

# Can Teaching Be Taught? Improving Teachers' Pedagogical Skills at Scale in Rural Peru<sup>1</sup>

<b>Juan F. Castro</b>	<b>Paul Glewwe</b>	<b>Alexandra Heredia-Mayo</b>	<b>Stephanie Majerowicz Nieto</b>	<b>Ricardo Montero</b>
Universidad del Pacifico	University of Minnesota	Universidad del Pacifico	Universidad de los Andes	University of Minnesota

July, 2023

## Abstract

We evaluate the impact of a large-scale teacher coaching program in Peru, a context with high teacher turnover, on teachers' pedagogical skills and student learning. Previous studies find that small-scale coaching programs can improve teaching of reading and science in developing countries. However, scaling up can reduce programs' effectiveness, and teacher turnover can erode compliance and cause spillovers onto non-program schools. We develop a framework that defines different treatment effects when teacher turnover is present, and explains which effects can be estimated. We evaluate this teacher coaching program, exploiting random assignment of that program's expansion to 3,797 rural schools in 2016. After two years, teachers assigned to the program increased their aggregate pedagogical skills by 0.20 standard deviations. The program also increased student learning; after one year, Grade 2 students' mathematics and reading scores increased by 0.10 and 0.07 standard deviations (of the distributions of those test scores) respectively. After three years, the cumulative effect increases slightly, to 0.11 and 0.10, respectively. One reason why these impacts are low is that some uncoached teachers moved into treated schools in years 2 and 3. Following our framework, we estimate that the impacts on students of having a "fully" coached teacher for all three years are 0.18 and 0.16 standard deviations for mathematics and reading comprehension, respectively.

Keywords: education, teacher coaching, pedagogical skill, student learning, teacher turnover.  
JEL Codes: I21, O15.

---

<sup>1</sup> We would like to thank seminar participants at the Department of Applied Economics of the University of Minnesota, the Department of Economics of Universidad del Rosario, the LACEA 2019 Annual Meeting, the Department of Agricultural and Consumer Economics at the University of Illinois, and the Department of Agricultural Economics and Rural Development at Seoul National University for their valuable comments. We are also grateful to Hugo Fernández for excellent research assistance. Any remaining errors are ours alone.

# 1. Introduction

Teacher quality is an essential determinant of student learning (Das et al. 2007, Clotfelter et al. 2010, Chetty et al. 2014). Yet many teachers lack mastery in the subjects they teach, or lack the pedagogical skills to teach them effectively. This is especially true for teachers in developing countries (World Bank, 2018). Can these teachers' skills be improved?

Every year, developing countries spend over \$1 billion on teacher training (Loyalka et al., 2019). Popova et al. (2016) find that about two thirds of the World Bank educational projects between 2000 and 2012 included in-service teacher training. Such training is attractive because it can be centrally designed and coordinated by the Ministry of Education and is usually supported by teachers' unions (Evans and Popova, 2016).

In this study, we evaluate the impact of a large-scale teacher coaching program, operating in a context of high teacher turnover, on teachers' pedagogical skills as well as on student learning outcomes. Evidence on the impacts of in-service training in developing countries is mixed, and programs vary widely in form and content. A survey by Evans and Popova (2016) found that programs with face-to-face training, follow-up visits, engagement of teachers to obtain their ideas, and adaptation to local context, tend to have larger effects on student learning. Coaching programs often have these features as they involve school visits, classroom observations, and personalized feedback for teachers by trained peers or coaches. Thus, coaching programs are a promising alternative to traditional in-service training that offers intensive sessions to large numbers of teachers at a centralized venue.

When programs are offered at the school level but are intended to operate through teachers, and teachers can move between schools, estimates of the average treatment effect (ATE) of the program based on a randomized control trial may be biased. In particular, movement of teachers across schools may lead to spillovers that will introduce biases when comparing treated and control schools, even when all schools comply with their random assignment and there are no biases due to the selection or attrition of students.

Education interventions that operate through teachers often have all teachers in a school share treatment status (i.e., all teachers are either treated or untreated). Most studies of the effectiveness of these types of interventions focus on student outcomes and compare treatment schools with control schools, and some of them evaluate results after enough time has passed for teachers to switch schools (Lucas et. al. 2014, Jukes et. al. 2017, Cilliers et. al. 2020). These studies usually address potential biases due to student attrition, yet they rarely mention the possibility of teacher turnover or the potential bias it may induce.

This risk of bias may occur not only for education interventions but also for any estimation of treatment effects in cluster randomized control trials (RCTs) with movement of service providers or program beneficiaries across clusters. Indeed, high turnover is reported for many non-education contexts. For example, Kovner et al. (2014) report that 17.5% of new nurses in the U.S. leave their jobs within one year of starting, and Banerjee et al. (2021) find, in their control sample, that one-third of police officers in India changed stations over an 18-month period. Despite its frequency, turnover is usually ignored in program evaluations. For example, Georgiadis and Pitellis (2016) compare treated and control enterprises (clusters) in a job training program but do not discuss the possibility of workers moving across firms.

We make a methodological contribution by developing a framework that clarifies the assumptions and data needed to obtain unbiased estimates of average treatment effects (ATE), intent to treat effects (ITT), and average causal response (ACR, an extension of local average treatment effects (LATE)) in a clustered RCT with movement of service providers across clusters. In our context, this framework explains how treatment effects differ, depending on whether one focuses on a particular set of teachers, following them if they move to other schools (in which case the outcome variables are those teachers' skills), or on the teachers and students in particular schools (in which case the outcome variables are the skills of these schools' teachers and the learning progress of these schools' students). Both sets of treatment effects are highly relevant from a policy perspective. The first set is relevant for policies that focuses on improving the skills of a particular group of teachers, such as teachers whose pedagogical skills are thought to be deficient. The second set is relevant for policies aimed at improving the teaching skills and learning progress, respectively, of the teachers and students in a particular group of schools, such as schools where students' academic performance is particularly low. We show how the latter set of effects depends not only on the direct effect of the program on participating teachers' skills but also on the indirect effect of the program on teacher composition: which teachers stay in these schools, which teachers leave these schools, and which teachers move into these schools. Previous research based on cluster RCTs where service providers move across clusters has ignored these composition effects.

We show that, in general, it is not possible to estimate average treatment effects (ATEs) for teacher skill and student learning, although under certain conditions lower bounds for ATEs can be estimated. We also show that comparisons of teachers in treated and control schools after turnover has occurred will, in general, lead to biased estimates of intent to treat (ITT) effects for teachers *in the program schools when the program started*. However, it is possible to estimate these ITT effects if one has a sample of teachers that follows them when

they change schools, or using the data of teachers in treated and control schools after turnover has occurred *if* that turnover is unrelated to the program. This last result is important because following teachers who change schools and, more generally, following service providers who change clusters, can be difficult, which raises the risk of attrition bias in ITT estimates.

We estimate the effects on teachers' pedagogical skills and on student learning of a teacher coaching program implemented in rural multigrade schools in Peru. Trained coaches visit classrooms and give specific advice to teachers on their pedagogical practices, providing customized strategies to improve them. Identification exploits random assignment of 6,218 schools to treatment and control groups when the program expanded in 2016. Teachers' skills were measured in late 2017 (after nearly two years of treatment) by observing teacher-student interactions and a broad range of instructional practices in a randomly selected subsample of 166 treated and 174 control schools. Student skills were tested in grades 2 and 4 in late 2016 and late 2018, respectively, for all public schools with five or more students in those grades, which provide student test score data for 2,567 of the 6,218 randomly assigned schools.

As in many developing countries, Peru's rural schools have very high rates of teacher turnover;<sup>2</sup> of the teachers in the subsample of 340 schools with teacher skills data, about 43% had moved between 2016 and the start of 2017. Importantly, classroom observation data were collected not only in these 340 schools, but also in many (but not all) of the schools that received the teachers who moved from these schools to other schools between 2016 and 2017.

Our main findings are as follows. For the teachers who, after turnover occurred (i.e. in 2017), were teaching in the schools assigned to the program, we find that the ITT effect of two years of coaching on their pedagogical skills is 0.20 standard deviations (s.d.) of the distribution of those skills. This is also our preferred estimate of the ITT effect on the skills of the teachers in the program schools when the program began, some of whom left those schools in the next year. We also show that this ITT estimate is, under plausible assumptions, a lower bound of the ATEs for both sets of teachers. Turning to specific skills, the largest ITT effects are for lesson planning and, to a lesser extent, encouraging students' critical thinking.

We also estimated treatment effects of the program on student learning after one and three years (we have no data for the second year). After one year, the program increased learning among the Grade 2 students who took the 2016 National Student Evaluation by 0.10 s.d. in mathematics and 0.07 s.d. in reading comprehension (of the distributions of those test scores). These are both ITT and ATE effects, since all teachers followed their random assign-

---

<sup>2</sup> High teacher turnover is common in developing countries: Zeitlin (2021) reports turnover of about 20% per year in Rwanda, and Schaffner, Glewwe and Sharma (2021) report 18-21% turnover per year for teachers in Nepal.

ment in the first year. After three years of exposure, the ITT effect increases only slightly, to 0.11 s.d. for mathematics and 0.10 s.d. for reading comprehension; these estimates, which are lower bounds for ATE (which cannot be estimated in year 3), reflect the fact that many teachers in program schools in year 3 did not have three full years of coaching, and some teachers who had moved to control schools by year 3 had been coached in previous years. The average causal response (ACR) estimates after three years, which adjust the ITT estimates to estimate the impact of three years of exposure to teachers who were coached in all three years, are 0.18 and 0.16 s.d. for mathematics and reading comprehension, respectively.

Our estimates for the effect of coaching on pedagogical skills are smaller than those found in developed countries (0.49 s.d. on instructional practices, see Kraft et al., 2018). This may reflect the scale of the program, and Peru's high rate of teacher turnover. Yet we address two unresolved questions on coaching's impact on teachers' pedagogical skills in developing countries. We show that: (i) A program implemented at scale, even with high teacher turnover, can still exhibit positive impacts; and (ii) *General* pedagogical skills can be increased.

Furthermore, while our estimated effects on student learning may seem small, they are similar, and in one sense larger, than those typically found in developing countries. Evans and Yuan (2022) surveyed 224 education studies and found that the median effect on learning outcomes is 0.10 s.d., and these effect sizes decrease with the size of the study. For large studies, those with over 5,000 students, the median effect is only 0.05 s.d.

To our knowledge, no previous study has evaluated the effects on pedagogy and student learning of a large-scale teacher coaching program in a developing country. Most in-service training programs evaluated in the developing world are small-scale pilots or efficacy trials run by researchers or NGOs (Evans and Popova, 2016). For example, Cilliers et al. (2020) compared the impacts of coaching and centralized teacher training on student reading skills in 180 public schools in South Africa, and Albornoz et al. (2020) estimated the impact of teacher coaching to improve student learning of science in 70 public schools in Argentina. In contrast, we evaluate a program randomly implemented in 3,797 rural schools in Peru.

The issue of scale is relevant for coaching programs' effectiveness because of two features of this type of in-service training. First, the program's success depends on the supply of qualified coaches. If these skills are scarce, expanding the program likely will reduce its quality, and thus its effectiveness. Second, classroom observation and personalized feedback requires coaches to travel to several schools. This can be costly and can complicate program delivery if scaling-up implies serving schools in very remote areas. This is very likely for rural schools in developing countries, whose teachers often require additional training.

Teacher turnover not only complicates identification of program effects, as discussed above, but may also make coaching less effective by reducing compliance. Teachers who leave a school before the program ends may not receive the full “dose” of coaching, and program schools that receive new teachers may have staff who are only partially trained.

We know of only one other study that considered teacher turnover when evaluating a teacher training program. Clare et al. (2010) estimated the effect of a literacy coaching program in 32 elementary schools in Texas. Stressing how such turnover can thwart schools’ efforts to improve instruction through teacher training, the authors estimated the program’s effect on the reading skills of the students of teachers recruited to replace those who left their school in the first year of the program. They found a positive association between teachers’ program participation and their students’ reading skills. However, the non-random composition of their sample (recruited teachers in program and non-program schools may not be comparable) casts doubt on the causal interpretation of their results.

Finally, the literature thus far does not provide a clear indication as to whether coaching can improve *general* pedagogical skills. Most evaluations of coaching programs focus on pedagogy for a specific topic or course. For example, Albornoz et al. (2020) focus on improving teaching of science, and Cilliers et al. (2020) focus on reading. Kraft et al. (2018) highlight a lack of causal evidence on the effect of coaching for subjects other than reading or literacy. Some papers measure the effect of training on teacher time allocation (Bruns et al. 2018) or on using specific types of teaching (Kotze et al. 2019), but not on their teaching skills. The pedagogical skills of public-school teachers in developing countries are generally low, and a key policy question is whether coaching can improve a broad set of teaching skills.

The rest of the paper is organized as follows. Section 2 describes the program and explains the evaluation design. Section 3 presents our analytical framework, defines several treatment effects, and explains which can be estimated. Sections 4 and 5 present estimates of the program’s impact on teachers’ pedagogical skills and on student learning, respectively. Section 6 provides concluding remarks, policy implications, and advice for future research.

## **2. The Coaching Program and its Evaluation Design**

### **2.1 Teacher Hiring and Movement in Peru**

There are two types of teachers in the Peruvian school system: Tenured (civil servant) teachers (*nombrados*), who have a permanent position in a particular school, and contract

teachers (*contratados*) on temporary one-year contracts who are filling in for tenured teachers who are temporarily absent or for unfilled vacancies in particular schools. In the schools we consider – multigrade and monolingual – most (70-75%) of teachers are tenured.

Teachers become tenured through a selection process with two stages. The first stage consists of a nationally administered exam that covers reading comprehension, logical reasoning, and knowledge of pedagogical practices. Teachers with the minimum passing grade on the exam proceed to a second stage that is carried out by regional education offices and includes an interview and in-classroom observation of teaching practices.

Teachers who do not reach a minimum passing grade in exam in the first stage of the selection process, or who receive a passing grade but are unsuccessful at the second stage, can fill temporary teaching positions as contract teachers (and can continue trying to obtain tenure). Contract teachers have annual contracts: at the end of each school year they must apply for either a renewed contract at their current school or for a contract position at another school. When applying to new schools, contract teachers can apply to as many schools as they want within one region. They are then ranked according to their scores on the latest exam, and teachers with the best scores get their top priority of schools. Teachers can maximize their probability of getting placed by ranking as many schools as they are willing to go to, and by selecting less popular schools (for example, schools located in remote rural areas).

Tenured teachers tend to move less frequently given their permanent position in their schools, but they can request a transfer to another tenured position. In order to do this, they must meet three requirements: have been in a tenured position for at least three years, have been in the *current* tenured position for at least 2 years, and cannot move to another school within the same school district (Peru has about 250 school districts).

## **2.2 The Coaching Program**

In 2010, the Peruvian government initiated coaching programs to improve public primary school teachers' pedagogical practices. As per Ministry of Education guidelines, the school district authority (UGEL) hires coaches for teachers in the schools targeted by the program, who are selected from top-performing teachers. Coaches must have a pedagogical college or university degree, five or more years of primary school teaching experience, and at least one year of experience training or providing support to teachers. Administrative data show that coaches rank much higher than other teachers in the Ministry's teacher evaluations. Coaches were paid the equivalent of US\$1,200 per month, about double the average teacher's wage.

The Ministry of Education sets the standards for hiring coaches, and for the general program design, but the UGELs select and hire the coaches. Each coach works with eight teachers, and UGELs decide how to match coaches to teachers. Coaches are hired annually. About 20% continue for another year, but only 5% stay in the same school the next year.

The coaching program is a substantial investment by Peru's government, costing over US\$ 130 million per year.<sup>3</sup> By 2016, teachers in over 14,000 public schools with more than 900,000 students were being coached under several coaching programs. Over 90% of these schools are primary schools. There are three versions of the program for primary schools: (i) bilingual coaching (for schools where most students speak a Peruvian indigenous language); (ii) monolingual multigrade coaching (for schools where most students speak Spanish and there are fewer teachers than grades taught); and (iii) monolingual full-teacher coaching (for schools large enough to have one teacher per grade and where most students speak Spanish).

This paper evaluates the second type of coaching program,<sup>4</sup> which operates primarily in rural areas.<sup>5</sup> Over 90% of Peru's rural public primary schools are multigrade, which typically have two teachers and about 30 students. Rural multigrade schools are the majority of schools with coaching programs. The monolingual multigrade program is particularly expensive because the target schools tend to be very far apart, so the program requires a large number of coaches and significant travel expenses. This version of the program alone, called *Acompañamiento Pedagógico Multigrado* (APM) in Spanish, cost the government about 40 million US\$ in 2016 and served 174,000 students. This implies an annual cost of 228 US\$ per student, which is over 20% of the total expenditure per student in Peru's primary schools (in 2015, average spending per primary school student was 2,800 soles, or about 940 US\$).

A coach's work consists of several tasks. First, the coach meets the school principal and gathers information about the educational context. Then, the coach attends all teachers' class sessions (one teacher per day) to observe their classroom performance and make an initial diagnostic assessment. The coach uses this assessment to identify the competencies that the teachers must improve and develops an improvement plan with each teacher. During the school year, the coach observes eight more of each teacher's class sessions at regular intervals. The program is usually implemented for three consecutive years. After each classroom observation, the coach and the teacher meet to discuss the progress made in terms

---

<sup>3</sup> It was not implemented in 2021 and 2022 due to Covid-19, after which it was restarted, but on a smaller scale.

<sup>4</sup> Although the three types of coaching programs have some differences (such as the teacher-to-coach ratio or the bilingual certification of coaches), what happens during the coaching sessions is very similar in all three versions.

<sup>5</sup> About 95% of the 6,218 schools in our study are located in rural areas.



of the improvement plan. The coach sends monthly and quarterly reports to the UGEL, and to the school principal, on each teacher's progress and on areas for improvement. At the end of the year, the coach provides a final feedback session for each teacher, collecting his or her impressions of the process, and then writes a final report for each teacher on the achievements, actions, and areas requiring further effort, referencing the initial improvement plan.

In addition to the classroom observations, each coach organizes eight workshops per year for his or her teachers to discuss pedagogical practices and encourage the exchange of ideas. In the workshops, all the teachers for a given coach gather with the coach to discuss a particular pedagogical topic of interest. The coach encourages and guides the exchange of ideas and successful practices among teachers and provides theoretical support on the chosen subject. At the end of each workshop, the group chooses a new topic for the next gathering.

Instead of content knowledge of the material, the program focuses on strengthening pedagogical skills and on developing the ability of teachers to periodically reflect on their own strengths and weaknesses and adjust their behavior accordingly:

“The pedagogical coaching promotes the development and strengthening of skills related to understanding the student in her context, curricular planning, guiding learning, ensuring a safe school environment, and evaluating student learning. In addition, it promotes the development of critical thinking skills like self-reflection and analysis, through exercises that seek reflection and critical analysis of the teacher's own performance.” (APM Manual)

APM uses a cascade system. Each coach is trained, supported and monitored by a pedagogical specialist. Each specialist is required to monitor each coach at least twice per year during the coach's classroom visits. The specialist also provides two workshops per year directly to teachers. Coaches and specialists follow the “Framework for Good Teaching Performance” developed by the Ministry of Education to guide their training. The framework specifies nine skills that teachers should master; the program focuses on seven of these skills:

- Knowledge and comprehension of the students' characteristics and backgrounds.
- Collaborative class preparation with other teachers in the same school.
- Fostering an environment that promotes learning, democratic values, and diversity.
- Guiding the learning process through mastery of the curricular content and the use of effective pedagogical strategies and resources.
- Permanent evaluation of the learning process and provision of feedback to students.
- Active participation in the school management.
- Fostering relationships of respect and collaboration with school community members.

## 2.3 Evaluation Design

In 2016, the APM program was expanded in a way that involved random assignment. All schools that started APM before 2016, and had not yet completed the full three years of the program, continued to participate in APM and were not part of the experimental sample. Monolingual multigrade schools that had low scores on Peru's Grade 2 national student evaluation and had not yet participated in APM were randomized into treatment and control groups. Of the 6,218 eligible schools, which we call the randomized expansion schools, 3,797 were randomly assigned to the treatment group and started the APM program in February of 2016 (Peru's school year runs from March to November). Henceforth, we call these schools APM schools. The other 2,421 schools, the control group, which we call non-APM schools, did not participate in any coaching program in 2016, 2017 and 2018. This randomization, shown in Figure 1, was stratified at the region (department) level, Peru's highest level of political division (Peru has 26 regions).<sup>6</sup>

The sample size is reduced by the availability of outcome data. Standardized tests scores are available only for the 2,567 schools with five or more students in the grade being tested, and the pedagogical skills were measured only for a stratified (at the region level) random subsample (340 schools, 166 APM, 174 non-APM) of the full experimental sample.<sup>7</sup> Appendix Table A1 provides summary statistics for the 6,218 randomized expansion schools (Column (3)), as well as for the wider population of all Peruvian public primary schools (Column (1)), and all monolingual multigrade public primary schools (Column 2), which is the target population of the APM program. Understandably, the randomized expansion schools tend to be smaller and more rural than the average public primary school, which includes large schools in main cities. They also differ in access to the internet and to computers, and in the quality of school infrastructure. Differences are much smaller in access to textbooks and workbooks: more than 70% of schools across our samples receive textbooks, and just above 65% receive workbooks. Finally, all schools have a similarly high percentage of teachers with degrees (97%) and similar teacher-student ratios and school-day lengths.

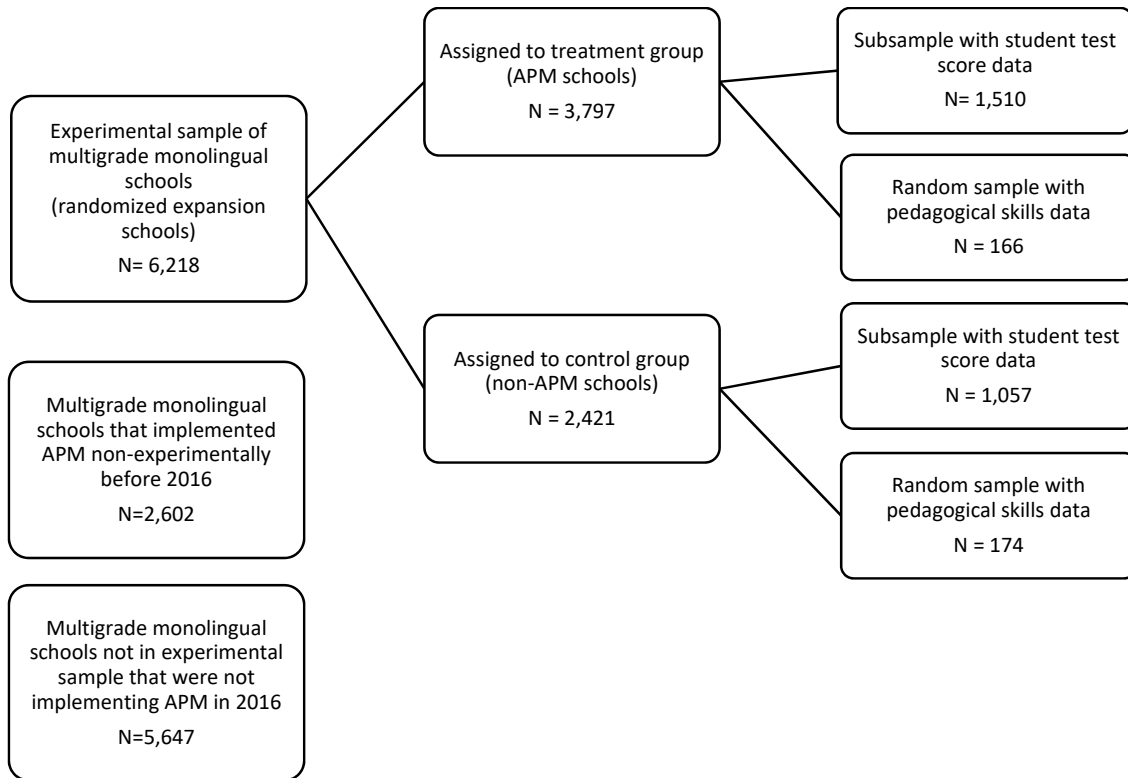
Our sample (Column (3)) is much more similar to the average Peruvian monolingual multigrade school (the target population, Column (2)), with a similar average size of about 29

---

<sup>6</sup> Some regions did not have enough eligible schools to provide equal numbers of APM and non-APM schools, which generated variation in the proportion of APM and non-APM schools across regions. Also, in two regions the number of eligible schools was less than or equal to the quota that they had to fill, which left them without any control schools. Since there was no random assignment in these regions, we exclude them from the analysis.

<sup>7</sup> The initial plan drew a random subsample of 364 schools (182 APM and 182 non-APM), but 24 of these schools could not be reached due to their very remote location.

**Figure 1. Samples from Random Assignment to APM and non-APM Schools**



students and 2 teachers per school, similar access to the internet (8%) and computers, and quality of school infrastructure. We conclude that our sample is very similar to the target population of the APM program, which are monolingual (Spanish-speaking) multigrade schools, but somewhat different from urban primary schools. This should be kept in mind when considering the external validity of our study.

Column (4) of Table A1 shows descriptive statistics for the subset of schools that have test scores. Only schools with five or more students in the tested grade level are eligible to take the national evaluation test, so this subsample includes slightly larger schools than our full randomized expansion sample in Column (3); the average school has 46 students and 2.6 teachers. Teacher-student ratios are slightly lower in the test score subsample, and several indicators of infrastructure quality display some small differences, which suggest that schools in this subsample are slightly better off than, but generally similar to, the full sample. On the other hand, the randomized expansion sample has lower baseline test scores than the average monolingual multigrade school, which reflects that it was targeted towards lower performing schools. Overall, while there are small differences between our test-score subsample and both the full sample and the average monolingual multigrade school, the differences may be small

enough to allow one to extrapolate from our sample to all monolingual multigrade schools. Lastly, Column (5) shows means for the subsample of 340 schools where teachers' pedagogical practices were measured. They closely resemble the overall randomized expansion sample (Column (3)), as expected since they are a random sample drawn from those schools.

*Timeline.* The random assignment was done in late 2015. APM schools began the program in early 2016 and operated it for three consecutive years. The school year begins in March, and the standardized tests are taken in November. We look at effects on students' 2016 and 2018 test scores, one year and three years after the program started. Unfortunately, standardized tests were not administered in 2017 due to a national strike. Our measurement of pedagogical skills took place near the end of 2017, two years after the program's implementation.

*Outcomes.* The measure of student learning outcomes is the National Student Evaluation (henceforth, ECE, its Spanish acronym) primary school exam that assesses students' mathematics and reading comprehension skills. It has been implemented annually since 2007 and is comparable across years.<sup>8</sup> All schools with five or more students in the tested grade take the exam; this means that some schools move in and out of the testing sample over time. Initially, the ECE tested students at the end of the second grade of primary school but, starting in 2018, it was shifted to fourth grade. This implies that, for our cohort of students, we have test score data at the student level first in 2016 when they were in second grade, and again in 2018 in fourth grade. Table A2 in the Appendix shows descriptive statistics of the exam. The ECE scores are reported both as levels of subject mastery and as a Rasch score with a nationally standardized mean of 500 and standard deviation of 100. Table A2 shows that average subject mastery is low; a large proportion of students (especially in the experimental sample) are ranked in the lowest category. For example, in 2015 only 22% of all students (and 14% of randomized expansion sample students) met the learning expectations for their grade in math.

Teachers' pedagogical practices were observed in the subsample of 340 schools at the end of the 2017 school year. In addition, many teachers who had left these 340 schools to go to other schools in 2017 were followed and observed in their new schools. The observers assessed eight pedagogical skills of these teachers (see Table 1). These measures of pedagogical skills, and the rubric used, were designed by experts at Peru's Ministry of Education.

---

<sup>8</sup> The standardized exam was continuously implemented from 2007 until 2016; it was discontinued in 2017 for one year due to a Ministerial decision in that year in response to a prolonged nation-wide teachers' strike. Students had missed several weeks of class and allegedly were not up to date with the subjects that the ECE covered, prompting the decision to cancel it. It was reinstated in 2018. The methodology for the ECE did not change after the gap and the 2018 results are comparable to the ECE exams taken before 2017.

**Table 1: Description of the Pedagogical Skills on which Teachers Were Observed**

Pedagogical Skill	Description
Lesson Planning	The session's purpose is stated explicitly, in a way that students can understand. Activities are planned and aligned with the stated purpose. The session is closed referring to its purpose.
Time Management	Almost all time is allocated to pedagogical activities. Routines, transitions, and interruptions are well managed. Students know the routines and require little teacher assistance to do them.
Promotion of Students' Critical Thinking	The activities promote analysis and reasoning. Most of the questions are open ended and students are given time to delve into them.
Promotion of Students' Participation	The teacher succeeds in getting students involved and actively participating, incorporating their opinions, ideas, and interests into the session. Students can influence the class dynamics.
Provision of Oral Feedback	The teacher pays attention to the difficulties, doubts, and errors of the students, encouraging them to develop their own answers (through questions or hints), helping them to improve their understanding of the subject and advancing in their learning process. The teacher gathers evidence of the students' progress.
Provision of Written Feedback	The teacher assesses the students' work, helping them to see how to achieve what is expected of them.
Quality of Relations between Teacher and Students	Relationships in the classroom are respectful. The class sessions possess a warm environment.
Management of Students' Behavior	The teacher employs positive strategies to promote and reinforce good behavior of students, who autoregulate. An environment that promotes learning is facilitated. Bad behavior is very rare.

Several studies have found that these pedagogical skills predict student's academic success.<sup>9</sup> We also construct an overall index by standardizing and then averaging these eight skills.

### 3. Framework and Treatment Effects

This section presents the empirical framework used in this paper. It starts by defining treatment effects for teacher skills and student learning for the context of the APM program. It then explains which treatment effects can be estimated with the available data, and then provides lower bounds for those that cannot be estimated. For details, see Appendix B.

#### 3.1 Four Types of Teachers

The Angrist, Imbens and Rubin (1996) framework divides the population of interest into *always takers*, who can always obtain the treatment, *never takers*, who can always avoid the treatment, and *compliers*, who follow their assigned treatment. Strictly speaking, these classifications are based on behavior, and do not imply any assumptions about preferences.

<sup>9</sup> For example, Akpur (2020) finds a link between student learning and promotion of critical thinking, Stronge, Ward and Grant (2011) find a similar link for effective use of class time (which is related to lesson planning), Gage et. al. (2018) and Wisniewski, Zierer and Hattie (2020) identify provision of feedback as a mediator on student learning, Allen et. al (2013) and Fauth et. al. (2019) provide evidence in favor of promoting a positive emotional climate. Stronge, Ward and Grant (2011) and Fauth et. al. (2019) highlight the benefits of monitoring and managing student behavior. Table A9 shows that our index is positively correlated with student learning.

In the APM context, changes in treatment status occur via turnover (teachers switching schools). Part of this turnover may be driven by the presence of the program, but part may also occur for reasons other than APM. If turnover is in part due to the program, it is reasonable to assume that such teachers have preferences regarding APM. We propose a framework that allows differences in preferences for APM to explain at least some teacher turnover, but we do not want turnover to be explained only by these preferences; teachers may switch schools for reasons completely unrelated to APM.

This requires changing the “traditional” classification of the population. For example, the traditional Angrist, Imbens and Rubin (1996) framework classifies a teacher moving from an APM school to a non-APM school as a *never taker*. If we assume that this is driven by a strong preference against APM, and ascribe that preference to *never takers*, we exclude the possibility that this move would have occurred even in the absence of APM.

To allow for teacher turnover that is unrelated to APM, we divide the population of teachers in the multigrade monolingual schools into four groups. First, we divide teachers into those who are relatively indifferent to APM and those with strong preferences for or against it. We further separate the latter group into *likers* (L) and *dislikers* (D). Likers are those teachers who like the program enough so that they would make a strong enough effort to secure a position in an APM school if they are not already working in one. According to the application process described in subsection 2.1, this can be done, for example, by giving a high ranking to schools that are remote or otherwise unappealing to most teachers, but that have the APM program. Conversely, dislikers are those teachers who dislike the program enough so that they would make a strong enough effort to secure a position in a non-APM school if they are not already working in one. Finally, we divide teachers indifferent to APM into those who have a strong enough preference for moving, but for reasons unrelated to the program, that they exert sufficient effort to move to a new school, whom we call *movers* (M), and those who remain in their schools, whom we call *remainers* (R).<sup>10</sup> We allow the impact of APM to differ by teachers, so we define  $\delta^M$ ,  $\delta^R$ ,  $\delta^L$  and  $\delta^D$  as the average effects on teacher skills of one year of APM on movers, remainers, likers and dislikers, respectively.

Since all 6,218 randomized expansion schools followed their random assignment in 2016, all teachers in those schools had no choice regarding participation in APM in year 1.<sup>11</sup>

---

<sup>10</sup> As almost all other studies do, we assume that there are no “defiers”. Such teachers would move to a non-APM school in year 2 if they were assigned to an APM school in year 1, or move to an APM school in year 2 if assigned to a non-APM school in year 1, because they want to defy their random assignment.

<sup>11</sup> When teachers learned of their random assignment for 2016 it was too late to switch schools in that year.

We assume that teachers' behavior in the following years is characterized as follows: (i) by definition, all likers assigned to non-APM schools in year 1 move to an APM school in year 2, and all dislikers assigned to APM schools in year 1 move to a non-APM school in year 2; (ii) a fixed proportion of likers switch from one APM school to another APM school every year; (iii) a fixed proportion of dislikers switch from a non-APM to another non-APM school every year; (iv) the number of teacher positions in APM and non-APM schools is fixed; and (v) our 6,218 multigrade monolingual schools are a representative sample of a larger system of multigrade monolingual schools within which most of the teachers remain, and teacher transitions in and out this system do not affect the proportions of likers, dislikers, movers and remainers in this system.<sup>12</sup> As explained below, these proportions are a function of the scale of the program since teachers compete for teaching positions, and the competition to move to, or move out of, an APM school depends on the proportion of schools that are APM schools.

Comparing our four groups of teachers with the “traditional” classification above, *likers* and *dislikers* would be classified as *always takers* and *never takers*, respectively, and *remainers* can be classified as *compliers*. The key difference is the addition of *movers*, whose behavior is consistent with that of any of these three traditional groups. If a mover does not change treatment status after changing schools, he or she could be seen as a *complier* according to the traditional classification. Yet if this teacher had moved from an APM school to a non-APM school he or she would be classified as a *never taker*, and if he or she had moved from a non-APM school to an APM school, he or she would be considered an *always taker*. Moreover, since movers do not take APM into account when changing schools, they always have a probability between 0 and 1 (and never equal to 0 or 1) of moving to an APM (or non-APM) school after year 1, which is not the case for any of the three “traditional” groups.

