski, R. Glennerster, S. Grantham-McGregor, K. Muralidharan, B. Piper, S. Ruto, J. Saavedra, S. Schmelkes, and H. Yoshikawa (2023). *2023 cost-effective approaches to improve global learning. What does recent evidence tell us are “smart buys” for improving learning in low- and middle-income countries?* London, UK; Washington, DC; New York, NY: Foreign, Commonwealth & Development Office (FCDO), World Bank, United Nations International Children’s Emergency Fund (UNICEF), United States Agency for International Development (USAID).
- Albornoz, F., M. V. Anuati, M. Furman, M. Luzuriaga, M. E. Podestá, and I. Taylor (2020). Training to teach science: Experimental evidence from Argentina. *World Bank Economic Review* 34(2), 393–417.
- Alderman, H., J. R. Behrman, D. R. Ross, and R. H. Sabot (1996). The returns to endogenous human capital in Pakistan’s rural wage labour market. *Oxford Bulletin of Economics and Statistics* 58(1), 29–55.
- Altonji, J. G. and C. R. Pierret (2001). Employer learning and statistical discrimination. *Quarterly Journal of Economics* 116(1), 313–350.
- Angrist, J. D. and V. Lavy (1997). The effect of a change in language of instruction on the returns to schooling in Morocco. *Journal of Labor Economics* 15(1), S48–S76.
- Araujo, M. C., P. M. Carneiro, Y. Cruz-Aguayo, and N. Schady (2016). Teacher quality and learning outcomes in kindergarten. *Quarterly Journal of Economics* 131(3), 1415–1453.
- Bando, R., E. Naslund-Hadley, and P. Gertler (2019). Effect of inquiry and problem based pedagogy on learning: Evidence from 10 field experiments in four countries. (NBER Working Paper No. 26280). Cambridge, MA: National Bureau of Economic Research (NBER).
- Banerjee, A. V., S. Cole, E. Duflo, and L. L. Linden (2007). Remedying education: Evidence from two randomized experiments in India. *Quarterly Journal of Economics* 122(3), 1235–1264.
- Barro, R. J. (2001). Human capital and growth. *American Economic Review* 91(2), 12–17.
- Beg, S. A., A. M. Lucas, W. Halim, and U. Saif (2022). Engaging teachers with technology increased achievement, bypassing teachers did not. *American Economic Journal: Economic Policy* 14(2), 61–90.

- Behrman, J. R., D. R. Ross, and R. H. Sabot (2008). Improving quality versus increasing the quantity of schooling: Estimates of rates of return from rural Pakistan. *Journal of Development Economics* 85(1-2), 94–104.
- Berlinski, S. and M. Busso (2017). Challenges in educational reform: An experiment on active learning in mathematics. *Economic Letters* 156, 172–175.
- Berry, J. and P. Mukherjee (2019). Pricing of private education in urban India: Demand, use, and impact. *Unpublished manuscript*. Athens, GA: University of Georgia.
- Beuermann, D. W., E. Naslund-Hadley, I. J. Ruprah, and J. Thompson (2013). The pedagogy of science and environment: Experimental evidence from Peru. *The Journal of Development Studies* 49(5), 719–736.
- Bhattacharjea, S., W. Wadhwa, and R. Banerji (2011). *Inside primary schools: A study of teaching and learning in rural India*. Mumbai, Maharashtra: Pratham.
- Bianchi, N. and M. Giorcelli (2020). Scientific education and innovation: From technical diplomas to university stem degrees. *Journal of the European Economic Association* 18(5), 2608–2646.
- Bishop, J. H. (1989). Is the test score decline responsible for the productivity growth decline? *American Economic Review* 79(1), 178–197.
- Boissiere, M., J. B. Knight, and R. H. Sabot (1985). Earnings, schooling, ability, and cognitive skills. *American Economic Review* 75(5), 1016–1030.
- Bold, T., D. Filmer, G. Martin, E. Molina, B. Stacy, C. Rockmore, J. Svensson, and W. Wane (2017). Enrollment without learning: Teacher effort, knowledge, and skill in primary schools in Africa. *Journal of Economic Perspectives* 31(4), 185–204.
- Bruns, B. and J. Luque (2014). *Great teachers: How to raise student learning in Latin America and the Caribbean*. Washington, DC: The World Bank.
- Chaudhury, N., J. Hammer, M. Kremer, K. Muralidharan, and F. H. Rogers (2006). Missing in action: Teacher and health worker absence in developing countries. *The Journal of Economic Perspectives* 20(1), 91–116.
- Cilliers, J., B. Fleisch, J. Kotze, N. Mohohlwane, S. Taylor, and T. Thulare (2022). Can virtual replace in-person coaching? Experimental evidence on teacher professional development and student learning in South Africa. *Journal of Development Economics* 155, 102815.
- Cilliers, J., B. Fleisch, C. Prinsloo, and S. Taylor (2019). How to improve teaching practice? an experimental comparison of centralized training and in-classroom coaching. *Journal of Human Resources*, 0618–9538R1.
- Cruz-Aguayo, Y., P. Ibararán, and N. Schady (2017). Do tests applied to teachers predict their effectiveness? *Economics Letters* 159, 108–111.
- de Barros, A. (2022). Explaining the productivity paradox: Experimental evidence from educational technology. *Unpublished manuscript*. Cambridge, MA: Massachusetts Institute of Technology (MIT).

- Dresser, R. (2012). The impact of scripted literacy instruction on teachers and students. *Issues in Teacher Education* 21(1), 71–87.
- Eble, A., C. Frost, A. Camara, B. Bouy, M. Bah, M. Sivaraman, P.-T. J. Hsieh, C. Jayanty, T. Brady, and P. Gawron (2021). How much can we remedy very low learning levels in rural parts of low-income countries? impact and generalizability of a multi-pronged para-teacher intervention from a cluster-randomized trial in the gambia. *Journal of Development Economics* 148, 102539.
- Fabregas, R. (2019). Broadcasting human capital? The long-term effects of Mexico’s Telesecundarias. *Unpublished manuscript*. Austin, TX: LBJ School of Public Affairs, The University of Texas at Austin.
- Gitomer, D. H., G. Phelps, B. H. Weren, H. Howell, and A. J. Croft (2014). *Evidence on the validity of content knowledge for teaching assessments*. San Francisco, CA: Jossey-Bass.
- Glewwe, P. (1996). The relevance of standard estimates of rates of return to schooling for education policy: A critical assessment. *Journal of Development Economics* 51(2), 267–290.
- Gray-Lobe, G., A. Keats, M. Kremer, I. Mbiti, and O. W. Ozier (2022). Can education be standardized? Evidence from Kenya. (Working Paper No. 2022-68). Chicago, IL: Becker Friedman Institute For Economics.
- Grogger, J. and E. Eide (1995). Changes in college skills and the rise in the college wage premium. *Journal of Human Resources*, 280–310.
- Hanushek, E. A. (2002). *Publicly provided education*, pp. 2045–2141. Amsterdam, the Netherlands; London, UK; and New York, NY: Elsevier Science, North Holland.
- Hanushek, E. A. and D. D. Kimko (2000). Schooling, labor-force quality, and the growth of nations. *American Economic Review* 90(5), 1184–1208.
- Hanushek, E. A., M. Piopiunik, and S. Wiederhold (2019). The value of smarter teachers: International evidence on teacher cognitive skills and student performance. *Journal of Human Resources* 54(4), 857–899.
- Hanushek, E. A. and L. Woessmann (2008). The role of cognitive skills in economic development. *Journal of Economic Literature* 46(3), 607–668.
- Hill, H. C., B. Rowan, and D. L. Ball (2005). Effects of teachers’ mathematical knowledge for teaching on student achievement. *American Educational Research Journal* 42(2), 371–406.
- IEA (2015). PIRLS 2016: Assessment framework. TIMSS & PIRLS International Study Center. Lynch School of Education, Boston College & International Association for the Evaluation of Educational Achievement (IEA).
- IEA (2017). TIMSS 2019: Assessment frameworks. Edited by Mullis, I. V. S. & Martin, M. O. TIMSS & PIRLS International Study Center. Lynch School of Education, Boston College & International Association for the Evaluation of Educational Achievement (IEA).

- Jain, T., A. Mukhopadhyay, N. Prakash, and R. Rakesh (2022). Science education and labor market outcomes in a developing economy. *Economic Inquiry* 60(2), 741–763.
- Jamison, E. A., D. T. Jamison, and E. A. Hanushek (2007). The effects of education quality on income growth and mortality decline. *Economics of Education Review* 26(6), 771–788.
- Jolliffe, D. (1998). Skills, schooling, and household income in Ghana. *The World Bank Economic Review* 12(1), 81–104.
- Kane, T. J., D. F. McCaffrey, T. Miller, and D. O. Staiger (2013). Have we identified effective teachers? Validating measures of effective teaching using random assignment. *Measures of Effective Teaching Project*. Seattle, WA: Bill and Melinda Gates Foundation.
- Kane, T. J., J. E. Rockoff, and D. O. Staiger (2008). What does certification tell us about teacher effectiveness? Evidence from New York City. *Economics of Education Review* 27(6), 615–631.
- Kane, T. J. and D. O. Staiger (2008). Estimating teacher impacts on student achievement: An experimental evaluation. (NBER Working Paper No. 14607). Cambridge, MA: National Bureau of Economic Research (NBER).
- Kane, T. J. and D. O. Staiger (2012). Gathering feedback for teaching: Combining high-quality observations with student surveys and achievement gains. Seattle, WA: Bill and Melinda Gates Foundation.
- Kremer, M., N. Chaudhury, F. H. Rogers, K. Muralidharan, and J. Hammer (2005). Teacher absence in India: A snapshot. *Journal of the European Economic Association* 3(2-3), 658–667.
- Lazear, E. P. (2003). Teacher incentives. *Swedish Economic Policy Review* 10(3), 179–214.
- List, J. A. (2022). *The voltage effect: How to make good ideas great and great ideas scale*. New York, NY: Currency.
- Loyalka, P., A. Popova, G. Li, and Z. Shi (2019). Does teacher training actually work? Evidence from a large-scale randomized evaluation of a national teacher training program. *American Economic Journal: Applied Economics* 11(3), 128–154.
- Mbiti, I. M., K. Muralidharan, M. Romero, Y. Schipper, C. Manda, and R. Rajani (2019). Inputs, incentives, and complementarities in education: Experimental evidence from Tanzania. *The Quarterly Journal of Economics* 134(3), 1627–1673.
- Metzler, J. and L. 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.
- MoE (2022). All India survey on higher education 2020-2021. New Delhi, India: Department of Higher Education, Ministry of Education, Government of India.
- Molina, E., S. F. Fatima, A. D. Ho, C. Melo, T. Wilichowski, and A. Pushparatnam (2020). Measuring the quality of teaching practices in primary schools: Assessing the validity of the Teach observation tool in Punjab, Pakistan. *Teaching and Teacher Education* 96, 103171.



- Mulligan, C. B. (1999). Galton versus the human capital approach to inheritance. *Journal of Political Economy* 107(SG), S184–S224.
- Muralidharan, K., J. Das, A. Holla, and A. Mohpal (2017). The fiscal cost of weak governance: Evidence from teacher absence in India. *Journal of Public Economics* 145(C), 116–135.
- Murnane, R. J., J. B. Willett, M. J. Braatz, and Y. Duhaldeborde (2001). Do different dimensions of male high school students’ skills predict labor market success a decade later? *Economics of Education Review* 20(4), 311–320.
- Murnane, R. J., J. B. Willett, and F. Levy (1995). The growing importance of cognitive skills in wage determination. *Review of Economics and Statistics*, 251–266.
- Naslund-Hadley, E., S. W. Parker, and J. M. Hernandez-Agramonte (2014). Fostering early math comprehension: Experimental evidence from Paraguay. *Global Education Review* 1, 135–154.
- Navarro-Sola, L. (2021). Secondary schools with televised lessons: The labor market returns of the Mexican Telesecundaria. *Unpublished manuscript*. Stockholm, Sweden: Institute for International Economic Studies, Stockholm University.
- NCERT (2018). National achievement survey 2017: District report cards. New Delhi, India: National Council of Educational Research and Training (NCERT).
- NCES (2018). Digest of education statistics 2018. U.S. Department of Education, National Center for Education Statistics, Common Core of Data. URL: [https://nces.ed.gov/programs/digest/d18/tables/dt18\\_215.30.asp?current=yes](https://nces.ed.gov/programs/digest/d18/tables/dt18_215.30.asp?current=yes). Last accessed: March 7, 2020.
- Neal, D. A. and W. R. Johnson (1996). The role of premarket factors in Black-White wage differences. *Journal of Political Economy* 104(5), 869–895.
- NIEPA (2017). Elementary education in India: Where do we stand? (District report cards 2016-2017, Vol. I). New Delhi, India: National Institute of Educational Planning and Administration (NIEPA).
- OECD (2019). TALIS 2018 results (Vol. I): Teachers and school leaders as lifelong learners. Paris, France: Organisation for Economic Co-operation and Deveopment (OECD).
- OECD (2022). Education at a Glance 2022: OECD Indicators. Paris, France: Organisation for Economic Co-operation and Deveopment.
- O’Neill, J. (1990). The role of human capital in earnings differences between black and white men. *Journal of Economic Perspectives* 4(4), 25–45.
- Piper, B., S. S. Zuilkowski, M. Dubeck, E. Jepkemei, and S. J. King (2018). Identifying the essential ingredients to literacy and numeracy improvement: Teacher professional development and coaching, student textbooks, and structured teachers’ guides. *World Development* 106, 324–336.