### 3.2 Treatment Effects for Teacher Skills.

The skill of teacher  $j$  at the end of year  $t$ , denoted by  $y_j^t$ , is assumed to be a linear function of his or her skills in the previous year ( $y_j^{t-1}$ ), the skill gained from one more year of experience ( $\lambda_j$ ), and whether he or she is treated (coached) in year  $t$  ( $T_j^t$ ). The average treatment impact can vary by the four teachers types (remainers (R), likers (L), dislikers (D) and movers (M)). General depreciation of teaching skills can be included in  $\lambda_j$ . Equation (1) provides the general expression of  $y_j^t$  for year  $t$ , and equation (2) shows the specific expression for year 1:

---

<sup>12</sup> Administrative data show that, in any given year, about 10% of teachers move out of multigrade monolingual schools into other schools, and another 10% leave the public education system.

$$y_j^t = y_j^{t-1} + \lambda_j + \delta^k T_j^t, \quad \text{for } k = R, L, D, M \quad (1)$$

$$y_j^1 = y_j^0 + \lambda_j + \delta^k T_j^1 = \theta_j^1 + \delta^k T_j^1, \quad \text{for } k = R, L, D, M \quad (2)$$

where  $\theta_j^1$  is convenient notation for  $y_j^0 + \lambda_j$ .<sup>13</sup>

In year 2, there may be interactions (denoted by  $\gamma_{1,2}^k$ ) of the coaching in years 1 and 2:

$$y_j^2 = y_j^1 + \lambda_j + \delta^k T_j^2 + \gamma_{1,2}^k T_j^1 T_j^2, \quad \text{for } k = R, L, D, M \quad (3)$$

$$= \theta_j^2 + \delta^k (T_j^1 + T_j^2) + \gamma_{1,2}^k T_j^1 T_j^2$$

The second line substitutes out  $y_j^1$  using (2), and  $\theta_j^2$  denotes  $\theta_j^1 + \lambda_j = y_j^0 + 2\lambda_j$ . If, for example, the second year's impact of coaching is less than that of the first year, then the interaction term  $\gamma_{1,2}^k$  is  $< 0$ . Also,  $\gamma_{1,2}^k$  can include depreciation of teacher skills produced by the program.

For year 3, further interaction effects are needed. The equation for  $y_j^3$  is:

$$y_j^3 = y_j^2 + \lambda_j + \delta^k T_j^3 + \gamma_{1,2}^k (T_j^1 T_j^2 + T_j^1 T_j^3 + T_j^2 T_j^3) + \gamma_{1,2,3}^k T_j^1 T_j^2 T_j^3, \quad \text{for } k = R, L, D, M \quad (4)$$

$$= \theta_j^3 + \delta^k (T_j^1 + T_j^2 + T_j^3) + \gamma_{1,2}^k (T_j^1 T_j^2 + T_j^1 T_j^3 + T_j^2 T_j^3) + \gamma_{1,2,3}^k T_j^1 T_j^2 T_j^3$$

where the second line uses (3) to substitute out  $y_j^2$ , and  $\theta_j^3 = \theta_j^2 + \lambda_j = y_j^0 + 3\lambda_j$ . Note that the interaction effect for any combination of two years of training is assumed to be the same, regardless of which two years they are; allowing for different interaction effects for each possible pair of years would do little beyond complicating the notation. The triple interaction  $\gamma_{1,2,3}^k$  can include depreciation of the skills of teachers who are coached for all three years.

For the APM program, three standard treatment effects can be defined for teacher skills. The first is the average treatment effect (ATE), APM's impact on the average teacher (when all teachers are treated, i.e. receive coaching). The counterfactual is that no teachers are treated, or equivalently that the program does not exist. ATE for year  $t$  is defined as:

$$\text{ATE}_{\text{chr}}(t) \equiv E[y_1^t - y_0^t] = E[y_1^t] - E[y^t | \text{No program exists}] \quad (5)$$

---

<sup>13</sup> An implicit assumption in equation (1), and thus of equations (2) – (4), is that there are no peer effects: the skills of teacher  $j$  are not affected by whether fellow teachers have been coached. It is not possible to check this assumption with our data, yet there are three reasons why it is unlikely that a coached teacher will have sizeable impacts on the skills of other teachers in the same school. First, about 20% of the schools in the test score data, and 49% in the teacher skill data, have only one teacher; for these schools peer effects are not possible. Second, almost all schools that have more than one teacher have only two or three teachers, and they all teach different grades. For example, one teacher teaches grade 1-3 and another teaches grades 4-6. Third, coaching is generally teacher-specific, addressing the pedagogical weaknesses of a specific teacher and the needs of that teacher's students; other teachers are likely to have different pedagogical weaknesses and students with different needs; this further reduces opportunities for peer effects. If peer effects do occur, such a SUTVA violation would lead to underestimation of ITT effects, so our ITT estimates would be lower bounds for the true ITT parameters.



where the “tchr” subscript indicates that the treatment effect refers to teachers’ skills. For  $y$ , the superscript is still years since the program started, but subscripts indicate potential outcomes (1 = treated, 0 = not treated). Implicit in this definition is that the two potential outcomes in year  $t$  ( $y_1^t$  and  $y_0^t$ ) maintain the same potential outcome status (treated or not treated) since year 1, so a teacher who is treated in year 1 is treated for all years between 1 and  $t$ , and a teacher who is not treated in year 1 is not treated for all years between 1 and  $t$ . The population of teachers for which this treatment effect is defined is all teachers who were teaching in multigrade monolingual schools in Peru in year 1.

A more specific example of equation (5) is for year 2 ( $t = 2$ ), which is the only year for which teacher skill data are available. This can be expressed as:

$$ATE_{tchr}(2) \equiv E[y_1^2 - y_0^2] = 2\bar{\delta} + \bar{\gamma}_{1,2} \quad (5')$$

where  $\bar{\delta} = \delta^R p^R + \delta^L p^L + \delta^D p^D + \delta^M p^M$ ,  $\bar{\gamma}_{1,2} = \gamma_{1,2}^R p^R + \gamma_{1,2}^L p^L + \gamma_{1,2}^D p^D + \gamma_{1,2}^M p^M$ , and  $p^k$  is the proportion of type  $k$  teachers. Appendix B gives expressions for  $ATE_{tchr}(1)$  and  $ATE_{tchr}(3)$ .

Next consider the intention to treat (ITT) effect. This is the program’s impact on skills in year  $t$  of teachers randomly assigned to APM schools in year 1, regardless of the school they were in (APM or non-APM) in later years. The counterfactual is random assignment to a non-APM school in year 1, regardless where they taught in later years. It is defined as:

$$ITT_{tchr}(t) \equiv E[y^t | R_{tchr, year 1} = 1] - E[y^t | R_{tchr, year 1} = 0] \quad (6)$$

$R_{tchr, year 1}$  refers to the teacher’s school in year 1, which can differ from his or her school in year  $t$ . An example of equation (6) is for year 2, the year with teacher skill data:<sup>14</sup>

$$\begin{aligned} ITT_{tchr}(2) &\equiv E[y^2 | R_{tchr, year 1} = 1] - E[y^2 | R_{tchr, year 1} = 0] \quad (6') \\ &= \bar{\delta} + p^R(\delta^R + \gamma_{1,2}^R) + p^L \gamma_{1,2}^L + p^M \tau \gamma_{1,2}^M \end{aligned}$$

where  $\tau$  is the proportion of teacher positions in the population of all monolingual multigrade schools. The intuition is that  $\bar{\delta}$  is the effect of the first year, when all teachers follow their random assignment, and the other terms are the effects on the teachers treated in the second year (remainders, likers, and the movers who randomly end up in APM schools in year 2). The counterfactual for remainders is being in a non-APM schools for both years, while the

---

<sup>14</sup> Note a slight abuse of notation: “R” is used in two different ways. If it is “normal” size (not a superscript) it indicates a school’s *random* assignment, but if it is a superscript it denotes *remainder* teachers.

counterfactual for likers and movers who randomly (with probability  $\tau$ ) end up in an APM school in year 2, is being in a non-APM school in year 1 and an APM school in year 2.

A final important point is that, unlike  $ATE_{\text{tchr}}(2)$ ,  $ITT_{\text{tchr}}(2)$  depends on  $\tau$ . In a “small-scale” RCT,  $\tau$  would be almost zero and so could be ignored, but in an “at-scale” RCT  $\tau$  will be larger and will affect  $ITT_{\text{tchr}}(2)$ . The intuition is that a proportion  $\tau$  of movers in APM schools in year 1 will also be in APM schools in year 2, which “turns on” the interaction effect from two years of coaching; if the proportion of APM schools had been very small, very few movers who moved into APM schools in year 2 would have been treated in year 1.

In addition, there is a more subtle impact of  $\tau$  on  $ITT_{\text{tchr}}(2)$ : it determines the level of competition among “potential likers” to move into APM schools, and similarly the extent of competition among “potential dislikers” to get into non-APM schools. This will ultimately determine the proportions of teachers who are actual likers and dislikers, and thus the proportions of teachers who are remainers and movers. However, if there are no likers or dislikers, then the value of  $\tau$  would not affect the proportions of remainers and movers.

Another treatment effect that is often estimated for randomized control trials is a local average treatment effect (LATE).<sup>15</sup> It is defined only for a binary treatment variable, but the APM treatment variable can have more than two values since teachers can switch schools: the treatment can be 0, 1, 2 or 3 years. Angrist and Imbens (1995) extended LATE to non-binary treatments, which they call an average causal response (ACR). The general definition is:

$$ACR_{\text{tchr}}(t) \equiv \sum_{s=1}^t E[y_s^t - y_{s-1}^t | T_1^t \geq s > T_0^t] \frac{\text{Prob}[T_1^t \geq s > T_0^t]}{\sum_{r=1}^t \text{Prob}[T_1^t \geq r > T_0^t]} \quad (7)$$

where  $T_0^t$  is the (potential) number of years of training in year  $t$  for a teacher assigned to a non-APM school in year 1, and  $T_1^t$  is the (potential) number of years of training in year  $t$  for a teacher assigned to an APM school in year 1.<sup>16</sup> The subscripts on  $y$  indicate the value of  $y$  given a (potential) number of *years* of treatment (which varies from 0 to 3), not the value of  $y$  given a binary “treatment or no treatment” variable, in contrast to the definition of  $ATE_{\text{tchr}}(t)$ .

Consider equation (7) for year 2, the only year with teacher skill data:

$$ACR_{\text{tchr}}(2) \equiv E[y_1^2 - y_0^2 | T_1^2 \geq 1 > T_0^2] \frac{\text{Prob}[T_1^2 \geq 1 > T_0^2]}{\text{Prob}[T_1^2 \geq 1 > T_0^2] + \text{Prob}[T_1^2 = 2 > T_0^2]} \quad (7')$$

<sup>15</sup> For the APM context, there is no ATT (average treatment effect on the treated) because ATT requires that some teachers assigned to the treatment ( $R_{\text{tchr, year 1}} = 1$ ) are never treated. Such teachers do not exist in the APM context because all the teachers who were randomly assigned to the APM schools were treated in year 1.

<sup>16</sup> For the general case, possible values for both  $T_0^t$  and  $T_1^t$  are integers from 0 to  $t$ . Yet, for the APM program, all teachers followed their random assignment in year 1, so possible values for  $T_0^t$  are 0 to  $t-1$ , and for  $T_1^t$  are 1 to  $t$ .

$$\begin{aligned}
& + E[y_2^2 - y_1^2 | T_1^2 = 2 > T_0^2] \frac{\text{Prob}[T_1^2=2>T_0^2]}{\text{Prob}[T_1^2 \geq 1 > T_0^2] + \text{Prob}[T_1^2=2>T_0^2]} \\
& = [\bar{\delta} + p^R(\delta^R + \gamma_{1,2}^R) + p^L\gamma_{1,2}^L + p^M\tau\gamma_{1,2}^M]/[1 + p^R] = \text{ITT}_{\text{chr}}(2)/[1 + p^R]
\end{aligned}$$

The intuition behind equation (7.2) is as follows. The term  $E[y_1^2 - y_0^2 | T_1^2 \geq 1 > T_0^2]$  is the impact on teacher skills of receiving one year of treatment, relative to having zero years of treatment, as indicated by the subscripts on the  $y$  terms, for teachers who would have had at least one year of treatment in year 2 if assigned to an APM school in year 1 ( $T_1^2 \geq 1$ ), but would not have been treated in year 2 if assigned to a non-APM school in year 1 ( $T_0^2 < 1$ ). Of the four teacher types, this includes all remainers and dislikers, and movers who randomly switched to a non-APM school in year 2 (for whom  $T_0^2 = 0$  and  $T_1^2 = 1$ ). The term  $E[y_2^2 - y_1^2 | T_1^2 = 2 > T_0^2]$  is the impact on teacher skills of receiving a *second* year of the treatment, *relative to having one year of treatment*, as indicated by the subscripts on the  $y$  terms, for teachers who would have had two years of treatment in year 2 if assigned to an APM school in year 1 but only zero or one year of treatment in year 2 if assigned to a non-APM school in year 1. This includes all remainers, all likers, and movers who randomly switched to APM schools in year 2 (for whom  $T_0^2 = 1$  and  $T_1^2 = 2$ ). Turning to the sum of the probabilities in the denominator,  $\text{Prob}[T_1^2 \geq 1 > T_0^2]$  is the probability that a teacher is a remainder, a disliker, or a mover who randomly switches to a non-APM school in year 2, and  $\text{Prob}[T_1^2 = 2 > T_0^2]$  is the probability that a teacher is a remainder, a liker, or a mover who randomly switches to an APM school in year 2. Their sum is greater than 1; remainers are “counted twice” since they are included in both probabilities. Likers, dislikers, and movers are “counted” only once.

In effect,  $\text{ACR}_{\text{chr}}(2)$  is an average of: a) the (average) impact on teacher skills of going from no treatment to one year of treatment for remainers, dislikers, and those movers who randomly move to a non-APM school in year 2; and b) the (average) impact on those skills of going from one to two years of treatment for remainers, likers, and the movers who randomly move to APM schools in year 2. Thus,  $\text{ACR}_{\text{chr}}(2)$  is the average of the impact on teacher skills for each additional year of treatment due to random assignment in year 1 to an APM school, with remainers getting “double weight” since that assignment raises their years of treatment by two years, but for all others random assignment increases years of treatment by only one year. Importantly, note that, for any  $t$ ,  $\text{ACR}_{\text{chr}}(t)$  is a *per year* (not a cumulative) impact that averages only over years of treatment induced by random assignment to an APM school in year 1. To obtain a cumulative impact over  $t$  years, multiply  $\text{ACR}_{\text{chr}}(t)$  by  $t$ .

A final notable characteristics of  $ACR_{\text{tchr}}(2)$  is that, like  $ITT_{\text{tchr}}(2)$ , it is also a function of  $\tau$ , which is to be expected since the numerator of  $ACR_{\text{tchr}}(2)$  equals  $ITT_{\text{tchr}}(2)$ .

The three treatment effects discussed so far focus on particular teachers, and so they follow teachers who move to other schools. But many teacher training or coaching programs focus on particular schools, so it is useful to define treatment effects for the teachers currently in the schools that implemented APM.

There are two possibilities for treatment effects that focus on schools.<sup>17</sup> The first is an average treatment effect (ATE) on teacher skills for those schools, where the counterfactual is no program at all, which we denote as  $ATE_{\text{sch}}$ . This is defined as follows for year  $t$ :

$$ATE_{\text{sch}}(t) \equiv E[y^i | R = 1] - E[y^i | \text{Program does not exist}] \quad (8)$$

As above, consider again the specific case of year 2, the only year with teacher skill data:

$$ATE_{\text{sch}}(2) = (2\delta^R + \gamma_{1,2}^R)p^R + (\delta^L(1+\tau) + \gamma_{1,2}^L\tau)(p^L/\tau) + (\delta^M(1+\tau) + \gamma_{1,2}^M\tau)p^M(\mu/\tau) \quad (8')$$

$$+ \bar{\theta}^{2,L}p^L((1-\tau)/\tau) - \bar{\theta}^{2,D}p^D + \bar{\theta}^{2,M}p^M((\mu/\tau) - 1)$$

where  $\mu$  is the proportion of all movers who move to an APM school in years 2 or 3, and the  $\bar{\theta}^{2,k}$  terms are averages of  $\theta_j^2$  for year 2 for type  $k$  teachers.<sup>18</sup> The first line of (8.2) is the “direct” treatment effect and the second is a “composition” effect, which accounts for differences in average  $\theta$  between likers, who move into APM schools in year 2, and dislikers, who move out of APM schools in year 2 (and also accounts for changes in the distribution of movers across the two types of schools, who compete with likers to get into APM schools and with dislikers to get into non-APM schools). Note that  $ATE_{\text{sch}}(2)$ , and more generally  $ATE_{\text{sch}}(t)$  with  $t \geq 2$ , also depends on  $\tau$ . Intuitively,  $\tau$  determines the proportions of likers and movers in APM schools (and of dislikers and movers in non-APM schools), yet this is no longer the case if there are no likers or dislikers, as explained below.

The second treatment effect for teacher skills that focuses on schools is  $ITT_{\text{sch}}$ ; it is similar to  $ATE_{\text{sch}}$  except that the counterfactual is the skills of teachers in non-APM schools:

$$ITT_{\text{sch}}(t) \equiv E[y^i | R = 1] - E[y^i | R = 0] \quad (9)$$

<sup>17</sup>  $ACR_{\text{sch}}(t)$  is not well defined since teachers who move into the 6,218 schools have no instrumental variable.

<sup>18</sup> To see where the  $\mu/\tau$  term comes from, recall that the number of teaching positions in a school rarely changes. If the number of those positions is fixed in all schools, this definition of  $\mu$  (where  $\mu$  is determined by the application process that also determines the proportions of teachers who are likers, dislikers, movers and remainers; see subsection 2.1), implies that, among all teachers in APM and non-APM schools, the proportion who are movers in APM schools in year 2 or 3 is  $\mu p^M$ . Focusing on APM schools, this proportion must be divided by  $\tau$ , yielding  $(\mu/\tau)p^M$ . Similar derivations show the proportion of movers in non-APM schools in year 2 or 3 is  $[(1-\mu)/(1-\tau)]p^M$ .

For year 2, this is:

$$\begin{aligned} \text{ITT}_{\text{sch}}(2) = & (2\delta^R + \gamma_{1,2}^R)p^R + (\delta^L(1+\tau) + \gamma_{1,2}^L\tau)(p^L/\tau) + (\delta^M(1+\tau) + \gamma_{1,2}^M\tau)p^M(\mu/\tau) \quad (9') \\ & - [\delta^D p^D(\tau/(1-\tau)) + \delta^M \tau p^M((1-\mu)/(1-\tau))] \\ & + \bar{\theta}^{2,L}(p^L/\tau) + \bar{\theta}^{2,M}p^M(\mu/\tau) - [\bar{\theta}^{2,D}(p^D/(1-\tau)) + \bar{\theta}^{2,M}p^M((1-\mu)/(1-\tau))] \end{aligned}$$

The first two lines are the (net) treatment effect; the last is the composition effect. As with  $\text{ATE}_{\text{sch}}(t)$ ,  $\text{ITT}_{\text{sch}}(t)$  depends on the proportion of schools that are treated ( $\tau$ ) when  $t \geq 2$ .

### 3.3 Treatment Effects for Student Learning

Next, consider treatment effects on student skills. Assume that the skill (measured by a test score) of student  $i$  at the end of year  $t$ , denoted by  $s_i^t$ , is determined by his or her skill at the end of the previous year ( $s_i^{t-1}$ ) and the skills of his or her teacher in year  $t$  ( $y_j^t$ ), where  $j$  is the teacher that student  $i$  had in year  $t$ , and  $\pi$  is the impact of teacher skill on student skills:

$$s_i^t = \sigma s_i^{t-1} + \pi y_j^t$$

Each school is randomly assigned to be either an APM ( $R = 1$ ) or non-APM ( $R = 0$ ) school, an assignment that is fixed over time. Analysis of student skills is simplified by the fact that few students change schools (see subsection 4.2), and each school follows its random assignment.

We define three treatment effects for student skills. The first two,  $\text{ATE}_{\text{stud}}$  and  $\text{ITT}_{\text{stud}}$ , are analogous to the two treatment effects defined for their schools ( $\text{ATE}_{\text{sch}}$  and  $\text{ITT}_{\text{sch}}$ ). All three treatment effects for years 2 and 3 are complex due to several possible “histories” for students’ teachers in those years. For example, in year 2 a student’s teacher in an APM school could be a liker who was in an APM school in years 1 and 2, or a liker who was in a non-APM school in year 1 but in an APM school in year 2. Another example is a student in an APM school in year 3; if he or she was taught by a treated teacher in year 1 (this is certain as the student was in an APM school in year 1), and by a teacher in year 2 who had APM in year 2 but not year 1, and by a teacher in year 3 who had APM in years 2 and 3 but not year 1, he or she was exposed to four years of teacher coaching, and the cumulative learning gain from this exposure is averaged over the four years. The general definition of  $\text{ATE}_{\text{stud}}$  for year  $t$  is:

$$\text{ATE}_{\text{stud}}(t) \equiv E[s^t | R = 1] - E[s^t | \text{Program does not exist}] \quad (10)$$

Applying this definition to year 1 yields  $\text{ATE}_{\text{stud}}(1) = \pi \bar{\delta}$ . Applying it to year 3 (recall that test score data exist only for 2016 and 2018) yields (see Appendix B for the derivations):

$$\begin{aligned}
ATE_{stud}(3) &= \sigma ATE_{stud}(2) + \pi ATE_{sch}(3) = \sigma(\sigma ATE_{stud}(1) + \pi ATE_{sch}(2)) + \pi ATE_{sch}(3) \quad (10') \\
&= \sigma^2 \pi \bar{\delta} + \sigma \pi [(2\delta^R + \gamma_{1,2}^R) p^R + (\delta^L(1+\tau) + \gamma_{1,2}^L \tau) (p^L/\tau) + (\delta^M(1+\tau) + \gamma_{1,2}^M \tau) p^M(\mu/\tau)] \\
&+ \pi [(3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) p^R + (\delta^L(2+\tau) + \gamma_{1,2}^L(2\tau+1) + \gamma_{1,2,3}^L) (p^L/\tau) + (\delta^M(1+2\tau) + \gamma_{1,2}^M \tau(2+\tau) + \tau^2 \gamma_{1,2,3}^M) p^M(\mu/\tau)] \\
&+ \sigma \pi [\bar{\theta}^{2,L} p^L((1-\tau)/\tau) - \bar{\theta}^{2,D} p^D + \bar{\theta}^{2,M} p^M((\mu/\tau) - 1)] + \pi [\bar{\theta}^{3,L} p^L((1-\tau)/\tau) - \bar{\theta}^{3,D} p^D + \bar{\theta}^{3,M} p^M((\mu/\tau) - 1)]
\end{aligned}$$

For  $ATE_{stud}(3)$ , the first two lines are the treatment effect, and the last line is the composition effect. Again, for  $t = 2$  or  $3$ ,  $ATE_{stud}(t)$  depends on  $\tau$ .

Turn next to ITT. The general definition for year  $t$  is:

$$ITT_{stud}(t) \equiv E[s^t | R = 1] - E[s^t | R = 0] \quad (11)$$

For year 1,  $ITT_{stud}(1) = ATE_{stud}(1) = \pi \bar{\delta}$ , as all teachers follow their schools' random assignment in year 1. For year 3, applying the general definition yields (Appendix B gives details):

$$\begin{aligned}
ITT_{stud}(3) &= \sigma ITT_{stud}(2) + \pi ITT_{sch}(3) = \sigma(\sigma ITT_{stud}(1) + \pi ITT_{sch}(2)) + \pi ITT_{sch}(3) \quad (11') \\
&= \sigma^2 \pi \bar{\delta} + \sigma \pi [(2\delta^R + \gamma_{1,2}^R) p^R + (\delta^L(1+\tau) + \gamma_{1,2}^L \tau) (p^L/\tau) + (\delta^M(1+\tau) + \gamma_{1,2}^M \tau) p^M(\mu/\tau) - [\delta^D p^D(\tau/1-\tau) + \delta^M \tau p^M((1-\mu)/(1-\tau))]] \\
&+ \pi [(3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) p^R + (\delta^L(2+\tau) + \gamma_{1,2}^L(2\tau+1) + \gamma_{1,2,3}^L) (p^L/\tau) + (\delta^M(1+2\tau) + \gamma_{1,2}^M \tau(2+\tau) + \tau^2 \gamma_{1,2,3}^M) p^M(\mu/\tau) \\
&\quad - \pi [\delta^D p^D(\tau/1-\tau) + (\delta^M 2\tau + \tau^2 \gamma_{1,2}^M) p^M((1-\mu)/(1-\tau))] \\
&\quad + \sigma \pi [\bar{\theta}^{2,L} (p^L/\tau) + \bar{\theta}^{2,M} p^M(\mu/\tau) - [\bar{\theta}^{2,D} (p^D/(1-\tau)) + \bar{\theta}^{2,M} p^M((1-\mu)/(1-\tau))]] \\
&\quad + \pi [\bar{\theta}^{3,L} (p^L/\tau) + \bar{\theta}^{3,M} p^M(\mu/\tau) - [\bar{\theta}^{3,D} (p^D/(1-\tau)) + \bar{\theta}^{3,M} p^M((1-\mu)/(1-\tau))]]
\end{aligned}$$

The first three lines are the (net) treatment effect, and the last two are the composition effect. Note again that, for  $t = 2$  or  $3$ , that  $ITT_{stud}(t)$  depends on  $\tau$ .

The third treatment effect for students is the (average) impact of an additional year of teacher training on student learning, averaged over all additional years of that training that a student experiences. In effect, this is a transfer of the  $ACR_{tchr}$  treatment effects on teacher skill onto student learning, which is complicated by the many different “histories” a student can have of treated teachers in years 2 and 3. We call these treatment effects  $ACR_{stud}$ , though they differ from  $ACR_{tchr}$  (and so differ from the Angrist and Imbens ACR effects) since students are not *directly* treated but instead are *indirectly* treated by exposure to treated teachers.

The general definition of  $ACR_{students}$  in year  $t$  (1, 2 or 3) is:

$$ACR_{stud}(t) \equiv \frac{E[s^t | R=1] - E[s^t | R=0]}{E[h_{tchr}(t) | R=1] - E[h_{tchr}(t) | R=0]} \quad (12)$$

where  $h_{\text{chr}}(t)$  is the cumulative “history” from year 1 to year  $t$  of a student’s exposure to teachers with APM coaching. For example, a student in a treated school in year 2 had a coached teacher in year 1, but in year 2 the teacher could have one or two years of coaching (e.g. one for a teacher in a non-APM school in year 1), so the student’s  $h_{\text{chr}}(2)$  could be 2 or 3. The expected value of  $h_{\text{chr}}(t)$  averages over the types of teachers in the school from year 1 to year  $t$ .

For year 1,  $\text{ACR}_{\text{stud}}(1) = \text{ATT}_{\text{stud}}(1) = \text{ITT}_{\text{stud}}(1)$  since all teachers follow their random assignment in year 1, so  $\text{ACR}_{\text{stud}}(1) = \pi\bar{\delta}$ . For year 3, applying the definition in (12) yields:

$$\begin{aligned} \text{ACR}_{\text{stud}}(3) &= \frac{E[s^3|R=1] - E[s^3|R=0]}{E[h_{\text{chr}}(3)|R=1] - E[h_{\text{chr}}(3)|R=0]} \quad (12') \\ &= \pi \frac{\sigma^2\bar{\delta} + ((3+2\sigma)\delta^R + (3+\sigma)\gamma_{1,2}^R + \gamma_{1,2,3}^R)p^R + (\delta^L(\sigma(1+\tau)+2+\tau) + \gamma_{1,2}^L(\sigma\tau+2\tau+1) + \tau\gamma_{1,2,3}^L)(p^L/\tau) + (\delta^M(\sigma(1+\tau)+2\tau+1) + \gamma_{1,2}^M\tau(\sigma+2+\tau) + \gamma_{1,2,3}^M\tau^2)p^M(\mu/\tau)}{[1+5p^R + (3+2\tau)p^L/\tau + (2+3\tau)p^M(\mu/\tau)] - [2\tau p^D/(1-\tau) + 3\tau p^M((1-\mu)/(1-\tau))]} \\ &\quad - \pi \frac{\delta^D p^D(\tau/(1-\tau))(\sigma+1) + (\tau(\sigma+2)\delta^M + \tau^2\gamma_{1,2}^M)p^M((1-\mu)/(1-\tau))}{[1+5p^R + (3+2\tau)p^L/\tau + (2+3\tau)p^M(\mu/\tau)] - [2\tau p^D/(1-\tau) + 3\tau p^M((1-\mu)/(1-\tau))]} \\ &\quad + \pi \frac{[(\sigma\bar{\theta}^{2,L} + \bar{\theta}^{3,L})(p^L/\tau) + (\sigma\bar{\theta}^{2,M} + \bar{\theta}^{3,M})p^M(\mu/\tau)] - [(\sigma\bar{\theta}^{2,D} + \bar{\theta}^{3,D})p^D/(1-\tau) + (\sigma\bar{\theta}^{2,M} + \bar{\theta}^{3,M})p^M((1-\mu)/(1-\tau))]}{[1+5p^R + (3+2\tau)p^L/\tau + (2+3\tau)p^M(\mu/\tau)] - [2\tau p^D/(1-\tau) + 3\tau p^M((1-\mu)/(1-\tau))]} \\ &= \frac{\text{ITT}_{\text{stud}}(3)}{[1+5p^R + (3+2\tau)p^L/\tau + (2+3\tau)p^M(\mu/\tau)] - [2\tau p^D/(1-\tau) + 3\tau p^M((1-\mu)/(1-\tau))]} \end{aligned}$$

To understand this derivation, note that the numerator is  $\text{ITT}_{\text{stud}}(3)$ . The first expression in brackets in the denominator,  $1 + (5p^R + (3+2\tau)p^L/\tau + (2+3\tau)p^M(\mu/\tau))$ , is  $E[h_{\text{chr}}(3) | R = 1]$ , is the average cumulative exposure to years of teacher coaching of a student in an APM school in year 3. The “1 +” term is exposure to a coached teacher in year 1. In years 2 and 3, the probability of getting a remainder teacher is  $p^R$ , and the probabilities of getting a liker or mover teacher are  $p^L/\tau$  and  $p^M(\mu/\tau)$ , respectively. If a student gets a remainder teacher in year 2, he or she is exposed to two more years of accumulated coaching since that teacher has had two years of coaching by year 2, and if the student gets a remainder teacher in year three he or she will get three more years of accumulated coaching, for a total of five additional years (beyond year 1). If the student gets a liker teacher in year 2, the average liker teacher will have had  $(1+\tau)$  years of coaching (one in year 2 and one more for a proportion  $\tau$  of those teachers in year 1), and if the student gets a liker teacher in year 3, that will expose him or her to an additional  $2+\tau$  years of accumulated coaching, so overall exposure to liker teachers will provide  $3 + 2\tau$  years of accumulated coaching. Finally, exposure to a mover teacher in year 2 leads to  $1+\tau$  additional years of accumulated coaching, and exposure to a mover in year 3 adds another  $1+ 2\tau$  (since movers move randomly every year). Similar calculations for students who randomly end up in non-APM schools in year 1 (and the next two years) lead to accumulated coaching from dislikers and movers of  $2\tau p^D/(1-\tau) + 3\tau p^M((1-\mu)/(1-\tau))$ .

Note that, as with  $ACR_{tchr}(2)$ ,  $ACR_{stud}(3)$  is a *per year* effect; to obtain a cumulative effect for exposure to fully coached teachers in all three years, multiply  $ACR_{stud}(3)$  by 6.

### 3.4 Treatment Effects if No Likers or Dislikers.

The treatment effects for years 2 and 3 in subsections 3.2 and 3.3 are much simpler if there are no likers or dislikers, leaving only remainers and movers. The equations simplify to:<sup>19</sup>

$$ATE_{tchr}(2) = 2\bar{\delta} + \bar{\gamma}_{1,2}, \text{ where } \bar{\delta} = \delta^R p^R + \delta^M p^M \text{ and } \bar{\gamma}_{1,2} = \gamma_{1,2}^R p^R + \gamma_{1,2}^M p^M \quad (5'')$$

$$ITT_{tchr}(2) = \bar{\delta} + p^R(\delta^R + \gamma_{1,2}^R) + p^M \tau \gamma_{1,2}^M \quad (6'')$$

$$ACR_{tchr}(2) = [\bar{\delta} + p^R(\delta^R + \gamma_{1,2}^R) + p^M \tau \gamma_{1,2}^M] / [2p^R + p^M] = ITT_{tchr}(2) / [1 + p^R] \quad (7'')$$

$$ATE_{sch}(2) = (2\delta^R + \gamma_{1,2}^R) p^R + (\delta^M(1+\tau) + \gamma_{1,2}^M \tau) p^M \quad (8'')$$

$$ITT_{sch}(2) = \bar{\delta} + (\delta^R + \gamma_{1,2}^R) p^R + p^M \tau \gamma_{1,2}^M \quad (9'')$$

Note that  $ITT_{sch}(2) = ITT_{tchr}(2)$ , but  $ATE_{sch}(2) \neq ATE_{tchr}(2)$ .

$$ATE_{stud}(3) = \sigma^2 \pi \bar{\delta} + \sigma \pi [(2\delta^R + \gamma_{1,2}^R) p^R + (\delta^M(1+\tau) + \gamma_{1,2}^M \tau) p^M] \quad (10.3')$$

$$+ \pi [(3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) p^R + (\delta^M(1+2\tau) + \gamma_{1,2}^M \tau(2+\tau) + \tau^2 \gamma_{1,2,3}^M) p^M]$$

$$ITT_{stud}(3) = \sigma^2 \pi \bar{\delta} + \sigma \pi [(2\delta^R + \gamma_{1,2}^R) p^R + (\delta^M + \gamma_{1,2}^M \tau) p^M] \quad (11.3')$$

$$+ \pi [(3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) p^R + (\delta^M + \gamma_{1,2}^M \tau + \tau^2 \gamma_{1,2,3}^M) p^M]$$

$$ACR_{stud}(3) = \pi \frac{\sigma^2 \bar{\delta} + ((3+2\sigma)\delta^R + (3+\sigma)\gamma_{1,2}^R + \gamma_{1,2,3}^R) p^R + (\delta^M(\sigma+1) + \gamma_{1,2}^M \tau(\sigma+2) + \gamma_{1,2,3}^M \tau^2) p^M}{1+5p^R+2p^M} \quad (12.3')$$

$$= \frac{ITT_{stud}(3)}{1+5p^R+2p^M}$$

Note that there are no composition effects for  $ATE_{sch}(2)$ ,  $ITT_{sch}(2)$ ,  $ATE_{stud}(3)$ , and  $ITT_{stud}(3)$ . Also, the absence of likers and dislikers ( $p^L = p^D = 0$ ) implies that there are only remainers and movers, and that  $\mu = \tau$  (movers will be equally distributed over APM and non-APM schools since they do not need to compete with likers and dislikers to move into an APM or non-APM school).