- Piper, B., S. S. Zuilkowski, and A. Mugenda (2014). Improving reading outcomes in Kenya: First-year effects of the PRIMR Initiative. *International Journal of Educational Development* 37, 11–21.
- Popova, A., D. K. Evans, M. E. Breeding, and V. Arancibia (2022). Teacher professional development around the world: The gap between evidence and practice. *The World Bank Research Observer* 37(1), 107–136.
- Rivkin, S. G., E. A. Hanushek, and J. F. Kain (2005). Teachers, schools, and academic achievement. *Econometrica*, 417–458.
- Rockoff, J. E. (2004). The impact of individual teachers on student achievement: Evidence from panel data. *American Economic Review*, 247–252.
- Rockoff, J. E., B. A. Jacob, T. J. Kane, and D. O. Staiger (2011). Can you recognize an effective teacher when you recruit one? *Education Finance and Policy* 6(1), 43–74.
- Romero, M., J. Sandefur, and W. A. Sandholtz (2020). Outsourcing education: Experimental evidence from Liberia. *American Economic Review* 110(2), 364–400.
- Sankar, D. and T. Linden (2014). How much and what kind of teaching is there in elementary education in India? Evidence from three states. Washington, DC: The World Bank.
- Santibanez, L. (2006). Why we should care if teachers get A’s: Teacher test scores and student achievement in Mexico. *Economics of Education Review* 25(5), 510–520.
- Schmidt, W. H., M. T. Tatto, K. Bankov, S. Blömeke, T. Cedillo, L. Cogan, S. I. Han, R. Houang, F. J. Hsieh, L. Paine, M. Santillan, and J. Schille (2007). *Teacher education for middle school mathematics in six countries*. East Lansing, MI: Center for Research in Mathematics and Science Education, Michigan State University.
- Stallings, J. A. (1977). *Learning to look: A handbook on classroom observation and teaching models*. Wadsworth Pub. Co., Belmont.
- Stallings, J. A., S. L. Knight, and D. Markham (2014). Using the Stallings observation system to investigate time on task in four countries. *Unpublished manuscript*. Washington, DC: The World Bank.
- Tatto, M. T., R. Peck, J. Schille, K. Bankov, S. L. Senk, M. Rodriguez, L. Ingvarson, M. Reckase, and G. Rowley (2012). *Policy, practice, and readiness to teach primary and secondary mathematics in 17 countries: Findings from the IEA Teacher Education and Development Study in Mathematics (TEDS-M)*. Amsterdam, the Netherlands: International Evaluation Association (IEA).
- UN (2019). World population prospects, 2019 revision. New York, NY: United Nations Department of Economic and Social Affairs, Population Dynamics. URL: <https://population.un.org/wpp/Publications/>. Last accessed: March 7, 2020.
- UNESCO (2022). Global education monitoring report 2021/2. Non-state actors in education. Who chooses? Who loses? Paris, France: United Nations Educational, Scientific, and Cultural Organization (UNESCO).

- Valencia, S. W., N. A. Place, S. D. Martin, and P. L. Grossman (2006). Curriculum materials for elementary reading: Shackles and scaffolds for four beginning teachers. *The Elementary School Journal* 107(1), 93–120.
- von Hippel, P. T. and L. Bellows (2018). How much does teacher quality vary across teacher preparation programs? reanalyses from six states. *Economics of Education Review* 64, 298–312.
- Wachanga, S. W. and J. G. Mwangi (2004). Effects of the cooperative class experiment teaching method on secondary school students’ chemistry achievement in Kenya’s Nakuru district. *International Education Journal* 5(1), 26–36.
- Woessmann, L. (2003). Specifying human capital. *Journal of Economic Surveys* 17(3), 239–270.
- World Bank (2017). What is happening inside classrooms in Indian secondary schools? A time on task study in Madhya Pradesh and Tamil Nadu. Washington, DC: The World Bank and Educational Initiatives.
- World Bank (2023). Education statistics (Edstats). <https://datatopics.worldbank.org/education/>. Retrieved: June 8, 2023.
- Yoshikawa, H., D. Leyva, C. E. Snow, E. Treviño, M. C. Arbour, M. C. Barata, C. Weiland, C. Gómez, L. Moreno, A. Rolla, and N. D’Sa (2015). Experimental impacts of a teacher professional development program in chile on preschool classroom quality and child outcomes. *Journal of Developmental Psychology* 51, 309–322.

Table 1: Differences between student characteristics (baseline)

	(1) Control	(2) Treatment	(3) Col. (2)-(1)
Female	0.506 [0.500]	0.542 [0.498]	0.036 (0.031)
Age	10.781 [1.011]	10.852 [1.038]	0.055 (0.048)
Speaks Marathi at home	0.690 [0.463]	0.690 [0.463]	-0.000 (0.023)
Speaks Hindi at home	0.169 [0.375]	0.164 [0.370]	-0.005 (0.018)
Speaks English at home	0.011 [0.102]	0.006 [0.080]	-0.004 (0.005)
Mother completed primary school	0.634 [0.482]	0.600 [0.490]	-0.034 (0.023)
Father completed primary school	0.768 [0.422]	0.778 [0.416]	0.009 (0.019)
Student has a desk to study	0.279 [0.449]	0.257 [0.437]	-0.022 (0.029)
Student has own room	0.172 [0.378]	0.204 [0.403]	0.032 (0.030)
Student has a computer	0.195 [0.396]	0.176 [0.381]	-0.019 (0.025)
Student has Internet	0.425 [0.495]	0.373 [0.484]	-0.053 (0.045)
Student has a TV	0.901 [0.299]	0.894 [0.307]	-0.006 (0.017)
Attends tuition in math	0.347 [0.476]	0.315 [0.465]	-0.032 (0.036)
Attends tuition in science	0.260 [0.439]	0.228 [0.420]	-0.032 (0.042)
Attends tuition in Marathi or Hindi	0.368 [0.482]	0.356 [0.479]	-0.012 (0.027)
Attends tuition in English	0.327 [0.469]	0.305 [0.461]	-0.022 (0.034)
Math (standardized score)	0.000 [1.000]	-0.131 [0.955]	-0.131* (0.068)
Science (standardized score)	-0.000 [1.000]	-0.050 [1.011]	-0.050 (0.066)
Language (standardized score)	0.000 [1.000]	-0.066 [1.016]	-0.065 (0.060)
Raven's matrices (standardized score)	0.000 [1.000]	0.045 [0.984]	0.045 (0.055)
N (students)	1,367	1,290	2,657

*Notes:* This table compares the students in the control and treatment groups at baseline. It shows the means and standard deviations for each group (columns 1-2) and tests for differences between groups including randomization-strata (i.e., grade) fixed effects (column 3). The sample includes all students observed at baseline. Standard deviations appear in brackets, and standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table 2: Differences between teacher characteristics (endline)

	(1)	(2)	(3)	(4)	(5)
	PMC teachers		Diff. between PMC teachers	SEI fellows	Diff. between PMC teachers and SEI fellows
	Control	Treatment	Col. (2)-(1)	Treatment	Col. (4)-(1)
Female	0.809	0.711	-0.101	0.633	-0.078
	[0.398]	[0.458]	(0.085)	[0.487]	(0.107)
Age	41.553	41.867	0.250	20.816	-20.989***
	[8.304]	[8.033]	(1.553)	[1.811]	(1.231)
Years of experience (total)	18.936	17.578	-1.472	1.184	-16.307***
	[8.578]	[8.286]	(1.601)	[0.486]	(1.231)
Years of experience (at this school)	5.191	5.089	-0.130	1.000	-4.085***
	[4.808]	[4.039]	(0.993)	[0.000]	(0.638)
Years of math experience (at this school)	1.511	2.533	1.007**	1.000	-1.524***
	[1.266]	[2.693]	(0.467)	[0.000]	(0.395)
Years of science experience (at this school)	1.489	2.533	1.029**	1.000	-1.524***
	[1.249]	[2.693]	(0.467)	[0.000]	(0.395)
Has bachelor's degree or higher	0.809	0.867	0.056	0.347	-0.520***
	[0.398]	[0.344]	(0.076)	[0.481]	(0.080)
Also teaches English	0.957	1.000	0.043	0.000	-1.000
	[0.204]	[0.000]	(0.030)	[0.000]	(0.000)
Also teaches Indian languages	0.894	1.000	0.107**	0.000	-1.000
	[0.312]	[0.000]	(0.045)	[0.000]	(0.000)
Teaches in English	0.064	0.022	-0.043	0.041	0.017
	[0.247]	[0.149]	(0.030)	[0.200]	(0.020)
Teaches in Hindi	0.043	0.178	0.136**	0.163	-0.015
	[0.204]	[0.387]	(0.052)	[0.373]	(0.012)
Teaches in Marathi	0.894	0.800	-0.093*	0.796	-0.002
	[0.312]	[0.405]	(0.051)	[0.407]	(0.021)
Teacher assessment (prop.-correct score)	0.643	0.621	-0.022	0.771	0.150***
	[0.099]	[0.096]	(0.019)	[0.077]	(0.017)
Content knowledge (prop.-correct score)	0.836	0.829	-0.007	0.931	0.102***
	[0.120]	[0.122]	(0.027)	[0.071]	(0.022)
Instructional practices (prop.-correct score)	0.308	0.280	-0.028	0.493	0.213***
	[0.165]	[0.179]	(0.032)	[0.163]	(0.038)
Student errors (prop.-correct score)	0.543	0.503	-0.041	0.689	0.187***
	[0.166]	[0.154]	(0.030)	[0.165]	(0.030)
N (instructors)	47	45	92	49	94

*Notes:* This table compares the Pune Municipal Corporation (PMC) teachers and Science Education Initiative (SEI) fellows at endline. It shows the means and standard deviations for each group and tests for differences between groups including randomization-strata (i.e., grade) fixed effects). The sample includes all teachers observed at endline. Standard deviations appear in brackets, and standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table 3: Impact on standardized test scores (endline, endline audit, and follow-up)

	Math		Science		Language	
	(1)	(2)	(3)	(4)	(5)	(6)
<i>A. Endline</i>						
Treatment	0.247*** (0.048)	0.340*** (0.048)	0.207*** (0.051)	0.216*** (0.049)	0.132** (0.052)	0.151*** (0.039)
Baseline score		0.646*** (0.022)		0.485*** (0.030)		0.555*** (0.026)
N (students)	2524	2307	2524	2307	2524	2307
<i>B. Endline audit</i>						
Treatment	0.046 (0.106)	0.091 (0.104)	0.367*** (0.116)	0.399*** (0.113)		
Baseline score		0.139** (0.061)		0.203*** (0.061)		
N (students)	553	551	553	551		
<i>C. Follow-up</i>						
Treatment	0.283*** (0.051)	0.356*** (0.046)	0.120** (0.046)	0.142** (0.053)	0.059 (0.042)	0.081** (0.037)
Baseline score		0.651*** (0.028)		0.321*** (0.028)		0.446*** (0.029)
N (students)	2369	2023	2369	2023	2369	2023

*Notes:* This table compares students' standardized scores in the control and treatment groups based on assessments administered at endline and endline "audit", about nine months after the rollout of the intervention (April 2018) and at follow-up, about 11 months after the rollout of the intervention (June 2018). Scores are standardized to have mean zero and standard deviation one in the control group. Standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table 4: Impact on proportion-correct repeated and non-repeated items (endline)

	Math		Science		Language	
	(1)	(2)	(3)	(4)	(5)	(6)
<i>A. All items</i>						
Treatment	0.049*** (0.010)	0.067*** (0.009)	0.036*** (0.009)	0.038*** (0.009)	0.031** (0.012)	0.035*** (0.009)
Baseline score		0.584*** (0.019)		0.509*** (0.031)		0.520*** (0.024)
Control mean						
N (students)	2524	2307	2524	2307	2524	2307
<i>B. Repeated items</i>						
Treatment	0.071*** (0.012)	0.094*** (0.012)	0.054*** (0.011)	0.058*** (0.011)	0.031** (0.013)	0.036*** (0.010)
Baseline score		0.749*** (0.024)		0.619*** (0.036)		0.522*** (0.027)
Control mean	0.542		0.462		0.539	
N (students)	2524	2307	2524	2307	2524	2307
<i>C. Non-repeated items</i>						
Treatment	0.028*** (0.008)	0.041*** (0.008)	0.018** (0.009)	0.017** (0.008)	0.030** (0.013)	0.034*** (0.010)
Baseline score		0.428*** (0.018)		0.393*** (0.030)		0.517*** (0.024)
Control mean	0.382		0.388		0.512	
N (students)	2524	2307	2524	2307	2524	2307
P-value of the difference	0.000	0.000	0.933	0.000	0.000	0.862

*Notes:* This table compares students' proportion-correct scores on items in both the baseline and endline assessments ("repeated") and on items only in the endline assessments ("non-repeated") across the control and treatment groups, about nine months after the rollout of the intervention (April 2018). Standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%. As Table 3 indicates, the mean standardized effects on all items without accounting for baseline performance are 0.247 SDs ( $p < 0.01$ ) in math, 0.207 SDs ( $p < 0.01$ ) in science, and 0.132 ( $p < 0.05$ ) in language. The corresponding effects for repeated items are 0.275 SDs ( $p < 0.01$ ) in math, 0.251 SDs ( $p < 0.01$ ) in science, and 0.125 SDs ( $p < 0.05$ ) in language. The corresponding effects for non-repeated items are 0.169 SDs ( $p < 0.01$ ) in math, 0.104 SDs ( $p < 0.05$ ) in science, and 0.124 SDs ( $p < 0.05$ ) in language. Standardizing with respect to the total percent-correct score, effects for repeated items would be 0.359 SDs ( $p < 0.01$ ) in math, 0.306 SDs ( $p < 0.01$ ) in science, and 0.134 SDs ( $p < 0.05$ ) in language; and those for non-repeated items are 0.142 SDs ( $p < 0.01$ ) in math, 0.101 SDs ( $p < 0.05$ ) in science, and 0.130 SDs ( $p < 0.05$ ) in language.

Table 5: Impact on proportion-correct repeated and non-repeated items (follow-up)

	Math		Science		Language	
	(1)	(2)	(3)	(4)	(5)	(6)
<i>A. All items</i>						
Treatment	0.058*** (0.010)	0.073*** (0.009)	0.019** (0.007)	0.022** (0.008)	0.006 (0.004)	0.008** (0.004)
Baseline score		0.609*** (0.026)		0.303*** (0.026)		0.186*** (0.012)
Control mean						
N (students)	2369	2023	2369	2023	2369	2023
<i>B. Repeated items</i>						
Treatment	0.070*** (0.013)	0.087*** (0.012)	0.027** (0.011)	0.032** (0.013)	0.018 (0.012)	0.023** (0.010)
Baseline score		0.738*** (0.033)		0.439*** (0.038)		0.461*** (0.030)
Control mean	0.531		0.380		0.476	
N (students)	2369	2023	2369	2023	2369	2023
<i>C. Non-repeated items</i>						
Treatment	0.041*** (0.010)	0.053*** (0.009)	0.015** (0.006)	0.017** (0.006)	0.005 (0.009)	0.011 (0.011)
Baseline score		0.439*** (0.023)		0.200*** (0.023)		0.302*** (0.024)
Control mean	0.350		0.280		0.413	
N (students)	2369	2023	2369	2023	2369	2023
P-value of the difference	0.012	0.315	0.292	0.027	0.259	0.426

*Notes:* This table compares students' proportion-correct scores on items in both the baseline and endline assessments ("repeated") and on items only in the endline assessments ("non-repeated") across the control and treatment groups, about nine months after the rollout of the intervention (April 2018). Standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%. As Table 3 indicates, the mean standardized effects on all items without accounting for baseline performance are 0.283 SDs ( $p < 0.01$ ) in math, 0.120 SDs ( $p < 0.05$ ) in science, and 0.059 ( $p > 0.10$ ) in language. The corresponding effects for repeated items are 0.274 SDs ( $p < 0.01$ ) in math, 0.110 SDs ( $p < 0.05$ ) in science, and 0.073 SDs ( $p > 0.10$ ) in language. The corresponding effects for non-repeated items are 0.227 SDs ( $p < 0.01$ ) in math, 0.112 SDs ( $p < 0.05$ ) in science, and 0.022 SDs ( $p > 0.10$ ) in language. Standardizing with respect to the total percent-correct score, effects for repeated items would be 0.345 SDs ( $p < 0.01$ ) in math, 0.169 SDs ( $p < 0.05$ ) in science, and 0.177 SDs ( $p > 0.10$ ) in language; and those for non-repeated items are 0.201 SDs ( $p < 0.01$ ) in math, 0.097 SDs ( $p < 0.05$ ) in science, and 0.051 SDs ( $p > 0.10$ ) in language.



Table 6: Impact on instructional practices (announced observations)

	(1)	(2)	(3)
	PMC teachers	SEI fellows	Difference between PMC teachers and SEI fellows
	Control	Treatment	Col. (2)-(1)
<i>A. Negative practices</i>			
Share of instructors who...			
...stayed in one spot during the lesson	0.188 [0.392]	0.060 [0.240]	-0.127** (0.060)
...sat down for a long period of time	0.042 [0.201]	0.020 [0.141]	-0.022 (0.029)
...used phone during lesson	0.052 [0.223]	0.080 [0.274]	0.028 (0.049)
...were upset at incorrect answers	0.062 [0.243]	0.180 [0.388]	0.117* (0.062)
...hit, pinched, or slapped students	0.021 [0.144]	0.000 [0.000]	-0.021 (0.014)
...shouted or used harsh language	0.031 [0.175]	0.020 [0.141]	-0.011 (0.027)
Composite index (std.)	-0.000 [1.000]	0.174 [0.919]	0.174 (0.172)
N (instructors)	96	50	146
<i>B. Positive practices</i>			
Share of instructors who...			
...made eye contact with nearly all students	0.750 [0.435]	0.740 [0.443]	-0.010 (0.062)
...called on nearly all students by name	0.635 [0.484]	0.640 [0.485]	0.005 (0.084)
...called on students from all seating rows	0.688 [0.466]	0.780 [0.418]	0.093 (0.098)
...asked both closed and open questions	0.625 [0.487]	0.780 [0.418]	0.155* (0.088)
...asked students to explain their answers	0.469 [0.502]	0.760 [0.431]	0.291*** (0.095)
...corrected student answers	0.719 [0.452]	0.940 [0.240]	0.221*** (0.070)
...allowed students to ask questions	0.385 [0.489]	0.580 [0.499]	0.195** (0.095)
...assigned classwork	0.635 [0.484]	0.740 [0.443]	0.105 (0.069)
...helped individual students	0.677 [0.470]	0.960 [0.198]	0.283*** (0.064)
...summarized lesson at the end	0.271 [0.447]	0.360 [0.485]	0.089 (0.091)
...assigned homework	0.479 [0.502]	0.640 [0.485]	0.161* (0.092)
...praise or encouraged students	0.760 [0.429]	0.960 [0.198]	0.200*** (0.058)
...smiled, joked, or laughed	0.167 [0.375]	0.280 [0.454]	0.113 (0.090)
Composite index (std.)	-0.000 [1.000]	0.730 [0.687]	0.730*** (0.175)
N (instructors)	96	50	146

*Notes:* This table compares the instructional practices of Pune Municipal Corporation (PMC) teachers in the control and treatment groups and of Science Education Initiative (SEI) fellows in the treatment group, based on announced observations about five months after the rollout of the intervention (December 2017). Observations of PMC teachers lasted 30 minutes and those of SEI fellows lasted 120 minutes. Panel A displays results expressed as a proportion of class time and Panel B shows results in terms of minutes per lesson. The composite indexes are the first principal components of the variables in each panel, standardized with respect to the control group. Standard deviations appear in brackets and standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table 7: Impact on attendance and punctuality (unannounced school visits)

	(1)	(2)	(3)	(4)	(5)
	PMC teachers		Impact on PMC teachers	SEI fellows	Difference between PMC teachers and SEI fellows
	Control	Treatment	Col. (2)-(1)	Treatment	Col. (4)-(1)
<i>A. Attendance and punctuality</i>					
Share of instructors who...					
...attended today	0.840 [0.370]	0.878 [0.331]	0.038 (0.078)	0.720 [0.454]	-0.120 (0.089)
...arrived on time today	0.740 [0.443]	0.714 [0.456]	-0.026 (0.087)	0.440 [0.501]	-0.300*** (0.084)
N (instructors)	50	49	99	50	100
<i>B. Location at the school</i>					
Share of present instructors who...					
...were in the classroom	0.048 [0.216]	0.047 [0.213]	-0.001 (0.031)	0.556 [0.504]	0.508*** (0.094)
...were in the principal's office	0.690 [0.468]	0.744 [0.441]	0.049 (0.079)	0.167 [0.378]	-0.523*** (0.095)
...were elsewhere in the school	0.262 [0.445]	0.209 [0.412]	-0.048 (0.074)	0.278 [0.454]	0.015 (0.098)
N (present instructors)	42	43	85	36	78

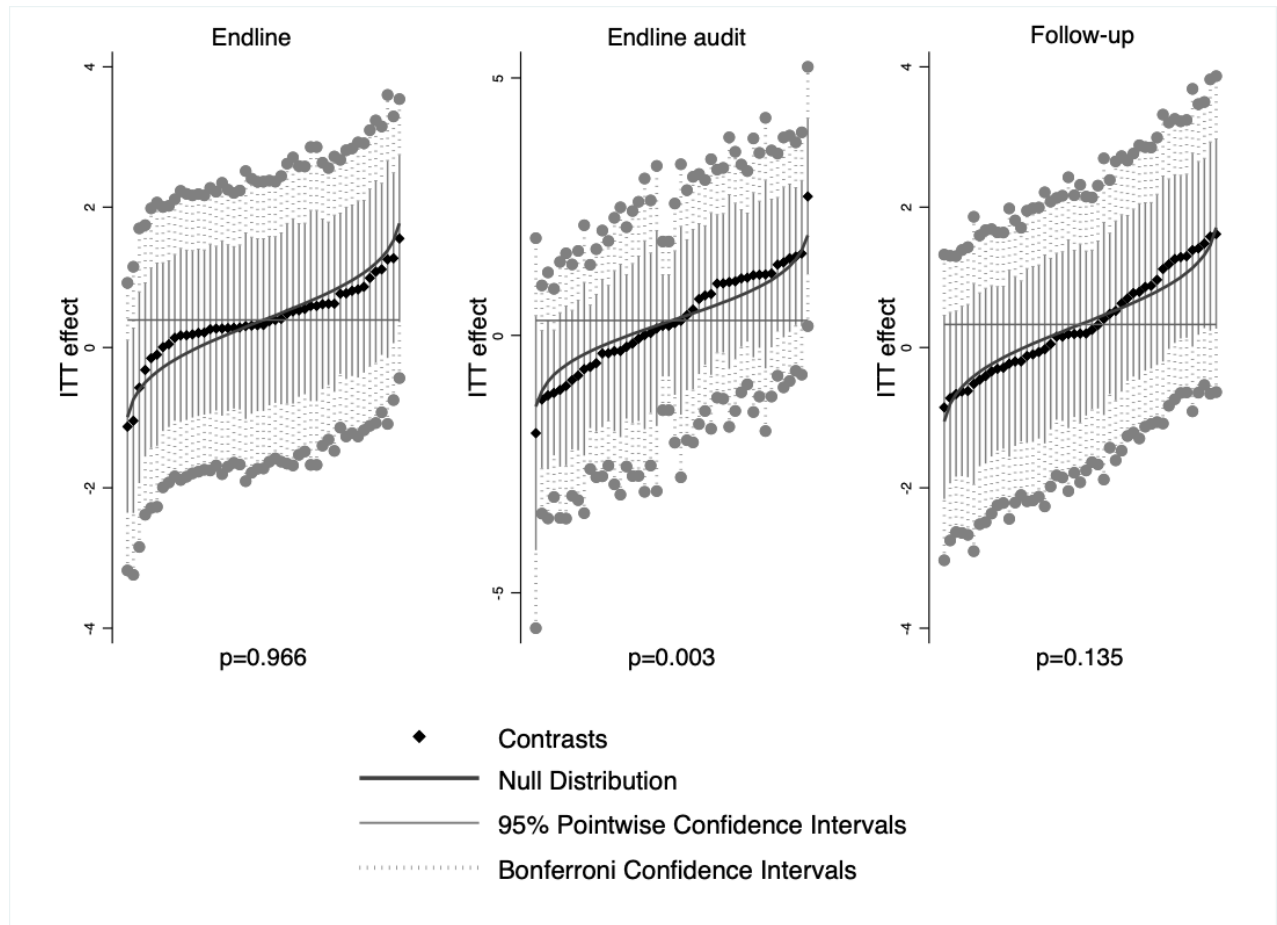
*Notes:* This table compares average attendance and punctuality of Pune Municipal Corporation (PMC) teachers in the control and treatment groups and of Science Education Initiative (SEI) fellows in the treatment group, based on unannounced visits about seven months after the rollout of the intervention (February 2018). Panel A displays results for all instructors and Panel B shows results only for instructors who were present on the day of the visit. Standard deviations appear in brackets and standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table 8: Impact on allocation of instructional time (announced observations)

	(1)	(2)	(3)	(4)	(5)
	PMC teachers		Impact on PMC teachers	SEI fellows	Difference between PMC teachers and SEI fellows
	Control	Treatment	Col. (2)-(1)	Treatment	Col. (4)-(1)
<i>A. Share of lesson</i>					
Share of lesson time...					
...on task	0.778 [0.199]	0.080 [0.134]	-0.698*** (0.026)	0.745 [0.156]	-0.034 (0.031)
...on class management	0.139 [0.134]	0.208 [0.264]	0.069* (0.041)	0.159 [0.126]	0.021 (0.021)
...off task	0.083 [0.135]	0.692 [0.342]	0.609*** (0.048)	0.096 [0.112]	0.013 (0.021)
N (instructors)	96	50	146	50	146
<i>B. Time in lesson</i>					
Minutes of lesson time...					
...on task	23.344 [5.956]	9.600 [16.081]	-13.744*** (2.295)	89.388 [18.665]	65.992*** (2.720)
...on class management	4.156 [4.022]	24.960 [31.688]	20.804*** (4.293)	19.102 [15.083]	14.967*** (2.161)
...off task	2.500 [4.052]	83.040 [41.060]	80.540*** (5.557)	11.510 [13.407]	9.041*** (1.904)
N (instructors)	96	50	146	50	146

*Notes:* This table compares the allocation of instructional time of Pune Municipal Corporation (PMC) teachers in the control and treatment groups and of Science Education Initiative (SEI) fellows in the treatment group, based on announced observations about five months after the rollout of the intervention (December 2017). Observations of PMC teachers lasted 30 minutes and those of SEI fellows lasted 120 minutes. Panel A displays results expressed as a proportion of lesson time and Panel B shows results in terms of minutes per lesson. Standard deviations appear in brackets and standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Figure 1: Heterogeneous impact on standardized test scores by school (endline, endline audit, and follow-up)



*Notes:* The figure provides a “caterpillar plot” (von Hippel and Bellows, 2018) of impacts on the first principal component of a principal-component analysis of standardized scores in math and science by school at endline, about nine months after the rollout of the intervention (April 2018); the endline “audit,” about nine months after the rollout of the intervention (April 2018); and follow-up, about 11 months after the rollout of the intervention (June 2018). Each black dot refers to the point estimate for a given school. Bonferroni confidence intervals adjust standard errors for multiple hypothesis testing. The black solid line shows the null distribution of “effects” that can be expected due to error. The p-values are for a test of the null hypothesis of no heterogeneity.

## Appendix A Additional tables and figures

Table A.1: Attrition rates in assessments (endline and follow-up)

	(1) Endline	(2) Follow-up
<i>A. Treatment</i>		
Treatment	-0.019 (0.014)	-0.004 (0.017)
N (students)	2657	2657
Control mean	0.141	0.243
<i>B. Treatment and baseline</i>		
Treatment	0.158 (0.171)	0.179 (0.196)
Female	-0.079** (0.035)	-0.064* (0.036)
Age	0.030** (0.011)	0.038** (0.016)
Composite score (at baseline)	-0.046*** (0.008)	-0.038*** (0.010)
Female $\times$ Treatment	0.059 (0.038)	0.051 (0.044)
Age $\times$ Treatment	-0.020 (0.016)	-0.021 (0.019)
Composite score $\times$ Treatment	-0.009 (0.009)	-0.008 (0.014)
N (students)	2404	2404
F-ratio (Interactions)	1.686	1.098
P-value	0.183	0.359

*Notes:* This table shows estimates from regressions predicting follow-up status in assessments administered at endline, about nine months after the rollout of the intervention (April 2018) and at follow-up, about 11 months after the rollout of the intervention (June 2018). Follow-up is defined as having an observed test score at endline. The sample includes students present at baseline. Panel A regresses follow-up on treatment status, and Panel B regresses follow-up on treatment status interacted with baseline characteristics. Both panels include randomization-strata fixed effects. Standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table A.2: Round of data collection for the study

(1)	(2)	(3)	(4)	(5)	(6)
			Student participation rates		
			Both		
Month	Event	Grades	groups	Control	Treatment
<i>A. 2017-2018 school year</i>					
June	School year starts				
<u>Baseline:</u>					
July	Student assessments	5, 6	100%	100%	100%
	Student surveys	5, 6	100%	100%	100%
<u>Midline:</u>					
December	Announced classroom observations	5, 6	99%	100%	98%
February	Unannounced school visits	5, 6	99%	100%	98%
	Administrative data on student attendance	5, 6	100%	100%	100%
<u>Endline:</u>					
March-	Student assessments	5, 6			
April	– Math and science		88%	87%	89%
	– Language		86%	85%	87%
	Student surveys	5, 6	80%	79%	81%
	Student “audit” assessments of math and science	5, 6	21%	22%	22%
April-	Teacher surveys	5, 6	95%	93%	98%
May	Teacher assessments	5, 6	94%	92%	98%
<i>B. 2018-2019 school year</i>					
June	School year starts				
<u>Follow-up:</u>					
June-	Student assessments	6, 7			
August	– Math		81%	81%	81%
	– Science		80%	80%	81%
	– Language		81%	80%	81%
	Administrative data on student attendance	6, 7			

*Notes:* The table shows the timeline for the interventions and rounds of data collection for the study, including the month in which each event occurred (column 1), a brief description of the event (column 2), the target grades (column 3), and the percentage of students that participated in each event by experimental group (columns 4-6). The student “audit” assessments only targeted 25% of the study sample.