### 3.5 What Do OLS and IV Regressions Estimate?

Most, but not all, of these treatment effects can be estimated by OLS or IV regression. We have two samples of teachers, one (imperfectly) follows the teachers who were in APM and non-APM schools in year 1 (Sample 1), and the other focuses on the teachers in the APM and

<sup>19</sup> They follow from the results in subsections 3.2 and 3.3:  $p^L = p^D = 0$  and  $\mu = \tau$  if there are no likers or dislikers.



non-APM schools in any given year (Sample 2). OLS regression of Sample 1 teachers' skills in year  $t$  on a constant term and a dummy variable for assignment to an APM school in year 1 yields an unbiased estimate of the  $ITT_{\text{tchr}}(t)$  treatment effect.<sup>20</sup> For example, consider year 2:

$$\begin{aligned}\hat{\beta}_1^y_{\text{OLS},t=2} &= E[y^2 | R_{\text{tchr}, \text{year } 1} = 1] - E[y^2 | R_{\text{tchr}, \text{year } 1} = 0] \quad (13) \\ &= \bar{\delta} + p^R(\delta^R + \gamma_{1,2}^R) + p^L\gamma_{1,2}^L + \tau p^M\gamma_{1,2}^M = ITT_{\text{tchr}}(2)\end{aligned}$$

The “1” subscript indicates Sample 1 teachers. Appendix B presents this derivation, as well as those for years 1 and 3. It also presents the derivations for the other OLS and IV estimators in this subsection, for all three years, and shows that OLS estimation applied to Sample 2 teachers estimates  $ITT_{\text{tchr}}(t)$  (recall that  $ITT_{\text{tchr}}(t) = ITT_{\text{sch}}(t)$  if there are no likers or dislikers).

Next, consider IV estimation using Sample 1 teachers. Let  $T^{\text{Tot},t}$  denote the number of years that a teacher has participated in the program up through year  $t$ . IV regression uses random assignment as an instrument for  $T^{\text{Tot},t}$  to estimate the (average) impact of a year of exposure to the program on teacher skills. This yields unbiased estimates of  $ACR_{\text{tchr}}(t)$ . For year 2:

$$\begin{aligned}\hat{\beta}_1^y_{\text{IV},t=2} &= \frac{E[y^2 | R_{\text{tchr}, \text{year } 1} = 1] - E[y^2 | R_{\text{tchr}, \text{year } 1} = 0]}{E[T^{\text{Tot},2} | R_{\text{tchr}, \text{year } 1} = 1] - E[T^{\text{Tot},2} | R_{\text{tchr}, \text{year } 1} = 0]} \quad (14) \\ &= (\bar{\delta} + p^R(\delta^R + \gamma_{1,2}^R) + p^L\gamma_{1,2}^L + p^M\tau\gamma_{1,2}^M) / (1 + p^R) = ACR_{\text{tchr}}(2)\end{aligned}$$

One can also apply OLS to Sample 2 teachers, the teachers who, in any given year, teach in the schools that were randomly assigned in year 1 to be APM or non-APM schools. An OLS regression of Sample 2 teachers' skills in year  $t$  on a constant and a dummy for teaching in an APM school in year  $t$  yields an unbiased estimate of  $ITT_{\text{sch}}(t)$ . So, for year 2:

$$\begin{aligned}\hat{\beta}_2^y_{\text{OLS},t=2} &= E[y^2 | R = 1] - E[y^2 | R = 0] \quad (15) \\ &= (2\delta^R + \gamma_{1,2}^R)p^R + (\delta^L(1+\tau) + \gamma_{1,2}^L\tau)(p^L/\tau) - [\delta^D p^D(\tau/(1-\tau)) + p^M[\delta^M(\mu-\tau^2)/(\tau-\tau^2) + \mu\gamma_{1,2}^M] \\ &\quad + \bar{\theta}^{2,L}(p^L/\tau) + \bar{\theta}^{2,M}p^M(\mu/\tau) - [\bar{\theta}^{2,D}(p^D/(1-\tau)) + \bar{\theta}^{2,M}p^M((1-\mu)/(1-\tau))]] = ITT_{\text{sch}}(2)\end{aligned}$$

In general, IV estimation cannot be used for Sample 2 teachers in year 2 since some of those teachers moved into both APM schools and non-APM schools that were not part of the initial random assignment, such as teachers working in monolingual multigrade schools in year 1 that had ECE scores above the threshold that determined eligibility for the randomized

<sup>20</sup> Almost all of the regressions in this paper have other explanatory variables, but since random assignment is by definition uncorrelated with these other variables the first line in (13) still holds by the Frisch-Waugh theorem. Regressions without these explanatory variables (e.g. Table 6) yields very similar results.

expansion (see subsection 2.3). These Sample 2 teachers have no instrument, so IV estimation cannot be done for Sample 2 teachers.

Next, consider OLS regression for student test scores, more specifically regressing those scores on a constant and a dummy indicating being in an APM school. OLS regression of students' test scores in year  $t$  on a constant and a dummy for being enrolled in an APM school in year  $t$  yields an unbiased estimate of  $ITT_{stud}(t)$ . For years 1 and 3 this implies that:

$$\hat{\beta}_{OLS,1}^S = E[s^1 | R = 1] - E[s^1 | R = 0] = \pi\bar{\delta} = ATE_{stud}(1) = ITT_{stud}(1) = ACR_{stud}(1) \quad (16)$$

$$\begin{aligned} \hat{\beta}_{OLS,t=3}^S &= E[s^3 | R = 1] - E[s^3 | R = 0] \quad (17) \\ &= \sigma^2 \pi e \gamma_{1,2}^R p^R + (\delta^L(1+\tau) + \gamma_{1,2}^L \tau)(p^L/\tau) + (\delta^M(1+\tau) + \gamma_{1,2}^M \tau)p^M(\mu/\tau) - [\delta^D p^D(\tau/(1-\tau)) + \delta^M \tau p^M((1-\mu)/(1-\tau))] \\ &+ \pi[(3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R)p^R + (\delta^L(2+\tau) + \gamma_{1,2}^L(2\tau+1) + \tau\gamma_{1,2,3}^L)(p^L/\tau) + (\delta^M(1+2\tau) + \gamma_{1,2}^M \tau(2+\tau) + \tau^2\gamma_{1,2,3}^M)p^M(\mu/\tau)] \\ &\quad - \pi[\delta^D p^D(\tau/(1-\tau)) + (\delta^M 2\tau + \tau^2\gamma_{1,2}^M)p^M((1-\mu)/(1-\tau))] \\ &\quad + \sigma\pi[\bar{\theta}^{2,L}(p^L/\tau) + \bar{\theta}^{2,M}p^M(\mu/\tau) - [\bar{\theta}^{2,D}(p^D/(1-\tau)) + \bar{\theta}^{2,M}p^M((1-\mu)/(1-\tau))]] \\ &\quad + \pi[\bar{\theta}^{3,L}(p^L/\tau) + \bar{\theta}^{3,M}p^M(\mu/\tau) - [\bar{\theta}^{3,D}(p^D/(1-\tau)) + \bar{\theta}^{3,M}p^M((1-\mu)/(1-\tau))]] = ITT_{stud}(3) \end{aligned}$$

Last, consider IV estimation for student test scores. The treatment for year  $t$  is the “history” from years 1 to  $t$  of students' exposure to treated teachers,  $h_{tchr}(t)$  (see subsection 3.3). Thus:

$$\hat{\beta}_{IV,t}^S \equiv \frac{E[s^t | R=1] - E[s^t | R=0]}{E[h_{tchr}(t) | R=1] - E[h_{tchr}(t) | R=0]} \quad (18)$$

This is an unbiased estimate of  $ACR_{stud}(t)$ . Applying this to year 1, it equals OLS since all teachers follow their random assignment in year 1:

$$\hat{\beta}_{IV,1}^S = \pi\bar{\delta} = ATE_{stud}(1) = ITT_{stud}(1) = ACR_{stud}(1) \quad (19)$$

For year 3, the IV estimate is:

$$\begin{aligned} \hat{\beta}_{IV,3}^S &= \frac{E[s^3 | R=1] - E[s^3 | R=0]}{E[h_{tchr}(3) | R=1] - E[h_{tchr}(3) | R=0]} \quad (20) \\ &= \pi \frac{\sigma^2 \bar{\delta} + ((3+2\sigma)\delta^R + (3+\sigma)\gamma_{1,2}^R + \gamma_{1,2,3}^R)p^R + (\delta^L(\sigma(1+\tau)+2\tau) + \gamma_{1,2}^L(\sigma\tau+2\tau+1) + \tau\gamma_{1,2,3}^L)(p^L/\tau) + (\delta^M(\sigma(1+\tau)+2\tau+1) + \gamma_{1,2}^M \tau(\sigma+2\tau) + \gamma_{1,2,3}^M \tau^2)p^M(\mu/\tau)}{[1+5p^R + (3+2\tau)p^L/\tau + (2+3\tau)p^M(\mu/\tau)] - [2\tau p^D/(1-\tau) + 3\tau p^M((1-\mu)/(1-\tau))]} \\ &\quad - \pi \frac{\delta^D p^D(\tau/(1-\tau))(\sigma+1) + (\tau(\sigma+2)\delta^M + \tau^2\gamma_{1,2}^M)p^M((1-\mu)/(1-\tau))}{[1+5p^R + (3+2\tau)p^L/\tau + (2+3\tau)p^M(\mu/\tau)] - [2\tau p^D/(1-\tau) + 3\tau p^M((1-\mu)/(1-\tau))]} \\ &\quad + \pi \frac{[(\sigma\bar{\theta}^{2,L} + \bar{\theta}^{3,L})(p^L/\tau) + (\sigma\bar{\theta}^{2,M} + \bar{\theta}^{3,M})p^M(\mu/\tau)] - [(\sigma\bar{\theta}^{2,D} + \bar{\theta}^{3,D})p^D/(1-\tau) + (\sigma\bar{\theta}^{2,M} + \bar{\theta}^{3,M})p^M((1-\mu)/(1-\tau))]}{[1+5p^R + (3+2\tau)p^L/\tau + (2+3\tau)p^M(\mu/\tau)] - [2\tau p^D/(1-\tau) + 3\tau p^M((1-\mu)/(1-\tau))]} \\ &= ACR_{stud}(3) = \frac{ITT_{stud}(3)}{[1+5p^R + (3+2\tau)p^L/\tau + (2+3\tau)p^M(\mu/\tau)] - [2\tau p^D/(1-\tau) + 3\tau p^M((1-\mu)/(1-\tau))]} \end{aligned}$$

### 3.6 Bounds on Treatment Effects for Years 2 and 3.

As explained above, for all  $t$  for which data are available,  $ITT_{tchr}(t)$  and  $ACR_{tchr}(t)$  can be estimated using Sample 1 teachers, and  $ITT_{stud}(t)$  and  $ACR_{stud}(t)$  can be estimated using student test scores from the APM and non-APM schools. In addition, for year 1  $ATE_{tchr}(1)$ , which equals  $ATE_{sch}(1)$ , and  $ATE_{stud}(1)$  can be estimated since they equal the corresponding ITT estimands. Unfortunately,  $ATE_{tchr}(t)$ ,  $ATE_{sch}(t)$  and  $ATE_{stud}(t)$  cannot be estimated for  $t \geq 2$ . Yet, under plausible assumptions it is possible to show that ITT estimands are lower bounds of these ATE treatment effects. Turn now to these results, focusing on ATEs for which we have data to estimate; derivations, and ATEs for which we have no data, are in Appendix B.

For year 2, consider  $ATE_{tchr}(2)$  and  $ITT_{tchr}(2)$ . Their difference is:

$$ATE_{tchr}(2) - ITT_{tchr}(2) = \delta^L p^L + (\delta^D + \gamma_{1,2}^D) p^D + (\delta^M + \gamma_{1,2}^M (1-\tau)) p^M \quad (21)$$

As long as the first year of the program does not have a negative effect on the skills of likers (i.e.  $\delta^L \geq 0$ ) and a second year does not have a negative effect on the skills of dislikers ( $\delta^D + \gamma_{1,2}^D \geq 0$ ) or movers ( $\delta^M + \gamma_{1,2}^M \geq 0$ ),  $ITT_{tchr}(2)$  will be a lower bound for  $ATE_{tchr}(2)$ .

It is less clear that  $ITT_{sch}(2)$  is a lower bound for  $ATE_{sch}(2)$ , because these treatment effects follow schools over time, as opposed to following teachers over time, and the composition of teachers in APM and non-APM schools can change over time. More specifically:

$$\begin{aligned} & ATE_{sch}(2) - ITT_{sch}(2) \quad (22) \\ &= \delta^D p^D (\tau / (1-\tau)) + \delta^M \tau p^M ((1-\mu) / (1-\tau)) - \bar{\theta}^{2,L} p^L + \bar{\theta}^{2,D} p^D (\tau / (1-\tau)) + \bar{\theta}^{2,M} p^M ((\tau-\mu) / (1-\tau)). \end{aligned}$$

It is reasonable to assume that the two  $\delta$  terms ( $\delta^D$  and  $\delta^M$ ) are  $\geq 0$ , but the sign of the combined effect of the  $\bar{\theta}^2$  terms, which reflect changes in teacher composition, is ambiguous, even though it is reasonable to assume that all the  $\bar{\theta}^2$  terms are  $> 0$ . Perhaps this combined effect is close to zero and, if negative, is smaller in absolute value than the (weighted) sum of the two  $\delta$  terms, so that  $ITT_{sch}(2)$  is a lower bound for  $ATE_{sch}(2)$ , but it could be that the sum of the composition terms is negative and larger in absolute value than the expression with the two  $\delta$  terms. However, if there are no likers or dislikers then there is no composition effect (since  $p^L = p^D = 0$  and  $\mu = \tau$ ) and so  $ITT_{sch}(2)$  is a lower bound for  $ATE_{sch}(2)$ . In particular,  $ATE_{sch}(2) - ITT_{sch}(2) = \delta^M \tau p^M$ , which is  $\geq 0$  as long as  $\delta^M \geq 0$ , which is plausible.

Finally, turn to student skills for year three to compare  $ITT_{stud}(3)$  with  $ATE_{stud}(3)$ :

$$ATE_{stud}(3) - ITT_{stud}(3) \quad (23)$$

$$\begin{aligned}
&= \sigma\pi[\delta^D p^D(\tau/(1-\tau)) + \delta^M \tau p^M((1-\mu)/(1-\tau))] + \pi[\delta^D p^D(\tau/(1-\tau)) + (\delta^M 2\tau + \tau^2 \gamma_{1,2}^M) p^M((1-\mu)/(1-\tau))] \\
&+ \sigma\pi[-\bar{\theta}^{2,L} p^L + \bar{\theta}^{2,D} p^D(\tau/(1-\tau)) + \bar{\theta}^{2,M} p^M((\tau-\mu)/(1-\tau))] + \pi[-\bar{\theta}^{3,L} p^L + \bar{\theta}^{3,D} p^D(\tau/(1-\tau)) + \bar{\theta}^{3,M} p^M((\tau-\mu)/(1-\tau))]
\end{aligned}$$

The first three  $\delta$  terms are  $\geq 0$  (assuming  $\delta^D$  and  $\delta^M$  are  $\geq 0$ ), and the  $\delta^M 2\tau + \tau^2 \gamma_{1,2}^M$  term is also  $\geq 0$  as long as two years of exposure to the program does not reduce the skills of movers (as long as  $2\delta^M + \gamma_{1,2}^M \geq 0$ ). Yet the sign of the combined effect of the  $\theta$  terms, which reflects changes in teacher composition, is ambiguous, even though it is reasonable to assume that all of the  $\bar{\theta}^3$  terms are  $> 0$ . Yet note that if there are no likers or dislikers then there is no composition effect (since  $p^L = p^D = 0$  and  $\mu = \tau$ ) and so  $ITT_{stud}(3)$  is a lower bound for  $ATE_{stud}(3)$ .

## 4. Fieldwork and Data

### 4.1 Baseline Balance

To verify that the randomization yielded balanced treatment and control groups we checked the baseline balance on several school characteristics. Table 2 shows the descriptive statistics and pairwise t-tests on the difference between control and treatment groups for those school characteristics. Our preferred specification has school district fixed effects, so balance regressions include school district fixed effects, but no other controls. (School districts are subdivisions of regions, so region fixed effects are redundant; in any case baseline characteristics are similarly balanced using region fixed effects instead of school district fixed effects.)

Table 2 shows that most covariates are balanced in the test score evaluation sample, the exceptions being the number of students and, consequently, the number of teachers (since teacher assignment depends on the number of students). While this is what one would roughly expect by chance (the joint F-test is insignificant, with a p-value of 0.152), and the teacher-student ratio (which may affect treatment outcomes directly) is balanced, we control for the numbers of students and teachers in our regressions to ensure that we do not wrongly attribute to the program any differences due to this imbalance. We also show that the results are generally robust to excluding these controls (see Table A3).

### 4.2 Attrition

As explained in Section 2, we use information on two sets of outcome variables: teacher-level outcomes (teachers' pedagogical skills measured for our research) and student-level outcomes (students' test scores gathered from the country's nationwide ECE assessments).

For the teacher-level outcomes, the planned pedagogical sample consisted of 364 schools from the 6,218 the randomized expansion schools: 182 were randomly selected from

**Table 2. Balance Table for Experimental Sample with Test Scores and Pedagogical Skills Measures**

Variable	Experimental sample with Test Scores in 2016					Experimental sample with Pedagogical Measures					
	Control		Treated		Pairwise t-test Conditional Diff.	Control		Treated		Pairwise t-test Conditional Diff.	
	N	Mean/(SD)	N	Mean/(SD)		N	Mean/(SD)	N	Mean/(SD)		
(1)	(2)	(3)	(4)	(5)	(6)						
Math Score (in 2015)	797	513.081 (93.721)	1120	513.439 (91.070)	0.566	81	494.397 (94.245)	69	501.125 (90.382)	2.630	
Language Score (in 2015)	797	517.649 (67.236)	1120	520.768 (64.580)	-0.113	81	508.422 (65.470)	69	514.273 (67.811)	1.816	
Number of Students	1,054	48.365 (28.040)	1509	44.283 (22.882)	5.420***	174	30.920 (25.637)	166	27.259 (21.687)	2.334	
Number of Teachers	1,054	2.644 (1.268)	1509	2.514 (1.192)	0.185***	174	1.983 (1.140)	166	1.801 (1.151)	0.156	
Number of Sections	1,054	5.824 (0.878)	1509	5.843 (0.763)	0.006	174	5.494 (0.942)	166	5.175 (1.427)	0.226	
Teacher-Student Ratio	1,036	0.063 (0.030)	1490	0.064 (0.031)	-0.001	174	0.088 (0.053)	161	0.098 (0.083)	-0.005	
Rurality scale	1,052	2.388 (0.774)	1509	2.199 (0.853)	-0.011	173	2.491 (0.687)	166	2.361 (0.756)	0.073	
Speak indigenous language (%)	1,054	4.350 (18.716)	1509	4.221 (18.972)	-0.139	174	2.244 (13.734)	166	2.487 (15.402)	-0.280	
Poverty Rates	1,030	64.653 (19.512)	1497	56.110 (23.456)	0.537	174	57.362 (21.635)	166	58.173 (20.872)	0.653	
Ceiling Material	1,019	5.717 (1.372)	1468	5.811 (1.543)	-0.063	173	5.543 (1.222)	160	5.631 (1.353)	-0.164	
Wall Material	1,019	6.068 (1.311)	1468	6.113 (1.386)	-0.088	173	5.983 (1.383)	160	6.088 (1.334)	-0.008	
Floor Material	1,019	2.856 (0.719)	1468	2.854 (0.705)	0.017	173	2.902 (0.826)	160	2.881 (0.648)	0.116	
% Teachers with degree	1,021	0.962 (0.122)	1480	0.962 (0.124)	0.004	169	0.984 (0.079)	159	0.979 (0.094)	0.005	
Internet Access	923	0.096 (0.295)	1327	0.102 (0.303)	0.016	158	0.070 (0.255)	150	0.093 (0.292)	-0.027	
Receives Textbooks	1,040	0.717 (0.451)	1489	0.758 (0.429)	0.006	174	0.661 (0.475)	161	0.720 (0.450)	-0.064	
Receives Notebooks	1,039	0.677 (0.468)	1489	0.701 (0.458)	0.016	174	0.598 (0.492)	161	0.640 (0.482)	-0.031	
School Day Length	1,039	8.208 (1.043)	1490	8.104 (0.737)	0.008	174	8.115 (0.930)	161	8.217 (0.871)	0.016	
Electricity	967	0.520 (0.500)	1401	0.560 (0.497)	0.016	155	0.503 (0.502)	153	0.549 (0.499)	-0.044	
Water	967	0.543 (0.498)	1401	0.519 (0.500)	0.015	155	0.484 (0.501)	153	0.497 (0.502)	0.019	
Sanitation	967	0.134 (0.341)	1401	0.148 (0.356)	-0.006	155	0.142 (0.350)	153	0.124 (0.331)	0.031	
Computers per Student	971	0.442 (3.135)	1405	0.460 (2.666)	0.071	155	0.194 (0.646)	153	0.307 (1.199)	-0.253**	
F-test of joint significance (F-stat)						1.318					
F-test, p-value						0.152					
F-test, number of observations						1,538					

Notes: This table presents the balance between treatment (APM) and control (non-APM) schools for the full experimental sample of schools with test scores, and for the subsample with measures of pedagogical practices. Columns 3 and 6 show the coefficient of regressing the treatment on each variable, including school district fixed effects. Significance: \*\*\*p<.01, \*\*p<.05, \*p<.1. Rurality is a categorical variable that takes values 0 for Urban schools, and 1, 2 and 3 for increasingly rural schools. Ceiling, wall and floor materials are categorical variables that take values up to 7, with higher values implying better materials.

the 3,797 schools randomly assigned to APM and 182 were randomly selected from the 2,421 schools randomly not assigned to APM. These 364 schools were selected to observe, in the third quarter of 2017, the pedagogical practices of the teachers who: (i) had worked in one of the 364 evaluation sample schools in 2016 (Sample 1); and (ii) had worked in an evaluation sample school in 2017 (Sample 2). The former required visiting schools not in the 364 sub-sample in year 2 because many Sample 1 teachers changed schools between 2016 and 2017.

**Table 3: Attrition of Sample 1 and Sample 2 Teachers and Evaluation Sample Schools in Year 2 (2017)**

	Sample 1 teachers			Sample 2 teachers			Evaluation sample schools		
	APM (1)	Non-APM (2)	Total (3)	APM (4)	Non-APM (5)	Total (6)	APM (7)	Non-APM (8)	Total (9)
<b>Original (2016)</b>	321	341	662	355	384	739	182	182	364
<b>Observed (2017)</b>	219	236	455	299	341	640	166	174	340
<b>Attrition rate (%)</b>	0.318	0.301	0.312	0.158	0.112	0.134	0.088	0.044	0.066
<b>Difference in attrition rates</b>		0.017 (0.036)			0.046* (0.025)			0.044* (0.026)	

APM is the treated group and Non-APM is the control group. \*  $p < 0.1$ ; \*\*  $p < 0.05$ ; \*\*\*  $p < 0.01$

It was not possible to observe the pedagogical skills of all Sample 1 teachers (see columns (1) – (3) in Table 3). In fact, Sample 1 attrition is high. This is mainly due to outdated information on teachers’ locations when the fieldwork was planned (in March of 2017, the start of Peru’s school year). The teacher location information at that time indicated that, to observe all Sample 1 teachers who were still teaching, 406 schools needed to be visited, including 104 that were not in the 364 pedagogical sample schools. During fieldwork, 91.6% (372) of these 406 schools were visited (34 in very remote areas could not be visited), but outdated information often led to situations where the teachers had moved to other schools, and by the time this was discovered it was logistically impossible to go to the schools where those teachers were working. As seen in Table 3, only 68.8% (455 out of 662) of the original Sample 1 teachers were observed in 2017. Of the 207 unobserved Sample 1 teachers, 50 (7.6% of the 662) had stopped teaching in public schools, 28 (4.2%) were in one of the 34 schools that were not visited, and 129 (19.5%) had moved to a public school that was not in the planned sample of 406 schools. Turning to Sample 2 teachers (those in the 364 evaluation sample schools in year 2), 86.6% (640 out of 739) were observed in year 2 (see columns (4) – (6) of Table 2). In this sample, attrition is due mainly to 24 pedagogical sample schools in hard-to-reach areas that could not be visited in year 2 (see columns (7) – (9) of Table 2).

Non-random attrition can lead to biased estimates, especially for Sample 1 teachers, given their high rate of attrition. Yet if the average characteristics of the missing teachers are similar for APM and non-APM teachers, which for Sample 1 would be the case if the data (on where teachers who moved were working) were outdated primarily due to random factors, then this attrition will not yield biased estimates.

To check for possible bias, we do three things. First, we compare the attrition rates of the APM and non-APM groups. Table 3 shows that the differences in attrition rates are 1.7 (Sample 1) and 4.6 (Sample 2) percentage points. Neither difference is statistically significant at the 5% level, although the Sample 2 difference is significant at the 10% level and the 4.6 percentage point difference may be a concern since it is a 41% higher (15.8% vs. 11.2%) rate.

Second, we compare several observable characteristics of (non-attrited) APM and non-APM schools and teachers. Random assignment to the program in 2016 should ensure that, before any attrition occurred, the teacher characteristics were balanced for the teachers working in the APM and non-APM schools in that year. Random assignment should also ensure that the baseline characteristics of the 364 schools in the teacher skills evaluation sample are balanced. If attrition among Sample 1 teachers is random, the characteristics of the 455 teachers in Table 2 who were observed in 2017 should be similar between those who were working in APM schools and those who were working non-APM schools in 2016.

Figures 2 and 3 show that the treatment and control groups are similar in terms of the observed characteristics of: (i) the original 662 Sample 1 teachers in year 1 (2016); (ii) the subsample of 455 Sample 1 teachers who remained in the sample in year 2 (2017); (iii) the original 364 evaluation sample schools in year 1 (2016); and (iv) the subsample of 340 schools visited in year 2 (2017). Importantly, none of the (standardized) differences is very large, and none is statistically significant at the 5% level.<sup>21</sup>

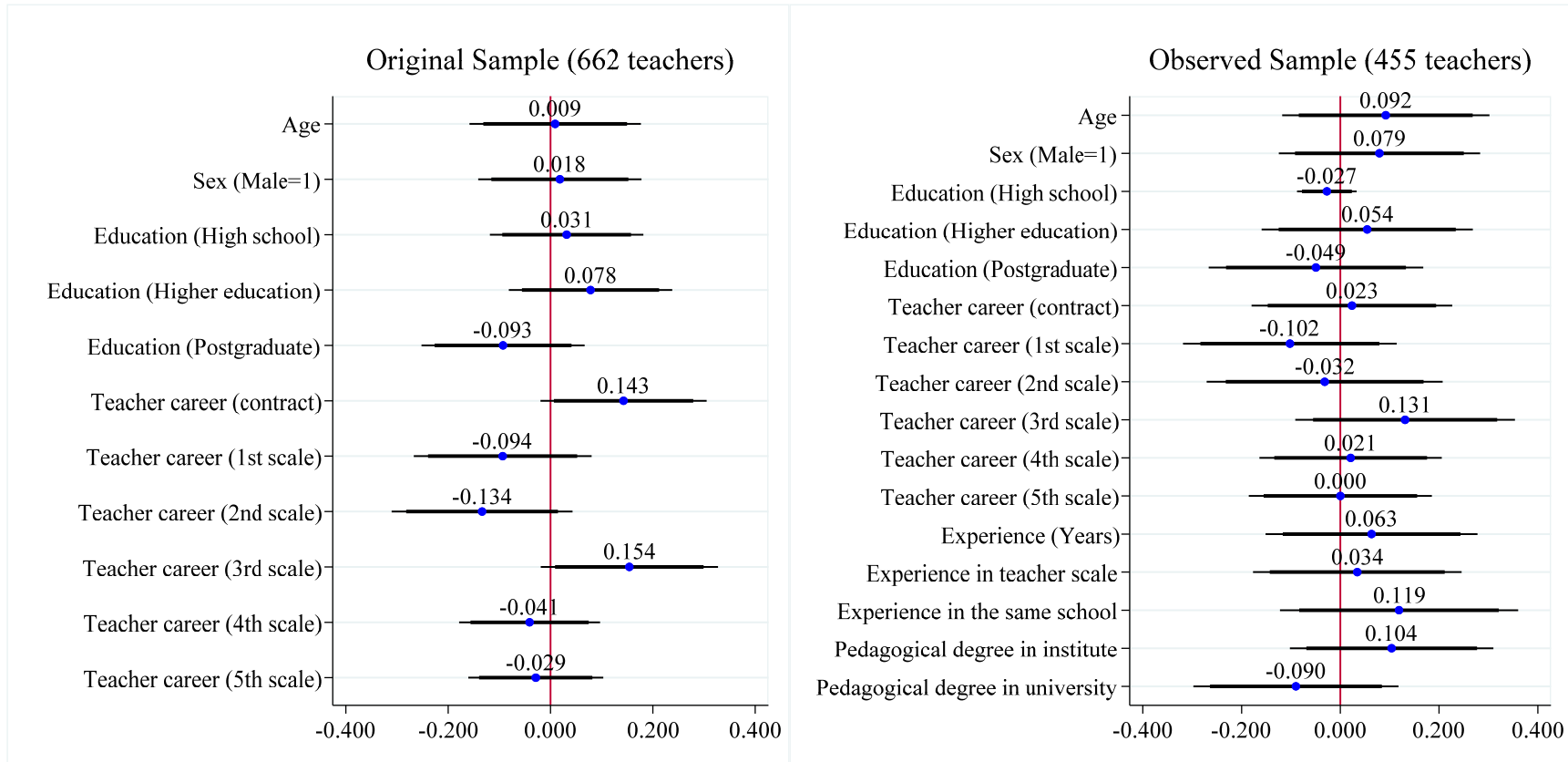
We do not compare Sample 2 teachers' baseline characteristics in 2017 (year 2) between APM and non-APM schools to check for balance at baseline because random assignment of schools in 2016 (year 1) does not ensure such balance across these two groups of schools in year 2. In particular, if certain types of teachers self-select into APM or non-APM schools in year 2, Sample 2 teachers' baseline characteristics may be correlated with the treatment status of the schools where they worked in year 2.

---

<sup>21</sup> Appendix A presents further evidence that attrition is uncorrelated with treatment assignment. Table A4 shows that teachers' pre-treatment characteristics do not predict assignment to an APM school. Table A7 shows that assignment to an APM school does not predict being observed at the end of 2017.

**Figure 2**

**Balance in Teacher Characteristics for the Original and Observed in Year 2 Teachers Who Worked in an Evaluation Sample School in 2016 (Sample 1)**



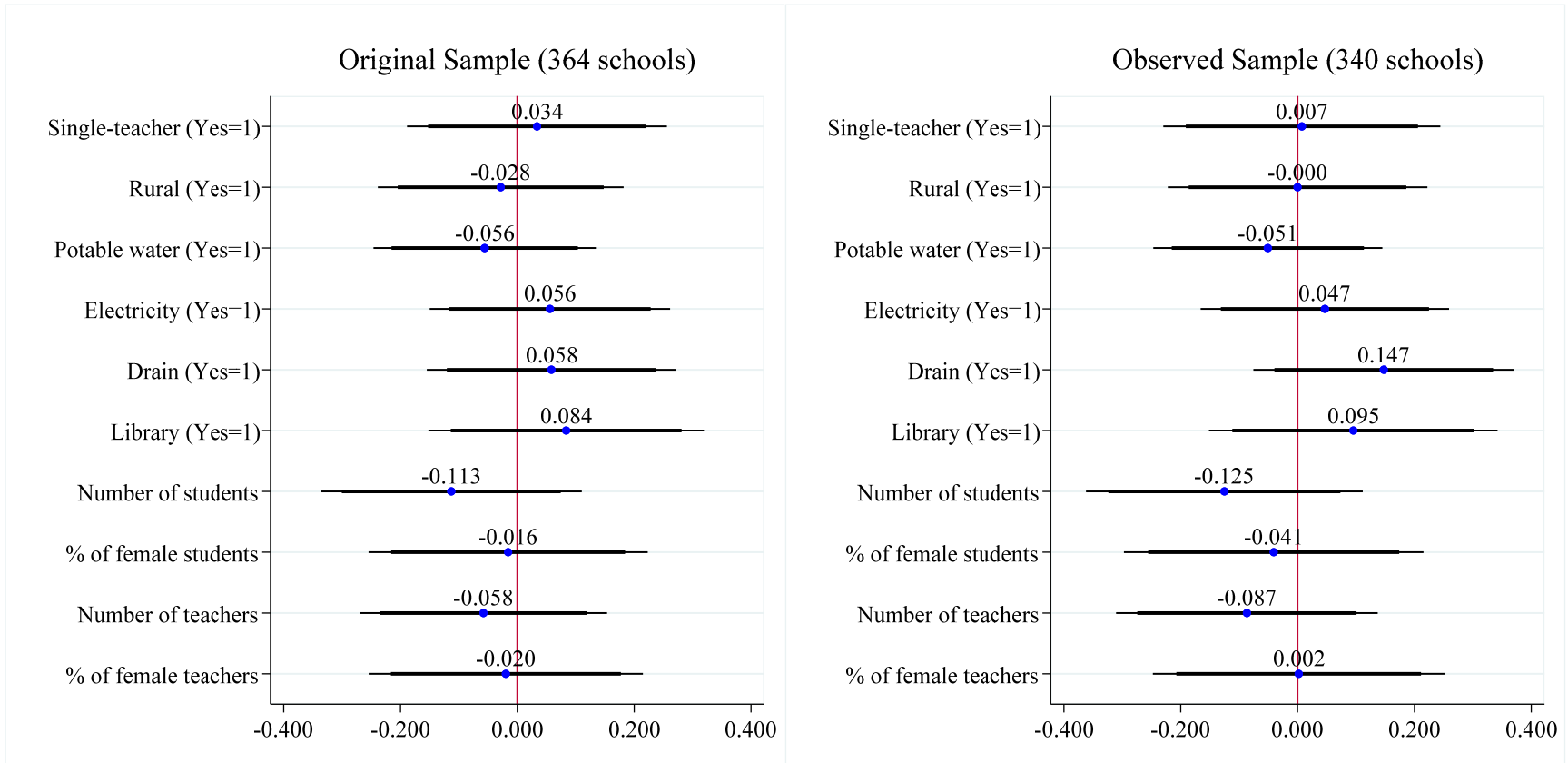
All regressions include UGEL fixed effects. Standard errors clustered at the school level.

Estimates indicate differences in the standardized characteristics of control and treatment groups. Thick and thin lines indicate 90% and 95% confidence intervals, respectively.

We do not present the differences in teacher experience and pedagogical degree for the original sample because we do not have information on those variables for the teachers that were not observed at the end of year 2.



**Figure 3**  
**Balance in School Characteristics in the Original and Observed Evaluation Sample Schools**



All regressions include UGEL fixed effects.

Estimates indicate differences in the standardized characteristics of control and treatment groups. Thick and thin lines indicate 90% and 95% confidence intervals, respectively.

Third, we used data from exams given to teachers in 2014 and 2015 that were used as part of the process by which contract teachers could become permanent civil service teachers and civil service teachers apply for promotion. We found that teachers who scored higher on those exams were less likely to move from other schools in Peru to either an APM school or a non-APM school in our 6,218 randomized expansion schools, and also that teachers who scored higher were less likely to move out of the 6,218 randomized expansion schools to other schools in Peru (see Appendix Table A10). Most importantly, there is no relationship between these test scores and whether the teachers moved to an APM or a non-APM school, which shows that there is no systematic movement of better (or worse) teachers into APM or non-APM schools.

For the student-level data, there is little attrition. Using administrative data on enrollment, we found almost all the students who started in our sample in 2016 (year 1). Student turnover, unlike teacher turnover, is relatively rare in rural primary schools, especially among those targeted by APM since almost 95% are in rural areas where there are very few schools to choose from. Excluding students in their final year of primary school, and averaging over the years 2013 to 2016, only 6.9% of the students in our 6,218 primary schools in a given year were not in the same school in the next year.

### 4.3 Teacher Turnover and the Proportions of the Four Types of Teacher

We use administrative data on the location of teachers as well as the framework established in Section 3 to examine teacher turnover and the proportions of the four types of teacher in the sample.<sup>22</sup> Table 4 shows the 2016-2017 turnover behavior of Sample 1 teachers (i.e. the 12,189<sup>23</sup> teachers in the 6,218 randomized schools in 2016).

**Table 4: Distribution of Year 1 Teachers by Their Destination School in Year 2**

Treatment Arm in 2016	2016-2017 Turnover	Teachers	Percent
APM school	Stayed in the same school	4,222	63.2
	Moved to an APM school	806	12.1
	Moved to a non-APM school	1,649	24.7
	Total	6,677	100.0
Non-APM school	Stayed in the same school	2,847	62.4
	Moved to an APM school	440	9.7
	Moved to a non-APM school	1,274	27.9
	Total	4,561	100.0

<sup>22</sup> Table A5 in the Appendix shows where teachers assigned to APM and non-APM schools in the randomization year end up in each type of school one year later according to their type and initial sorting.

<sup>23</sup> Table 4 excludes 951 teachers (7.8% of the 12,189 teachers) in the 2016 randomization sample who were not found in the administrative data in 2017; they most likely left the public education system.

By comparing the proportions of teachers in APM and non-APM schools in year 1 who moved to an APM school in year 2 (the difference between equations (A4) and (A1) in Appendix Table A5), we estimate that  $\sigma p^L = -0.024$ , where  $\sigma$  is the proportion of likers in an APM school in a given year (e.g. year 1) who remain in the same school in the next year (e.g. year 2), rather than moving to a different APM school.<sup>24</sup> Similarly, by comparing the proportions of teachers in APM and non-APM schools who moved to a non-APM school from year 1 to year 2 (the difference between equations A5 and A2 in Appendix Table A5), we estimate that  $v p^D$  equals  $-0.032$ , where  $v$  is the proportion of dislikers in a non-APM school in a given year (e.g. year 1) who remain in the same school in the next year (e.g. year 2), rather than moving to a different non-APM school.

Both  $\sigma p^L$  and  $v p^D$  are very close to 0. For  $\sigma p^L$  to equal 0, either  $\sigma$  or  $p^L$  (or both) must equal 0. If  $\sigma = 0$ , then all likers change from one APM school to another APM school in the following year. Similarly,  $v = 0$  implies that all dislikers already in a non-APM school in a given year move to another non-APM school the next year. Such turnover seems very unlikely since most teachers (63%) remained in the same school even before the randomized expansion of the APM program (see Table 5). By definition, likers and dislikers have strong incentives to move between schools if, in year 1, they find themselves in a school that is the opposite of their preference (likers starting in a non-APM school or dislikers starting in an APM school), but when they are placed in the school of their preferred type, we would expect turnover to be similar to what was observed in the sample before the program started, 36.6%, not 100%. Therefore, both  $\sigma = 0$  and  $v = 0$  seem very unlikely. The other option, which we consider the most realistic, is that  $p^L$  and  $p^D$  are equal to 0: there are no likers or dislikers.

The conclusion that there are no likers or dislikers is a strong claim, so we offer two additional pieces of supporting evidence. First, we analyze how teacher turnover changed over time. If there are likers and dislikers, we would expect an increased movement of teachers in the first year after the randomized expansion of APM as likers and dislikers move to the schools of their preferred type. Since schools stick to their random

---

<sup>24</sup> To see how this was calculated, this definition of  $\sigma$  implies that the proportion of likers who move to another APM school is  $1 - \sigma$ . Recall that  $\mu$  is the proportion of movers in any school who (randomly) move to an APM school in the following year. Thus, of all teachers in an APM school in year 1,  $p^L(1 - \sigma) + p^M\mu$  is the proportion who move to other APM schools in year 2, and our data show that this proportion is 0.121 (see Appendix Table A5). Similarly, the proportion of teachers in non-APM schools in year 1 who move to an APM school in year 2 is  $p^L + p^M\mu$ , and this proportion equals 0.097 in our data. The difference between these two proportions equals  $\sigma p^L$ , which is  $-0.024$  in our data. Note that this difference includes the estimates for the mentioned parameters as well as random differences in proportions that arise due to sampling. Thus small negative estimates are possible if a parameter equals 0.

assignment in later years, we would expect that most of this extra turnover would occur in year 2 (2017), although some could occur in later years if some “potential” likers and dislikers are unable to move to their preferred schools in year 2. Therefore, if there are likers or dislikers, there should be a large spike in the number of teachers moving across treatment arms between 2016 and 2017, followed by a gradual return to regular levels of movement (from movers randomly moving between APM and non-APM schools, and likers and dislikers moving to another school of their preferred type). Table 5 shows the evolution of teacher movement across treatment arms from 2015 to 2019. There is no spike in the movement from APM to non-APM schools from 2016 to 2017; it remains at 14%, the same rate as from 2015 to 2016, and slightly less than from 2017 to 2018. A similar pattern holds for movement from non-APM to APM schools, which from 2016 to 2017 increased slightly to 12% (from 11% from 2015 to 2016) and remained at 12% from 2017 to 2018. These trends are consistent with the claim of no likers or dislikers.

**Table 5: Teacher Turnover Between APM and non-APM schools**

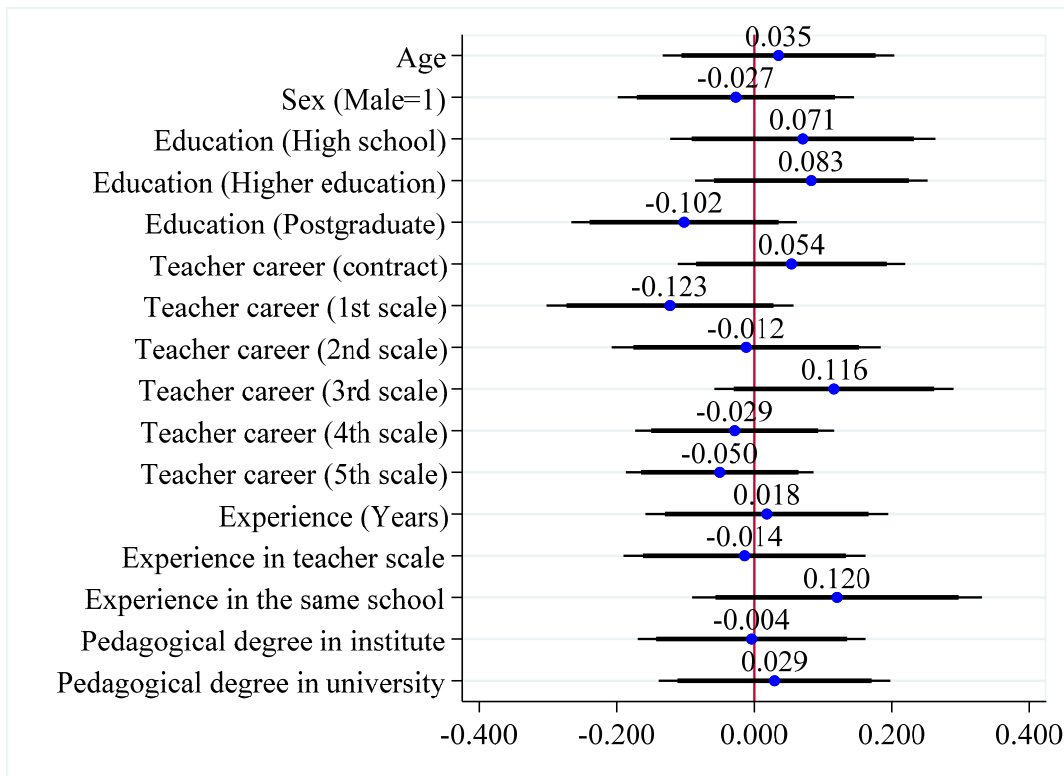
APM Schools	2015 to 2016	2016 to 2017	2017 to 2018	2018 to 2019
Stayed in Same School	63%	65%	62%	65%
Moved to an APM School	8%	10%	10%	9%
Moved to a Non-APM School	<b>14%</b>	<b>14%</b>	<b>15%</b>	<b>12%</b>
Moved Out of Target Schools	15%	12%	14%	13%
Non-APM Schools	2015	2016	2017	2018
Stayed in Same School	63%	66%	62%	65%
Moved to an APM School	<b>11%</b>	<b>12%</b>	<b>12%</b>	<b>11%</b>
Moved to a Non-APM School	11%	10%	11%	10%
Moved Out of Target Schools	15%	12%	15%	14%
All Schools	2015	2016	2017	2018
Stayed in Same School	63%	65%	62%	65%
Moved to an APM School	10%	11%	11%	10%
Moved to a Non-APM School	12%	12%	13%	11%
Moved Out of Target Schools	15%	12%	14%	14%

Note: This table shows the year-to-year turnover status of teachers who started each two-year period in a school within one of the 6,218 randomized expansion schools.

A second piece of additional evidence for the claim of no likers or dislikers is comparisons of the characteristics of teachers who worked in the randomized pedagogical skill sample in 2017 (Sample 2). If there were likers or dislikers one would expect the characteristics of teachers to differ between APM and non-APM schools after turnover, as likers would be only in APM schools while dislikers would be only in non-

APM schools. Figure 4 shows estimates of treatment effects of APM on a wide set of teacher characteristics in the randomized expansion sample in 2017. We find no effect for any of the characteristics, suggesting that there was no systematic selection of teachers into either APM or non-APM schools, further supporting the claim of no likers or dislikers.

**Figure 4: Treatment Effects on the Composition of Teacher Characteristics among the Teachers in Randomized Pedagogical Skill Sample Schools in 2017 (Sample 2)**



All regressions include UGEL fixed effects. Standard errors are clustered at the school level.

Estimates indicate differences in the standardized characteristics of control and treatment groups.

Thick and thin lines indicate 90% and 95% confidence intervals, respectively.

## 5. The Treatment Effects of APM

### 5.1 Teacher Skills

*Overall teacher skills.* This subsection presents estimates of  $E[y^2 | R_{tchr, year 1} = 1] - E[y^2 | R_{tchr, year 1} = 0]$ , that is estimates of  $ITT_{tchr}(2)$  in equation (6'), and  $E[y^2 | R = 1] - E[y^2 | R = 0]$ , estimates of  $ITT_{sch}(2)$  in equation (9'), using OLS regressions for the

455 Sample 1 teachers and the 640 Sample 2 teachers (see Table 3), respectively. We also present the estimates obtained by regressing  $y^2$  on the predicted years of treatment, instrumented by random assignment in year 1, using Sample 1 teachers. As explained in subsection 3.5, this IV approach provides a consistent estimate of the  $ACR_{\text{tchr}}(2)$  treatment effect. For all estimates, the dependent variable,  $y^2$ , is an index of pedagogical skills that averages the standardized scores of the eight indicators obtained from classroom observations (see subsection 2.3). We present estimates with and without teacher characteristics as covariates when using Sample 1.<sup>25</sup> Table 6 presents these results.

**Table 6: Aggregate Skill: Ordinary Least Squares (OLS) Estimates and IV Estimates**

	Ordinary Least Squares Estimates			IV Estimates	
	Sample 1		Sample 2	Sample 1	
	(1)	(2)	(3)	(4)	(5)
<b>Treatment</b>	<b>0.287***</b> <b>(0.108)</b>	<b>0.314***</b> <b>(0.102)</b>	<b>0.195**</b> <b>(0.097)</b>	<b>0.159***</b> <b>(0.054)</b>	<b>0.174***</b> <b>(0.050)</b>
Experience	--	0.000 (0.009)	--	--	-0.000 (0.008)
Contract teacher	--	0.152 (0.162)	--	--	0.145 (0.145)
Teacher career Level	--	0.114** (0.046)	--	--	0.113*** (0.041)
Sex (men = 1)	--	-0.313*** (0.099)	--	--	-0.315*** (0.089)
Age	--	-0.029*** (0.009)	--	--	-0.028*** (0.008)
R <sup>2</sup>	0.29	0.37	0.23	0.29	0.37
Sample Size	455	455	640	455	455

\*  $p < 0.1$ ; \*\*  $p < 0.05$ ; \*\*\*  $p < 0.01$

All regressions include UGEL fixed effects. Standard errors clustered at the school level are in parentheses.

Before discussing the results, recall the claim (subsection 4.3) that our population of teachers has no likers or dislikers. Recall also (subsection 3.4) that, if there are no likers or dislikers, both  $\hat{\beta}_1^y_{\text{OLS},t=2}$  and  $\hat{\beta}_2^y_{\text{OLS},t=2}$  estimate  $ITT_{\text{tchr}}(2)$ , which equals  $ITT_{\text{sch}}(2)$ . Thus, all OLS estimates in Table 6 consistently estimate the same parameter.

<sup>25</sup> The use of teacher characteristics as covariates is appropriate only for Sample 1 because characteristics of Sample 2 teachers can be affected by the treatment. In Table A6 in the Appendix, we test for interactions between the treatment status and the characteristics of Sample 1 teachers. We find no evidence of heterogeneity by teacher experience, type of contract, position in the teacher career or sex. These results are important as they support the linearity assumption for the teacher skills production function in equation (1).

The first and second columns of Table 6 present estimates of  $ITT_{tchr}(2)$ . The estimate in Column (1), which does not control for teacher characteristics, indicates that offering APM for two years increases teachers' pedagogical skills by 0.29 standard deviations (s.d.). The estimate in Column (2), when teacher characteristics are added as covariates, is very similar: 0.31 s.d. The estimate for  $ITT_{sch}(2)$  in Column (3), 0.20 s.d., is somewhat lower, even though  $ITT_{sch}(2)$  should equal  $ITT_{tchr}(2)$ . Recall that Sample 1 teachers had high rates of attrition due to difficulties finding teachers who moved; this implies that remainers are very likely overrepresented in Sample 1. In contrast, the proportions of remainers and movers in Sample 2 should correspond to their proportions in the population of teachers in the 6,218 randomized expansion schools. Thus, the Column (3) estimate is our preferred estimate of  $ITT_{tchr}(2)$ , which also equals  $ITT_{sch}(2)$ ; the effect after two years on teachers' aggregate pedagogical skill of *assigning* them to an APM school in year 1 is a 0.20 s.d. increase in those skills

Our estimate that  $ITT_{tchr}(2) = ITT_{sch}(2) = 0.20$  sheds some light on other parameters of interest. Recall that, in general,  $ATE_{tchr}(2) \geq ITT_{tchr}(2)$ , and if there are no likers and dislikers then  $ATE_{sch}(2) \geq ITT_{sch}(2)$ . Thus the effect of two years of APM coaching on the aggregate pedagogical practice of the average teacher,  $ATE_{tchr}(2)$ , and the effect of APM on the aggregate pedagogical practice of the teachers in APM schools in year 2,  $ATE_{sch}(2)$ , are at least as large as, and likely larger than, 0.2 s.d.

Columns (4) and (5) in Table 6 present our IV estimates of  $ACR_{tchr}(2)$  using Sample 1 teachers. They show that, averaging over all years of coaching received, an additional year of coaching increases by 0.16 to 0.17 s.d. the average pedagogical skill of all teachers, but this average gives remainers a "double weight" because random assignment to an APM school induces them to obtain two years of coaching. Consistent with the fact that  $ACR_{tchr}(2)$  equals  $ITT_{tchr}(2)/(1+p^R)$ , this IV estimate, which is a per year estimate, is somewhat larger than (half of) the Sample 1 estimate of  $ITT_{tchr}(2)$ , an estimate of cumulative impact over two years, in column (2).

*Specific Pedagogical Skills.* The discussion thus far has focused on the aggregate index of pedagogical skills, but one can also estimate  $ITT_{tchr}(2)$  for each of the eight more specific pedagogical skills shown in Table 1. Table 7 shows these results. To minimize spurious statistical significance due to multiple hypothesis testing, Table 7 also presents adjusted p-values, using the Romano and Wolf (2016) stepdown method to account for multiple hypothesis testing; these are in brackets below the standard errors.

The estimates in Table 7 indicate that the biggest impact of assigning teachers to the APM program, in terms of both the size and the statistical significance of the estimated parameters, is on teachers' lesson planning; the point estimates are 0.34 s.d. for Sample 1 and 0.38 s.d. for Sample 2. There is also evidence that APM raises teachers' pedagogical skills in developing their students' critical thinking, although the statistical significance is at best only marginal after controlling for multiple hypothesis testing.

## 5.2 Student Learning

This subsection explores the impact of the APM coaching program on student learning, as measured by the National Student Evaluation (ECE) taken one and three years after the program began (that is, 2016 and 2018). We compare student test scores in the APM and non-APM schools in the much larger student test score sample. This sample is not restricted to the 340 schools with pedagogical practices data, but it is restricted to those schools that participated in the 2016 ECE and the 2018 ECE. As explained earlier, only schools with five or more students in the relevant grade take the ECE, so we have test scores for only 2,567 of the 6,218 randomized expansion schools.

Table 8 presents estimates of the APM coaching program's treatment effects on average ECE scores for the sample of 2,567 schools in 2016 and 2018, after one and three years of coaching. The ECE is taken at the end of the school year (which is also the end of the calendar year), so the 2016 ECE yields estimates of the APM program's impact after one year for students in grade 2. All teachers complied with their random assignment in 2016, so this is an estimate of  $ATE_{stud}(1)$ , the average treatment effect of one year of APM on student learning. In 2018, the ECE was conducted again, but this time it was done in grade 4, which in general contains the same students who were tested in 2016 in grade 2, except that it excludes students who repeated grade 2 or 3 (about 7-8% repeat each year). The 2018 ECE allows us to test for the impact of the program after three full years of implementation. Students almost always comply with treatment assignment, yet many teachers switched schools between 2016 and 2018, so we cannot estimate the average treatment effect,  $ATE_{stud}(3)$  for three years. Rather, we estimate  $ITT_{stud}(3)$ , which is a lower bound of  $ATE_{stud}(3)$  if there are no likers or dislikers.



**Table 7**  
**Disaggregated Skills: Ordinary Least Squares Estimates**

	(1) Lesson Planning	(2) Time Management	(3) Critical Thinking	(4) Student Participation	(5) Class Feedback	(6) Written Feedback	(7) Classroom Relationships	(8) Behavior Management
<b>Panel A. Sample 1</b>								
Treatment	0.335 (0.105)*** [0.018]**	0.0811 (0.105) [0.701]	0.268 (0.099)*** [0.073]*	0.168 (0.107) [0.488]	0.193 (0.104)* [0.348]	0.138 (0.098) [0.511]	0.0719 (0.116) [0.701]	0.123 (0.105) [0.580]
N	448	450	450	450	450	448	450	450
R-squared	0.307	0.221	0.281	0.364	0.371	0.332	0.263	0.277
<b>Panel B. Sample 2</b>								
Treatment	0.375 (0.088)*** [0.002]***	-0.0673 (0.092) [0.926]	0.194 (0.094)** [0.284]	0.0627 (0.090) [0.926]	0.0881 (0.096) [0.891]	0.175 (0.098)* [0.422]	0.0225 (0.095) [0.967]	0.0190 (0.089) [0.967]
N	633	633	633	632	633	631	633	632
R-squared	0.245	0.171	0.200	0.260	0.277	0.236	0.209	0.238

\*  $p < 0.1$ ; \*\*  $p < 0.05$ ; \*\*\*  $p < 0.01$ .

Note: Effects are measured in standard deviations. Regressions of Panel A include the following control variables: experience, contract teacher, teacher career level, sex and age. All regressions include UGEL fixed effects. Standard errors clustered at the school level are reported in parentheses and adjusted p-values for multiple hypotheses testing are reported in brackets. We calculate the adjusted p-values using the stepdown method of Romano and Wolf (2016).

*5.2.1 Results after one year.* Table 8 presents estimates of the program’s treatment effects on standardized test scores for mathematics and reading comprehension.<sup>26</sup> Columns (1) and (4) show estimates of  $ATE_{stud}(1)$  after one year of implementation, while columns (2), (3) (5) and (6) show ITT and ACR estimates after three years of the program. While the program was designed by the Ministry of Education, it was implemented by each local school district (UGEL),<sup>27</sup> so our preferred specification, shown in this table, includes school district fixed effects, which also control for any differences in actual program implementation within each region. All Table 8 regressions also control for school size (number of teachers and students), which was slightly unbalanced at baseline. We cluster standard errors at the school level in all regressions, following Abadie et al. (2017), since the treatment is assigned at the school level.

The APM coaching program has significantly positive impacts on student learning. After one year, average test scores increase by 0.10 and 0.07 standard deviations (s.d.) in math and reading comprehension, respectively. These are average treatment effects,  $ATE_{stud}(1)$ , and they suggest that coaching that provides regular, individualized support to teachers can be an effective policy to increase student learning. For perspective, note that the effect after one year is similar in magnitude to the median effect on learning outcomes of 234 education studies in low and middle income countries reviewed by Evans and Yuan (2022). And when compared to the median for large studies (those with over 5,000 students), the effect of the APM program after only one year is almost double that median effect (0.05 s.d.).

Table A3 in the appendix shows how estimates change when using regional, rather than school district, fixed effects, and when excluding controls. Both of those changes reduce the size of the coefficient slightly, but the results are generally robust to these changes.<sup>28</sup> Table A3 includes another specification, column (4), that adds to the analysis the panel data available from 2010 to 2018 and adds school-level fixed effects and state-specific time trends, without any controls; its results are very close to those of main OLS specification in Table 8.

*5.2.2 Results after three years.* Columns (2) and (5) of Table 8 show the effects of the APM program in 2018, after three years. Recall that in 2018 the standardized test is for grade 4, so

---

<sup>26</sup> Recall that ECE scores exist only for schools with five or more students in a given grade; this greatly reduces the number of schools in the student test score sample. Table A1 in the appendix shows that almost all characteristics of the schools with test scores are very similar to those for the 6,218 randomized expansion schools. The baseline balance in Table 2 is for this smaller subsample of schools, which is the relevant sample for analysis.

<sup>27</sup> Peru’s 225 school districts (UGELs) are managed by school boards, which implement education policies in their districts. Each UGEL is overseen by its Regional Education Board (Dirección de Educación Regional).

<sup>28</sup> The exception is reading comprehension scores after one year of APM; they are significant only if controls are included. Yet the treatment effects after three years, however, are robust to excluding controls for both subjects.

**Table 8: Results on Student Learning After One and Three Years of Coaching**

	Mathematics			Reading		
	1 Year		3 Years	1 Year		3 Years
	OLS (1)	OLS (2)	IV (3)	OLS (4)	OLS (5)	IV (6)
Treatment	0.100*** (0.032)	0.113*** (0.032)		0.072** (0.030)	0.099*** (0.031)	
Cumulative Years Treated			0.177*** (0.052)			0.158*** (0.049)
Coefficient on random assignment in first-stage regression			0.623*** (0.008)			0.623*** (0.008)
F-stat			6,126			6,126
Control Mean	-0.590	-0.734	-0.734	-0.861	-0.812	-0.812
Observations	22,308	18,357	18,261	22,309	18,371	18,275
Schools	2,566	2,066	2,053	2,566	2,066	2,053
R <sup>2</sup>	0.142	0.182	0.184	0.162	0.168	0.169

Note: This table shows the average treatment effect of the coaching program on standardized student test scores. Columns 1 and 4 show the effect after one year of treatment in 2016, while columns 2 and 5 show the effects after three years of treatment in 2018. Finally, columns 3 and 6 present 2SLS estimates using the random treatment assignment as an instrument for the proportion of teachers effectively treated for three years. All specifications include school district fixed effects and control for school size (number of teachers and students), which is not balanced at baseline (See Table A3 for additional specifications). All results use standardized exam scores and can be interpreted as standard deviations. Regressions are run at the student level, with robust standard errors clustered by school presented in parentheses. \*\*\* p-value <0.01, \*\* p-value <0.05, \* p-value <0.1.

that, except for repeaters, we follow the same students observed in 2016 in grade 2 after two more years of exposure to APM. The estimated program effects, which are now ITT effects ( $ITT_{stud}(3)$ ) and so are lower bounds for ATE ( $ATE_{stud}(3)$ ), remain positive after three years of the program and are slightly higher: 0.11 s.d. for math, 0.10 s.d. for reading comprehension.

These ITT results show the average effect on students learning after three years for schools that were randomly assigned to the APM program in 2016. Yet the exposure of students to treated teachers, and therefore the effective treatment dose, differs widely among APM schools as a result of teacher turnover. To estimate the impact on students of being exposed to one more year of teacher coaching, we use random assignment in 2016 to instrument students' exposure to coached teachers in each school. We have data on teachers' school assignment, so we constructed a variable that captures the intensity of coaching for the teachers present in each year (since the program started) in a given school. This incorporates the coaching history of all teachers that the students had over the course of three years.<sup>29</sup>

We created a variable that measures the variation in the intensity of coaching received by teachers, who received either one, two or three years of coaching in the past three years;

<sup>29</sup> Strictly speaking, we construct and average "history" over all teachers in a given school in a given year, since we cannot match students to individual teachers. Note, however, that 20% of the schools in our student test score sample had only one teacher, so for these schools we are matching students to their specific teacher.

students, in turn, were exposed, over those three years, to teachers with varying years of treatment. Our constructed variable is based on the total years of coaching that each teacher received and calculates for all students the average intensity of coaching that the teachers in their school had received, for each of the three years that a student was in his or her school.

For example, students in a school A that was randomly assigned to be an APM school would have been exposed to teachers with one year of coaching in year 1 (since all schools complied with their random assignment and teachers had not yet been able to switch schools). If all teachers remain in that school the students in school A would be exposed in their second year to a teacher coached for two years, bringing their total coaching exposure to three years (one in the first year and two in the second), and similarly (if there were no teacher turnover) would be exposed to teachers with an average of three years of coaching in year 3, bringing students total exposure to six years of coaching over the three years. In contrast, students in an APM school B that experiences full teacher turnover each year would be exposed to one year of coaching in year 1, another year of coaching in year 2 (assuming that all teachers left and all new teachers had not been coached in their first year), and another year of coaching in year 3 (if all teachers once again left and all new teachers had not been coached in years 1 and 2), for a total of 3 years of coaching exposure. Students in non-APM schools that receive treated teachers can also receive exposure to some years of coaching, depending on the extent of coaching their teachers received previously.

We created a variable for student exposure to coached teachers up through year 3 that ranges from 0 to 6 years. To simplify the interpretation, we divide this dosage variable by 6 to generate a variable that varies from 0 to 1, where 1 indicates students who were exposed to a full dosage of coaching, meaning that all the teachers to whom these students were exposed had been coached in the year of student exposure and in all previous years. The coefficient on this instrumented proportion thus measures the effect on student test scores of being exposed to a set of “fully coached” teachers for three years. Thus, it estimates  $ACR_{stud}(3)$ , multiplied by 6 (recall that  $ACR_{stud}(3)$  is a per-year-of-coaching effect). The average value of this (rescaled) exposure variable in 2018 is 0.68 for APM schools and 0.05 for non-APM schools.

As with the OLS estimates, our preferred specifications for both stages of the IV estimates include school district fixed effects and controls for number of teachers and number of students. Standard errors are clustered at the school level.

Columns (3) and (6) of Table 8 present the IV estimates of the APM program on student learning, instrumenting the proportion of teachers’ accumulated coaching by schools’ random assignment. As expected, they are larger than the OLS estimates in columns (2) and

(5) because they measure the impact of “full” student exposure to APM coaching over three years. These estimates of  $ACR_{stud}(3)$  are 0.177 for mathematics and 0.158 for reading comprehension. This is consistent with the fact that the first stage coefficient is only about 0.623, which means that being assigned to an APM school increases the proportion of treated teacher-years by 62 percentage points relative to non-APM schools, which on average have only 0.05 of their teacher-years treated. This high rate of teachers switching schools, which means that the average teacher in a school treated in year 3 (2018) had been coached for only two years, could explain why the ITT effects after three years of treatment (columns (2) and (5)) are not significantly higher than the ITT effects after one year (columns (1) and (4)). Even so, once we account for teacher turnover, the IV estimates in Table 8 suggest decreasing but positive marginal effects to the second and third years of coaching.

*5.2.3 Student Skill Distribution.* Teachers could respond to the treatment in various ways: for example, they could focus their efforts on the lower end of the student learning distribution to help the weaker students, or given the high stakes nature of the tests<sup>30</sup> they could focus on top students by shifting resources and attention away from those who struggle, or they could acquire skills that help them engage with students across the entire student skill distribution. To test which part of the student grade distribution is shifting in response to the treatment we run a quantile regression, taking advantage of the availability of individual student test scores.

Figure 5 shows the results for quantile regressions for each decile of the student test score distribution, using our preferred specification that has school district fixed effects, after three years of the APM program (in 2018). We find that the program raised student test scores along the entire student skill distribution and we cannot reject that treatment effects are constant across all deciles. This suggests that the program, which focuses on individual teacher weaknesses, helps teachers to deal with the particular challenges their students face regardless of those students’ position in the student skill distribution.

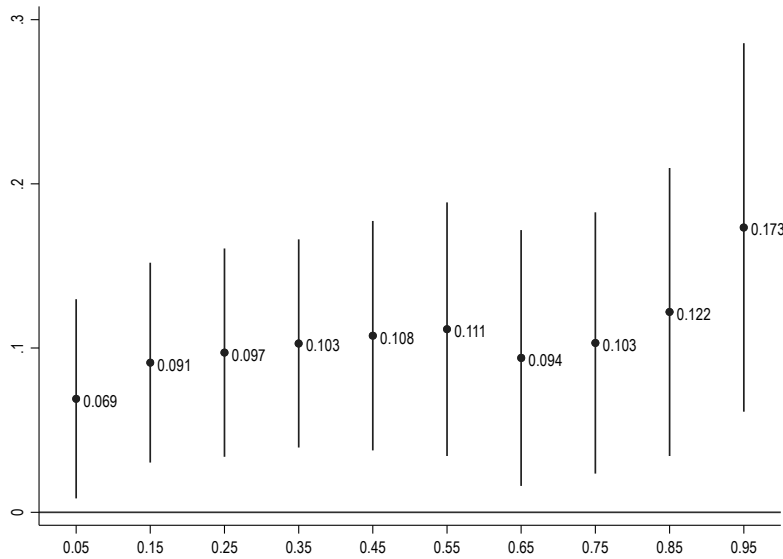
## 6. Concluding Remarks

Teacher quality plays a key role in student learning, but the quality of teachers is often low, especially in developing countries. Given the vital role that human capital plays in economic growth and individuals’ income and well-being, a key policy priority is to develop policies

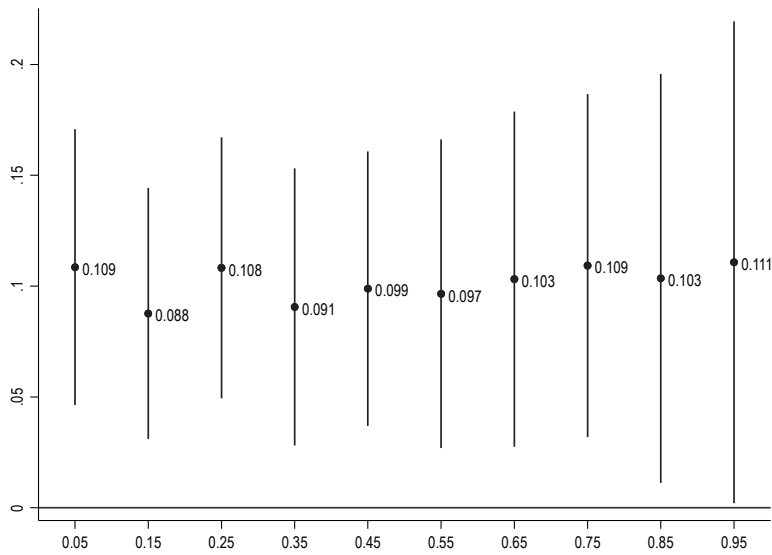
---

<sup>30</sup> There are some incentive payment schemes that pay teachers bonuses according to their schools’ performance on these tests. They should not affect our estimates since they apply to both APM and non-APM schools.

**Figure 5: Quantile Regression Results After Three Years of Implementation (2018)**



a) Mathematics



b) Reading Comprehension

Note: These figures show the quantile regression coefficients for the effect of the program on standardized test scores after three years of implementation (2018) for each decile of the distribution of student test scores. 95% C.I. shown with standard errors clustered by school. All specifications include school district fixed effects and control for school size.

that increase teacher quality. The success of teacher training programs in raising teacher quality is, at best, mixed, but teacher coaching programs are a promising policy option.

We have estimated the effect of a large-scale teacher coaching program operating in a context of high teacher turnover in rural Peru on a broad range of pedagogical skills, and on

student learning. Our study contributes to the literature on teacher training and pedagogy by addressing the issues of scale and teacher turnover as potential threats to the effectiveness of coaching, and by presenting evidence that the general pedagogical skills of the current stock of teachers can be improved. This research also contributes to the literature by developing an analytical framework that defines different types of treatment effects when teacher turnover is present and explains which treatment effects can be estimated.

When teacher turnover is present, the success of teacher training or coaching programs can be judged from two perspectives, the impact on the teachers who were initially offered the treatment, regardless of whether they stay in their schools or move to a different school, and the impact on the teachers and students in treated schools after turnover has occurred. It is possible to estimate intent to treat (ITT) effects for the first perspective if one has a sample of teachers that follows them when they change schools, or if one has data on teacher skills from the schools that were randomly assigned to treatment and control groups *and* teacher turnover is unrelated to the program, and for the second perspective using data on teachers' skills and student learning in treated and control schools after turnover has occurred. We also show that, unfortunately, average treatment effects (ATE) cannot be estimated without bias even when turnover is unrelated to the program. However, we show that ITT estimates serve as a lower bound for ATE for the teachers who were initially offered the treatment (the first perspective). Yet from the second perspective ITTs are a lower bound for ATEs only if teacher turnover is unrelated to the program. We believe that this framework can be useful for future education evaluations carried out in contexts of high teacher turnover or, more generally, in any evaluations where treatments are offered at a cluster-level and service providers can change clusters while the intervention is still in progress.

We find that, after two years, the program has an (average) intent to treat (ITT) effect that increases teachers' pedagogical skills by 0.20 s.d.: this estimate applies to both perspectives. This effect is concentrated on two dimensions of the pedagogical practices: lesson planning and, to a lesser extent, encouraging students' critical thinking. We also estimated the ITT effect of the program on student learning and found positive effects after one year (0.07-0.10 s.d.) and after three years (0.10-0.11 s.d.) of coaching.