Table A.3: Script adherence among SEI fellows (announced observations)

	(1) All fellows	(2) High- scoring	(3) Low- scoring
Share of SEI fellows who...			
...wrote on blackboard as indicated in the script	0.818 [0.390]	0.800 [0.408]	0.842 [0.375]
...asked students questions in the script	0.750 [0.438]	0.800 [0.408]	0.684 [0.478]
...used materials as indicated in the script	0.523 [0.505]	0.480 [0.510]	0.579 [0.507]
...assigned students activities in the script	0.500 [0.506]	0.600 [0.500]	0.368 [0.496]
...changed/rephrased parts of the script	0.341 [0.479]	0.280 [0.458]	0.421 [0.507]
...added to or expanded on parts of the script	0.341 [0.479]	0.240 [0.436]	0.474 [0.513]
...excluded parts of the script	0.182 [0.390]	0.160 [0.374]	0.211 [0.419]
N (fellows)	50	28	22

*Notes:* This table displays the prevalence of script adherence among Science Education Initiative (SEI) fellows in treatment grades, based on announced observations about five months after the rollout of the intervention (December 2017). Standard deviations appear in brackets. High-scoring fellows are those who scored above the fellow-specific median on a test of content knowledge, instructional practices, and understanding of students' misconceptions (described in section 3); low-scoring fellows are those below that median.

Table A.4: Impact on allocation of time on task (announced observations)

	(1)	(2)	(3)	(4)	(5)
	PMC teachers		Impact on PMC teachers	SEI fellows	Difference between PMC teachers and SEI fellows
	Control	Treatment	Col. (2)-(1)	Treatment	Col. (4)-(1)
<i>A. Share of lesson</i>					
Share of lesson time...					
...reading aloud	0.049 [0.122]	0.000 [0.000]	-0.049*** (0.013)	0.014 [0.041]	-0.035** (0.014)
...on explanation/lecture	0.343 [0.208]	0.020 [0.061]	-0.323*** (0.022)	0.337 [0.191]	-0.007 (0.030)
...on interactive demo	0.036 [0.070]	0.000 [0.000]	-0.036*** (0.007)	0.008 [0.028]	-0.028*** (0.009)
...on question and answers	0.181 [0.153]	0.002 [0.014]	-0.179*** (0.016)	0.153 [0.126]	-0.028 (0.021)
...on practice and drill	0.007 [0.030]	0.002 [0.014]	-0.005* (0.003)	0.004 [0.020]	-0.003 (0.005)
...on classwork	0.129 [0.147]	0.002 [0.014]	-0.127*** (0.016)	0.157 [0.117]	0.028 (0.021)
...copying	0.032 [0.067]	0.000 [0.000]	-0.032*** (0.007)	0.057 [0.076]	0.025** (0.012)
N (instructors)	96	50	146	50	146
<i>B. Time in lesson</i>					
Minutes of lesson time...					
...reading aloud	1.469 [3.668]	0.000 [0.000]	-1.469*** (0.395)	1.714 [4.899]	0.245 (0.765)
...on explanation/lecture	10.281 [6.228]	2.400 [7.273]	-7.881*** (1.213)	40.408 [22.938]	30.067*** (3.194)
...on interactive demo	1.094 [2.093]	0.000 [0.000]	-1.094*** (0.218)	0.980 [3.320]	-0.116 (0.544)
...on question and answers	5.438 [4.592]	0.240 [1.697]	-5.198*** (0.511)	18.367 [15.120]	12.943*** (2.009)
...on practice and drill	0.219 [0.897]	0.240 [1.697]	0.021 (0.199)	0.490 [2.399]	0.273 (0.364)
...on classwork	3.875 [4.416]	0.240 [1.697]	-3.635*** (0.525)	18.857 [14.071]	14.969*** (1.971)
...copying	0.969 [2.018]	0.000 [0.000]	-0.969*** (0.209)	6.857 [9.165]	5.892*** (1.270)
N (instructors)	96	50	146	50	146

*Notes:* This table compares the allocation of time on task of Pune Municipal Corporation (PMC) teachers in the control and treatment groups and of Science Education Initiative (SEI) fellows in the treatment group, measured in announced observations about five months after the rollout of the intervention (December 2017). All lessons taught by PMC teachers lasted 30 minutes and all lessons taught by SEI fellows lasted 120 minutes. Panel A displays results expressed as a proportion of lesson time and Panel B shows results in terms of minutes per lesson. Standard deviations appear in brackets and standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.



Table A.5: Impact on allocation of time on classroom management (announced observations)

	(1)	(2)	(3)	(4)	(5)
	PMC teachers		Impact on PMC teachers	SEI fellows	Difference between PMC teachers and SEI fellows
	Control	Treatment	Col. (2)-(1)	Treatment	Col. (4)-(1)
<i>A. Share of lesson</i>					
Share of lesson time...					
...on instructions	0.011 [0.035]	0.006 [0.024]	-0.005 (0.005)	0.043 [0.076]	0.031*** (0.010)
...on discipline	0.015 [0.038]	0.010 [0.036]	-0.005 (0.006)	0.016 [0.051]	0.002 (0.008)
...on class mgmt. w/students	0.057 [0.084]	0.008 [0.027]	-0.049*** (0.011)	0.082 [0.075]	0.024* (0.014)
...on class mgmt. alone	0.055 [0.099]	0.184 [0.245]	0.129*** (0.035)	0.018 [0.049]	-0.037*** (0.013)
N (instructors)	96	50	146	50	146
<i>B. Time in lesson</i>					
Minutes of lesson time...					
...on instructions	0.344 [1.055]	0.720 [2.879]	0.376 (0.410)	5.143 [9.165]	4.799*** (1.262)
...on discipline	0.438 [1.150]	1.200 [4.371]	0.762 (0.595)	1.959 [6.171]	1.523* (0.874)
...on class mgmt. w/students	1.719 [2.529]	0.960 [3.289]	-0.759 (0.582)	9.796 [9.058]	8.088*** (1.312)
...on class mgmt. alone	1.656 [2.980]	22.080 [29.430]	20.424*** (3.930)	2.204 [5.834]	0.558 (0.853)
N (instructors)	96	50	146	50	146

*Notes:* This table compares the allocation of time on classroom management of Pune Municipal Corporation (PMC) teachers in the control and treatment groups and of Science Education Initiative (SEI) fellows in the treatment group, measured in announced observations about five months after the rollout of the intervention (December 2017). All lessons taught by PMC teachers lasted 30 minutes and all lessons taught by SEI fellows lasted 120 minutes. Panel A displays results expressed as a proportion of class time and Panel B shows results in terms of minutes per lesson. Standard deviations appear in brackets and standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table A.6: Impact on allocation of time off task (announced observations)

	(1)	(2)	(3)	(4)	(5)
	PMC teachers		Impact on PMC teachers	SEI fellows	Difference between PMC teachers and SEI fellows
	Control	Treatment	Col. (2)-(1)	Treatment	Col. (4)-(1)
<i>A. Share of lesson</i>					
Share of lesson time...					
...on social interaction w/students	0.007 [0.030]	0.002 [0.014]	-0.005 (0.004)	0.002 [0.014]	-0.005 (0.004)
...on social interaction w/adults	0.015 [0.046]	0.016 [0.037]	0.001 (0.008)	0.006 [0.024]	-0.008 (0.006)
...uninvolved	0.038 [0.094]	0.102 [0.235]	0.064* (0.034)	0.018 [0.049]	-0.019 (0.014)
...out of room	0.024 [0.086]	0.572 [0.395]	0.548*** (0.053)	0.069 [0.102]	0.046*** (0.016)
N (instructors)	96	50	146	50	146
<i>B. Time in lesson</i>					
Minutes of lesson time...					
...on social interaction w/students	0.219 [0.897]	0.240 [1.697]	0.021 (0.263)	0.245 [1.714]	0.030 (0.261)
...on social interaction w/adults	0.438 [1.375]	1.920 [4.444]	1.482** (0.660)	0.735 [2.907]	0.298 (0.419)
...uninvolved	1.125 [2.829]	12.240 [28.220]	11.115*** (3.925)	2.204 [5.834]	1.084 (0.927)
...out of the room	0.719 [2.566]	68.640 [47.448]	67.921*** (6.274)	8.327 [12.297]	7.628*** (1.717)
N (instructors)	96	50	146	50	146

*Notes:* This table compares the allocation of time off task of Pune Municipal Corporation (PMC) teachers in the control and treatment groups and of Science Education Initiative (SEI) fellows in the treatment group, measured in announced observations about five months after the rollout of the intervention (December 2017). All lessons taught by PMC teachers lasted 30 minutes and all lessons taught by SEI fellows lasted 120 minutes. Panel A displays results expressed as a proportion of class time and Panel B shows results in terms of minutes per lesson. Standard deviations appear in brackets and standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table A.7: Impact on allocation of time engaging students (announced observations)

	(1)	(2)	(3)	(4)	(5)
	PMC teachers		Impact on PMC teachers	SEI fellows	Difference between PMC teachers and SEI fellows
	Control	Treatment	Col. (2)-(1)	Treatment	Col. (4)-(1)
<i>A. Share of lesson</i>					
Share of lesson time engaging...					
...no students	0.091 [0.143]	0.746 [0.296]	0.655*** (0.042)	0.106 [0.121]	0.016 (0.022)
...one student	0.119 [0.127]	0.014 [0.040]	-0.105*** (0.015)	0.129 [0.122]	0.010 (0.023)
...a small group of students	0.044 [0.084]	0.000 [0.000]	-0.044*** (0.009)	0.045 [0.074]	0.001 (0.014)
...a large group of students	0.232 [0.199]	0.016 [0.042]	-0.216*** (0.024)	0.298 [0.230]	0.066 (0.045)
...all students	0.451 [0.248]	0.024 [0.059]	-0.427*** (0.031)	0.404 [0.230]	-0.048 (0.041)
N (instructors)	96	50	146	50	146
<i>B. Time in lesson</i>					
Number of minutes engaging...					
...no students	2.719 [4.289]	89.520 [35.566]	86.801*** (4.775)	12.735 [14.576]	10.049*** (2.024)
...one student	3.562 [3.803]	1.680 [4.855]	-1.883** (0.765)	15.429 [14.697]	11.850*** (2.142)
...a small group of students	1.312 [2.531]	0.000 [0.000]	-1.312*** (0.272)	5.388 [8.853]	4.070*** (1.291)
...a large group of students	6.969 [5.976]	1.920 [5.062]	-5.049*** (0.952)	35.755 [27.549]	28.789*** (4.302)
...all students	13.531 [7.425]	2.880 [7.093]	-10.651*** (1.472)	48.490 [27.600]	34.938*** (4.203)
N (instructors)	96	50	146	50	146

*Notes:* This table compares the allocation of time engaging students of Pune Municipal Corporation (PMC) teachers in the control and treatment groups and of Science Education Initiative (SEI) fellows in the treatment group, measured in announced observations about five months after the rollout of the intervention (December 2017). All lessons taught by PMC teachers lasted 30 minutes and all lessons taught by SEI fellows lasted 120 minutes. Panel A displays results expressed as a proportion of class time and Panel B shows results in terms of minutes per lesson. Standard deviations appear in brackets and standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table A.8: Impact on proportion-correct test scores (endline, endline audit, and follow-up)

	Math		Science		Language	
	(1)	(2)	(3)	(4)	(5)	(6)
<i>A. Endline</i>						
Treatment	0.049*** (0.010)	0.067*** (0.009)	0.036*** (0.009)	0.038*** (0.009)	0.031** (0.012)	0.035*** (0.009)
Baseline score		0.584*** (0.019)		0.509*** (0.031)		0.520*** (0.024)
N (students)	2524	2307	2524	2307	2524	2307
<i>B. Endline audit</i>						
Treatment	0.011 (0.026)	0.067*** (0.009)	0.092*** (0.029)	0.038*** (0.009)		
Baseline score		0.584*** (0.019)		0.509*** (0.031)		
N (students)	553	2307	553	2307		
<i>C. Follow-up</i>						
Treatment	0.058*** (0.010)	0.073*** (0.009)	0.019** (0.007)	0.022** (0.008)	0.006 (0.004)	0.008** (0.004)
Baseline score		0.609*** (0.026)		0.303*** (0.026)		0.186*** (0.012)
N (students)	2369	2023	2369	2023	2369	2023

*Notes:* This table compares students' proportion-correct scores in the control and treatment groups based on assessments administered at endline and endline "audit", about nine months after the rollout of the intervention (April 2018) and at follow-up, about 11 months after the rollout of the intervention (June 2018). Scores are standardized to have mean zero and standard deviation one in the control group. Standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table A.9: Impact on IRT-scaled test scores (endline, endline audit, and follow-up)

	Math		Science		Language	
	(1)	(2)	(3)	(4)	(5)	(6)
<i>A. Endline</i>						
Treatment	0.266*** (0.047)	0.347*** (0.043)	0.192*** (0.053)	0.192*** (0.048)	0.115** (0.048)	0.159*** (0.035)
Baseline score		0.703*** (0.024)		0.508*** (0.031)		0.598*** (0.026)
N (students)	2504	2288	2497	2284	2432	2231
<i>B. Endline audit</i>						
Treatment	0.098 (0.085)	0.134 (0.085)	0.222*** (0.067)	0.236*** (0.063)		
Baseline score		0.140*** (0.043)		0.119*** (0.032)		
N (students)	553	551	553	551		
<i>C. Follow-up</i>						
Treatment	0.293*** (0.055)	0.353*** (0.048)	0.089* (0.044)	0.106** (0.052)	0.033 (0.038)	0.055 (0.034)
Baseline score		0.672*** (0.032)		0.284*** (0.028)		0.354*** (0.029)
N (students)	2270	1925	2271	1926	2267	1924

*Notes:* This table compares students' Item Response Theory (IRT)-scaled scores in the control and treatment groups based on assessments administered at endline and endline "audit", about nine months after the rollout of the intervention (April 2018) and at follow-up, about 11 months after the rollout of the intervention (June 2018). Scores are standardized to have mean zero and standard deviation one in the control group. Standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table A.10: Heterogeneous impact on standardized test scores by students' baseline score (endline, endline audit, follow-up)