This research also contributes to the discussion about how to improve the pedagogical skill of teachers serving rural schools in ways that are most cost-effective. Rural schools are often located in hard-to-reach areas that tend to be avoided by teachers if they are given a choice. One potential way to improve pedagogical skills and student learning in rural schools is to offer incentives to attract more talented teachers. The rural bonus scheme in Peru pursues

this objective by offering a 30% salary increase to those teachers who accept a position in a rural school. This bonus has had a small positive effect on the probability of filling a teacher vacancy but has shown no effects on learning outcomes (Castro and Esposito, 2022).

The cost of the coaching program evaluated in this study is around US\$ 3,000 per teacher, per year. This is about 30% of the average annual salary of a primary school teacher in Peru, and it is similar to the wage premium offered by the bonus program, with two important differences: coaching is only a three-year investment (not a permanent salary increase), and we have shown that it is effective for increasing student learning.

Another policy to increase pedagogical skills and student learning in rural schools is to offer incentives for (current) teachers to increase their productivity. Recent studies have shown that expensive policies based on large unconditional salary increases can reduce the number of teachers taking second jobs but have no effects on teacher productivity (de Ree et al., 2018). Pay-for-performance programs offer another alternative to improve teachers' productivity. The impact of these types of incentives has been examined in several low and middle-income countries, with mixed results. Very few studies, however, have estimated the effect of such programs in the context of a nation-wide intervention. A recent study by Bellés-Obrero and Lombardi (2019) evaluated the effect of a national pay-for-performance program implemented in 2015 in public secondary schools in Peru. The program, *Bono Escuela*, offers an additional monthly salary to the principal and teachers of the schools that rank in the top 20% of the national 8<sup>th</sup> grade student evaluation within their school district. The authors found no effect on student learning, as well as evidence that this lack of effect was due to teachers' uncertainty regarding which pedagogical practices raise student learning.

Our results show that a large-scale coaching program can be an effective policy to improve the performance of existing teachers at a reasonable cost. Rather than offering incentives for teachers to devote more time and effort to the task (something that might not be effective if teachers lack the requisite pedagogical skills), the results of this paper suggest that it is more effective to directly intervene to enhance their teaching skills.



## References

- Abadie, A., Athey, S., Imbens, G. W., and Wooldridge, J. (2017). "When should you adjust standard errors for clustering?" Technical Report 24003, National Bureau of Economic Research.
- Akpur, U. (2020). "Critical, reflective, creative thinking and their reflections on academic achievement." *Thinking Skills and Creativity*, 37, 100683.
- Albornoz, F., Anauati, M., Furman, M., Luzuriaga, M., Podesta, M., & Taylor, I. (2020). "Training to Teach Science: Experimental Evidence from Argentina." *The World Bank Economic Review*, 34(2), 393-417.
- Allen, J., Gregory, A., Mikami, A., Lun, J., Hamre, B., & Pianta, R. (2013). "Observations of Effective Teacher-Student Interactions in Secondary School Classrooms: Predicting Student Achievement With the Classroom Assessment Scoring System-Secondary." *School Psychology Review*, 42(1), 76-98.
- Angrist, J., and Imbens, G. (1995). "Two-Stage Least Squares Estimation of Average Causal Effects in Models with Variable Treatment Intensity." *Journal of the American Statistical Association* 90(430), 431-442.
- Angrist, J., G. Imbens & Rubin, D. (1996). "Identification of Causal Effects Using Instrumental Variables." *Journal of the American Statistical Association* 91(434), 444-455.
- Banerjee, A., Chattopadhyay, R., Duflo, E., Keniston, D., & Singh, N. (2021). "Improving Police Performance in Rajasthan, India: Experimental Evidence on Incentives, Managerial Autonomy, and Training." *American Economic Journal: Economic Policy*, 13(1), 36-66.
- Bellés-Obrero, C., & Lombardi, M. (2019). "Teacher Performance Pay and Student Learning: Evidence from a Nationwide Program in Peru." *IZA Discussion Paper No. 12600*.
- Bennett, D., Naqvi, A., & Schmidt, W. P. (2018). "Learning, Hygiene and Traditional Medicine." *The Economic Journal*, 128(612), F545-F574.
- Bruns, B., Costa, L., & Cunha, N. (2018). "Through the looking glass: Can classroom observation and coaching improve teacher performance in Brazil?" *Economics of Education Review*, 64(1), 214-250.
- Castro, J., & Esposito, B. (2022). "The Effect of Bonuses on Teacher Retention and Student Learning in Rural Schools: A Story of Spillovers." *Education, Finance and Policy*, 17(4), 693-718.
- Chetty, R., Friedman, J., & Rockoff, J. (2014). "Measuring the Impacts of Teachers I: Evaluating Bias in Teacher Value-Added Estimates." *American Economic Review*, 104(9), 2593-2632.
- Cilliers, J., Fleisch, B., Prinsloo, C., & Taylor, S. (2020). "How to Improve Teaching Practice? An Experimental Comparison of Centralized Training and In-Classroom Coaching." *Journal of Human Resources*, 55(3):926-962.
- Clare, L., Garnier, H., Junker, B., & Correnti, R. (2010). "Investigating the Effectiveness of a Comprehensive Literacy Coaching Program in Schools with High Teacher Mobility." *The Elementary School Journal*, 111(1), 35-62.

- Clotfelter, C., Ladd, H., & Vigdor, J. (2010). "Teacher Credentials and Student Achievement in High School: A Cross Subject Analysis with Fixed Effects." *Journal of Human Resources*, 45(3), 655-681.
- Das, J., Dercon, S., Habyarimana, J., & Krishnan, P. (2007). "Teacher Shocks and Student Learning. Evidence from Zambia." *Journal of Human Resources*, 42(4), 820-862.
- de Ree, J., Muralidharan, K., Pradhan, M., & Rogers, H. (2018). "Double for Nothing? Experimental Evidence on an Unconditional Teacher Salary Increase in Indonesia." *Quarterly Journal of Economics*, 133(2), 993-1039.
- Evans, D., & Popova, A. (2016). "What Really Works to Improve Learning in Developing Countries? An Analysis of Divergent Findings in Systematic Reviews." *The World Bank Research Observer*, 31(2), 242-270.
- Evans, D. K. and Yuan, F. (2022). "How Big Are Effect Sizes in International Education Studies?" *Educational Evaluation and Policy Analysis*, 44(3), 532–540.
- Fauth, B., Decristan, J., Decker, A. T., Büttner, G., Hardy, I., Klieme, E., & Kunter, M. (2019). "The effects of teacher competence on student outcomes in elementary science education: The mediating role of teaching quality." *Teaching and Teacher Education*, 86, 102882.
- Gage, N. A., Scott, T., Hirn, R., & MacSuga-Gage, A. S. (2018). The Relationship Between Teachers' Implementation of Classroom Management Practices and Student Behavior in Elementary School. *Behavioral Disorders*, 43(2), 302–315.  
<https://doi.org/10.1177/0198742917714809>.
- Georgiadis, A., & Pitelis, C. N. (2016). "The Impact of Employees' and Managers' Training on the Performance of Small-and Medium-Sized Enterprises: Evidence from a Randomized Natural Experiment in the UK Service Sector." *British Journal of Industrial Relations*, 54(2), 409-421.
- Jukes, M., Turner, E., Dubeck, M., Halliday, K., Inyega, H., Wolf, S., and Brooker, S. (2017). "Improving literacy instruction in Kenya through teacher professional development and text messages support: A cluster randomized trial". *Journal of Research on Educational Effectiveness*, 10(3):449-481.
- Kotze, J., Fleisch, B., & Taylor, S. (2019). "Alternative forms of early grade instructional coaching: Emerging evidence from field experiments in South Africa." *International Journal of Educational Development*, 66, 203-213.
- Kovner, C. T., Brewer, C. S., Fatehi, F., & Jun, J. (2014). "What Does Nurse Turnover Rate Mean and What is the Rate?" *Policy, Politics, & Nursing Practice*, 15(3-4), 64-71.
- Kraft, M., Blazar, D., & Hogan, D. (2018). "The Effect of Teacher Coaching on Instruction and Achievement: A Meta-Analysis of the Causal Evidence." *Review of Educational Research*, 88(4), 547-588.
- Loyalka, Prashant, Anna Popova, Guirong Li, Chengfang Liu, and Henry 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.
- Lucas, Adrienne, Patrick McEwan, Moses Ngware and Moses Oketch. 2014. "Improving Early-grade Literacy in East Africa: Experimental Evidence from Kenya and Uganda". *Journal of Policy Analysis and Management* 33(4): 950-976.

- Popova, A., Evans, D., & Arancibia, V. (2016). "Training Teachers on the Job: What Works and How to Measure It." Policy Research Working Paper 7834. The World Bank: Washington, DC.
- Romano, J. P., and Wolf M. (2016). "Efficient Computation of Adjusted P-Values for Resampling-Based Stepdown Multiple Testing." *Statistics & Probability Letters*, 113, 38-40.
- Schaffner, Julie, Paul Glewwe and Uttam Sharma (2021). "Why Programs Fail: Lessons for Improving Public Service Quality from a Mixed Methods Evaluation of an Unsuccessful Teacher Training Program." Tufts University and University of Minnesota.
- Stronge, J. H., Ward, T. J., & Grant, L. W. (2011). "What Makes Good Teachers Good? A Cross-Case Analysis of the Connection Between Teacher Effectiveness and Student Achievement." *Journal of Teacher Education*, 62(4), 339–355.  
<https://doi.org/10.1177/0022487111404241>.
- Wisniewski, B., Zierer, K., & Hattie, J. (2020). "The power of feedback revisited: A meta-analysis of educational feedback research." *Frontiers in Psychology*, 10, 3087.
- World Bank. (2018). *World Development Report: Learning to Realize Education's Promise*. The World Bank: Washington DC.
- Zeitlin, Andrew (2021). "Teacher turnover in Rwanda." *Journal of African Economies*, 30:1, 81-102.

## Appendix A: Additional Tables

### Table A1. Descriptive Statistics

	Public Primary Schools		Randomized Expansion Sample		
	All (1)	Monolingual Multigrade (2)	All (3)	Test Score Sample (4)	Pedagogical Sample (5)
Math score (2015)	550.92 (93.37)	526.70 (97.67)	512.17 (93.45)	513.29 (92.16)	497.49 (92.24)
Language score (2015)	552.39 (69.62)	529.19 (69.13)	519.02 (66.78)	519.47 (65.70)	511.11 (66.40)
Number of students	88.07 (161.70)	28.78 (24.06)	28.66 (23.48)	45.96 (25.21)	29.13 (23.83)
Number of teachers	4.74 (6.79)	1.91 (1.18)	1.90 (1.14)	2.57 (1.23)	1.89 (1.15)
Number of sections (classes)	6.82 (4.81)	5.33 (1.26)	5.36 (1.19)	5.83 (0.81)	5.34 (1.21)
Teacher-Student Ratio	0.09 (0.09)	0.10 (0.09)	0.10 (0.07)	0.06 (0.03)	0.09 (0.07)
Rurality	1.90 (1.05)	2.25 (0.79)	2.35 (0.78)	2.28 (0.83)	2.43 (0.72)
% students indigenous mother tongue	23.70 (41.60)	3.22 (16.95)	5.02 (20.86)	4.27 (18.86)	2.36 (14.55)
Poverty rates	55.16 (22.96)	55.67 (20.83)	56.79 (22.83)	59.59 (22.33)	57.76 (21.24)
Ceiling material	5.90 (1.63)	5.66 (1.41)	5.64 (1.42)	5.77 (1.48)	5.59 (1.29)
Wall material	6.21 (1.20)	6.10 (1.27)	6.03 (1.37)	6.09 (1.36)	6.03 (1.36)
Floor material	2.81 (0.78)	2.86 (0.71)	2.83 (0.76)	2.85 (0.71)	2.89 (0.74)
% teachers with degree	0.97 (0.11)	0.98 (0.10)	0.98 (0.10)	0.96 (0.12)	0.98 (0.09)
Internet access	0.20 (0.40)	0.08 (0.27)	0.08 (0.28)	0.10 (0.30)	0.08 (0.27)
Receives textbooks	0.77 (0.42)	0.76 (0.42)	0.73 (0.45)	0.74 (0.44)	0.69 (0.46)
Receives workbooks	0.70 (0.46)	0.69 (0.46)	0.66 (0.47)	0.69 (0.46)	0.62 (0.49)
School-day length	8.16 (0.86)	8.22 (0.98)	8.16 (0.83)	8.15 (0.88)	8.16 (0.90)
Electricity	0.66 (0.47)	0.54 (0.50)	0.53 (0.50)	0.54 (0.50)	0.53 (0.50)
Water	0.59 (0.49)	0.50 (0.50)	0.48 (0.50)	0.53 (0.50)	0.49 (0.50)
Sanitation	0.25 (0.44)	0.13 (0.33)	0.12 (0.33)	0.14 (0.35)	0.13 (0.34)
Computers per Student	0.63 (4.24)	0.39 (2.50)	0.43 (2.58)	0.45 (2.87)	0.25 (0.96)
Schools	30,539	14,467	6,218	2,567	340

Note: This table shows the descriptive statistics for the experimental sample compared to all Peruvian public primary schools and multigrade schools. Column 1 shows the average characteristics for all public primary schools in Peru, while Column 2 restricts the sample to monolingual multigrade primary schools, the target population of the coaching program. Columns 3 and 4 show descriptive statistics for the experimental sample, with Column 3 including all schools in the sample and Column 4 restricting the sample to the subset of schools with test scores in 2016. Finally, Column 5 shows the subsample for which we measure the pedagogical skills of teachers, which is missing data from 24 very remote schools that could not be surveyed. Rurality is a categorical variable that takes values 0 for Urban schools, and 1, 2 and 3 for increasingly rural schools. Ceiling, Wall and Floor Materials are categorical variables that take values up to 7, with higher values implying better materials. Standard deviations are shown in parentheses.

**Table A2. Descriptive Statistics for Outcomes - Baseline**

	Public Primary Schools			Experimental Sample	
	All	Not in Sample	In Sample	Control	Treated
Math					
Rasch	550.18	555.49	512.15	510.29	513.44
Level 1 (beginning)	0.40	0.38	0.54	0.55	0.53
Level 2 (in process)	0.38	0.39	0.32	0.31	0.33
Level 3 (satisfactory)	0.22	0.23	0.14	0.14	0.14
Reading					
Rasch	566.45	573.08	518.96	516.09	520.95
Level 1 (beginning)	0.11	0.09	0.23	0.24	0.23
Level 2 (in process)	0.50	0.49	0.58	0.58	0.58
Level 3 (satisfactory)	0.39	0.41	0.19	0.18	0.20
N	21,396	18,774	2,622	1,070	1,552

**Table A.3. Robustness Checks on Treatment Effects on Students' Standardized Test Scores**

	Preferred Specification (1)	Without Controls (2)	Region Fixed Effects (3)	School Fixed Effects (4)
<i>Panel A. Mathematics Test Scores 2016</i>				
Treatment	0.100*** (0.032)	0.071** (0.033)	0.085*** (0.033)	0.099*** (0.033)
Observations	22308	22317	22308	132016
Schools	2566	2567	2566	2567
R <sup>2</sup>	0.142	0.133	0.089	0.285
<i>Panel B. Mathematics Test Scores 2018</i>				
Treatment	0.113*** (0.032)	0.080** (0.033)	0.091*** (0.033)	0.109*** (0.028)
Observations	18357	18357	18357	150373
Schools	2066	2066	2066	2567
R <sup>2</sup>	0.182	0.176	0.124	0.277
<i>Panel C. Reading Comprehension Test Scores 2016</i>				
Treatment	0.072** (0.030)	0.039 (0.031)	0.060* (0.031)	0.073** (0.029)
Observations	22309	22318	22309	132027
Schools	2566	2567	2566	2567
R <sup>2</sup>	0.162	0.152	0.118	0.31
<i>Panel D. Reading Comprehension Test Scores 2018</i>				
Treatment	0.099*** (0.031)	0.068** (0.032)	0.082*** (0.032)	0.090*** (0.026)
Observations	18371	18371	18371	150397
Schools	2066	2066	2066	2567
R <sup>2</sup>	0.168	0.162	0.117	0.297
Fixed Effects Controls	School District Yes	School District No	Region Yes	School No

Note: This table shows the average treatment effect of the coaching program on standardized student test scores for a number of different specifications for robustness. Panels A and C show the effect after one year of treatment in 2016, while Panels B and D show the effects after three years of treatment in 2018. Column 1 shows the preferred specification that we showed in the main text which include school district fixed effects and two unbalanced controls (number of students and teachers). Column 2 removes the controls, and Column 3 includes region fixed effects instead of school fixed effects. Finally, Column 4 takes advantage of the panel data from 2010-2018 in order to include school fixed effects, without any additional controls. All results are over standardized exam scores and can be interpreted as standard deviations. Regressions are run at the student level, with robust standard errors clustered by school presented in parentheses. \*\*\* p-value <0.01, \*\* p-value <0.05, \* p-value <0.1.

**Table A4**  
**Selective Attrition Test: Regression of Treatment Status on Pre-treatment Characteristics**

	(1) All teachers	(2) Observed teachers at the end of 2017	(3) Non-observed teachers at the end of 2017
Age	0.001 (0.003)	0.003 (0.003)	-0.000 (0.005)
Male	-0.001 (0.038)	0.022 (0.051)	-0.107 (0.090)
<i>Education</i>			
Higher education	-0.059 (0.183)	0.069 (0.095)	-0.302 (0.353)
Postgraduate	-0.102 (0.194)	0.027 (0.127)	---
<i>Teacher career</i>			
Contract	0.105 (0.066)	0.059 (0.082)	0.158 (0.130)
2nd scale	-0.017 (0.059)	0.024 (0.075)	-0.066 (0.117)
3rd scale	0.110 (0.072)	0.088 (0.092)	0.180 (0.141)
4th scale	-0.015 (0.091)	0.048 (0.123)	-0.213 (0.258)
5th scale	-0.037 (0.159)	0.040 (0.197)	---
Joint Significance Test - pvalue	0.507	0.942	0.136
N	646	444	202
R-squared	0.212	0.244	0.440

Note: All regressions include UGEL fixed effects. Standard errors clustered at the school level in parentheses.  
\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

**Table A5: Move Decisions of Sample 1 Teachers after One Year**

	Move decision (observed)	Likers	Movers	Remainers	Dislikers	Row Sum (observed)	Equation
Assigned to APM school	move to APM school	$p^L(1-\sigma)$	$p^M\mu$	0	0	0.1207	(A1)
	move to non-APM school	0	$p^M(1-\mu)$	0	$p^D$	0.2470	(A2)
	stayed in same school	$p^L\sigma$	0	$p^R$	0	0.6323	(A3)
Assigned to non-APM school	move to APM school	$p^L$	$p^M\mu$	0	0	0.0965	(A4)
	move to non-APM school	0	$p^M(1-\mu)$	0	$p^D(1-v)$	0.2793	(A5)
	stayed in same school	0	0	$p^R$	$p^Dv$	0.6242	(A6)

The fraction of likers who stay in the same school is  $\sigma$ , and the fraction of dislikers who stay in the same school is  $v$ . See Appendix C for how these relationships are derived.

**Table A6 Pedagogical Skills Heterogeneous Treatment Effects (Sample 1)**

	(1)	(2)	(3)	(4)	(5)
Treatment	0.314*** (0.102)	0.213 (0.240)	0.314*** (0.117)	0.273* (0.159)	0.236 (0.147)
Experience	0.000 (0.009)	-0.002 (0.011)	0.000 (0.009)	0.000 (0.009)	0.000 (0.009)
Contract Teacher	0.152 (0.162)	0.153 (0.163)	0.154 (0.226)	0.154 (0.163)	0.140 (0.162)
Magisterial Level	0.114** (0.046)	0.115** (0.046)	0.114** (0.046)	0.102* (0.060)	0.114** (0.046)
Sex (Men=1)	-0.313*** (0.099)	-0.315*** (0.099)	-0.313*** (0.099)	-0.313*** (0.099)	-0.396*** (0.147)
Age	-0.029*** (0.009)	-0.029*** (0.009)	-0.029*** (0.009)	-0.029*** (0.009)	-0.029*** (0.009)
Treatment #Experience		0.005 (0.011)			
Treatment #Contract			-0.004 (0.247)		
Treatment #M. Level				0.025 (0.081)	
Treatment #Sex					0.170 (0.188)
$R^2$	0.37	0.37	0.37	0.37	0.37
$N$	455	455	455	455	455

\*  $p < 0.1$ ; \*\*  $p < 0.05$ ; \*\*\*  $p < 0.01$

Note: All regressions include UGEL fixed effects. Standard errors clustered at the school level are reported in parenthesis.



**Table A7**  
**Selective Attrition Test: Regression of Observed Indicator on Treatment Status,  
Pre-treatment Characteristics, and Interactions of Both Variables**

	Observed in Year 2 (Yes=1)	
Treatment	-0.706*	(0.408)
Treatment x Age	0.007	(0.005)
Treatment x Male	-0.040	(0.082)
Treatment x Higher Education	0.403	(0.338)
Treatment x Postgraduate	0.271	(0.350)
Treatment x Contract	-0.029	(0.115)
Treatment x 2nd scale	0.075	(0.100)
Treatment x 3rd scale	-0.054	(0.121)
Treatment x 4th scale	-0.007	(0.189)
Treatment x 5th scale	0.318	(0.247)
Age	-0.003	(0.003)
Male	0.008	(0.058)
Higher education	0.049	(0.295)
Postgraduate	0.344	(0.303)
Contract Teacher	-0.153*	(0.086)
2nd scale	-0.068	(0.070)
3rd scale	-0.023	(0.082)
4th scale	-0.077	(0.154)
5th scale	0.016	(0.121)
Joint Significance Test: p-value (treatment and interactions between treatment and pre-treatment characteristics)		0.373
N	646	
R-squared	0.255	

Note: The regression includes UGEL fixed effects. Standard errors clustered at the school level in parentheses.  
\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

**Table A8: Teacher Skills Estimates using Regional Fixed Effects**

	Sample 1		Sample 2
	(1)	(2)	(3)
Treatment	0.182* (0.096)	0.204** (0.094)	0.218** (0.110)
Experience		-0.000 (0.008)	
Contract Teacher		0.095 (0.142)	
Teacher Career Level		0.087** (0.042)	
Sex (men=1)		-0.291*** (0.090)	
Age		-0.025*** (0.009)	
$R^2$	0.15	0.22	0.18
$N$	455	455	347

\*  $p < 0.1$ ; \*\*  $p < 0.05$ ; \*\*\*  $p < 0.01$

All regressions include regional fixed effects. Standard errors clustered at the school level are in parentheses.

**Table A9: Regression of 2016 Students' ECE Scores on Teacher Pedagogical Practices Index (2016)**

	ECE Scores			
	Math		Reading	
	(1)	(2)	(3)	(4)
Teacher Skill Index	0.042*** (0.015)	0.060 (0.042)	0.043*** (0.015)	0.072* (0.042)
$R^2$	0.19	0.19	0.24	0.24
$N$	1,487	298	1,487	298
Grades	All	2nd	All	2nd

\*  $p < 0.1$ ; \*\*  $p < 0.05$ ; \*\*\*  $p < 0.01$

Note: These regressions measure the correlation between the teacher skill index and students' test scores in the schools those teachers are. The measurement of teachers' practices uses the same rubrics and methodology as in our sample, corresponds to a random sample of primary schools not related to APM or its randomized expansion. Columns (1) and (3) include all observed teachers in the school, while columns (2) and (4) only uses information on teachers in second grade (which is the grade evaluated in the ECE). These observational results support the idea that the skills measured are correlated with student success in standardized tests. The point estimates are higher when sample is limited to teachers in the corresponding grade, but precision is lost due to a much smaller sample. All regressions include regional fixed effects. Huber-White heteroskedasticity robust standard errors in parentheses.

**Table A10: Differences in Teacher Test Scores across Movement Status**

	Teachers who move into randomized sample					Teachers who move out of randomized sample				
	Into Control Schools		Into Treatment Schools		Pairwise t-test	Out of Control Schools		Out of Treatment Schools		Pairwise t-test
	N	Mean/ (SD)	N	Mean. (SD)	Coef/ (P-value)	N	Mean/ (SD)	N	Mean/ (SD)	Coef/ (P-value)
Standardized teacher exam score in 2014 or 2015		(1)		(2)	(1)-(2)		(3)		(4)	(3)-(4)
<i>Year 2016</i>	1099	-0.130 (0.936)	1523	-0.112 (0.929)	-0.018 (0.635)	1630	0.014 (0.948)	2351	0.046 (0.992)	-0.032 (0.311)
<i>Year 2017</i>	1332	-0.313 (0.897)	1823	-0.298 (0.896)	-0.015 (0.636)	1413	-0.339 (0.899)	2009	-0.314 (0.865)	-0.025 (0.412)
<i>Year 2018</i>	1019	-0.323 (0.942)	1395	-0.355 (0.920)	0.032 (0.396)	1570	-0.191 (0.934)	2213	-0.155 (0.948)	-0.036 (0.249)

Note: This table shows the difference in the test scores obtained by teachers in exams by whether they move into or out of the 6,218 randomized sample schools. Teacher exam scores come from tests designed to measure teachers' pedagogical skills and content knowledge taken in order to get into the civil service track or to get a promotion within it. These test were all taken in 2014 and 2015 prior to treatment and scores are standardized over the entire sample of teachers with test scores. Columns (1) and (2) show the average score of teachers who moved into schools in the randomized sample, regardless of where they were coming from, while columns (3) and (4) show the average score for those who moved out of the randomized sample school, regardless of where they moved to. The columns show the mean and standard deviation for the test scores of those teachers, and the following column shows the coefficient and the p-value on the pairwise t-test of the difference. For columns (1) and (2), moving in year 2016 means that teachers moved from 2015 to 2016 (and therefore were in an APM school in 2016) while in columns (3) and (4) show teachers who were in APM schools in that year and moved at the end of the year. Therefore, in both cases, teachers were in the APM schools in the year shown. \* p<0.01, \*\* p<0.05 \*\*\* p<0.01

## Appendix B: Definitions and Derivations for All the Treatment Effects for APM Program

### I. The Basic Set-Up

The test score of student  $i$  at the end of year  $t$  ( $s_i^t$ ) is determined by his or her test score at the end of the previous year ( $s_i^{t-1}$ ) and the skills of the teacher that this student had in the current year ( $y_j^t$ ), where  $j$  denotes the particular teacher that student  $i$  had in year  $t$ :

$$s_i^t = \sigma s_i^{t-1} + \pi y_j^t \quad (\text{B.1})$$

where  $\sigma$  is the impact of the previous year's skills and  $\pi$  is the impact of teacher skill (which, for simplicity, is assumed to be unidimensional).

Schools are randomly assigned in year 1 to be either APM schools ( $R = 1$ ) or non-APM schools ( $R = 0$ ). This school assignment does not change over time, and  $R$  always refers to assignment of schools.

The skill of teacher  $j$  at the end of year  $t$ ,  $y_j^t$ , is assumed to be a linear function of his or her skills in the previous year ( $y_j^{t-1}$ ), the skill gained from one more year of experience ( $\lambda_j$ ), and whether he or she is treated in the current year ( $T_j^t$ ). The treatment impact,  $\delta^k$ , can vary by the type of teacher (remainder (R), liker (L), dislike (D) and mover (M)). Depreciation of teaching skills that are unrelated to the coaching program can be included in  $\lambda_j$ . Equation (B.2) provides the general expression of  $y_j^t$  for year  $t$ :

$$y_j^t = y_j^{t-1} + \lambda_j + \delta^k T_j^t, \quad \text{for } k = R, L, D, M \quad (\text{B.2})$$

Applying (B.2) to year 1 yields the following expression for the skills of teacher  $j$  in that year:

$$\begin{aligned} y_j^1 &= y_j^0 + \lambda_j + \delta^k T_j^1, \quad \text{for } k = R, L, D, M \quad (\text{B.3}) \\ &= \theta_j^1 + \delta^k T_j^1 \end{aligned}$$

where  $\theta_j^1$  is convenient notation for  $y_j^0 + \lambda_j$ .<sup>31</sup>

For year 2, we need to allow for interaction effects of treatments in different years; for example, the impact of a second year of exposure to the program on teachers' skills could be smaller than the impact for the first year of exposure. The equation for  $y_j^2$  is:

$$\begin{aligned} y_j^2 &= y_j^1 + \lambda_j + \delta^k T_j^2 + \gamma_{1,2}^k T_j^1 T_j^2, \quad \text{for } k = R, L, D, M \quad (\text{B.4}) \\ &= \theta_j^1 + \delta^k T_j^1 + \lambda_j + \delta^k T_j^2 + \gamma_{1,2}^k T_j^1 T_j^2 \\ &= \theta_j^2 + \delta^k (T_j^1 + T_j^2) + \gamma_{1,2}^k T_j^1 T_j^2 \end{aligned}$$

---

<sup>31</sup> Note that "R" is used in two different ways, to indicate randomization of a school and to denote "remainder" teachers. When "R" is in normal size text, it indicates the former, and when it is a superscript it indicates the latter.

where  $\theta_j^2$  is convenient notation for  $\theta_j^1 + \lambda_j = y_j^0 + 2\lambda_j$ . If the impact of a second year of exposure to the program is smaller than that of the first year of exposure, then  $\gamma_{1,2}^k$  would be  $< 0$ . Note also that  $\gamma_{1,2}^k$  can include depreciation of teacher skills produced by the program.

For year 3, we need to allow for further interaction effects. The equation for  $y_j^3$  is:

$$\begin{aligned} y_j^3 &= y_j^2 + \lambda_j + \delta^k T_j^3 + \gamma_{1,2}^k (T_j^1 T_j^3 + T_j^2 T_j^3) + \gamma_{1,2,3}^k T_j^1 T_j^2 T_j^3, \quad \text{for } k = R, L, D, M \quad (\text{B.5}) \\ &= \theta_j^2 + \delta^k (T_j^1 + T_j^2) + \gamma_{1,2}^k T_j^1 T_j^2 + \lambda_j + \delta^k T_j^3 + \gamma_{1,2}^k (T_j^1 T_j^3 + T_j^2 T_j^3) + \gamma_{1,2,3}^k T_j^1 T_j^2 T_j^3 \\ &= \theta_j^3 + \delta^k (T_j^1 + T_j^2 + T_j^3) + \gamma_{1,2}^k (T_j^1 T_j^2 + T_j^1 T_j^3 + T_j^2 T_j^3) + \gamma_{1,2,3}^k T_j^1 T_j^2 T_j^3 \end{aligned}$$

where  $\theta_j^3$  is convenient notation for  $\theta_j^2 + \lambda_j = y_j^0 + 3\lambda_j$ . Note that the interaction effect for any combination of two years of training is assumed to be the same, regardless of whether the two years are year 1 and year 2, or year 1 and year 3, or year 2 and year 3. Allowing for different interaction effects for each possible pair of years would do little beyond complicating the notation.

## II. Definitions of Treatment Effects for Teachers' Skills (denoted by $y$ )

We first define three standard treatment effects for teacher skills, then we explain how to define two of these three definitions when the focus is not on set of teachers but on the teachers in a set of schools, in a context where teachers often move between schools.

### 1. Average Treatment Effect (ATE).

We begin with the average treatment effect (ATE). This is defined as what happens when all teachers are treated, and the counterfactual is that no teachers are treated (the program does not exist). In the notation below, superscripts on  $y$  refer to number of years of the program, subscripts on  $y$  refer to the potential outcome (1 = treated, 0 = not treated), and the  $p^k$  terms ( $k = R, L, D$  and  $M$ ) refer to the proportion of teachers in the population who are of type  $k$ . ATE for teacher skills in year  $t$  is defined as:

$$\text{ATE}_{\text{tchr}}(t) \equiv E[y_1^t - y_0^t] = E[y_1^t] - E[y^t | \text{No program exists}] \quad (\text{B.6})$$

Applying this definition to years 1, 2 and 3 yields the following more specific definitions:

$$\text{ATE}_{\text{tchr}}(1) \equiv E[y_1^1 - y_0^1] = \bar{\delta}, \quad \text{where } \bar{\delta} = \delta^R p^R + \delta^L p^L + \delta^D p^D + \delta^M p^M \quad (\text{B.7})$$

$$\text{ATE}_{\text{tchr}}(2) \equiv E[y_1^2 - y_0^2] = 2\bar{\delta} + \bar{\gamma}_{1,2}, \quad \text{where } \bar{\gamma}_{1,2} = \gamma_{1,2}^R p^R + \gamma_{1,2}^L p^L + \gamma_{1,2}^D p^D + \gamma_{1,2}^M p^M \quad (\text{B.8})$$

$$\text{ATE}_{\text{tchr}}(3) \equiv E[y_1^3 - y_0^3] = 3\bar{\delta} + 3\bar{\gamma}_{1,2} + \bar{\gamma}_{1,2,3}, \quad \text{where } \bar{\gamma}_{1,2,3} = \gamma_{1,2,3}^R p^R + \gamma_{1,2,3}^L p^L + \gamma_{1,2,3}^D p^D + \gamma_{1,2,3}^M p^M \quad (\text{B.9})$$

## 2. Intention to Treat Effect (ITT)

Next, consider the intent to treat (ITT) effect. This is defined as the impact on teacher skills for the teachers who were offered the treatment in year 1, that is, the teachers who were in the treated schools in year 1. The counterfactual is being assigned to a control school in year 1. The general definition for ITT in year  $t$  is:

$$ITT_{\text{tchr}}(t) \equiv E[y^t | R_{\text{tchr, year 1}} = 1] - E[y^t | R_{\text{tchr, year 1}} = 0] \quad (\text{B.10})$$

$R_{\text{tchr, year 1}}$  refers to the teacher's school in year 1, which can differ from his or her school in year  $t$ . Applying this definition to years 1, 2 and 3 yields the following, more specific, definitions:

$$ITT_{\text{tchr}}(1) \equiv E[y^1 | R_{\text{tchr, year 1}} = 1] - E[y^1 | R_{\text{tchr, year 1}} = 0] = \bar{\delta} \quad (\text{B.11})$$

$$\begin{aligned} ITT_{\text{tchr}}(2) &\equiv E[y^2 | R_{\text{tchr, year 1}} = 1] - E[y^2 | R_{\text{tchr, year 1}} = 0] && (\text{B.12}) \\ &= \bar{\theta}^2 + [p^R(2\delta^R + \gamma_{1,2}^R) + p^L(2\delta^L + \gamma_{1,2}^L) + p^D\delta^D + p^M(\tau(2\delta^M + \gamma_{1,2}^M) + (1-\tau)\delta^M)] - [\bar{\theta}^2 + p^L\delta^L + p^M\tau\delta^M] \\ &= p^R(2\delta^R + \gamma_{1,2}^R) + p^L(\delta^L + \gamma_{1,2}^L) + p^D\delta^D + p^M(\delta^M + \tau\gamma_{1,2}^M) \\ &= \bar{\delta} + p^R(\delta^R + \gamma_{1,2}^R) + p^L\gamma_{1,2}^L + p^M\tau\gamma_{1,2}^M \end{aligned}$$

$$ITT_{\text{tchr}}(3) \equiv E[y^3 | R_{\text{tchr, year 1}} = 1] - E[y^3 | R_{\text{tchr, year 1}} = 0] \quad (\text{B.13})$$

$$\begin{aligned} &= \bar{\theta}^3 + p^R(3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) + p^L(3\delta^L + 3\gamma_{1,2}^L + \gamma_{1,2,3}^L) + p^D\delta^D \\ &+ p^M[\tau^2(3\delta^M + 3\gamma_{1,2}^M + \gamma_{1,2,3}^M) + 2\tau(1-\tau)(2\delta^M + \gamma_{1,2}^M) + (1-\tau)^2\delta^M] \\ &- [\bar{\theta}^3 + p^L(2\delta^L + \gamma_{1,2}^L) + p^M[\tau^2(2\delta^M + \gamma_{1,2}^M) + 2\tau(1-\tau)\delta^M]] \\ &= p^R(3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) + p^L(\delta^L + 2\gamma_{1,2}^L + \gamma_{1,2,3}^L) + p^D\delta^D + p^M(\delta^M + 2\tau\gamma_{1,2}^M + \tau^2\gamma_{1,2,3}^M) \\ &= \bar{\delta} + p^R(2\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) + p^L(2\gamma_{1,2}^L + \gamma_{1,2,3}^L) + p^M(2\tau\gamma_{1,2}^M + \tau^2\gamma_{1,2,3}^M) \end{aligned}$$

where  $\bar{\theta}^2 = p^R\bar{\theta}^{2,R} + p^L\bar{\theta}^{2,L} + p^D\bar{\theta}^{2,D} + p^M\bar{\theta}^{2,M}$  and  $\bar{\theta}^3 = p^R\bar{\theta}^{3,R} + p^L\bar{\theta}^{3,L} + p^D\bar{\theta}^{3,D} + p^M\bar{\theta}^{3,M}$ . The intuition behind the  $p^M(\tau(2\delta^M + \gamma_{1,2}^M) + (1-\tau)\delta^M)$  term in the second line of (B.12) is that, of the movers in APM schools in year 2, a proportion  $\tau$  were also in an APM school in year 1, so the combined effect of two years of treatment for them is  $2\delta^M + \gamma_{1,2}^M$ , and a proportion  $1-\tau$  were in non-APM schools in year 1, so the effect of treatment for one year for them is  $\delta^M$ .

## 3. Average Causal Response (ACR)

The third treatment effect for teacher skills is similar to a local average treatment effect (LATE), but it differs from LATE because treatment can vary by 1, 2 or 3 years. Angrist and Imbens

(1995) extended LATE to this case, and they called this treatment effect the average causal response (ACR). The general definition after  $t$  years is:

$$ACR_{tchr}(t) \equiv \sum_{s=1}^t E[y_s^t - y_{s-1}^t | T_1^t \geq s > T_0^t] \frac{\text{Prob}[T_1^t \geq s > T_0^t]}{\sum_{r=1}^t \text{Prob}[T_1^t \geq r > T_0^t]} \quad (\text{B.14})$$

where  $T_0^t$  is the (potential) number of years of training in year  $t$  for a teacher who was assigned to a non-APM school in year 1, and  $T_1^t$  is the (potential) number of years of training in year  $t$  for a teacher assigned to an APM school in year 1.<sup>32</sup> The subscripts on  $y$  indicate the value of  $y$  given a (potential) number of *years* of treatment (which varies from 0 to 3), not the value of  $y$  given “treatment or no treatment” (a binary variable), as was the case for the definition for  $ATE_{tchr}(t)$ .

Applying this general definition to years 1, 2 and 3 yields:

$$\begin{aligned} ACR_{tchr}(1) &\equiv E[y_1^1 - y_0^1 | T_1^1 \geq 1 > T_0^1] \frac{\text{Prob}[T_1^1 \geq 1 > T_0^1]}{\text{Prob}[T_1^1 \geq 1 > T_0^1]} \quad (\text{B.15}) \\ &= E[y_1^1 - y_0^1 | T_1^1 \geq 1 > T_0^1] \\ &= E[y^1 | R = 1] - E[y^1 | R = 0] = \bar{\delta} \end{aligned}$$

$$\begin{aligned} ACR_{tchr}(2) &\equiv E[y_1^2 - y_0^2 | T_1^2 \geq 1 > T_0^2] \frac{\text{Prob}[T_1^2 \geq 1 > T_0^2]}{\text{Prob}[T_1^2 \geq 1 > T_0^2] + \text{Prob}[T_1^2 = 2 > T_0^2]} \quad (\text{B.16}) \\ &+ E[y_2^2 - y_1^2 | T_1^2 = 2 > T_0^2] \frac{\text{Prob}[T_1^2 = 2 > T_0^2]}{\text{Prob}[T_1^2 \geq 1 > T_0^2] + \text{Prob}[T_1^2 = 2 > T_0^2]} \\ &= \frac{\delta^R p^R + \delta^D p^D + \delta^M (1-\tau) p^M}{p^R + p^D + (1-\tau) p^M} \times \frac{p^R + p^D + (1-\tau) p^M}{1 + p^R} \\ &+ \frac{(\delta^R + \gamma_{1,2}^R) p^R + (\delta^L + \gamma_{1,2}^L) p^L + (\delta^M + \gamma_{1,2}^M) \tau p^M}{p^R + p^L + \tau p^M} \times \frac{p^R + p^L + \tau p^M}{1 + p^R} \\ &= [\bar{\delta} + p^R(\delta^R + \gamma_{1,2}^R) + p^L \gamma_{1,2}^L + p^M \tau \gamma_{1,2}^M] / [1 + p^R] = ITT_{tchr}(2) / [1 + p^R] \end{aligned}$$

$$\begin{aligned} ACR_{tchr}(3) &\equiv E[y_1^3 - y_0^3 | T_1^3 \geq 1 > T_0^3] \frac{\text{Prob}[T_1^3 \geq 1 > T_0^3]}{\text{Prob}[T_1^3 \geq 1 > T_0^3] + \text{Prob}[T_1^3 = 2 > T_0^3] + \text{Prob}[T_1^3 = 3 > T_0^3]} \quad (\text{B.17}) \\ &+ E[y_2^3 - y_1^3 | T_1^3 \geq 2 > T_0^3] \frac{\text{Prob}[T_1^3 \geq 2 > T_0^3]}{\text{Prob}[T_1^3 \geq 1 > T_0^3] + \text{Prob}[T_1^3 = 2 > T_0^3] + \text{Prob}[T_1^3 = 3 > T_0^3]} \\ &+ E[y_3^3 - y_2^3 | T_1^3 \geq 3 > T_0^3] \frac{\text{Prob}[T_1^3 \geq 3 > T_0^3]}{\text{Prob}[T_1^3 \geq 1 > T_0^3] + \text{Prob}[T_1^3 = 2 > T_0^3] + \text{Prob}[T_1^3 = 3 > T_0^3]} \end{aligned}$$

---

<sup>32</sup> For the general case, possible values for both  $T_0^t$  and  $T_1^t$  are integers from 0 to  $t$ . However, for the APM program all teachers followed their random assignment in year 1, so possible values for  $T_0^t$  are 0 to  $t-1$ , and for  $T_1^t$  are 1 to  $t$ .

$$\begin{aligned}
&= \frac{\delta^R p^R + \delta^D p^D + \delta^M (1-\tau)^2 p^M}{p^R + p^D + (1-\tau)^2 p^M} \times \frac{p^R + p^D + (1-\tau)^2 p^M}{1 + 2p^R} \\
&+ \frac{(\delta^R + \gamma_{1,2}^R) p^R + (\delta^M + \gamma_{1,2}^M) 2\tau(1-\tau) p^M}{p^R + 2\tau(1-\tau) p^M} \times \frac{p^R + 2\tau(1-\tau) p^M}{1 + 2p^R} \\
&+ \frac{(\delta^R + 2\gamma_{1,2}^R + \gamma_{1,2,3}^R) p^R + (\delta^L + 2\gamma_{1,2}^L + \gamma_{1,2,3}^L) p^L + \tau^2 (\delta^M + 2\gamma_{1,2}^M + \gamma_{1,2,3}^M) p^M}{p^R + p^L + \tau^2 p^M} \times \frac{p^R + p^L + \tau^2 p^M}{1 + 2p^R} \\
&= [\bar{\delta} + p^R(2\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) + p^L(2\gamma_{1,2}^L + \gamma_{1,2,3}^L) + p^M(2\tau\gamma_{1,2}^M + \tau^2\gamma_{1,2,3}^M)]/[1 + 2p^R] \\
&= ITT_{\text{tchr}}(3)/[1 + 2p^R]
\end{aligned}$$

To see the intuition behind  $ACR_{\text{tchr}}(t)$ , consider  $ACR_{\text{tchr}}(2)$  given in equation (B.16). The term  $E[y_1^2 - y_0^2 | T_1^2 \geq 1 > T_0^2]$  is the impact on teacher skills of receiving one year of treatment, relative to having zero years of treatment, as indicated by the subscripts on the  $y$  terms, for teachers who would have had one or two years of treatment in year 2 if assigned to an APM school in year 1 ( $T_1^2 \geq 1$ ), but would not have been treated in year 2 if assigned to a non-APM school in year 1 ( $T_0^2 < 1$ ). Of the four teacher types, this includes all remainders and dislikers, and movers who randomly switched to a non-APM school in year 2 (for whom  $T_0^2 = 0$  and  $T_1^2 = 1$ ). The term  $E[y_2^2 - y_1^2 | T_1^2 = 2 > T_0^2]$  is the impact on teacher skills of receiving a *second* year of the treatment, *relative to having one year of treatment*, as indicated by the subscripts on the  $y$  terms, for teachers who would have had two years of treatment in year 2 if assigned to an APM school in year 1 but only zero or one year of in year 2 if assigned to a non-APM school in year 1. This includes all remainders, all likers, and movers who randomly switched to APM schools in year 2 (for whom  $T_0^2 = 1$  and  $T_1^2 = 2$ ). Turning to the sum of the two probability terms in the denominator,  $\text{Prob}[T_1^2 \geq 1 > T_0^2]$  is the probability that a teacher is a remainder, a disliker, or a mover who randomly switches to a non-APM school in year 2, and  $\text{Prob}[T_1^2 = 2 > T_0^2]$  is the probability that a teacher is a remainder, a liker, or a mover who randomly switches to an APM school in year 2. Their sum is greater than 1; remainders are “counted twice” since they are included in both probabilities. Likers, dislikers, and movers are “counted” only once.

In effect,  $ACR_{\text{tchr}}(2)$  is an average of: a) the (average) impact on teacher skills of going from no treatment to one year of treatment for remainders, dislikers, and those movers who randomly move to a non-APM school in year 2; and b) the (average) impact on those skills of going from one to two years of treatment for remainders, likers, and the movers who randomly move to APM schools in year 2. Thus,  $ACR_{\text{tchr}}(2)$  is the average of the impact on teacher skills for each additional year of treatment brought about by random assignment to an APM school in year 1, with remainders getting “double weight” since that assignment raises their years of treatment by two years, while for all others random assignment increases years of treatment by only one year. Importantly, note that, for any  $t$ ,  $ACR_{\text{tchr}}(t)$  is a *per year* (not a cumulative) impact that averages only over years of treatment induced by random assignment to an APM school in year 1. To obtain a cumulative impact over  $t$  years, multiply  $ACR_{\text{tchr}}(t)$  by  $t$ .

#### 4. Average Treatment Effect (on teacher skills) for teachers in treated schools ( $ATE_{\text{sch}}$ )

The three treatment effects discussed so far, in effect, follow teachers who move to other schools. But many teacher training or coaching programs focus on particular schools, so it is



useful to define treatment effects for the teachers in the schools that are implementing the APM program. As explained above, the number of teaching positions in a given school rarely changes. If the number of teaching positions in all schools is fixed, the proportion of teachers in treated schools in years 2 and 3 who are movers is  $(\mu/\tau)p^M$ , where  $\mu$  is the proportion of all movers who move to an APM school in year 2 or year 3 (which is determined by the application process that also determines the proportions of teachers who are remainers, likers, dislikers and movers), and the proportion of teachers in control schools in years 2 and 3 who are movers is  $[(1-\mu)/(1-\tau)]p^M$ .<sup>33</sup>

There are two possibilities for treatment effects that focus on schools. The first is an average treatment effect (ATE) on teacher skills for those schools, where the counterfactual is no program at all, which we denote as  $ATE_{sch}$ . This is defined as follows for year  $t$ :

$$ATE_{sch}(t) \equiv E[y^t | R = 1] - E[y^t | \text{Program does not exist}] \quad (B.18)$$

Applying this general definition for years 1, 2 and 3 yields:

$$\begin{aligned} ATE_{sch}(1) &\equiv E[y^1 | R = 1] - E[y^1 | \text{Program does not exist}] \\ &= \bar{\theta}^1 + \bar{\delta} - \bar{\theta}^1 = \bar{\delta} \end{aligned}$$

$$ATE_{sch}(2) \equiv E[y^2 | R = 1] - E[y^2 | \text{Program does not exist}] \quad (B.19)$$

$$\begin{aligned} &= (\bar{\theta}^{2,R} + 2\delta^R + \gamma_{1,2}^R)p^R + (\bar{\theta}^{2,L} + 2\delta^L + \gamma_{1,2}^L)p^L + (\bar{\theta}^{2,L} + \delta^L)p^L((1-\tau)/\tau) \\ &+ (\bar{\theta}^{2,M} + 2\delta^M + \gamma_{1,2}^M)\tau(\mu/\tau)p^M + (\bar{\theta}^{2,M} + \delta^M)(1-\tau)(\mu/\tau)p^M - [\bar{\theta}^{2,R}p^R + \bar{\theta}^{2,L}p^L + \bar{\theta}^{2,D}p^D + \bar{\theta}^{2,M}p^M] \\ &= (2\delta^R + \gamma_{1,2}^R)p^R + (\delta^L(1+\tau) + \gamma_{1,2}^L\tau)p^L/\tau + [(1+\tau)\delta^M + \tau\gamma_{1,2}^M](\mu/\tau)p^M \\ &\quad + \bar{\theta}^{2,L}p^L((1-\tau)/\tau) - \bar{\theta}^{2,D}p^D + \bar{\theta}^{2,M}p^M((\mu/\tau) - 1) \end{aligned}$$

$$ATE_{sch}(3) \equiv E[y^3 | R = 1] - E[y^3 | \text{Program does not exist}] \quad (B.20)$$

$$\begin{aligned} &= (\bar{\theta}^{3,R} + 3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R)p^R + (\bar{\theta}^{3,L} + 3\delta^L + 3\gamma_{1,2}^L + \gamma_{1,2,3}^L)p^L + (\bar{\theta}^{3,L} + 2\delta^L + \gamma_{1,2}^L)p^L((1-\tau)/\tau) \\ &+ (\bar{\theta}^{3,M} + 3\delta^M + 3\gamma_{1,2}^M + \gamma_{1,2,3}^M)\tau^2(\mu/\tau)p^M + (\bar{\theta}^{3,M} + 2\delta^M + \gamma_{1,2}^M)2\tau(1-\tau)(\mu/\tau)p^M + (\bar{\theta}^{3,M} + \delta^M)(1-\tau)^2(\mu/\tau)p^M \\ &\quad - [\bar{\theta}^{3,R}p^R + \bar{\theta}^{3,L}p^L + \bar{\theta}^{3,D}p^D + \bar{\theta}^{3,M}p^M] \\ &= (3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R)p^R + (\delta^L(2+\tau) + (1+2\tau)\gamma_{1,2}^L + \tau\gamma_{1,2,3}^L)p^L/\tau + [(1+2\tau)\delta^M + \tau(2+\tau)\gamma_{1,2}^M + \gamma_{1,2,3}^M\tau^2](\mu/\tau)p^M \\ &\quad + \bar{\theta}^{3,L}p^L((1-\tau)/\tau) - \bar{\theta}^{3,D}p^D + \bar{\theta}^{3,M}p^M((\mu/\tau) - 1) \end{aligned}$$

---

<sup>33</sup> This definition of  $\mu$  implies that, among all teachers in APM and non-APM schools, the proportion who are movers in APM schools in year 2 or 3 is  $\mu p^M$ . Focusing on APM schools only, this proportion must be divided by  $\tau$ , yielding  $(\mu/\tau)p^M$ . A similar derivation shows that the proportion of movers in non-APM schools is  $[(1-\mu)/(1-\tau)]p^M$ .

The first line of the final expressions for  $ATE_{sch}(2)$  and  $ATE_{sch}(3)$  are the treatment effect, and the last line is the composition effect.

### 5. Intent to Treat Effect (on teacher skills) for teachers in treated schools ( $ITT_{sch}$ )

The second treatment effect for teacher skills that focuses on schools is an ITT effect; it is similar to  $ATE_{sch}$  except that the counterfactual is the skills of teachers in non-APM schools:

$$ITT_{sch}(t) \equiv E[y^t | R = 1] - E[y^t | R = 0] \quad (B.21)$$

Applying this general definition for years 1, 2 and 3 yields:

$$\begin{aligned} ITT_{sch}(1) &\equiv E[y^1 | R = 1] - E[y^1 | R = 0] & (B.22) \\ &= \bar{\theta}^1 + \bar{\delta} - \bar{\theta}^1 = \bar{\delta} \end{aligned}$$

$$\begin{aligned} ITT_{sch}(2) &\equiv E[y^2 | R = 1] - E[y^2 | R = 0] & (B.23) \\ &= (\bar{\theta}^{2,R} + 2\delta^R + \gamma_{1,2}^R)p^R + (\bar{\theta}^{2,L} + 2\delta^L + \gamma_{1,2}^L)p^L + (\bar{\theta}^{2,L} + \delta^L)p^L((1-\tau)/\tau) \\ &\quad + (\bar{\theta}^{2,M} + 2\delta^M + \gamma_{1,2}^M)\tau(\mu/\tau)p^M + (\bar{\theta}^{2,M} + \delta^M)(1-\tau)(\mu/\tau)p^M \\ &\quad - [\bar{\theta}^{2,R}p^R + (\delta^D\tau + \bar{\theta}^{2,D})p^D/(1-\tau) + (\delta^M\tau + \bar{\theta}^{2,M}p^M)((1-\mu)/(1-\tau))] \\ &= (2\delta^R + \gamma_{1,2}^R)p^R + (\delta^L(1+\tau) + \tau\gamma_{1,2}^L)(p^L/\tau) - \delta^D\tau(p^D/(1-\tau)) + ((\delta^M(1+\tau) + \gamma_{1,2}^M\tau)(\mu/\tau) - \delta^M\tau(1-\mu)/(1-\tau))p^M \\ &\quad + \bar{\theta}^{2,L}(p^L/\tau) - \bar{\theta}^{2,D}(p^D/(1-\tau)) + \bar{\theta}^{2,M}p^M((\mu/\tau) - (1-\mu)/(1-\tau)) \end{aligned}$$

$$\begin{aligned} ITT_{sch}(3) &\equiv E[y^3 | R = 1] - E[y^3 | R = 0] & (B.24) \\ &= (\bar{\theta}^{3,R} + 3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R)p^R + (\bar{\theta}^{3,L} + 3\delta^L + 3\gamma_{1,2}^L + \gamma_{1,2,3}^L)p^L + (\bar{\theta}^{3,L} + 2\delta^L + \gamma_{1,2}^L)p^L((1-\tau)/\tau) \\ &\quad + (\bar{\theta}^{3,M} + 3\delta^M + 3\gamma_{1,2}^M + \gamma_{1,2,3}^M)\tau^2(\mu/\tau)p^M + (\bar{\theta}^{3,M} + 2\delta^M + \gamma_{1,2}^M)2\tau(1-\tau)(\mu/\tau)p^M + (\bar{\theta}^{3,M} + \delta^M)(1-\tau)^2(\mu/\tau)p^M \\ &\quad - [\bar{\theta}^{3,R}p^R + (\delta^D\tau + \bar{\theta}^{3,D})p^D/(1-\tau) + ((2\delta^M + \gamma_{1,2}^M)\tau^2 + 2\tau(1-\tau)\delta^M + \bar{\theta}^{3,M})p^M((1-\mu)/(1-\tau))] \\ &= (3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R)p^R + (\delta^L(2+\tau) + \gamma_{1,2}^L(1+2\tau) + \gamma_{1,2,3}^L\tau)(p^L/\tau) - \delta^D\tau(p^D/(1-\tau)) \\ &\quad + [(\mu/\tau)(\delta^M(1+2\tau) + \gamma_{1,2}^M\tau(2+\tau) + \gamma_{1,2,3}^M\tau^2) - ((1-\mu)/(1-\tau)(\delta^M2\tau + \gamma_{1,2}^M\tau^2))]p^M \\ &\quad + \bar{\theta}^{3,L}(p^L/\tau) - \bar{\theta}^{3,D}(p^D/(1-\tau)) + \bar{\theta}^{3,M}p^M((\mu/\tau) - (1-\mu)/(1-\tau)) \end{aligned}$$

The first line of  $ITT_{sch}(2)$  is (first two lines of  $ITT_{sch}(3)$  are) the (net) treatment effect, and the last line is the composition effect.

### III. Definitions of Treatment Effects for Students' Skills (denoted by s)

Focusing on students' skills is simplified by the fact that students are assumed not to change schools, and that the schools they are in always follow their (the schools') random assignment.

We define three treatment effects for student skills. The first two,  $ATE_{stud}$  and  $ITT_{stud}$ , correspond to the two treatment effects defined for their schools ( $ATE_{sch}$  and  $ITT_{sch}$ ). These treatment effects for years 2 and 3 are complex because there are several possible "histories" for students' teachers in those years. For example, in year 2 a student's teacher in a treated school could be a liker who was in an APM school in years 1 and 2, or a liker who was in a non-APM school in year 1 but in an APM school in year 2. Another example is a student in a treated school in year 3; if he or she was taught by treated teacher in year 1 (who by definition had had one year of APM at that time), then by a teacher in year 2 who had APM in year 2 but not year 1, and by a teacher in year 3 who had APM in years 2 and 3 but not in year 1, he or she has been exposed to four years of teacher treatment, and his or her cumulative gain in learning from exposure to those teachers will be averaged over the four years. The general definition of  $ATE_{stud}$  for year t is:

$$ATE_{stud}(t) \equiv E[s^t | R = 1] - E[s^t | \text{Program does not exist}] \quad (\text{B.25})$$

Applying this to years 1, 2 and 3 yield the specific treatment effects for those years:

$$ATE_{stud}(1) \equiv E[s^1 | R = 1] - E[s^1 | \text{Program does not exist}] \quad (\text{B.26})$$

$$\begin{aligned} &= E[\sigma s^0 + \pi y^1 | R = 1] - E[\sigma s^0 + \pi y^1 | \text{Program does not exist}] \\ &= E[\pi y^1 | R = 1] - E[\pi y^1 | \text{Program does not exist}] \\ &= \pi E[y^1 | R = 1] - \pi E[y^1 | \text{Program does not exist}] \\ &= \pi(\bar{\theta}^1 + \bar{\delta}) - \pi\bar{\theta}^1 \\ &= \pi\bar{\delta}, \text{ where } \bar{\delta} = \delta^R p^R + \delta^L p^L + \delta^D p^D + \delta^M p^M \end{aligned}$$

To obtain  $ATE_{stud}(2)$ , one can use the results for  $ATE_{sch}(2)$  in equation (B.19):

$$\begin{aligned} &ATE_{stud}(2) \equiv E[s^2 | R = 1] - E[s^2 | \text{Program does not exist}] \quad (\text{B.26}) \\ &= E[\sigma s^1 + \pi y^2 | R = 1] - E[\sigma s^1 + \pi y^2 | \text{Program does not exist}] \\ &= \sigma(E[s^1 | R = 1] - E[s^1 | \text{Program does not exist}]) \\ &\quad + \pi(E[y_j^2 | R = 1] - E[y_j^2 | \text{Program does not exist}]) \\ &= \sigma\pi\bar{\delta} + \pi[(2\delta^R + \gamma_{1,2}^R)p^R + (\delta^L(1+\tau) + \gamma_{1,2}^L\tau)p^L/\tau + ((1+\tau)\delta^M + \tau\gamma_{1,2}^M)(\mu/\tau)p^M] \\ &\quad + \pi[\bar{\theta}^{2,L}p^L((1-\tau)/\tau) - \bar{\theta}^{2,D}p^D + \bar{\theta}^{2,M}p^M((\mu/\tau) - 1)] \end{aligned}$$

The first line is the treatment effect and the second line is the composition effect.

Year 3 is slightly more complicated since movers continue to move but likers and dislikers (and remainers) do not move between years 2 and 3. Using the results for  $ATE_{sch}(3)$  from (B.20).

$$\begin{aligned}
ATE_{stud}(3) &\equiv E[s^3 | R = 1] - E[s^3 | \text{Program does not exist}] \quad (B.27) \\
&= E[\sigma s^2 + \pi y^3 | R = 1] - E[\sigma s^2 + \pi y^3 | \text{Program does not exist}] \\
&= \sigma(E[s^2 | R = 1] - E[s^2 | \text{Program does not exist}]) \\
&\quad + \pi(E[y^3 | R = 1] - E[y^3 | \text{Program does not exist}]) \\
&= \sigma^2 \pi \bar{\delta} + \sigma \pi [(2\delta^R + \gamma_{1,2}^R) p^R + (\delta^L(1+\tau) + \gamma_{1,2}^L \tau)(p^L/\tau) + ((1+\tau)\delta^M + \tau \gamma_{1,2}^M)(\mu/\tau) p^M] \\
&\quad + \sigma \pi [\bar{\theta}^{2,L} p^L((1-\tau)/\tau) - \bar{\theta}^{2,D} p^D + \bar{\theta}^{2,M} p^M((\mu/\tau) - 1)] \\
&+ \pi [(3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) p^R + (\delta^L(2+\tau) + (1+2\tau)\gamma_{1,2}^L + \tau \gamma_{1,2,3}^L)(p^L/\tau) + ((1+2\tau)\delta^M + \tau(2+\tau)\gamma_{1,2}^M + \gamma_{1,2,3}^M \tau^2)(\mu/\tau) p^M] \\
&\quad + \pi [\bar{\theta}^{3,L} p^L((1-\tau)/\tau) - \bar{\theta}^{3,D} p^D + \bar{\theta}^{3,M} p^M((\mu/\tau) - 1)] \\
&= \sigma^2 \pi \bar{\delta} + \sigma \pi [(2\delta^R + \gamma_{1,2}^R) p^R + (\delta^L(1+\tau) + \gamma_{1,2}^L \tau)(p^L/\tau) + ((1+\tau)\delta^M + \tau \gamma_{1,2}^M)(\mu/\tau) p^M] \\
&+ \pi [(3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) p^R + (\delta^L(2+\tau) + (1+2\tau)\gamma_{1,2}^L + \tau \gamma_{1,2,3}^L)(p^L/\tau) + ((1+2\tau)\delta^M + \tau(2+\tau)\gamma_{1,2}^M + \gamma_{1,2,3}^M \tau^2)(\mu/\tau) p^M] \\
&\quad + \sigma \pi [\bar{\theta}^{2,L} p^L((1-\tau)/\tau) - \bar{\theta}^{2,D} p^D + \bar{\theta}^{2,M} p^M((\mu/\tau) - 1)] + \pi [\bar{\theta}^{3,L} p^L((1-\tau)/\tau) - \bar{\theta}^{3,D} p^D + \bar{\theta}^{3,M} p^M((\mu/\tau) - 1)]
\end{aligned}$$

The first two lines are the treatment effect and the last line is the composition effect.

Next, turn to  $ITT_{stud}$ . The general definition is:

$$ITT_{stud}(t) \equiv E[s^t | R = 1] - E[s^t | R = 0] \quad (B.28)$$

Applying this to years 1, 2 and 3 yield the specific treatment effects for those years:

$$\begin{aligned}
ITT_{stud}(1) &\equiv E[s^1 | R = 1] - E[s^1 | R = 0] \quad (B.29) \\
&= E[\sigma s^0 + \pi y^1 | R = 1] - E[\sigma s^0 + \pi y^1 | R = 0] \\
&= E[\pi y^1 | R = 1] - E[\pi y^1 | R = 0] \\
&= \pi E[y^1 | R = 1] - \pi E[y^1 | R = 0] \\
&= \pi(\bar{\theta}^1 + \bar{\delta}) - \pi \bar{\theta}^1 = \pi \bar{\delta}
\end{aligned}$$

$$ITT_{\text{stud}}(2) \equiv E[s^2 | R = 1] - E[s^2 | R = 0] \quad (\text{B.30})$$

$$\begin{aligned} &= E[\sigma s^1 + \pi y^2 | R = 1] - E[\sigma s^1 + \pi y^2 | R = 0] \\ &= \sigma(E[s^1 | R = 1] - E[s^1 | R = 0]) + \pi(E[y^2 | R = 1] - E[y^2 | R = 0]) \\ &= \sigma\pi\bar{\delta} + \pi[(2\delta^R + \gamma_{1,2}^R)p^R + (\delta^L(1+\tau) + \tau\gamma_{1,2}^L)(p^L/\tau) - \delta^D\tau(p^D/(1-\tau)) + ((\delta^M(1+\tau) + \gamma_{1,2}^M\tau)(\mu/\tau) - \delta^M\tau(1-\mu)/(1-\tau))p^M] \\ &\quad + \pi[\bar{\theta}^{2,L}(p^L/\tau) - \bar{\theta}^{2,D}(p^D/(1-\tau)) + \bar{\theta}^{2,M}p^M((\mu/\tau) - (1-\mu)/(1-\tau))] \end{aligned}$$

The first line is the (net) treatment effect and the second line is the composition effect.

$$ITT_{\text{stud}}(3) \equiv E[s^3 | R = 1] - E[s^3 | R = 0] \quad (\text{B.31})$$

$$\begin{aligned} &= E[\sigma s_i^2 + \pi y_j^3 | R = 1] - E[\sigma s_i^2 + \pi y_j^3 | R = 0] \\ &= \sigma(E[s^2 | R = 1] - E[s^2 | R = 0]) + \pi(E[y^3 | R = 1] - E[y^3 | R = 0]) \\ &= \sigma^2\pi\bar{\delta} + \sigma\pi[(2\delta^R + \gamma_{1,2}^R)p^R + (\delta^L(1+\tau) + \tau\gamma_{1,2}^L)(p^L/\tau) - \delta^D\tau(p^D/(1-\tau)) + ((\delta^M(1+\tau) + \gamma_{1,2}^M\tau)(\mu/\tau) - \delta^M\tau(1-\mu)/(1-\tau))p^M] \\ &\quad + \sigma\pi[\bar{\theta}^{2,L}(p^L/\tau) - \bar{\theta}^{2,D}(p^D/(1-\tau)) + \bar{\theta}^{2,M}p^M((\mu/\tau) - (1-\mu)/(1-\tau))] \\ &\quad + \pi[(3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R)p^R + (\delta^L(2+\tau) + \gamma_{1,2}^L(1+2\tau) + \gamma_{1,2,3}^L\tau)(p^L/\tau) - \delta^D\tau(p^D/(1-\tau))] \\ &\quad + \pi[(\mu/\tau)(\delta^M(1+2\tau) + \gamma_{1,2}^M\tau(2+\tau) + \gamma_{1,2,3}^M\tau^2) - ((1-\mu)/(1-\tau)(\delta^M2\tau + \gamma_{1,2}^M\tau^2))]p^M \\ &\quad + \pi[\bar{\theta}^{3,L}(p^L/\tau) - \bar{\theta}^{3,D}(p^D/(1-\tau)) + \bar{\theta}^{3,M}p^M((\mu/\tau) - (1-\mu)/(1-\tau))] \\ &= \sigma^2\pi\bar{\delta} + \sigma\pi[(2\delta^R + \gamma_{1,2}^R)p^R + (\delta^L(1+\tau) + \tau\gamma_{1,2}^L)(p^L/\tau) - \delta^D\tau(p^D/(1-\tau)) + (\delta^M((1+\tau) + \gamma_{1,2}^M\tau)(\mu/\tau) - \delta^M\tau(1-\mu)/(1-\tau))p^M] \\ &\quad + \pi[(3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R)p^R + (\delta^L(2+\tau) + \gamma_{1,2}^L(1+2\tau) + \gamma_{1,2,3}^L\tau)(p^L/\tau) - \delta^D\tau(p^D/(1-\tau))] \\ &\quad + \pi[(\mu/\tau)(\delta^M(1+2\tau) + \gamma_{1,2}^M\tau(2+\tau) + \gamma_{1,2,3}^M\tau^2) - ((1-\mu)/(1-\tau)(\delta^M2\tau + \gamma_{1,2}^M\tau^2))]p^M \\ &\quad + \sigma\pi[\bar{\theta}^{2,L}(p^L/\tau) - \bar{\theta}^{2,D}(p^D/(1-\tau)) + \bar{\theta}^{2,M}p^M((\mu/\tau) - (1-\mu)/(1-\tau))] \\ &\quad + \pi[\bar{\theta}^{3,L}(p^L/\tau) - \bar{\theta}^{3,D}(p^D/(1-\tau)) + \bar{\theta}^{3,M}p^M((\mu/\tau) - (1-\mu)/(1-\tau))] \end{aligned}$$

The first three lines are the (net) treatment effect and the last two lines are the composition effect.

The third treatment effect for students is the (average) impact of an additional year of teacher training on student learning, averaged over all additional years of that training that a student experiences. In effect, this is a transfer of the  $ACR_{\text{tchr}}$  treatment effects on teacher skill onto student learning, which is complicated by the many different “histories” a student can have

in terms of treated teachers in years 2 and 3. We call these treatment effects  $ACR_{stud}$  effects, though they differ from  $ACR_{tchr}$  (and so differ from the ACR effects of Angrist and Imbens, 1995) since students are not *directly* treated but instead are *indirectly* treated by exposure to treated teachers.

The general definition of  $ACR_{students}$  in year  $t$  (1, 2 or 3) is:

$$ACR_{stud}(t) \equiv \frac{E[s^t|R=1] - E[s^t|R=0]}{E[h_{tchr}(t)|R=1] - E[h_{tchr}(t)|R=0]} \quad (B.32)$$

where  $h_{tchr}(t)$  is the cumulative “history” from year 1 to year  $t$  of a student’s exposure to teachers with APM coaching. For example, a student in a treated school in year 2 had a coached teacher in year 1, but in year 2 the teacher could have one or two years of coaching (i.e. one if the teacher was in a non-APM school in year 1), so the student could have  $h_{tchr}(2)$  of either 2 or 3. The expected value of  $h_{tchr}(t)$  averages over the types of teachers in the school from year 1 to year  $t$ .