	Math	Science	Language
	(1)	(2)	(3)
<i>A. Endline</i>			
Treatment	0.339*** (0.048)	0.214*** (0.050)	0.154*** (0.038)
Baseline score	0.627*** (0.030)	0.463*** (0.037)	0.585*** (0.029)
Treat $\times$ Baseline score	0.040 (0.038)	0.044 (0.040)	-0.063 (0.043)
P-value of the sum	0.000	0.000	0.000
N (students)	2307	2307	2307
<i>B. Endline audit</i>			
Treatment	0.091 (0.104)	0.399*** (0.113)	
Baseline score	0.133* (0.078)	0.206** (0.078)	
Treat $\times$ Baseline score	0.014 (0.121)	-0.006 (0.105)	
P-value of the sum	0.129	0.021	
N (students)	551	551	
<i>C. Follow-up</i>			
Treatment	0.356*** (0.046)	0.138** (0.052)	0.081** (0.036)
Baseline score	0.654*** (0.039)	0.291*** (0.035)	0.448*** (0.037)
Treat $\times$ Baseline score	-0.006 (0.048)	0.061 (0.050)	-0.004 (0.043)
P-value of the sum	0.000	0.000	0.000
N (students)	2023	2023	2023

*Notes:* This table shows the impact of the intervention on assessments of math, science, and language administered at endline and endline “audit”, about nine months after the rollout of the intervention (April 2018), and at follow-up, about 11 months after the rollout of the intervention (June 2018) by students' baseline score. Endline scores are standardized to have mean zero and standard deviation one in the control group. The p-value of the sum is from a test that the sum of the main effect of treatment and its interaction is equal to zero. Standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table A.11: Heterogeneous impact on standardized test scores by students' socio-economic status (endline, endline audit, follow-up)

	Math		Science		Language	
	(1)	(2)	(3)	(4)	(5)	(6)
<i>A. Endline</i>						
Treatment	0.261*** (0.051)	0.344*** (0.048)	0.191*** (0.052)	0.210*** (0.050)	0.113** (0.052)	0.138*** (0.041)
SES index	-0.030 (0.033)	-0.039* (0.021)	-0.069* (0.037)	-0.096*** (0.032)	-0.084** (0.037)	-0.074** (0.029)
Treat $\times$ SES index	0.112** (0.045)	0.056* (0.032)	0.165*** (0.053)	0.144*** (0.045)	0.111* (0.058)	0.057 (0.046)
Baseline score		0.654*** (0.023)		0.502*** (0.032)		0.574*** (0.026)
P-value of the sum	0.017	0.464	0.010	0.173	0.539	0.648
N (students)	2255	2255	2255	2255	2255	2255
<i>B. Endline audit</i>						
Treatment	0.059 (0.111)	0.098 (0.111)	0.392*** (0.121)	0.424*** (0.115)		
SES index	0.016 (0.080)	0.003 (0.079)	-0.082 (0.067)	-0.105 (0.066)		
Treat $\times$ SES index	-0.093 (0.104)	-0.091 (0.100)	0.106 (0.108)	0.104 (0.108)		
Baseline score		0.154** (0.064)		0.224*** (0.063)		
P-value of the sum	0.210	0.117	0.751	0.985		
N (students)	532	532	532	532		
<i>C. Follow-up</i>						
Treatment	0.274*** (0.052)	0.355*** (0.045)	0.133*** (0.049)	0.152*** (0.053)	0.059 (0.048)	0.077** (0.036)
SES index	-0.001 (0.033)	-0.012 (0.024)	-0.048 (0.030)	-0.068** (0.029)	-0.035 (0.032)	-0.032 (0.026)
Treat $\times$ SES index	0.083* (0.044)	0.036 (0.033)	0.138*** (0.046)	0.125*** (0.041)	0.089** (0.042)	0.049 (0.032)
Baseline score		0.658*** (0.030)		0.340*** (0.029)		0.451*** (0.029)
P-value of the sum	0.031	0.405	0.013	0.057	0.139	0.555
N (students)	1972	1972	1972	1972	1972	1972

Notes: This table shows the impact of the intervention on assessments of math, science, and language administered at endline and endline "audit", about nine months after the rollout of the intervention (April 2018), and at follow-up, about 11 months after the rollout of the intervention (June 2018) by students' socio-economic status. Endline scores are standardized to have mean zero and standard deviation one in the control group. The p-value of the sum is from a test that the sum of the main effect of treatment and its interaction is equal to zero. Standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table A.12: Heterogeneous impact on standardized test scores by students' sex (endline, endline audit, follow-up)

	Math		Science		Language	
	(1)	(2)	(3)	(4)	(5)	(6)
<i>A. Endline</i>						
Treatment	0.218** (0.083)	0.267*** (0.065)	0.131 (0.080)	0.148** (0.070)	0.025 (0.074)	0.095 (0.064)
Female	-0.016 (0.067)	-0.004 (0.047)	0.064 (0.087)	0.053 (0.073)	0.161** (0.067)	0.147*** (0.053)
Treat $\times$ Female	0.085 (0.104)	0.145* (0.081)	0.114 (0.114)	0.117 (0.105)	0.163* (0.089)	0.080 (0.088)
Baseline score		0.657*** (0.022)		0.501*** (0.032)		0.568*** (0.028)
P-value of the sum	0.433	0.047	0.066	0.038	0.000	0.002
N (students)	2258	2258	2258	2258	2258	2258
<i>B. Endline audit</i>						
Treatment	-0.005 (0.136)	0.028 (0.141)	0.452*** (0.124)	0.487*** (0.126)		
Female	0.155 (0.171)	0.158 (0.167)	0.154 (0.131)	0.162 (0.124)		
Treat $\times$ Female	0.121 (0.219)	0.130 (0.221)	-0.142 (0.202)	-0.149 (0.199)		
Baseline score		0.155** (0.061)		0.220*** (0.062)		
P-value of the sum	0.121	0.103	0.950	0.943		
N (students)	533	533	533	533		
<i>C. Follow-up</i>						
Treatment	0.262*** (0.096)	0.310*** (0.066)	0.096 (0.074)	0.104 (0.077)	0.052 (0.080)	0.114* (0.061)
Female	0.068 (0.083)	0.071 (0.057)	0.103 (0.073)	0.087 (0.069)	0.134* (0.072)	0.115* (0.059)
Treat $\times$ Female	0.022 (0.122)	0.083 (0.081)	0.065 (0.099)	0.086 (0.103)	0.014 (0.103)	-0.066 (0.086)
Baseline score		0.662*** (0.030)		0.340*** (0.028)		0.450*** (0.029)
P-value of the sum	0.305	0.012	0.029	0.020	0.069	0.435
N (students)	1974	1974	1974	1974	1974	1974

*Notes:* This table shows the impact of the intervention on assessments of math, science, and language administered at endline and endline “audit”, about nine months after the rollout of the intervention (April 2018), and at follow-up, about 11 months after the rollout of the intervention (June 2018) by students' sex. Endline scores are standardized to have mean zero and standard deviation one in the control group. The p-value of the sum is from a test that the sum of the main effect of treatment and its interaction is equal to zero. Standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.



Table A.13: Heterogeneous impact on standardized test scores by students' caste (endline, endline audit, follow-up)

	Math		Science		Language	
	(1)	(2)	(3)	(4)	(5)	(6)
<i>A. Endline</i>						
Treatment	0.337*** (0.055)	0.389*** (0.045)	0.223*** (0.067)	0.205*** (0.065)	0.108 (0.066)	0.077 (0.067)
Scheduled caste/tribe	-0.111 (0.078)	-0.017 (0.047)	-0.101 (0.084)	-0.071 (0.065)	-0.132 (0.081)	-0.122** (0.060)
Treat $\times$ SC/ST	-0.106 (0.088)	-0.092 (0.062)	-0.018 (0.101)	0.022 (0.087)	0.079 (0.091)	0.160* (0.083)
Baseline score		0.647*** (0.021)		0.494*** (0.027)		0.559*** (0.025)
P-value of the sum	0.003	0.061	0.094	0.405	0.538	0.602
N (students)	2236	2236	2236	2236	2236	2236
<i>B. Endline audit</i>						
Treatment	0.206 (0.165)	0.232 (0.165)	0.517*** (0.155)	0.523*** (0.154)		
Scheduled caste/tribe	0.019 (0.139)	0.043 (0.143)	0.048 (0.102)	0.038 (0.100)		
Treat $\times$ SC/ST	-0.274 (0.197)	-0.276 (0.198)	-0.188 (0.178)	-0.169 (0.181)		
Baseline score		0.129** (0.064)		0.174*** (0.061)		
P-value of the sum	0.042	0.063	0.369	0.411		
N (students)	537	537	537	537		
<i>C. Follow-up</i>						
Treatment	0.307*** (0.055)	0.366*** (0.051)	0.178** (0.068)	0.158** (0.065)	0.136** (0.056)	0.096* (0.053)
Scheduled caste/tribe	-0.082 (0.071)	0.013 (0.043)	-0.057 (0.063)	-0.048 (0.057)	-0.001 (0.064)	-0.009 (0.061)
Treat $\times$ SC/ST	-0.025 (0.086)	-0.023 (0.061)	-0.090 (0.093)	-0.041 (0.085)	-0.090 (0.086)	-0.006 (0.085)
Baseline score		0.665*** (0.026)		0.323*** (0.029)		0.448*** (0.030)
P-value of the sum	0.202	0.866	0.037	0.153	0.173	0.780
N (students)	1955	1955	1955	1955	1955	1955

*Notes:* This table shows the impact of the intervention on assessments of math, science, and language administered at endline and endline "audit", about nine months after the rollout of the intervention (April 2018), and at follow-up, about 11 months after the rollout of the intervention (June 2018) by students' caste. Endline scores are standardized to have mean zero and standard deviation one in the control group. The p-value of the sum is from a test that the sum of the main effect of treatment and its interaction is equal to zero. Standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table A.14: Heterogeneous impact on standardized test scores by student-instructor sex match (endline, endline audit, follow-up)

	Math		Science		Language	
	(1)	(2)	(3)	(4)	(5)	(6)
<i>A. Endline</i>						
Treatment	0.312*** (0.068)	0.367*** (0.047)	0.220*** (0.071)	0.242*** (0.061)	0.148** (0.066)	0.161*** (0.059)
Sex match	0.084 (0.074)	0.042 (0.049)	0.147 (0.097)	0.124 (0.076)	0.228*** (0.075)	0.137** (0.057)
Treat $\times$ Sex match	-0.037 (0.101)	0.029 (0.087)	0.014 (0.122)	-0.010 (0.109)	-0.022 (0.116)	0.003 (0.105)
Baseline score		0.656*** (0.021)		0.505*** (0.030)		0.566*** (0.028)
P-value of the sum	0.585	0.327	0.107	0.203	0.037	0.104
N (students)	2128	2128	2128	2128	2128	2128
<i>B. Endline audit</i>						
Treatment	-0.054 (0.137)	-0.016 (0.137)	0.466*** (0.152)	0.505*** (0.140)		
Sex match	-0.028 (0.170)	-0.032 (0.165)	0.146 (0.134)	0.165 (0.119)		
Treat $\times$ Sex match	0.344 (0.251)	0.346 (0.255)	-0.180 (0.244)	-0.207 (0.243)		
Baseline score		0.154** (0.065)		0.230*** (0.063)		
P-value of the sum	0.145	0.145	0.872	0.846		
N (students)	506	506	506	506		
<i>C. Follow-up</i>						
Treatment	0.337*** (0.081)	0.378*** (0.055)	0.121** (0.059)	0.128* (0.064)	0.101 (0.062)	0.100** (0.048)
Sex match	0.159* (0.084)	0.105* (0.058)	0.103 (0.065)	0.076 (0.055)	0.150* (0.079)	0.071 (0.064)
Treat $\times$ Sex match	-0.068 (0.127)	0.012 (0.093)	0.089 (0.094)	0.094 (0.093)	-0.039 (0.102)	-0.012 (0.079)
Baseline score		0.655*** (0.030)		0.342*** (0.030)		0.449*** (0.028)
P-value of the sum	0.318	0.076	0.024	0.039	0.128	0.288
N (students)	1870	1870	1870	1870	1870	1870

*Notes:* This table shows the impact of the intervention on assessments of math, science, and language administered at endline and endline “audit”, about nine months after the rollout of the intervention (April 2018), and at follow-up, about 11 months after the rollout of the intervention (June 2018) by whether the sex of the instructor matches the sex of the student. Endline scores are standardized to have mean zero and standard deviation one in the control group. The p-value of the sum is from a test that the sum of the main effect of treatment and its interaction is equal to zero. Standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table A.15: Heterogeneous impact on standardized test scores by instructors' assessment score (endline, endline audit, follow-up)

	Math		Science		Language	
	(1)	(2)	(3)	(4)	(5)	(6)
<i>A. Endline</i>						
Treatment	0.255** (0.101)	0.449*** (0.064)	0.162 (0.107)	0.211** (0.091)	0.127 (0.123)	0.242*** (0.082)
Instructor score	0.104* (0.058)	0.068 (0.046)	0.117** (0.053)	0.093* (0.050)	0.039 (0.075)	0.021 (0.052)
Treat $\times$ Instructor score	-0.083 (0.089)	-0.123** (0.053)	-0.064 (0.098)	-0.073 (0.074)	-0.026 (0.105)	-0.079 (0.074)
Baseline score		0.649*** (0.021)		0.486*** (0.029)		0.557*** (0.025)
P-value of the sum	0.723	0.119	0.444	0.687	0.873	0.395
N (students)	2374	2174	2374	2174	2374	2174
<i>B. Endline audit</i>						
Treatment	0.099 (0.236)	0.167 (0.242)	0.295 (0.219)	0.365* (0.214)		
Instructor score	-0.178* (0.097)	-0.196** (0.095)	0.014 (0.096)	-0.011 (0.094)		
Treat $\times$ Instructor score	0.119 (0.196)	0.123 (0.198)	0.046 (0.190)	0.037 (0.179)		
Baseline score		0.169*** (0.061)		0.212*** (0.063)		
P-value of the sum	0.704	0.641	0.689	0.849		
N (students)	523	521	523	521		
<i>C. Follow-up</i>						
Treatment	0.233** (0.114)	0.368*** (0.080)	0.052 (0.078)	0.106 (0.064)	0.008 (0.083)	0.090 (0.064)
Instructor score	0.118* (0.064)	0.072 (0.046)	0.028 (0.036)	0.004 (0.038)	0.051 (0.043)	-0.004 (0.032)
Treat $\times$ Instructor score	-0.059 (0.101)	-0.070 (0.069)	0.028 (0.067)	0.024 (0.058)	-0.004 (0.072)	0.001 (0.057)
Baseline score		0.649*** (0.029)		0.325*** (0.030)		0.446*** (0.027)
P-value of the sum	0.323	0.955	0.272	0.533	0.338	0.956
N (students)	2231	1909	2231	1909	2231	1909

*Notes:* This table shows the impact of the intervention on assessments of math, science, and language administered at endline and endline “audit”, about nine months after the rollout of the intervention (April 2018), and at follow-up, about 11 months after the rollout of the intervention (June 2018) by instructors' standardized scores on an assessment of instructional practice, knowledge of student errors, and content knowledge in math (calculated as the first principal component of a principal-component analysis of the scores on all three domains, standardized with respect to instructors in the control group). Endline scores are standardized to have mean zero and standard deviation one in the control group. The p-value of the sum is from a test that the sum of the main effect of treatment and its interaction is equal to zero. Standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table A.16: Impact on standardized test scores among high-scoring teachers (endline, endline audit, follow-up)

	Math		Science		Language	
	(1)	(2)	(3)	(4)	(5)	(6)
<i>A. Endline</i>						
Treatment	0.256*** (0.069)	0.390*** (0.060)	0.176** (0.085)	0.200** (0.079)	0.102 (0.091)	0.193*** (0.054)
Baseline score		0.634*** (0.032)		0.460*** (0.035)		0.554*** (0.028)
N (students)	1382	1276	1382	1276	1382	1276
<i>B. Endline audit</i>						
Treatment	0.252 (0.175)	0.331* (0.175)	0.305* (0.176)	0.337* (0.172)		
Baseline score		0.214*** (0.075)		0.144** (0.068)		
N (students)	308	307	308	307		
<i>C. Follow-up</i>						
Treatment	0.232** (0.093)	0.371*** (0.076)	0.118 (0.076)	0.160** (0.064)	0.039 (0.077)	0.127** (0.053)
Baseline score		0.675*** (0.042)		0.287*** (0.036)		0.440*** (0.031)
N (students)	1312	1130	1312	1130	1312	1130

*Notes:* This table shows the impact of the intervention on assessments of math, science, and language administered at endline and endline “audit”, about nine months after the rollout of the intervention (April 2018), and at follow-up, about 11 months after the rollout of the intervention (June 2018) among teachers who obtained a proportion-correct score between 60 and 80% on an assessment of instructional practice, knowledge of student errors, and content knowledge in math. Endline scores are standardized to have mean zero and standard deviation one in the control group. The p-value of the sum is from a test that the sum of the main effect of treatment and its interaction is equal to zero. Standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table A.17: Impact on student attitudes, mindsets, and aspirations (endline)

	(1) Control	(2) Treatment	(3) Difference	(4) N
<i>A. Math and science</i>				
Share of students who...				
...enjoy learning math	0.814 [0.390]	0.830 [0.376]	0.017 (0.017)	2,524
...enjoy learning science	0.740 [0.439]	0.773 [0.419]	0.033 (0.020)	2,524
...wish they did not have math	0.277 [0.448]	0.257 [0.437]	-0.019 (0.026)	2,524
...wish they did not have science	0.278 [0.448]	0.257 [0.437]	-0.020 (0.028)	2,524
...feel nervous about math	0.422 [0.494]	0.380 [0.486]	-0.041 (0.027)	2,524
...feel nervous about science	0.402 [0.491]	0.404 [0.491]	0.002 (0.028)	2,524
...find math useful for life	0.728 [0.445]	0.761 [0.427]	0.033* (0.020)	2,524
...find science useful for life	0.721 [0.449]	0.725 [0.446]	0.005 (0.020)	2,524
...stop trying when math gets hard	0.234 [0.424]	0.262 [0.440]	0.029 (0.025)	2,524
...stop trying when science gets hard	0.193 [0.395]	0.234 [0.424]	0.041* (0.022)	2,524
Composite index (std.)	-0.000 [1.000]	0.038 [0.992]	0.039 (0.041)	2,524
<i>B. Intelligence</i>				
Share of students who believe...				
...people cannot change their intelligence	0.538 [0.499]	0.519 [0.500]	-0.019 (0.028)	2,524
...people have a fixed amount of intelligence	0.461 [0.499]	0.447 [0.497]	-0.014 (0.028)	2,524
...only some people are people	0.439 [0.496]	0.424 [0.494]	-0.014 (0.025)	2,524
...boys are more intelligent than girls	0.385 [0.487]	0.365 [0.482]	-0.020 (0.025)	2,524
...boys are better at math and science	0.371 [0.483]	0.357 [0.479]	-0.014 (0.029)	2,524
Composite index (std.)	-0.000 [1.000]	-0.052 [1.015]	-0.050 (0.061)	2,524
<i>C. Aspirations</i>				
Share of students who...				
...want to continue studying after high school	0.675 [0.469]	0.696 [0.460]	0.021 (0.024)	2,246
...want to study a STEM subject in high school	0.756 [0.430]	0.800 [0.400]	0.044 (0.027)	2,239
...want a STEM-related job	0.333 [0.472]	0.334 [0.472]	0.000 (0.027)	2,231
Composite index (std.)	-0.000 [1.000]	0.044 [0.975]	0.042 (0.052)	2,216

*Notes:* This table compares students' attitudes, mindsets, and aspirations in the control and treatment groups based on surveys administered at endline about nine months after the rollout of the intervention (April 2018). It shows the means and standard deviations for each group (columns 1-2) and tests for differences between groups including randomization-strata (i.e., grade) fixed effects (column 3). Standard deviations appear in brackets, and standard errors (clustered by school) appear in parentheses. The composite indexes are the first principal components of the variables in each panel, standardized with respect to the control group. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table A.18: Heterogeneous impact on student beliefs by students' sex (endline)

	STEM	Intelligence	Aspirations
	(1)	(2)	(3)
<i>A. Endline</i>			
Treatment	0.067 (0.067)	-0.017 (0.068)	0.115 (0.076)
Female	0.080 (0.072)	-0.326*** (0.081)	0.220*** (0.079)
Treat $\times$ Female	-0.010 (0.091)	-0.022 (0.112)	-0.058 (0.096)
P-value of the sum	0.211	0.000	0.014
N (students)	2258	2258	2009

*Notes:* This table shows the impact of the intervention on surveys administered at endline about nine months after the rollout of the intervention (April 2018), by students' sex. Indexes are standardized to have mean zero and standard deviation one in the control group. The p-value of the sum is from a test that the sum of the main effect of treatment and its interaction is equal to zero. Standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table A.19: Heterogeneous impact on student beliefs by students' caste (endline)

	STEM	Intelligence	Aspirations
	(1)	(2)	(3)
<i>A. Endline</i>			
Treatment	0.060 (0.063)	-0.115* (0.068)	0.148** (0.065)
Scheduled caste/tribe	-0.090 (0.056)	-0.112** (0.048)	0.042 (0.067)
Treat $\times$ SC/ST	0.009 (0.081)	0.122 (0.085)	-0.114 (0.094)
P-value of the sum	0.158	0.882	0.248
N (students)	2236	2236	1989

*Notes:* This table shows the impact of the intervention on surveys administered at endline about nine months after the rollout of the intervention (April 2018), by students' socio-economic status. Indexes are standardized to have mean zero and standard deviation one in the control group. The p-value of the sum is from a test that the sum of the main effect of treatment and its interaction is equal to zero. Standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table A.20: Heterogeneous impact on student beliefs by students' socio-economic status (endline)

	STEM	Intelligence	Aspirations
	(1)	(2)	(3)
<i>A. Endline</i>			
Treatment	0.066 (0.043)	-0.032 (0.065)	0.084 (0.052)
SES index	0.007 (0.039)	0.040 (0.032)	0.038 (0.035)
Treat $\times$ SES index	-0.014 (0.048)	-0.082* (0.045)	0.025 (0.048)
P-value of the sum	0.806	0.176	0.102
N (students)	2255	2255	2006

*Notes:* This table shows the impact of the intervention on surveys administered at endline about nine months after the rollout of the intervention (April 2018), by students' socio-economic status. Indexes are standardized to have mean zero and standard deviation one in the control group. The p-value of the sum is from a test that the sum of the main effect of treatment and its interaction is equal to zero. Standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.



Table A.21: Heterogeneous impact on student beliefs by baseline score (endline)

	STEM	Intelligence	Aspirations
	(1)	(2)	(3)
<i>A. Endline</i>			
Treatment	0.063 (0.042)	-0.050 (0.064)	0.094* (0.055)
Baseline score	0.034 (0.022)	-0.042* (0.023)	0.111*** (0.023)
Treat $\times$ Baseline	-0.046 (0.028)	0.018 (0.037)	0.008 (0.027)
P-value of the sum	0.568	0.396	0.000
N (students)	2307	2307	2052

*Notes:* This table shows the impact of the intervention on surveys administered at endline about nine months after the rollout of the intervention (April 2018), by students' socio-economic status. Indexes are standardized to have mean zero and standard deviation one in the control group. The p-value of the sum is from a test that the sum of the main effect of treatment and its interaction is equal to zero. Standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table A.22: Impact on student attendance (unannounced visits, school registers, and endline)

	(1) Control	(2) Treatment	(3) Difference	(4) N
<i>A. Unannounced visits</i>				
Share of students observed at school	0.787 [0.410]	0.789 [0.408]	0.002 (0.028)	2,442
<i>B. School registers</i>				
Share of students marked as...				
...absent 0 days this week	0.551 [0.498]	0.544 [0.498]	-0.008 (0.025)	2,619
...absent 1-2 days this week	0.274 [0.446]	0.271 [0.445]	-0.003 (0.022)	2,619
...absent 3-4 days this week	0.074 [0.261]	0.083 [0.276]	0.009 (0.012)	2,619
...absent 5+ days this week	0.101 [0.301]	0.102 [0.303]	0.001 (0.017)	2,619
<i>C. Endline student survey</i>				
Share of students who reported being...				
...late 0 days this week	0.401 [0.490]	0.388 [0.487]	-0.014 (0.030)	2,219
...late 1-2 days this week	0.299 [0.458]	0.337 [0.473]	0.038** (0.018)	2,219
...late 3-4 days this week	0.165 [0.372]	0.140 [0.348]	-0.025 (0.020)	2,219
...late 5+ days this week	0.135 [0.341]	0.135 [0.342]	0.001 (0.018)	2,219
...absent 0 days this week	0.410 [0.492]	0.412 [0.492]	0.002 (0.031)	2,222
...absent 1-2 days this week	0.334 [0.472]	0.308 [0.462]	-0.026 (0.022)	2,222
...absent 3-4 days this week	0.128 [0.334]	0.148 [0.356]	0.021 (0.017)	2,222
...absent 5+ days this week	0.129 [0.335]	0.133 [0.339]	0.003 (0.017)	2,222

*Notes:* This table compares the attendance of students in the control and treatment groups based on various measures collected during unannounced visits to schools about seven months after the rollout of the intervention (February 2018) and at endline two months later (April 2018). It shows the means and standard deviations for each group (columns 1-2) and tests for differences between groups including randomization-strata (i.e., grade) fixed effects (column 3). Standard deviations appear in brackets, and standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table A.23: Impact on student demand for tuition (endline)

	(1) Control	(2) Treatment	(3) Difference	(4) N
<i>A. Subjects</i>				
Share of students who attend tuition...				
...in math	0.323 [0.468]	0.352 [0.478]	0.030 (0.026)	2,195
...in science	0.256 [0.437]	0.286 [0.452]	0.029 (0.026)	2,173
...in language	0.238 [0.426]	0.266 [0.442]	0.029 (0.028)	2,151
...in English	0.318 [0.466]	0.345 [0.476]	0.027 (0.028)	2,196
<i>B. Duration</i>				
Share of students who attend tuition...				
...less than 2 hours per week	0.036 [0.186]	0.038 [0.192]	0.002 (0.010)	2,237
...2-4 hours per week	0.097 [0.296]	0.099 [0.299]	0.003 (0.015)	2,237
...4-6 hours per week	0.048 [0.215]	0.046 [0.210]	-0.001 (0.011)	2,237
...more than 6 hours per week	0.216 [0.412]	0.252 [0.434]	0.036* (0.020)	2,237

*Notes:* This table compares students' demand for private tuition in the control and treatment groups based on surveys administered at endline about nine months after the rollout of the intervention (April 2018). It shows the means and standard deviations for each group (columns 1-2) and tests for differences between groups including randomization-strata (i.e., grade) fixed effects (column 3). Standard deviations appear in brackets, and standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table A.24: Impact on class materials (announced observations)

	(1) Control	(2) Treatment	(3) Col. (2)-(1)	(4) N
Class has blackboard	0.980 [0.141]	1.000 [0.000]	0.020 (0.020)	150
Class has whiteboard	0.020 [0.141]	0.000 [0.000]	-0.020 (0.020)	150
Class has chalk/markers	0.980 [0.141]	1.000 [0.000]	0.020 (0.019)	150
Class has textbook for teachers	0.820 [0.388]	0.300 [0.461]	-0.520*** (0.083)	150
Class has textbook for students	0.720 [0.454]	0.260 [0.441]	-0.460*** (0.098)	150
Class has laptop	0.060 [0.240]	0.040 [0.197]	-0.020 (0.045)	150
Class has digital whiteboard	0.040 [0.198]	0.040 [0.197]	0.000 (0.004)	150
Class has LCD projector	0.080 [0.274]	0.060 [0.239]	-0.020 (0.036)	150
Class has TV	0.160 [0.370]	0.060 [0.239]	-0.100 (0.066)	150
Class has science/math equipment	0.300 [0.463]	0.160 [0.368]	-0.140* (0.081)	150
Class has maps	0.240 [0.431]	0.120 [0.327]	-0.120 (0.078)	150
Class has charts/poster	0.760 [0.431]	0.780 [0.416]	0.020 (0.080)	150
Class has toys/games	0.080 [0.274]	0.040 [0.197]	-0.040 (0.040)	150
Composite index (std.)	0.000 [1.000]	-0.126 [0.711]	-0.126 (0.101)	150