For year 1,  $ACR_{stud}(1) = ATT_{stud}(t) = ITT_{stud}(t)$  since all teachers follow their random assignment in year 1, so:

$$ACR_{stud}(1) = \pi\bar{\delta} \quad (B.33)$$

For year 2, the definition in (B.32) gives (using the derivations in (B.30)):

$$\begin{aligned} ACR_{stud}(2) &\equiv \frac{E[s^2|R=1] - E[s^2|R=0]}{E[h_{tchr}(2)|R=1] - E[h_{tchr}(2)|R=0]} \quad (B.34) \\ &= \frac{\sigma\pi\bar{\delta} + \pi[(2\delta^R + \gamma_{1,2}^R)p^R + (\delta^L(1+\tau) + \tau\gamma_{1,2}^L)(p^L/\tau) - \delta^D\tau(p^D/(1-\tau)) + ((\delta^M(1+\tau) + \gamma_{1,2}^M\tau)(\mu/\tau) - \delta^M\tau(1-\mu)/(1-\tau))p^M]}{1 + 2p^R + (p^L/\tau)(\tau+1) + p^M(\mu/\tau)(1+\tau) - [\tau p^D/(1-\tau) + \tau p^M((1-\mu)/(1-\tau))]} \\ &\quad + \frac{\pi[\bar{\theta}^{2,L}(p^L/\tau) - \bar{\theta}^{2,D}(p^D/(1-\tau)) + \bar{\theta}^{2,M}p^M((\mu/\tau) - (1-\mu)/(1-\tau))]}{1 + 2p^R + (p^L/\tau)(\tau+1) + p^M(\mu/\tau)(1+\tau) - [\tau p^D/(1-\tau) + \tau p^M((1-\mu)/(1-\tau))]} \end{aligned}$$

For year 3, applying the definition in (B.32) yields (using the derivations in (B.31)):

$$\begin{aligned} ACR_{stud}(3) &\equiv \frac{E[s^3|R=1] - E[s^3|R=0]}{E[h_{tchr}(3)|R=1] - E[h_{tchr}(3)|R=0]} \quad (B.35) \\ &= \pi \frac{\sigma^2\bar{\delta} + ((3+2\sigma)\delta^R + (3+\sigma)\gamma_{1,2}^R + \gamma_{1,2,3}^R)p^R + (\delta^L(\sigma(1+\tau) + 2+\tau) + \gamma_{1,2}^L(\sigma\tau + 2\tau + 1) + \tau\gamma_{1,2,3}^L)(p^L/\tau) + (\delta^M(\sigma(1+\tau) + 2\tau + 1) + \gamma_{1,2}^M\tau(\sigma + 2 + \tau) + \gamma_{1,2,3}^M\tau^2)p^M(\mu/\tau)}{[1+5p^R + (3+2\tau)p^L/\tau + (2+3\tau)p^M(\mu/\tau)] - [2\tau p^D/(1-\tau) + 3\tau p^M((1-\mu)/(1-\tau))]} \\ &\quad - \pi \frac{\delta^D p^D (\tau/(1-\tau))(\sigma+1) + (\tau(\sigma+2)\delta^M + \tau^2\gamma_{1,2}^M)p^M((1-\mu)/(1-\tau))}{[1+5p^R + (3+2\tau)p^L/\tau + (2+3\tau)p^M(\mu/\tau)] - [2\tau p^D/(1-\tau) + 3\tau p^M((1-\mu)/(1-\tau))]} \\ &\quad + \pi \frac{[(\sigma\bar{\theta}^{2,L} + \bar{\theta}^{3,L})(p^L/\tau) + (\sigma\bar{\theta}^{2,M} + \bar{\theta}^{3,M})p^M(\mu/\tau)] - [(\sigma\bar{\theta}^{2,D} + \bar{\theta}^{3,D})p^D/(1-\tau) + (\sigma\bar{\theta}^{2,M} + \bar{\theta}^{3,M})p^M((1-\mu)/(1-\tau))]}{[1+5p^R + (3+2\tau)p^L/\tau + (2+3\tau)p^M(\mu/\tau)] - [2\tau p^D/(1-\tau) + 3\tau p^M((1-\mu)/(1-\tau))]} \\ &= \frac{ITT_{stud}(3)}{[1+5p^R + (3+2\tau)p^L/\tau + (2+3\tau)p^M(\mu/\tau)] - [2\tau p^D/(1-\tau) + 3\tau p^M((1-\mu)/(1-\tau))]} \end{aligned}$$

#### IV. How Does These Treatment Effects Simplify when There Are No Likers or Dislikers?

Basically,  $p^L$  and  $p^D$  both = 0, so  $p^R + p^M = 1$ .

##### 1. $ATE_{tchr}$

$$ATE_{tchr}(1) \equiv E[y_1^1 - y_0^1] = \bar{\delta}, \text{ where } \bar{\delta} = \delta^R p^R + \delta^M p^M \quad (B.36)$$

$$ATE_{tchr}(2) \equiv E[y_1^2 - y_0^2] = 2\bar{\delta} + \bar{\gamma}_{1,2}, \text{ where } \bar{\gamma}_{1,2} = \gamma_{1,2}^R p^R + \gamma_{1,2}^M p^M \quad (B.37)$$

$$ATE_{tchr}(3) \equiv E[y_1^3 - y_0^3] = 3\bar{\delta} + 3\bar{\gamma}_{1,2} + \bar{\gamma}_{1,2,3}, \text{ where } \bar{\gamma}_{1,2,3} = \gamma_{1,2,3}^R p^R + \gamma_{1,2,3}^M p^M \quad (B.38)$$

##### 2. $ITT_{tchr}$

$$ITT_{tchr}(1) \equiv E[y^1 | R_{tchr, year 1} = 1] - E[y^1 | R_{tchr, year 1} = 0] = \bar{\delta}, \text{ where } \bar{\delta} = \delta^R p^R + \delta^M p^M \quad (B.39)$$

$$ITT_{tchr}(2) \equiv E[y^2 | R_{tchr, year 1} = 1] - E[y^2 | R_{tchr, year 1} = 0] = \bar{\delta} + p^R(\delta^R + \gamma_{1,2}^R) + p^M \tau \gamma_{1,2}^M \quad (B.40)$$

$$ITT_{tchr}(3) \equiv E[y^3 | R_{tchr, year 1} = 1] - E[y^3 | R_{tchr, year 1} = 0] = \bar{\delta} + p^R(2\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) + p^M(2\tau \gamma_{1,2}^M + \tau^2 \gamma_{1,2,3}^M) \quad (B.41)$$

##### 3. Average Causal Response (ACR)/Local Average Treatment Effect (LATE)

$$ACR_{tchr}(1) \equiv E[y_1^1 - y_0^1 | T_1^1 \geq 1 > T_0^1] \frac{\text{Prob}[T_1^1 \geq 1 > T_0^1]}{\text{Prob}[T_1^1 \geq 1 > T_0^1]} = \bar{\delta}, \text{ where } \bar{\delta} = \delta^R p^R + \delta^M p^M \quad (B.42)$$

$$ACR_{tchr}(2) \equiv E[y_1^2 - y_0^2 | T_1^2 \geq 1 > T_0^2] \frac{\text{Prob}[T_1^2 \geq 1 > T_0^2]}{\text{Prob}[T_1^2 \geq 1 > T_0^2] + \text{Prob}[T_1^2 = 2 > T_0^2]} \quad (B.43)$$

$$+ E[y_1^2 - y_0^2 | T_1^2 = 2 > T_0^2] \frac{\text{Prob}[T_1^2 = 2 > T_0^2]}{\text{Prob}[T_1^2 \geq 1 > T_0^2] + \text{Prob}[T_1^2 = 2 > T_0^2]}$$

$$= [\bar{\delta} + p^R(\delta^R + \gamma_{1,2}^R) + p^M \tau \gamma_{1,2}^M] / (1 + p^R) = ITT_{tchr}(2) / (1 + p^R)$$

$$ACR_{teachers, 3 \text{ years}} \equiv E[y_1^3 - y_0^3 | T_1^3 \geq 1 > T_0^3] \frac{\text{Prob}[T_1^3 \geq 1 > T_0^3]}{\text{Prob}[T_1^3 \geq 1 > T_0^3] + \text{Prob}[T_1^3 = 2 > T_0^3] + \text{Prob}[T_1^3 = 3 > T_0^3]} \quad (B.44)$$

$$+ E[y_1^3 - y_0^3 | T_1^3 \geq 2 > T_0^3] \frac{\text{Prob}[T_1^3 \geq 2 > T_0^3]}{\text{Prob}[T_1^3 \geq 1 > T_0^3] + \text{Prob}[T_1^3 = 2 > T_0^3] + \text{Prob}[T_1^3 = 3 > T_0^3]}$$

$$+ E[y_1^3 - y_0^3 | T_1^3 \geq 3 > T_0^3] \frac{\text{Prob}[T_1^3 \geq 3 > T_0^3]}{\text{Prob}[T_1^3 \geq 1 > T_0^3] + \text{Prob}[T_1^3 = 2 > T_0^3] + \text{Prob}[T_1^3 = 3 > T_0^3]}$$

$$= [\bar{\delta} + p^R(2\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) + p^M(2\tau \gamma_{1,2}^M + \tau^2 \gamma_{1,2,3}^M)] / (1 + 2p^R) = ITT_{teachers, 3 \text{ years}} / (1 + 2p^R)$$

##### 4. $ATE_{sch}$

$$\begin{aligned} \text{ATE}_{\text{sch}}(1) &\equiv E[y^1 | R = 1] - E[y^1 | \text{Program does not exist}] \quad (\text{B.45}) \\ &= \bar{\delta}, \text{ where } \bar{\delta} = \delta^R p^R + \delta^M p^M \end{aligned}$$

$$\begin{aligned} \text{ATE}_{\text{sch}}(2) &\equiv E[y^2 | R = 1] - E[y^2 | \text{Program does not exist}] \quad (\text{B.46}) \\ &= (2\delta^R + \gamma_{1,2}^R) p^R + ((1+\tau)\delta^M + \tau\gamma_{1,2}^M) p^M \end{aligned}$$

Note that there is no composition effect, which also implies that  $\mu = \tau$ .

$$\begin{aligned} \text{ATE}_{\text{sch}}(3) &\equiv E[y^3 | R = 1] - E[y^3 | \text{Program does not exist}] \quad (\text{B.47}) \\ &= (3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) p^R + [(1+2\tau)\delta^M + \tau(2+\tau)\gamma_{1,2}^M + \gamma_{1,2,3}^M \tau^2] p^M \end{aligned}$$

Note that there is no composition effect, which also implies that  $\mu = \tau$ .

## 5. $\text{ITT}_{\text{sch}}$

$$\text{ITT}_{\text{sch}}(1) \equiv E[y^1 | R = 1] - E[y^1 | R = 0] = \bar{\delta}, \text{ where } \bar{\delta} = \delta^R p^R + \delta^M p^M \quad (\text{B.48})$$

$$\begin{aligned} \text{ITT}_{\text{sch}}(2) &\equiv E[y^2 | R = 1] - E[y^2 | R = 0] \quad (\text{B.49}) \\ &= (2\delta^R + \gamma_{1,2}^R) p^R + (\delta^M + \gamma_{1,2}^M \tau) p^M \end{aligned}$$

Again, there is no composition effect, which again implies that  $\mu = \tau$ .

$$\begin{aligned} \text{ITT}_{\text{sch}}(3) &\equiv E[y^3 | R = 1] - E[y^3 | R = 0] \quad (\text{B.50}) \\ &= (3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) p^R + (\delta^M + 2\tau\gamma_{1,2}^M + \tau^2\gamma_{1,2,3}^M) p^M \end{aligned}$$

Again, there is no composition effect, which again implies that  $\mu = \tau$ .

## 6. $\text{ATE}_{\text{stud}}$

$$\begin{aligned} \text{ATE}_{\text{stud}}(1) &\equiv E[s^1 | R = 1] - E[s^1 | \text{Program does not exist}] \quad (\text{B.51}) \\ &= \pi\bar{\delta}, \text{ where } \bar{\delta} = \delta^R p^R + \delta^M p^M \end{aligned}$$

$$\begin{aligned} \text{ATE}_{\text{stud}}(2) &\equiv E[s^2 | R = 1] - E[s^2 | \text{Program does not exist}] \quad (\text{B.52}) \\ &= \sigma\pi\bar{\delta} + \pi[(2\delta^R + \gamma_{1,2}^R) p^R + ((1+\tau)\delta^M + \tau\gamma_{1,2}^M) p^M] \end{aligned}$$

Again, there is no composition effect, which again implies that  $\mu = \tau$ .

$$\text{ATE}_{\text{stud}}(3) \equiv E[s^3 | R = 1] - E[s^3 | \text{Program does not exist}] \quad (\text{B.53})$$



$$\begin{aligned}
&= \sigma^2 \pi \bar{\delta} + \sigma \pi [(2\delta^R + \gamma_{1,2}^R) p^R + ((1+\tau)\delta^M + \tau \gamma_{1,2}^M) p^M] \\
&+ \pi [(3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) p^R + ((1+2\tau)\delta^M + \tau(2+\tau)\gamma_{1,2}^M + \gamma_{1,2,3}^M \tau^2) p^M]
\end{aligned}$$

Again, there is no composition effect, which again implies that  $\mu = \tau$ .

### 7. ITT<sub>stud</sub>

$$\text{ITT}_{\text{stud}}(1) \equiv E[s^1 | R = 1] - E[s^1 | R = 0] = \pi \bar{\delta}, \text{ where } \bar{\delta} = \delta^R p^R + \delta^M p^M \quad (\text{B.54})$$

$$\begin{aligned}
\text{ITT}_{\text{stud}}(2) &\equiv E[s^2 | R = 1] - E[s^2 | R = 0] \quad (\text{B.55}) \\
&= \sigma \pi \bar{\delta} + \pi [(2\delta^R + \gamma_{1,2}^R) p^R + (\delta^M + \gamma_{1,2}^M \tau) p^M]
\end{aligned}$$

Again, no composition effect, which again implies that  $\mu = \tau$ .

$$\begin{aligned}
\text{ITT}_{\text{stud}}(3) &\equiv E[s^3 | R = 1] - E[s^3 | R = 0] \quad (\text{B.56}) \\
&= \sigma^2 \pi \bar{\delta} + \sigma \pi [(2\delta^R + \gamma_{1,2}^R) p^R + (\delta^M + \gamma_{1,2}^M \tau) p^M] \\
&+ \pi [(3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) p^R + (\delta^M + 2\tau \gamma_{1,2}^M + \gamma_{1,2,3}^M \tau^2) p^M]
\end{aligned}$$

Again, there is no composition effect, which again implies that  $\mu = \tau$ .

### 8. ACR<sub>stud</sub>

$$\text{ACR}_{\text{stud}}(1) \equiv \frac{E[s^1 | R=1] - E[s^1 | R=0]}{E[h_{\text{tchr}}(1) | R=1] - E[h_{\text{tchr}}(1) | R=0]} = \pi \bar{\delta}, \text{ where } \bar{\delta} = \delta^R p^R + \delta^M p^M \quad (\text{B.57})$$

$$\begin{aligned}
\text{ACR}_{\text{stud}}(2) &\equiv \frac{E[s^2 | R=1] - E[s^2 | R=0]}{E[h_{\text{tchr}}(2) | R=1] - E[h_{\text{tchr}}(2) | R=0]} \quad (\text{B.58}) \\
&= \frac{\sigma \pi \bar{\delta} + \pi [(2\delta^R + \gamma_{1,2}^R) p^R + (\delta^M + \gamma_{1,2}^M \tau) p^M]}{1 + 2p^R + p^M} \\
&= \frac{\text{ITT}_{\text{stud}}(2)}{1 + 2p^R + p^M}
\end{aligned}$$

$$\begin{aligned}
\text{ACR}_{\text{stud}}(3) &\equiv \frac{E[s^3 | R=1] - E[s^3 | R=0]}{E[h_{\text{tchr}}(3) | R=1] - E[h_{\text{tchr}}(3) | R=0]} \quad (\text{B.59}) \\
&= \pi \frac{\sigma^2 \bar{\delta} + ((3+2\sigma)\delta^R + (3+\sigma)\gamma_{1,2}^R + \gamma_{1,2,3}^R) p^R + (\delta^M(\sigma+1) + \gamma_{1,2}^M \tau(\sigma+2) + \gamma_{1,2,3}^M \tau^2) p^M}{1 + 5p^R + 2p^M} \\
&= \frac{\text{ITT}_{\text{stud}}(3)}{1 + 5p^R + 2p^M}
\end{aligned}$$

## V. What Do OLS and IV Regressions Estimate?

Now turn to what we are able to estimate, starting with regressions based on the teacher skill data and then turning to student test scores.

### 1. Applied to Teacher Skill Variables

We have two samples of teachers, one that (imperfectly) follows the teachers who were in APM and non-APM schools in year 1 (Sample 1), and one that focuses on the teachers who are in the APM and non-APM schools in any given year (Sample 2). Start with the Sample 1 teachers.

#### *OLS Applied to Sample 1 Teachers*

The OLS estimate for Sample 1 applied to year  $t$ , can be denoted by  $\hat{\beta}_1^y_{OLS, year\ t}$ , where the “1” subscript indicates Sample 1 teachers. Regressing teacher skills on a constant and a variable for assignment to an APM school in year 1 yields (for any value of  $t$ ) the OLS estimate  $\hat{\beta}_1^y_{OLS, year\ t}$ :

$$\hat{\beta}_1^y_{OLS, year\ t} = E[y^t | R_{tchr, year\ 1} = 1] - E[y^t | R_{tchr, year\ 1} = 0] \quad (B.60)$$

Start with year 1 (we do not have the data, but we show for completeness). Using (B.11), and noting that teachers follow their random assignment in year 1, we have:

$$\hat{\beta}_1^y_{OLS, t=1} = E[y^1 | R_{tchr, year\ 1} = 1] - E[y^1 | R_{tchr, year\ 1} = 0] = \bar{\delta} \quad (B.61)$$

So, if we have data on teacher skills at the end of year 1 (which, unfortunately, we do not have), we can estimate all the teacher skill treatment effects that we described above for year 1 ( $ATE_{tchr}(1)$ ,  $ITT_{tchr}(1)$ ,  $ACR_{tchr}(1)$ ,  $ATE_{sch}(1)$ , and  $ITT_{sch}(1)$ ), because these are all equal to  $\bar{\delta}$  due to perfect compliance in year 1.

Next, turn to year 2. The OLS estimate for Sample 1 is:

$$\hat{\beta}_1^y_{OLS, t=2} = E[y^2 | R_{tchr, year\ 1} = 1] - E[y^2 | R_{tchr, year\ 1} = 0] \quad (B.62)$$

Using (B.12) we have:

$$\begin{aligned} \hat{\beta}_1^y_{OLS, t=2} &= E[y^2 | R_{tchr, year\ 1} = 1] - E[y^2 | R_{tchr, year\ 1} = 0] \quad (B.63) \\ &= \bar{\delta} + p^R(\delta^R + \gamma_{1,2}^R) + p^L\gamma_{1,2}^L + p^M\tau\gamma_{1,2}^M \end{aligned}$$

This equals to  $ITT_{tchr}(2)$ , so we can estimate  $ITT_{tchr}(2)$  by applying OLS to the Sample 1 teachers in year 2.

Next, turn to year 3. The OLS estimate for Sample 1 is:

$$\hat{\beta}_1^y_{OLS,t=3} = E[y^3 | R_{tchr, year 1} = 1] - E[y^3 | R_{tchr, year 1} = 0] \quad (B.64)$$

Again, for Sample 1, we are following the same teachers over time, so their proportions do not change. Using (B.13), we have:

$$\begin{aligned} \hat{\beta}_1^y_{OLS,t=3} &= E[y^3 | R_{tchr, year 1} = 1] - E[y^3 | R_{tchr, year 1} = 0] \quad (B.65) \\ &= \bar{\delta} + p^R(2\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) + p^L(2\gamma_{1,2}^L + \gamma_{1,2,3}^L) + p^M(2\tau\gamma_{1,2}^M + \tau^2\gamma_{1,2,3}^M) \end{aligned}$$

This equals to  $ITT_{tchr}(3)$ , so we can estimate  $ITT_{tchr}(3)$  by applying OLS to the Sample 1 teachers in year 3.

### ***OLS Applied to Sample 2 Teachers***

The OLS estimate of the impact of the APM program on the skills of Sample 2 teachers in year  $t$ , which can be denoted as  $\hat{\beta}_2^y_{OLS, year t}$ , is equal to  $E[y^t | R = 1] - E[y^t | R = 0]$ , where  $R$  refers to the random assignment (in year 1) of the school in which the teacher is in year  $t$ . That is, it compares the teachers who are in treated and control schools in year  $t$ , regardless of their random assignment (regardless of the schools in which they were teaching) in year 1.

For year 1, this is the same as OLS applied to Sample 1 teachers, since  $y^t = y^1$  and  $R = R_{tchr, year 1}$ , so there is no need to show this again.

Next, turn to year 2. The OLS estimate for Sample 2 is:

$$\hat{\beta}_2^y_{OLS,t=2} = E[y^2 | R = 1] - E[y^2 | R = 0] \quad (B.66)$$

For Sample 2, the proportions of teachers who are in the APM and non-APM schools will change, so we need to account for that. All likers will move to APM schools and all dislikers will move to non-APM schools. So the proportion of remainder, liker and mover teachers in APM schools will be  $p^R$  (no change),  $p^L/\tau$  and  $p^M(\mu/\tau)$ , where  $\mu$  is the proportion of all movers who end up in APM schools. Similarly, the proportion of remainder, disliker and mover teachers in non-APM schools will be  $p^R$  (no change),  $p^D/(1-\tau)$  and  $p^M((1-\mu)/(1-\tau))$ .

To calculate  $\hat{\beta}_2^y_{OLS,t=2}$ , start with  $E[y^2 | R = 1]$ :

$$\begin{aligned} E[y^2 | R = 1] &= \bar{\theta}^{2,R}p^R + \bar{\theta}^{2,L}(p^L/\tau) + \bar{\theta}^{2,M}p^M(\mu/\tau) \quad (B.67) \\ &+ p^R(2\delta^R + \gamma_{1,2}^R) + (p^L/\tau)(\delta^L(1+\tau) + \tau\gamma_{1,2}^L) + (p^M(\mu/\tau))(\delta^M(1+\tau) + \tau\gamma_{1,2}^M) \end{aligned}$$

Next, work out  $E[y^2 | R = 0]$ :

$$E[y^2 | R = 0] = \bar{\theta}^{2,R}p^R + \bar{\theta}^{2,D}(p^D/(1-\tau) + \bar{\theta}^{2,M}p^M((1-\mu)/(1-\tau))) + (p^D/(1-\tau))\tau\delta^D + p^M((1-\mu)/(1-\tau))\tau\delta^M \quad (B.68)$$

Equations (B.67) and (B.68) can then be used to obtain  $\hat{\beta}_2^y_{OLS,t=2}$ :

$$\hat{\beta}_2^y_{OLS,t=2} = E[y^2 | R = 1] - E[y^2 | R = 0] \quad (B.69)$$

$$\begin{aligned} &= p^R(2\delta^R + \gamma_{1,2}^R) + (p^L/\tau)(\delta^L(1+\tau) + \tau\gamma_{1,2}^L) + (p^M(\mu/\tau))(\delta^M(1+\tau) + \tau\gamma_{1,2}^M) - [(p^D/(1-\tau))\tau\delta^D + p^M((1-\mu)/(1-\tau))\tau\delta^M] \\ &\quad + \bar{\theta}^{2,L}(p^L/\tau) + \bar{\theta}^{2,M}p^M(\mu/\tau) - [\bar{\theta}^{2,D}(p^D/(1-\tau)) + \bar{\theta}^{2,M}p^M((1-\mu)/(1-\tau))] \\ &= p^R(2\delta^R + \gamma_{1,2}^R) + (p^L/\tau)(\delta^L(1+\tau) + \tau\gamma_{1,2}^L) - (p^D/(1-\tau))\tau\delta^D + p^M[\delta^M(\mu-\tau^2)/(\tau-\tau^2) + \mu\gamma_{1,2}^M] \\ &\quad + \bar{\theta}^{2,L}(p^L/\tau) - \bar{\theta}^{2,D}(p^D/(1-\tau)) + \bar{\theta}^{2,M}p^M[(\mu/\tau) - ((1-\mu)/(1-\tau))] \end{aligned}$$

This is equal to  $ITT_{sch}(2)$ , which means that we can estimate  $ITT_{sch}(2)$  by applying OLS to the Sample 2 teachers in year 2.

Finally, turn to year 3. The OLS estimate for Sample 2 teachers in year 3 is:

$$\hat{\beta}_2^y_{OLS,t=3} = E[y^3 | R = 1] - E[y^3 | R = 0] \quad (B.70)$$

This is similar to year 2, except we need to account for the fact that movers can move again between years 2 and 3. However, the proportions of the 4 types of teachers are the same as in year 3, we just have to adjust for 3 types of movers.

To calculate  $\hat{\beta}_2^y_{OLS,t=3}$ , start with  $E[y^3 | R = 1]$ :

$$E[y^3 | R = 1] = \bar{\theta}^{3,R}p^R + \bar{\theta}^{3,L}(p^L/\tau) + \bar{\theta}^{3,M}p^M(\mu/\tau) \quad (B.71)$$

$$+ p^R(3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) + (p^L/\tau)(\delta^L(2+\tau) + (2\tau+1)\gamma_{1,2}^L + \tau\gamma_{1,2,3}^L) + (p^M(\mu/\tau))(\delta^M(1+2\tau) + \tau(2+\tau)\gamma_{1,2}^M + \tau^2\gamma_{1,2,3}^M)$$

Then calculate  $E[y^3 | R = 0]$ :

$$E[y^3 | R = 0] = \bar{\theta}^{2,R}p^R + \bar{\theta}^{2,D}(p^D/(1-\tau)) + \bar{\theta}^{2,M}p^M((1-\mu)/(1-\tau)) \quad (B.72)$$

$$+ (p^D/(1-\tau))\tau\delta^D + p^M((1-\mu)/(1-\tau))(\delta^M2\tau + \tau^2\gamma_{1,2}^M)$$

Equations (B.71) and (B.72) can then be used to obtain  $\hat{\beta}_2^y_{OLS,t=3}$ :

$$\hat{\beta}_2^y_{OLS,t=3} = E[y^3 | R = 1] - E[y^3 | R = 0] \quad (B.73)$$

$$\begin{aligned} &= p^R(3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) + (p^L/\tau)(\delta^L(2+\tau) + (2\tau+1)\gamma_{1,2}^L + \tau\gamma_{1,2,3}^L) + (p^M(\mu/\tau))(\delta^M(1+2\tau) + \tau(2+\tau)\gamma_{1,2}^M + \tau^2\gamma_{1,2,3}^M) \\ &\quad - [(p^D/(1-\tau))\tau\delta^D + p^M((1-\mu)/(1-\tau))(\delta^M2\tau + \tau^2\gamma_{1,2}^M)] \\ &\quad + \bar{\theta}^{3,L}(p^L/\tau) + \bar{\theta}^{3,M}p^M(\mu/\tau) - [\bar{\theta}^{3,D}(p^D/(1-\tau)) + \bar{\theta}^{3,M}p^M((1-\mu)/(1-\tau))] \end{aligned}$$

The first two lines are the (net) treatment effect, and the last line is the composition effect.

This is equal to  $ITT_{sch}(3)$ , which means that we can estimate  $ITT_{sch}(3)$  by applying OLS to the Sample 2 teachers in year 3.

#### *IV Applied to Sample 1 Teachers*

Consider a simple IV estimation. The equation of interest is the impact of years of participation in APM on teacher skills. We can write this equation as follows, where  $T_j^{Tot,t}$  is the number of years that teacher  $j$  has been exposed to the program at time  $t$ :

$$y_j^t = \beta T_j^{Tot,t} + u_j \quad (B.74)$$

The first stage regression is random assignment to an APM school in year 1:

$$T_j^{Tot,t} = \alpha R_{tchr, year 1, j} + v_j \quad (B.75)$$

where  $R_{tchr, year 1, j}$  denotes teacher  $j$ 's random assignment in year 1. Simple IV regression of these equations estimates  $\beta$  as follows:

$$\hat{\beta}_1^y_{IV, year t} = \frac{Cov(y^t, R_{tchr, year 1})}{Cov(T^{Tot,t}, R_{tchr, year 1})} = \frac{E[y^t | R_{tchr, year 1} = 1] - E[y^t | R_{tchr, year 1} = 0]}{E[T^{Tot,t} | R_{tchr, year 1} = 1] - E[T^{Tot,t} | R_{tchr, year 1} = 0]} \quad (B.76)$$

where the second equality follows from the definition of covariance and the fact that  $R$  equals either 0 or 1. Note that  $\hat{\beta}_1^y_{IV, year t}$  estimates a ‘‘per year’’ effect of the treatment; to obtain the cumulative effect over all  $t$  years multiply  $\hat{\beta}_1^y_{IV, year t}$  by  $t$  (or, when running the regression, divide  $T_j^{Tot,t}$  by  $t$ ).

For year 1, applying this is straightforward. Since all teachers follow their random assignment in year 1, the denominator of  $\hat{\beta}_1^y_{IV, t=1}$  is 1. The numerator can be obtained from the derivations for ITT given in Section II, which implies that:

$$\hat{\beta}_1^y_{IV, t=1} = E[y^1 | R_{tchr, year 1} = 1] - E[y^1 | R_{tchr, year 1} = 0] = \bar{\delta} \quad (B.77)$$

This is equal to all the treatment effects defined in Section II, including  $ACR_{tchr}(1)$ .

For year 2, we need to estimate:

$$\hat{\beta}_1^y_{IV, t=2} = \frac{E[y^2 | R_{tchr, year 1} = 1] - E[y^2 | R_{tchr, year 1} = 0]}{E[T^{Tot,2} | R_{tchr, year 1} = 1] - E[T^{Tot,2} | R_{tchr, year 1} = 0]} \quad (B.78)$$

The numerator can be obtained from equation (B.12):

$$E[y^2 | R_{tchr, year 1} = 1] - E[y^2 | R_{tchr, year 1} = 0] = \bar{\delta} + p^R(\delta^R + \gamma_{1,2}^R) + p^L\gamma_{1,2}^L + p^M\tau\gamma_{1,2}^M \quad (B.79)$$

The denominator is straightforward to calculate:

$$\begin{aligned}
& E[T^{\text{Tot},2} | R_{\text{tchr, year 1}} = 1] - E[T^{\text{Tot},2} | R_{\text{tchr, year 1}} = 0] \quad (\text{B.80}) \\
& = (2p^R + 2p^L + p^D + p^M(2\tau + (1-\tau))) - (0p^R + p^L + 0p^D + \tau p^M) \\
& = 2p^R + p^L + p^D + p^M = 1 + p^R
\end{aligned}$$

Thus  $\hat{\beta}_1^y_{IV,t=2} = (\bar{\delta} + p^R(\delta^R + \gamma_{1,2}^R) + p^L\gamma_{1,2}^L + p^M\tau\gamma_{1,2}^M)/(1 + p^R)$ , and so it estimates  $ACR_{\text{tchr}}(2)$ .

For year 3, we need to estimate:

$$\hat{\beta}_1^y_{IV,t=3} = \frac{E[y^3 | R_{\text{tchr, year 1}} = 1] - E[y^3 | R_{\text{tchr, year 1}} = 0]}{E[T^{\text{Tot},3} | R_{\text{tchr, year 1}} = 1] - E[T^{\text{Tot},3} | R_{\text{tchr, year 1}} = 0]} \quad (\text{B.81})$$

The numerator can be obtained from equation (B.13):

$$\begin{aligned}
& E[y^3 | R_{\text{tchr, year 1}} = 1] - E[y^3 | R_{\text{tchr, year 1}} = 0] \quad (\text{B.82}) \\
& = \bar{\delta} + p^R(2\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) + p^L(2\gamma_{1,2}^L + \gamma_{1,2,3}^L) + p^M(2\tau\gamma_{1,2}^M + \tau^2\gamma_{1,2,3}^M)
\end{aligned}$$

The denominator is again straightforward to calculate:

$$\begin{aligned}
& E[T^{\text{Tot},3} | R_{\text{tchr, year 1}} = 1] - E[T^{\text{Tot},3} | R_{\text{tchr, year 1}} = 0] \quad (\text{B.83}) \\
& = (3p^R + 3p^L + p^D + p^M(3\tau^2 + 4\tau(1-\tau) + (1-\tau)^2)) - (0p^R + 2p^L + 0p^D + p^M(2\tau^2 + 2\tau(1-\tau) + 0(1-\tau)^2)) \\
& = 3p^R + p^L + p^D + p^M = 1 + 2p^R
\end{aligned}$$

Thus,  $\hat{\beta}_1^y_{IV,t=3} = (\bar{\delta} + p^R(2\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) + p^L(2\gamma_{1,2}^L + \gamma_{1,2,3}^L) + p^M(2\tau\gamma_{1,2}^M + \tau^2\gamma_{1,2,3}^M))/(1 + 2p^R)$ , so  $\hat{\beta}_1^y_{IV,t=3}$  estimates  $ACR_{\text{tchr}}(3)$ .

#### *IV Applied to Sample 2 Teachers*

If we have an “open” system, it is not possible to apply IV to Sample 2 teachers because some of them will come from outside of our randomization sample and thus the instrumental variable (random assignment) does not exist for some of those teachers. While one could argue that such teachers could be treated as not randomly assigned to the treatment groups, so that  $R^1 = 0$ , I do not think that this is correct, because teachers who want to move into our set of schools are a (self-)selected group of teachers who may, for example, be attracted by the APM program.

However, if our system is “closed”, so that we do have a valid IV for all Sample 2 teachers, it is possible to apply IV estimation to Sample 2. For year 1, as always all teachers follow their random assignment and this could again estimate  $\bar{\delta}$ , so  $\beta_{2,IV,\text{year 1}} = \bar{\delta}$ .

For year 2, we simply use the same approach as for Sample 1 teachers, except that  $R_{\text{chr, year 1}}$  is replaced by  $R$ ; that is, the focus is on the teachers in the APM and non-APM schools in year  $t$ ; not the teachers who were on those schools in year 1. For the teachers in Sample 2, IV estimates:

$$\beta_{2,IV,\text{year 2}} = \frac{E[y^2|R=1] - E[y^2|R=0]}{E[T^{\text{Tot},2}|R=1] - E[T^{\text{Tot},2}|R=0]} \quad (\text{B.84})$$

We simply need to work this out for the Sample 2 teachers. The defining characteristics of those teachers is where they were in year 2. Those in the treated schools in year 2 have  $T^2 = 1$ , and those in the control schools have  $T^2 = 0$ . Thus, the numerator of the above expression was derived in equation (B.69), which we show again:

$$\begin{aligned} & E[y^2 | R = 1] - E[y^2 | R = 0] \quad (\text{B.69}) \\ &= p^R(2\delta^R + \gamma_{1,2}^R) + (p^L/\tau)(\delta^L(1+\tau) + \tau\gamma_{1,2}^L) - (p^D/(1-\tau))\tau\delta^D + p^M[\delta^M(\mu-\tau^2)/(\tau-\tau^2) + \mu\gamma_{1,2}^M] \\ & \quad + \bar{\theta}^{2,L}(p^L/\tau) + \bar{\theta}^{2,M}p^M(\mu/\tau) - [\bar{\theta}^{2,D}(p^D/(1-\tau) + \bar{\theta}^{2,M}p^M((1-\mu)/(1-\tau)))] \end{aligned}$$

Next, consider the denominator for this IV estimate:

$$\begin{aligned} & E[T^{\text{Tot},2} | R = 1] - E[T^{\text{Tot},2} | R = 0] \quad (\text{B.85}) \\ &= p^R2 + (p^L/\tau)[2\tau + (1-\tau)] + (p^M(\mu/\tau))[2\tau + (1-\tau)] - [(p^D/(1-\tau))\tau + p^M((1-\mu)/(1-\tau))\tau] \\ &= 2p^R + (p^L/\tau)(1+\tau) + p^M(\mu-\tau^2)/(\tau-\tau^2) - (p^D/(1-\tau))\tau \end{aligned}$$

Combining the numerator and denominator for the IV estimate for Sample 2, for a closed system, yields:

$$\beta_{2,IV,\text{year 2}} \quad (\text{B.86})$$

$$= \frac{p^R(2\delta^R + \gamma_{1,2}^R) + (p^L/\tau)(\delta^L(1+\tau) + \tau\gamma_{1,2}^L) - (p^D/(1-\tau))\tau\delta^D + p^M[\delta^M(\mu-\tau^2)/(\tau-\tau^2) + \mu\gamma_{1,2}^M] + \bar{\theta}^{2,L}(p^L/\tau) + \bar{\theta}^{2,M}p^M(\mu-\tau)/(\tau(1-\tau)) - \bar{\theta}^{2,D}p^D/(1-\tau)}{2p^R + (p^L/\tau)(1+\tau) + p^M(\mu-\tau^2)/(\tau-\tau^2) - (p^D/(1-\tau))\tau}$$

Unlike OLS estimation for year 2 (for a “closed” or “open” system), IV estimation for Sample 2 teachers in year 2 for a “closed” system is exactly equal to IV estimation for Sample 1 teachers in year 2. The intuition is that, in a closed system, Sample 1 and Sample 2 teachers are the same population of teachers and are equally distributed in APM and non-APM schools when random assignment, the defining feature of IV estimation, was done in year 1. In contrast, OLS estimates for Sample 2 (but not Sample 1) teachers are defined in terms of where teachers were in year 2.