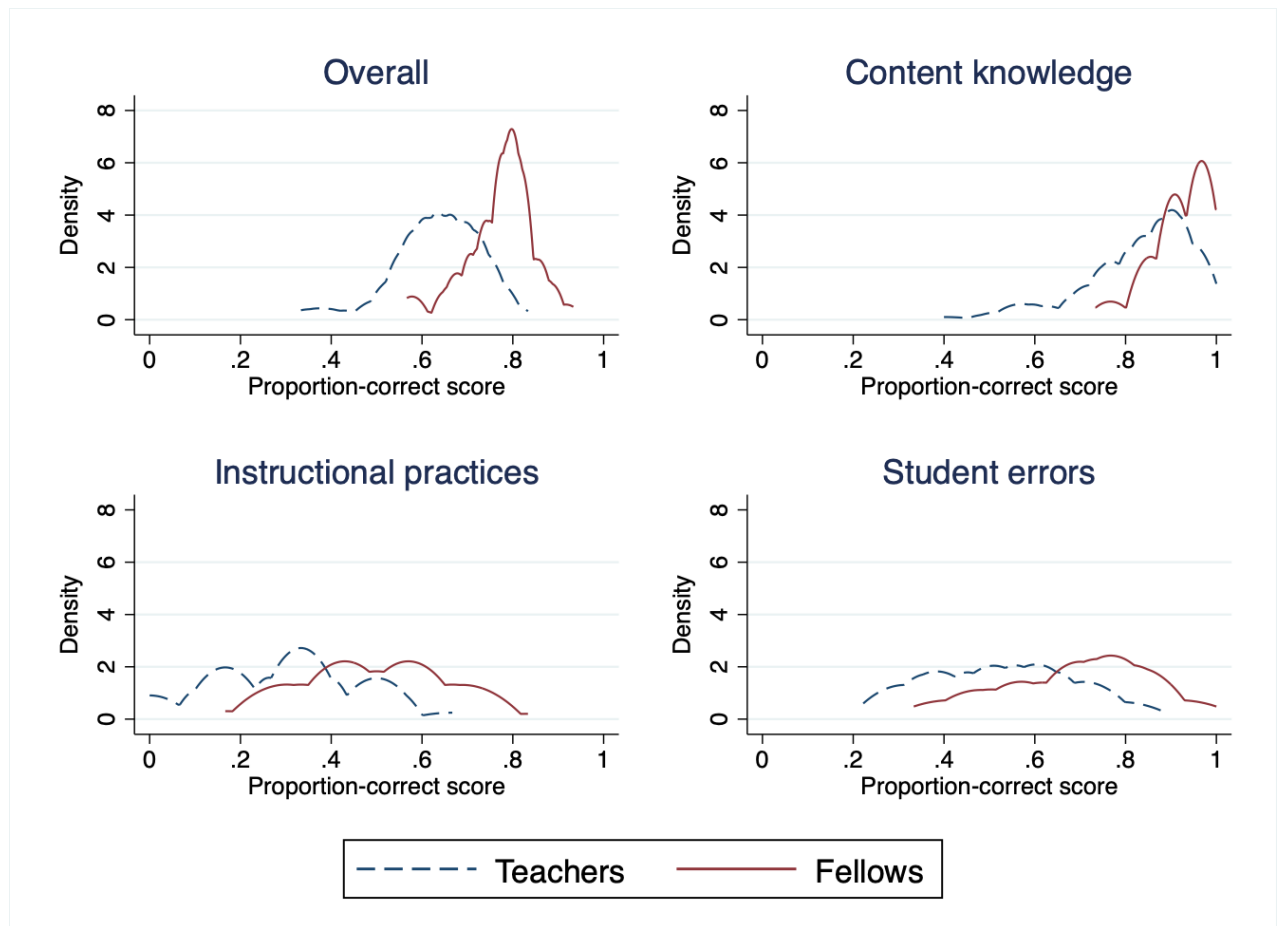
*Notes:* This table compares the availability of materials in control and treatment classes, based on announced observations about five months after the rollout of the intervention (December 2017). Standard deviations appear in brackets, and standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table A.25: Impact on variability of standardized test scores (endline, endline audit, and follow-up)

	Math		Science		Language	
	(1)	(2)	(3)	(4)	(5)	(6)
<i>A. Endline</i>						
Treatment	-0.001 (0.035)	0.014 (0.034)	0.015 (0.040)	0.030 (0.038)	-0.053* (0.030)	-0.061* (0.033)
Baseline score		0.287** (0.122)		0.330** (0.163)		0.419** (0.207)
N (students)	93	93	93	93	93	93
<i>B. Endline audit</i>						
Treatment	-0.119 (0.075)	-0.114 (0.069)	0.015 (0.064)	0.017 (0.066)		
Baseline score		0.080 (0.259)		0.039 (0.240)		
N (students)	91	91	91	91		
<i>C. Follow-up</i>						
Treatment	0.006 (0.039)	0.032 (0.038)	0.129*** (0.035)	0.136*** (0.037)	-0.000 (0.035)	-0.007 (0.034)
Baseline score		0.432*** (0.131)		0.157 (0.189)		0.437*** (0.153)
N (students)	94	94	94	94	94	94

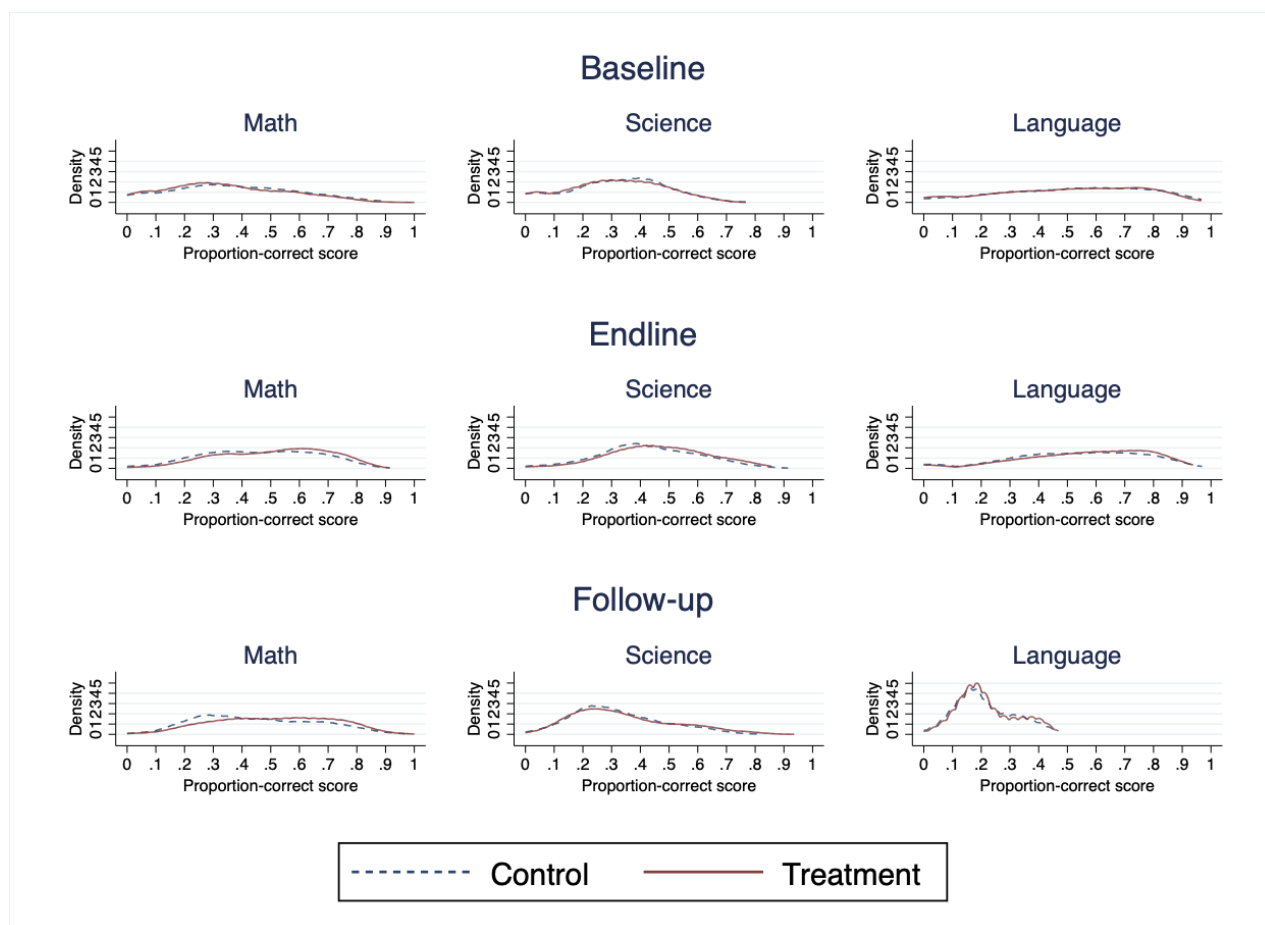
*Notes:* This table compares the variability in students' standardized scores in the control and treatment groups based on assessments administered at endline and endline "audit", about nine months after the rollout of the intervention (April 2018) and at follow-up, about 11 months after the rollout of the intervention (June 2018). Scores are standardized to have mean zero and standard deviation one in the control group. Standard errors (clustered by school) appear in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Figure A.1: Distributions of proportion-correct scores on teacher assessments



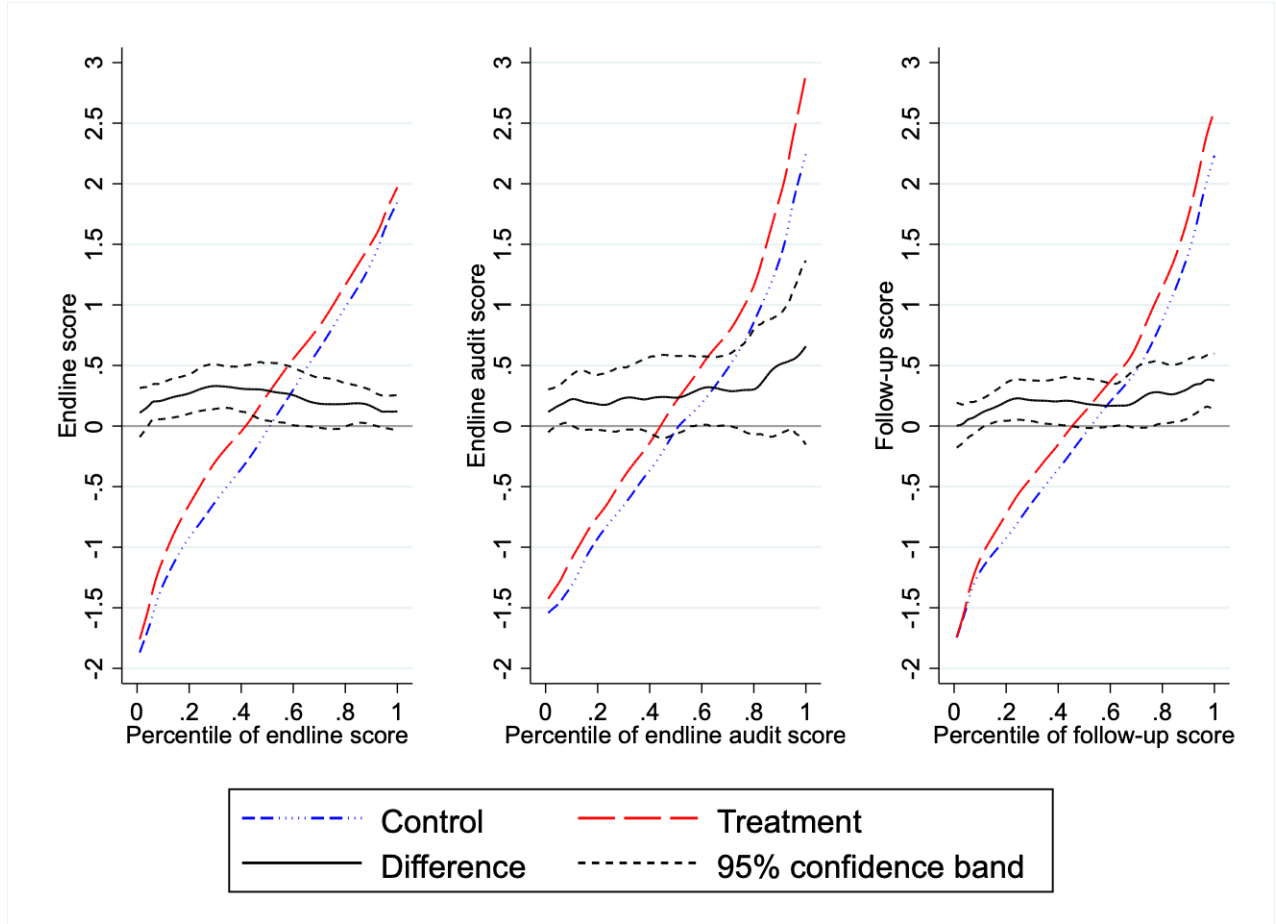
*Notes:* The figure shows the distribution of proportion-scaled scores on the teacher assessments. Proportion-correct scores indicate the proportion of items on each test answered correctly.

Figure A.2: Distributions of proportion-correct scores on student assessments



*Notes:* The figure shows the distribution of proportion-scaled scores on the student assessments. Proportion-correct scores indicate the proportion of items on each test answered correctly.

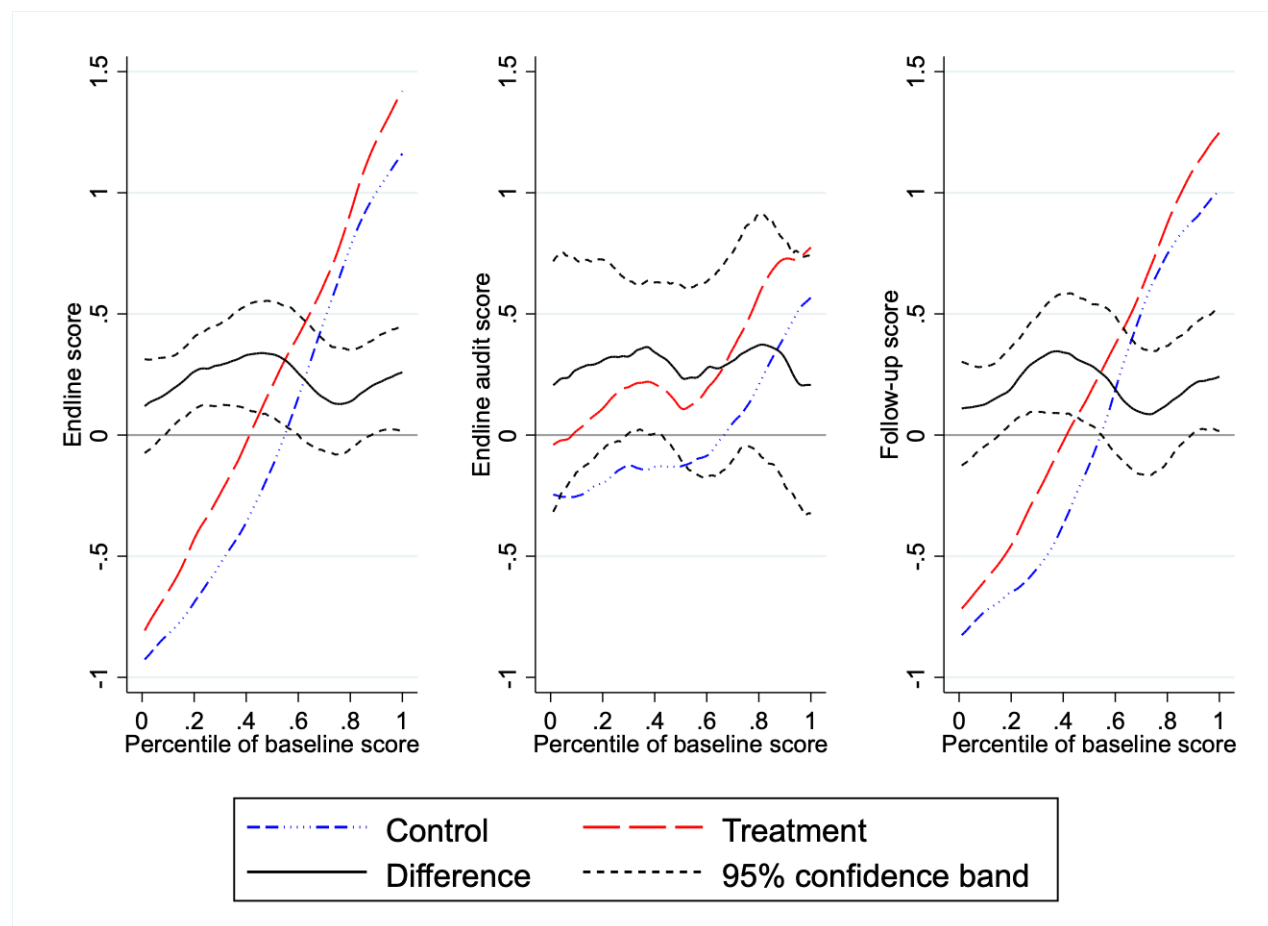
Figure A.3: Quantile treatment effects (endline, endline audit, and follow-up)



*Notes:* The figure shows quantiles of endline, endline audit, and follow-up composite scores (the first principal component of a principal component analysis of all subjects assessed at each round) for treatment and control students who participated in the baseline and each of these rounds, estimated by polynomial regressions of each round's scores on that round's percentiles separately by experimental group. Dashed black lines display bootstrapped 95% confidence intervals.

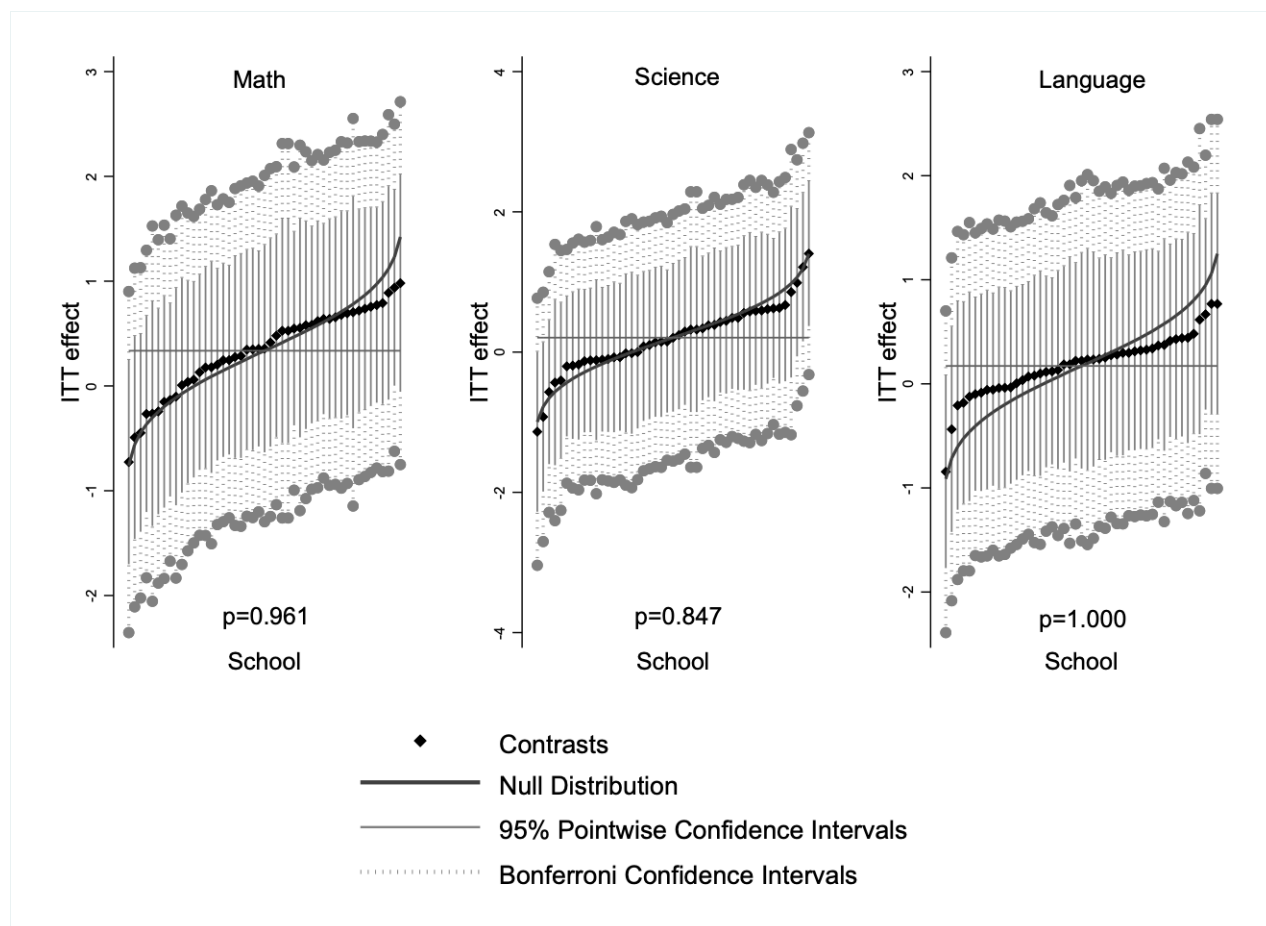


Figure A.4: Average treatment effects by baseline score (endline, endline audit, and follow-up)



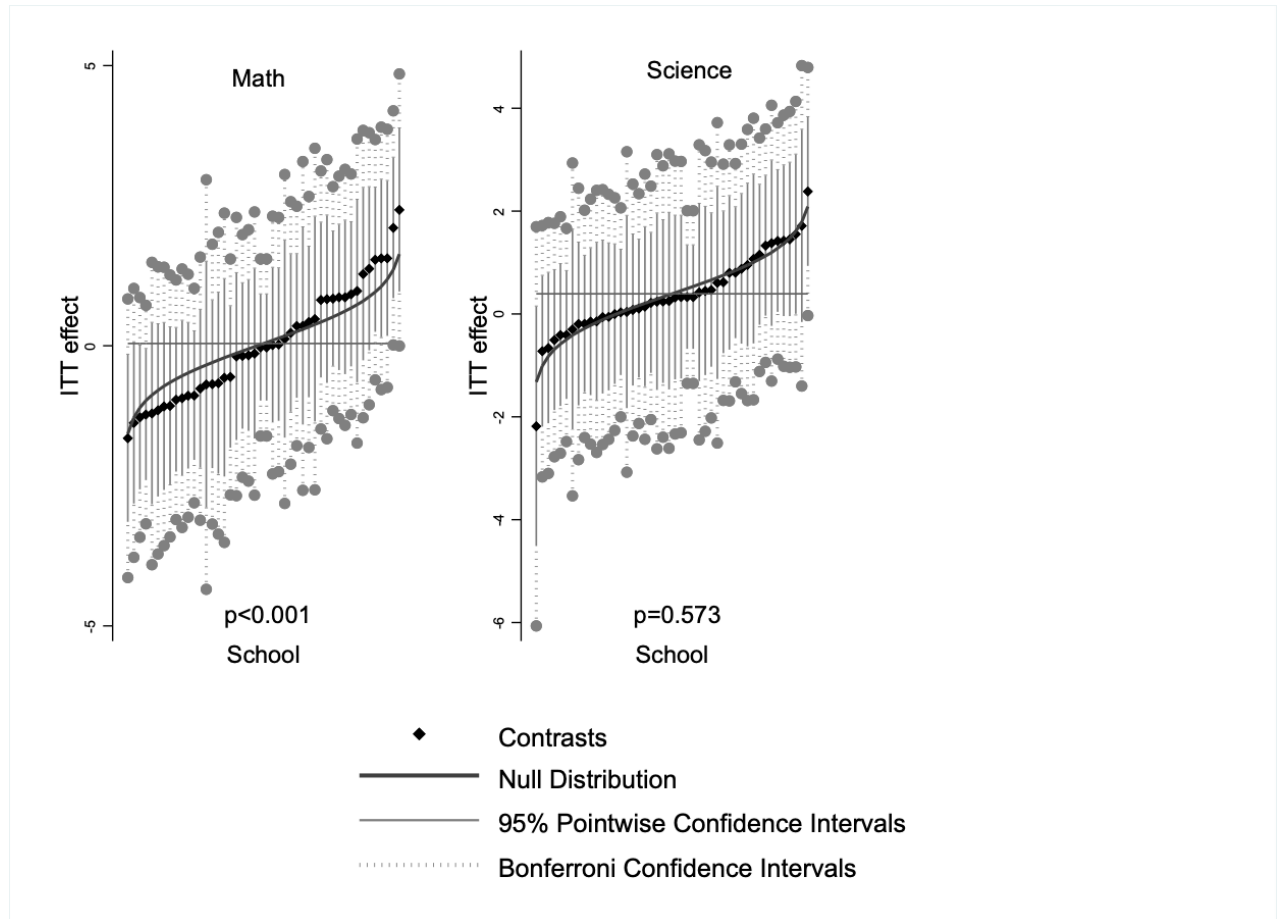
*Notes:* The figure shows estimates of average endline, endline audit, and follow-up composite scores (the first principal component of a principal component analysis of all subjects assessed at each round) and treatment effects at each percentile of baseline composite score for treatment and control students who participated in the baseline and each of these rounds, estimated by polynomial regression. Dashed black lines display bootstrapped 95% confidence intervals.

Figure A.5: Heterogeneous impact on standardized test scores by school (endline)



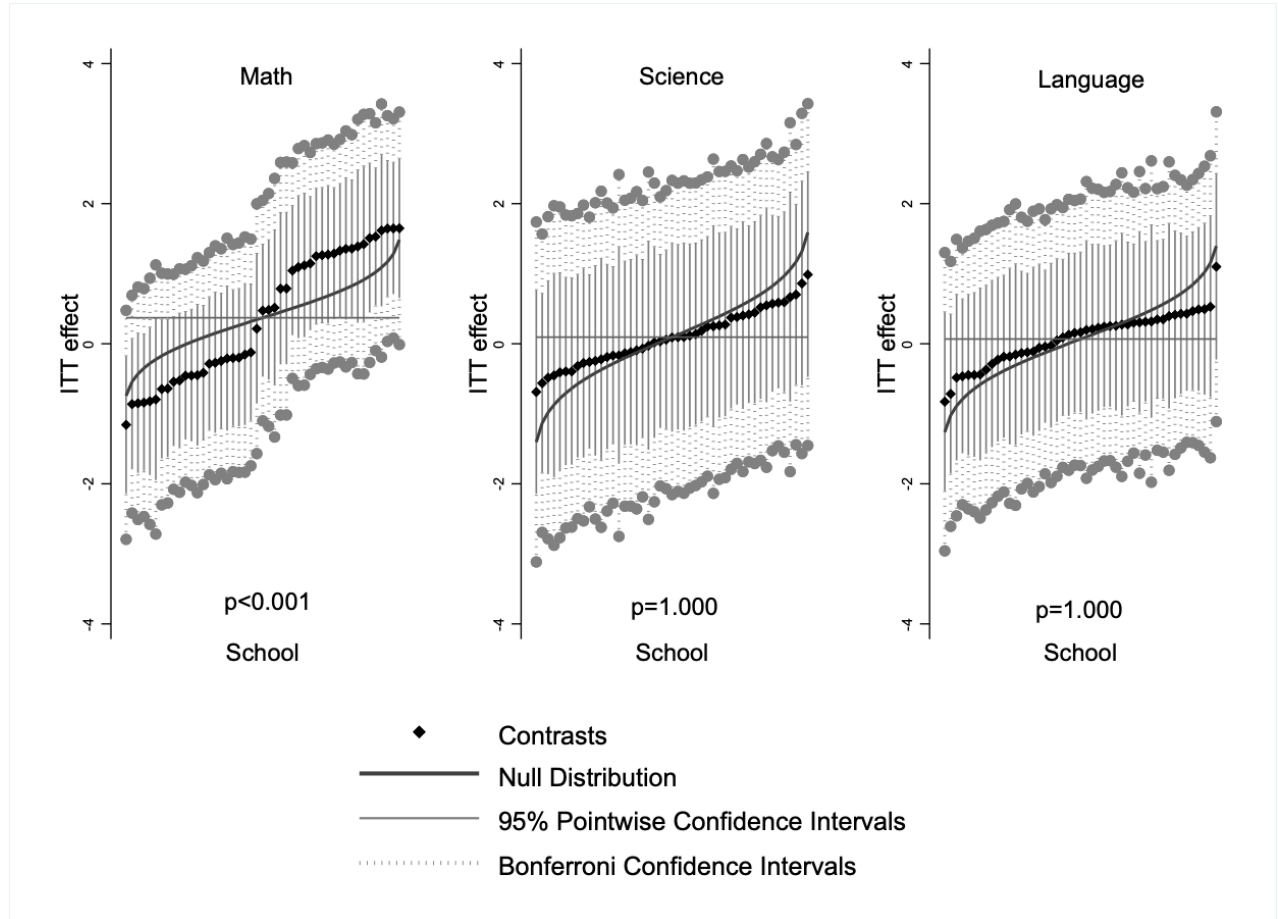
*Notes:* The figure provides a “caterpillar plot” (von Hippel and Bellows, 2018) of impacts on standardized scores by school at endline, about nine months after the rollout of the intervention (April 2018). Each black dot refers to the point estimate for a given school. Bonferroni confidence intervals adjust standard errors for multiple hypothesis testing. The black solid line shows the null distribution of “effects” that can be expected due to error. The p-values are for a test of the null hypothesis of no heterogeneity.

Figure A.6: Heterogeneous impact on standardized test scores by school (endline audit)



*Notes:* The figure provides a “caterpillar plot” (von Hippel and Bellows, 2018) of impacts on standardized scores by school at the endline “audit,” about nine months after the rollout of the intervention (April 2018). Each black dot refers to the point estimate for a given school. Bonferroni confidence intervals adjust standard errors for multiple hypothesis testing. The black solid line shows the null distribution of “effects” that can be expected due to error. The p-values are for a test of the null hypothesis of no heterogeneity.

Figure A.7: Heterogeneous impact on standardized test scores by school (follow-up)



*Notes:* The figure provides a “caterpillar plot” (von Hippel and Bellows, 2018) of impacts on standardized scores by school at follow-up, about 11 months after the rollout of the intervention (June 2018). Each black dot refers to the point estimate for a given school. Bonferroni confidence intervals adjust standard errors for multiple hypothesis testing. The black solid line shows the null distribution of “effects” that can be expected due to error. The p-values are for a test of the null hypothesis of no heterogeneity.