For year 3, the IV estimate for Sample 2 teachers for a closed system is:

$$\beta_{2,IV,\text{year 3}} = \frac{E[y^3|R=1] - E[y^3|R=0]}{E[T^{\text{Tot},3}|R=1] - E[T^{\text{Tot},3}|R=0]} \quad (\text{B.87})$$

The numerator is from equation (B.73), which is:

$$E[y^3 | R = 1] - E[y^3 | R = 0] \quad (\text{B.73})$$

$$\begin{aligned} &= p^R(3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) + (p^L/\tau)(\delta^L(2+\tau) + (2\tau+1)\gamma_{1,2}^L + \tau\gamma_{1,2,3}^L) + (p^M(\mu/\tau))(\delta^M(1+2\tau) + \tau(2+\tau)\gamma_{1,2}^M + \tau^2\gamma_{1,2,3}^M) \\ &\quad - [(p^D/(1-\tau))\tau\delta^D + p^M((1-\mu)/(1-\tau))(\delta^M 2\tau + \tau^2\gamma_{1,2}^M)] \\ &\quad \bar{\theta}^{3,L}(p^L/\tau) + \bar{\theta}^{3,M}p^M(\mu/\tau) - [\bar{\theta}^{3,D}(p^D/(1-\tau) + \bar{\theta}^{3,M}p^M((1-\mu)/(1-\tau)))] \end{aligned}$$

The denominator is:

$$E[T^{\text{Tot},2} | R = 1] - E[T^{\text{Tot},2} | R = 0] \quad (\text{B.88})$$

$$\begin{aligned} &= p^R 3 + (p^L/\tau)(3\tau + 2(1-\tau)) + (p^M(\mu/\tau))[3\tau^2 + 2 \times 2\tau(1-\tau) + (1-\tau)^2] - [(p^D/(1-\tau))\tau + p^M((1-\mu)/(1-\tau))2\tau] \\ &= 3p^R + (p^L/\tau)(2+\tau) + p^M(\mu+\mu\tau-2\tau^2)/(\tau-\tau^2) - (p^D/(1-\tau))\tau \end{aligned}$$

Combining the numerator and the denominator gives, for Sample 2 for a closed system:

$$\beta_{2,IV,\text{year } 3} \quad (\text{B.89})$$

$$\begin{aligned} &= \frac{p^R(3\delta^R+3\gamma_{1,2}^R+\gamma_{1,2,3}^R)+(p^L/\tau)(\delta^L(2+\tau)+(2\tau+1)\gamma_{1,2}^L+\tau\gamma_{1,2,3}^L) - (p^D/(1-\tau))\tau\delta^D + (p^M/(\tau-\tau^2))[\delta^M(\mu(1+\tau)-2\tau^2) + \gamma_{1,2}^M(\mu(2\tau-\tau^2)-\tau^3)]}{3p^R + (p^L/\tau)(2+\tau) + p^M(\mu+\mu\tau-2\tau^2)/(\tau-\tau^2) - (p^D/(1-\tau))\tau} \\ &\quad + \frac{\bar{\theta}^{3,L}(p^L/\tau) + \bar{\theta}^{3,M}p^M(\mu-\tau)/(\tau(1-\tau)) - \bar{\theta}^{3,D}p^D/(1-\tau)}{3p^R + (p^L/\tau)(2+\tau) + p^M(\mu+\mu\tau-2\tau^2)/(\tau-\tau^2) - (p^D/(1-\tau))\tau} \end{aligned}$$

## 2. Applied to Student Test Scores

Since students are assumed not to move between schools, there is only one sample of students, who are classified by the schools in which they are enrolled.

An OLS regression of students' test scores in year  $t$  on a constant term and the type of school (APM or non-APM) that the student is in that year will yield the following coefficient for the type of school, which we can denote by  $\hat{\beta}_{OLS, \text{year } t}^S$ :

$$\hat{\beta}_{OLS, \text{year } t}^S = E[s^t | R = 1] - E[s^t | R = 0] \quad (\text{B.90})$$

Consider  $\hat{\beta}_{OLS, \text{year } t}^S$  separately for years 1, 2 and 3.

Start with year 1. Applying equation (B.1) yields  $s_i^1 = \sigma s_i^0 + \pi y_j$ . Since teachers are unable to move in year 1, no teachers in non-APM schools are trained and all teachers in APM schools are trained. Thus:

$$E[s^1 | R = 1] = \sigma E[s^0 | R = 1] + \pi E[y^1 | R = 1] \quad (\text{B.91})$$



$$= \sigma E[s^0 | R = 1] + \pi[(\bar{\theta}^{1,R}p^R + \bar{\theta}^{1,L}p^L + \bar{\theta}^{1,D}p^D + \bar{\theta}^{1,M}p^M) + (\delta^R p^R + \delta^L p^L + \delta^D p^D + \delta^M p^M)]$$

$$E[s^1 | R = 0] = \sigma E[s^0 | R = 0] + \pi(\bar{\theta}^{1,R}p^R + \bar{\theta}^{1,L}p^L + \bar{\theta}^{1,D}p^D + \bar{\theta}^{1,M}p^M) \quad (B.92)$$

$$\begin{aligned} \hat{\beta}_{OLS, t=1}^s &= E[s^1 | R = 1] - E[s^1 | R = 0] = \pi(\delta^R p^R + \delta^L p^L + \delta^D p^D + \delta^M p^M) \quad (B.93) \\ &= \pi\bar{\delta} \end{aligned}$$

where  $\bar{\delta}$  is the population-weighted average of the four  $\delta$  terms. Note also that  $E[s^0 | R = 1] = E[s^0 | R = 0]$  since  $R$  was randomly assigned.

Thus, for year 1 OLS produces an unbiased estimate of both  $ATE_{stud}(1)$ , which also equals  $ITT_{stud}(1)$  and  $ACR_{stud}(1)$ , which is intuitively plausible since neither teachers nor students move in year 1.

Next, turn to year 2. The OLS estimate of  $\beta_2$ , which we denote by  $\hat{\beta}_{OLS, t=2}^s$ , is derived as follows, using equations (B.90), (B.93) and (B.23).

$$\hat{\beta}_{OLS, t=2}^s = E[s^2 | R = 1] - E[s^2 | R = 0] \quad (B.94)$$

$$= E[\sigma s^1 + \pi y^2 | R = 1] - E[\sigma s^1 + \pi y^2 | R = 0]$$

$$= \sigma \{E[s^1 | R = 1] - E[s^1 | R = 0]\} + \pi \{E[y^2 | R = 1] - E[y^2 | R = 0]\}$$

$$= \sigma\pi\bar{\delta}$$

$$\begin{aligned} &+ \pi[(2\delta^R + \gamma_{1,2}^R)p^R + (\delta^L(1+\tau) + \tau\gamma_{1,2}^L)(p^L/\tau) - \delta^D\tau(p^D/(1-\tau)) + ((\delta^M(1+\tau) + \gamma_{1,2}^M\tau)(\mu/\tau) - \delta^M\tau(1-\mu)/(1-\tau))p^M \\ &+ \bar{\theta}^{2,L}(p^L/\tau) - \bar{\theta}^{2,D}(p^D/(1-\tau)) + \bar{\theta}^{2,M}p^M((\mu/\tau) - (1-\mu)/(1-\tau))] \end{aligned}$$

The first and second lines are the direct effect on students while the third line is the composition effect due to teachers switching schools. This equals  $ITT_{stud}(2)$ .

Finally, turn to year 3. The OLS estimate of  $\beta_3$ , which we denote by  $\hat{\beta}_{OLS, t=3}^s$ , is derived as follows, using equations (B.90), (B.94) and (B.23).

$$\hat{\beta}_{OLS, t=3}^s = E[s^3 | R = 1] - E[s^3 | R = 0] \quad (B.95)$$

$$= E[\sigma s^2 + \pi y^3 | R = 1] - E[\sigma s^2 + \pi y^3 | R = 0]$$

$$= \sigma \{E[s^2 | R = 1] - E[s^2 | R = 0]\} + \pi \{E[y^3 | R = 1] - E[y^3 | R = 0]\}$$

$$= \sigma \{ \sigma\pi\bar{\delta} \}$$

$$\begin{aligned}
& + \pi[(2\delta^R + \gamma_{1,2}^R)p^R + (\delta^L(1+\tau) + \tau\gamma_{1,2}^L)(p^L/\tau) - \delta^D\tau(p^D/(1-\tau)) + ((\delta^M(1+\tau) + \gamma_{1,2}^M\tau)(\mu/\tau) - \delta^M\tau(1-\mu)/(1-\tau))p^M \\
& \quad + \bar{\theta}^{2,L}(p^L/\tau) - \bar{\theta}^{2,D}(p^D/(1-\tau)) + \bar{\theta}^{2,M}p^M((\mu/\tau) - (1-\mu)/(1-\tau))] \} \\
& + \pi\{(3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R)p^R + (\delta^L(2+\tau) + \gamma_{1,2}^L(1+2\tau) + \gamma_{1,2,3}^L\tau)(p^L/\tau) - \delta^D\tau(p^D/(1-\tau)) \\
& \quad + [(\mu/\tau)(\delta^M(1+2\tau) + \gamma_{1,2}^M\tau(2+\tau) + \gamma_{1,2,3}^M\tau^2) - ((1-\mu)/(1-\tau)(\delta^M2\tau + \gamma_{1,2}^M\tau^2))]p^M \\
& \quad + \bar{\theta}^{3,L}(p^L/\tau) - \bar{\theta}^{3,D}(p^D/(1-\tau)) + \bar{\theta}^{3,M}p^M((\mu/\tau) - (1-\mu)/(1-\tau))\}
\end{aligned}$$

The first, second, fourth and fifth lines are the direct effect on students, while the third and sixth lines are the composition effect due to teachers switching schools. This equals  $ITT_{stud}(3)$ .

Finally, consider IV estimates of student test scores. The variable being instrumented is the student's exposure to teachers with APM coaching. More specifically, as explained above, it is the "history" from year 1 to year  $t$  of students' exposure to treated teachers, which is denoted by  $h_{tchr}(t)$ . Random assignment (R) is the instrument. For simple IV estimation with a constant term and no other variables in the first stage and second stage equations, the IV estimate from year  $t$ , denoted by  $\hat{\beta}_{IV,t}^s$ , is:

$$\begin{aligned}
\hat{\beta}_{IV,t}^s & \equiv \frac{\text{Cov}(s^t, R)}{\text{Cov}(h_{tchr}(t), R)} \quad (B.96) \\
& = \frac{E[s^t|R=1] - E[s^t|R=0]}{E[h_{tchr}(t)|R=1] - E[h_{tchr}(t)|R=0]}
\end{aligned}$$

where the second line uses the fact that R is a binary variable.

For year 1, we have:

$$\hat{\beta}_{IV,1}^s \equiv \frac{E[s^1|R=1] - E[s^1|R=0]}{E[h_{tchr}(1)|R=1] - E[h_{tchr}(1)|R=0]} \quad (B.97)$$

All teachers follow their random assignment, so the denominator equals 1. The numerator is given in equation (B.29), which implies, as one would expect from full compliance in year 1:

$$\hat{\beta}_{IV,1}^s = \pi\bar{\delta} \quad (B.98)$$

For year 2, the definition in (B.107) gives:

$$\hat{\beta}_{IV,2}^s = \frac{E[s^2|R=1] - E[s^2|R=0]}{E[h_{tchr}(2)|R=1] - E[h_{tchr}(2)|R=0]} \quad (B.99)$$

The derivations above in equation (B.34) show that:

$$\hat{\beta}_{IV,2}^s = \frac{\sigma\pi\bar{\delta} + \pi[(2\delta^R + \gamma_{1,2}^R)p^R + (\delta^L(1+\tau) + \tau\gamma_{1,2}^L)(p^L/\tau) - \delta^D\tau(p^D/(1-\tau)) + ((\delta^M(1+\tau) + \gamma_{1,2}^M\tau)(\mu/\tau) - \delta^M\tau(1-\mu)/(1-\tau))p^M]}{1 + 2p^R + (p^L/\tau)(\tau+1) + p^M(\mu/\tau)(1+\tau) - [\tau p^D/(1-\tau) + \tau p^M((1-\mu)/(1-\tau))]} \quad (B.100)$$

$$+ \frac{\pi[\bar{\theta}^{2,L}(p^L/\tau) - \bar{\theta}^{2,D}(p^D/(1-\tau)) + \bar{\theta}^{2,M}p^M((\mu/\tau) - (1-\mu)/(1-\tau))]}{1 + 2p^R + (p^L/\tau)(\tau+1) + p^M(\mu/\tau)(1+\tau) - [\tau p^D/(1-\tau) + \tau p^M((1-\mu)/(1-\tau))]}$$

This is equal to  $ACR_{\text{stud}}(2)$ .

For year 3, applying the definition in (B.96) yields:

$$\hat{\beta}_{IV,3}^s = \frac{E[s^3|R=1] - E[s^3|R=0]}{E[h_{\text{tchr}}(3)|R=1] - E[h_{\text{tchr}}(3)|R=0]} \quad (\text{B.101})$$

The derivations in equation (B.35) show that

$$\hat{\beta}_{IV,3}^s = \quad (\text{B.102})$$

$$\begin{aligned} & \pi \frac{\sigma^2 \bar{\delta} + \left( (3+2\sigma)\delta^R + (3+\sigma)\gamma_{1,2}^R + \gamma_{1,2,3}^R \right) p^R + (\delta^L(\sigma(1+\tau) + 2\tau) + \gamma_{1,2}^L(\sigma\tau + 2\tau + 1) + \tau\gamma_{1,2,3}^L)(p^L/\tau) + (\delta^M(\sigma(1+\tau) + 2\tau + 1) + \gamma_{1,2}^M\tau(\sigma + 2\tau) + \gamma_{1,2,3}^M\tau^2)p^M(\mu/\tau)}{[1 + 5p^R + (3+2\tau)p^L/\tau + (2+3\tau)p^M(\mu/\tau)] - [2\tau p^D/(1-\tau) + 3\tau p^M((1-\mu)/(1-\tau))]} \\ & - \pi \frac{\delta^D p^D(\tau/(1-\tau))(\sigma+1) + (\tau(\sigma+2)\delta^M + \tau^2\gamma_{1,2}^M)p^M((1-\mu)/(1-\tau))}{[1 + 5p^R + (3+2\tau)p^L/\tau + (2+3\tau)p^M(\mu/\tau)] - [2\tau p^D/(1-\tau) + 3\tau p^M((1-\mu)/(1-\tau))]} \\ & + \pi \frac{[(\sigma\bar{\theta}^{2,L} + \bar{\theta}^{3,L})(p^L/\tau) + (\sigma\bar{\theta}^{2,M} + \bar{\theta}^{3,M})p^M(\mu/\tau)] - [(\sigma\bar{\theta}^{2,D} + \bar{\theta}^{3,D})p^D/(1-\tau) + (\sigma\bar{\theta}^{2,M} + \bar{\theta}^{3,M})p^M((1-\mu)/(1-\tau))]}{[1 + 5p^R + (3+2\tau)p^L/\tau + (2+3\tau)p^M(\mu/\tau)] - [2\tau p^D/(1-\tau) + 3\tau p^M((1-\mu)/(1-\tau))]} \end{aligned}$$

This is equal to  $ACR_{\text{stud}}(3)$ .

## VI. Bounds for ATE for Years 2 and 3

As shown above, we can estimate  $ITT_{\text{tchr}}(t)$  and  $ACR_{\text{tchr}}(t)$  using Sample 1 teachers,  $ITT_{\text{sch}}(t)$  using Sample 2 teachers, and  $ITT_{\text{stud}}(t)$  and  $ACR_{\text{stud}}(t)$  using student test scores in our school data. In addition, for year 1 we call also estimate  $ATE_{\text{tchr}}(1)$ , which also equals  $ATE_{\text{sch}}(1)$ , and  $ATE_{\text{stud}}(1)$  since in year 1 these are all equal to the corresponding ITT estimands. Unfortunately, we cannot estimate  $ATE_{\text{tchr}}(2)$ ,  $ATE_{\text{tchr}}(3)$ ,  $ATE_{\text{sch}}(2)$ ,  $ATE_{\text{sch}}(3)$ ,  $ATE_{\text{stud}}(2)$  or  $ATE_{\text{stud}}(3)$ . However, under plausible assumptions it is possible to show that ITT estimands are lower bounds on several ATE estimands.

Consider first  $ATE_{\text{tchr}}(2)$  and  $ITT_{\text{tchr}}(2)$ . Their difference is:

$$\begin{aligned} ATE_{\text{tchr}}(2) - ITT_{\text{tchr}}(2) &= 2\bar{\delta} + \bar{\gamma}_{1,2} - [\bar{\delta} + p^R(\delta^R + \gamma_{1,2}^R) + p^L\gamma_{1,2}^L + p^M\tau\gamma_{1,2}^M] \quad (\text{B.103}) \\ &= \delta^R p^R + \delta^L p^L + \delta^D p^D + \delta^M p^M + \gamma_{1,2}^R p^R + \gamma_{1,2}^L p^L + \gamma_{1,2}^D p^D + \gamma_{1,2}^M p^M - [p^R(\delta^R + \gamma_{1,2}^R) + p^L\gamma_{1,2}^L + p^M\tau\gamma_{1,2}^M] \\ &= \delta^L p^L + \delta^D p^D + \delta^M p^M + \gamma_{1,2}^D p^D + (1-\tau)\gamma_{1,2}^M p^M \\ &= \delta^L p^L + (\delta^D + \gamma_{1,2}^D) p^D + (\delta^M + \gamma_{1,2}^M(1-\tau)) p^M \end{aligned}$$

As long as the first year of the program does not have a negative effect on the skills of likers (i.e.  $\delta^L \geq 0$ ) and the second year does not have a negative effect on the skills of dislikers ( $\delta^D + \gamma_{1,2}^D \geq 0$ ) or movers ( $\delta^M + \gamma_{1,2}^M \geq 0$ ),  $ITT_{\text{chr}}(2)$  will be a lower bound for  $ATE_{\text{chr}}(2)$ .

It is less clear that  $ITT_{\text{sch}}(2)$  is a lower bound for  $ATE_{\text{sch}}(2)$ , because for these treatment effects follow schools over time, as opposed to following teachers over time, and the composition of teachers in APM and non-APM schools can change over time. More specifically:

$$ATE_{\text{sch}}(2) - ITT_{\text{sch}}(2) \quad (\text{B.104})$$

$$\begin{aligned} &= (2\delta^R + \gamma_{1,2}^R)p^R + (\delta^L(1+\tau) + \gamma_{1,2}^L\tau)(p^L/\tau) + (\delta^M(1+\tau) + \gamma_{1,2}^M\tau)p^M(\mu/\tau) + \bar{\theta}^{2,L}((1-\tau)/\tau) - \bar{\theta}^{2,D}p^D + \bar{\theta}^{2,M}p^M((\mu/\tau) - 1) \\ &\quad - [(2\delta^R + \gamma_{1,2}^R)p^R + (\delta^L(1+\tau) + \gamma_{1,2}^L\tau)(p^L/\tau) + (\delta^M(1+\tau) + \gamma_{1,2}^M\tau)p^M(\mu/\tau) - (\delta^D p^D(\tau/1-\tau)) + \delta^M \tau p^M((1-\mu)/(1-\tau))] \\ &\quad \quad - [\bar{\theta}^{2,L}(p^L/\tau) + \bar{\theta}^{2,M}p^M(\mu/\tau) - (\bar{\theta}^{2,D}(p^D/(1-\tau)) + \bar{\theta}^{2,M}p^M((1-\mu)/(1-\tau)))] \\ &\quad = \bar{\theta}^{2,L}((1-\tau)/\tau) - \bar{\theta}^{2,D}p^D + \bar{\theta}^{2,M}p^M((\mu/\tau) - 1) + \delta^D p^D(\tau/1-\tau) + \delta^M \tau p^M((1-\mu)/(1-\tau)) \\ &\quad \quad - [\bar{\theta}^{2,L}(p^L/\tau) + \bar{\theta}^{2,M}p^M(\mu/\tau) - (\bar{\theta}^{2,D}(p^D/(1-\tau)) + \bar{\theta}^{2,M}p^M((1-\mu)/(1-\tau)))] \\ &= \delta^D p^D(\tau/1-\tau) + \delta^M \tau p^M((1-\mu)/(1-\tau)) - \bar{\theta}^{2,L}p^L + \bar{\theta}^{2,D}p^D(\tau/(1-\tau)) + \bar{\theta}^{2,M}p^M((\tau-\mu)/(1-\tau)). \end{aligned}$$

The two  $\delta$  terms are  $\geq 0$  (assuming  $\delta^D$  and  $\delta^M$  are  $\geq 0$ ), but the sign of the combined effect of the  $\theta$  terms, which reflect changes in teacher composition, is ambiguous, even though it is reasonable to assume that all of the  $\bar{\theta}^2$  terms are  $> 0$ . One could argue that this combined effect is not far from zero and, if negative, is smaller in absolute value than the sum of the two  $\delta$  terms, and so  $ITT_{\text{sch}}(2)$  is a lower bound for  $ATE_{\text{sch}}(2)$ , but it is possible that the sum of the composition terms is negative and larger in absolute value than the two  $\delta$  terms. Note, however, that if there are no likers or dislikers then there is no composition effect (since  $p^L = p^D = 0$  and  $\mu = \tau$ ) and so  $ITT_{\text{sch}}(2)$  is a lower bound for  $ATE_{\text{sch}}(2)$ . In particular,  $ATE_{\text{sch}}(2) - ITT_{\text{sch}}(2) = \delta^M \tau p^M$ , which is  $\geq 0$  as long as  $\delta^M \geq 0$ , which is plausible.

Next, consider  $ATE_{\text{stud}}(2)$  and  $ITT_{\text{stud}}(2)$ , and more specifically their difference:

$$ATE_{\text{stud}}(2) - ITT_{\text{stud}}(2) \quad (\text{B.105})$$

$$\begin{aligned} &= \sigma\pi\bar{\delta} + \pi[(2\delta^R + \gamma_{1,2}^R)p^R + (\delta^L(1+\tau) + \gamma_{1,2}^L\tau)p^L/\tau + ((1+\tau)\delta^M + \tau\gamma_{1,2}^M)(\mu/\tau)p^M] \\ &\quad + \pi[\bar{\theta}^{2,L}p^L((1-\tau)/\tau) - \bar{\theta}^{2,D}p^D + \bar{\theta}^{2,M}p^M((\mu/\tau) - 1)] \\ &\quad - [\sigma\pi\bar{\delta} + \pi((2\delta^R + \gamma_{1,2}^R)p^R + (\delta^L(1+\tau) + \tau\gamma_{1,2}^L)(p^L/\tau) - \delta^D\tau(p^D/(1-\tau)) + ((\delta^M(1+\tau) + \gamma_{1,2}^M\tau)(\mu/\tau) - \delta^M\tau(1-\mu)/(1-\tau))p^M)] \\ &\quad \quad - \pi[\bar{\theta}^{2,L}(p^L/\tau) - \bar{\theta}^{2,D}(p^D/(1-\tau)) + \bar{\theta}^{2,M}p^M((\mu/\tau) - (1-\mu)/(1-\tau))] \\ &= \pi[\delta^D\tau(p^D/(1-\tau)) + \delta^M\tau(1-\mu)/(1-\tau)p^M - \bar{\theta}^{2,L}p^L + \bar{\theta}^{2,D}p^D(\tau/(1-\tau)) + \bar{\theta}^{2,M}p^M((\tau-\mu)/(1-\tau))] \end{aligned}$$

This is simply  $\pi$  multiplied by  $ATE_{sch}(2) - ITT_{sch}(2)$ , and so, as with  $ATE_{sch}(2) - ITT_{sch}(2)$ , the sign of this expression for the general case is ambiguous. yet again if there are no likers or dislikers then  $ITT_{stud}(2)$  is a lower bound for  $ATE_{stud}(2)$  as long as  $\delta^M \geq 0$ .

Next, consider  $ATE_{tchr}(3)$  and  $ITT_{tchr}(3)$ . Their difference is:

$$\begin{aligned} & ATE_{tchr}(3) - ITT_{tchr}(3) \quad (B.106) \\ &= 3\bar{\delta} + 3\bar{\gamma}_{1,2} + \bar{\gamma}_{1,2,3} - [\bar{\delta} + p^R(2\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) + p^L(2\gamma_{1,2}^L + \gamma_{1,2,3}^L) + p^M(2\tau\gamma_{1,2}^M + \tau^2\gamma_{1,2,3}^M)] \\ &= p^L(2\delta^L + \gamma_{1,2}^L) + p^D(2\delta^D + 3\gamma_{1,2}^D + \gamma_{1,2,3}^D) + p^M(2\delta^M + (3-2\tau)\gamma_{1,2}^M + (1-\tau^2)\gamma_{1,2,3}^M) \end{aligned}$$

This is plausibly  $\geq 0$ . The term  $p^L(2\delta^L + \gamma_{1,2}^L)$  is  $\geq 0$  as long as two years of exposure to the program does not reduce the skills of liker teachers. The term  $p^D(2\delta^D + 3\gamma_{1,2}^D + \gamma_{1,2,3}^D)$  is  $\geq 0$  as long as three years of exposure to the program does not reduce the skills of disliker teachers relative to the skills they would obtain from one year of exposure to the program. Finally,  $p^M(2\delta^M + (3-2\tau)\gamma_{1,2}^M + (1-\tau^2)\gamma_{1,2,3}^M) \geq 0$  as long as three years of exposure to the program does not reduce the skills of mover teachers relative to the skills they would obtain relative to one year of exposure to the program.

In contrast, it is less clear that  $ITT_{sch}(3)$  can serve as a lower bound for  $ATE_{sch}(3)$ . Their difference is:

$$\begin{aligned} & ATE_{sch}(3) - ITT_{sch}(3) \quad (B.107) \\ &= (3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R)p^R + (\delta^L(2+\tau) + (1+2\tau)\gamma_{1,2}^L + \tau\gamma_{1,2,3}^L)p^L/\tau + [(1+2\tau)\delta^M + \tau(2+\tau)\gamma_{1,2}^M + \gamma_{1,2,3}^M\tau^2](\mu/\tau)p^M \\ &\quad + \bar{\theta}^{3,L}p^L((1-\tau)/\tau) - \bar{\theta}^{3,D}p^D + \bar{\theta}^{3,M}p^M((\mu/\tau) - 1) \\ &\quad - [(3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R)p^R + (\delta^L(2+\tau) + \gamma_{1,2}^L(1+2\tau) + \gamma_{1,2,3}^L\tau)(p^L/\tau) - \delta^D\tau(p^D/(1-\tau))] \\ &\quad - [(\mu/\tau)(\delta^M(1+2\tau) + \gamma_{1,2}^M\tau(2+\tau) + \gamma_{1,2,3}^M\tau^2) - ((1-\mu)/(1-\tau)(\delta^M2\tau + \gamma_{1,2}^M\tau^2))]p^M \\ &\quad - [\bar{\theta}^{3,L}(p^L/\tau) - \bar{\theta}^{3,D}(p^D/(1-\tau)) + \bar{\theta}^{3,M}p^M((\mu/\tau) - (1-\mu)/(1-\tau))] \\ &= \delta^D\tau(p^D/(1-\tau)) + ((1-\mu)/(1-\tau)(\delta^M2\tau + \gamma_{1,2}^M\tau^2)p^M \\ &\quad - \bar{\theta}^{3,L}p^L + \bar{\theta}^{3,D}p^D(\tau/(1-\tau)) + \bar{\theta}^{3,M}p^M((\tau-\mu)/(1-\tau)) \end{aligned}$$

The two  $\delta$  terms are  $\geq 0$  (assuming  $\delta^D$  and  $\delta^M$  are  $\geq 0$ ) as long as two years of exposure to the program does not reduce the skills of movers (as long as  $2\delta^M + \gamma_{1,2}^M \geq 0$ ), but the sign of the combined effect of the  $\theta$  terms, which again reflects changes in teacher composition, is ambiguous, even though it is reasonable to assume that all of the  $\bar{\theta}^2$  terms are  $> 0$ . Note, however, that if there are no likers or dislikers then there is no composition effect (since  $p^L = p^D = 0$  and  $\mu = \tau$ ) and so  $ITT_{sch}(3)$  is a lower bound for  $ATE_{sch}(3)$ .

Finally, turn to student skills for year three and compare  $ITT_{stud(3)}$  with  $ATE_{stud(3)}$ :

$$\begin{aligned}
& ATE_{stud(3)} - ITT_{stud(3)} \quad (B.108) \\
&= \sigma^2 \pi \bar{\delta} + \sigma \pi [(2\delta^R + \gamma_{1,2}^R) p^R + (\delta^L(1+\tau) + \gamma_{1,2}^L \tau)(p^L/\tau) + ((1+\tau)\delta^M + \tau\gamma_{1,2}^M)(\mu/\tau)p^M] \\
&+ \pi [(3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) p^R + (\delta^L(2+\tau) + (1+2\tau)\gamma_{1,2}^L + \tau\gamma_{1,2,3}^L)(p^L/\tau) + ((1+2\tau)\delta^M + \tau(2+\tau)\gamma_{1,2}^M + \gamma_{1,2,3}^M \tau^2)(\mu/\tau)p^M] \\
&+ \sigma \pi [\bar{\theta}^{2,L}((1-\tau)/\tau) - \bar{\theta}^{2,D} p^D + \bar{\theta}^{2,M} p^M((\mu/\tau) - 1)] + \pi [\bar{\theta}^{3,L} p^L((1-\tau)/\tau) - \bar{\theta}^{3,D} p^D + \bar{\theta}^{3,M} p^M((\mu/\tau) - 1)] \\
&- \sigma^2 \pi \bar{\delta} + \sigma \pi [(2\delta^R + \gamma_{1,2}^R) p^R + (\delta^L(1+\tau) + \tau\gamma_{1,2}^L)(p^L/\tau) - \delta^D \tau(p^D/(1-\tau)) + (\delta^M((1+\tau) + \gamma_{1,2}^M \tau)(\mu/\tau) - \delta^M \tau(1-\mu)/(1-\tau)) p^M] \\
&- \pi [(3\delta^R + 3\gamma_{1,2}^R + \gamma_{1,2,3}^R) p^R + (\delta^L(2+\tau) + \gamma_{1,2}^L(1+2\tau) + \gamma_{1,2,3}^L \tau)(p^L/\tau) - \delta^D \tau(p^D/(1-\tau))] \\
&- \pi [(\mu/\tau)(\delta^M(1+2\tau) + \gamma_{1,2}^M \tau(2+\tau) + \gamma_{1,2,3}^M \tau^2) - ((1-\mu)/(1-\tau)(\delta^M 2\tau + \gamma_{1,2}^M \tau^2))] p^M \\
&- \sigma \pi [\bar{\theta}^{2,L}(p^L/\tau) - \bar{\theta}^{2,D}(p^D/(1-\tau)) + \bar{\theta}^{2,M} p^M((\mu/\tau) - (1-\mu)/(1-\tau))] \\
&- \pi [\bar{\theta}^{3,L}(p^L/\tau) - \bar{\theta}^{3,D}(p^D/(1-\tau)) + \bar{\theta}^{3,M} p^M((\mu/\tau) - (1-\mu)/(1-\tau))] \\
&= \sigma \pi [\delta^D p^D(\tau/1-\tau) + \delta^M \tau p^M((1-\mu)/(1-\tau))] + \pi [\delta^D p^D(\tau/1-\tau) + (\delta^M 2\tau + \tau^2 \gamma_{1,2}^M) p^M((1-\mu)/(1-\tau))] \\
&+ \sigma \pi [-\bar{\theta}^{2,L} p^L + \bar{\theta}^{2,D} p^D(\tau/(1-\tau)) + \bar{\theta}^{2,M} p^M((\tau-\mu)/(1-\tau))] + \pi [-\bar{\theta}^{3,L} p^L + \bar{\theta}^{3,D} p^D(\tau/(1-\tau)) + \bar{\theta}^{3,M} p^M((\tau-\mu)/(1-\tau))]
\end{aligned}$$

The first three  $\delta$  terms are  $\geq 0$  (assuming  $\delta^D$  and  $\delta^M$  are  $\geq 0$ ), and the  $\delta^M 2\tau + \tau^2 \gamma_{1,2}^M$  term is also  $\geq 0$  as long as two years of exposure to the program does not reduce the skills of movers (as long as  $2\delta^M + \gamma_{1,2}^M \geq 0$ ). Yet the sign of the combined effect of the  $\theta$  terms, which again reflects changes in teacher composition, is ambiguous, even though it is reasonable to assume that all of the  $\bar{\theta}^2$  terms are  $> 0$ . Note, however, that if there are no likers or dislikers then there is no composition effect (since  $p^L = p^D = 0$  and  $\mu = \tau$ ) and so  $ITT_{stud(3)}$  is a lower bound for  $ATE_{stud(3)}$ .

## Appendix C: Teacher Allocation during the Second Year of Treatment

Under the framework developed in Section 3 and in Appendix B, schools are randomly assigned to treatment and control arms during the first year of treatment. Teachers cannot change schools during that year and, therefore, are also randomly distributed between APM and non-APM schools. Therefore, we expect the proportion of each type of teacher to be equally distributed between treatment arms.

**Table C1: Teacher Distribution during Year One**

School	Proportion of each type of teachers				
Treatment Arm	Likers	Movers	Remainers	Dislikers	Total
Treatment	$p^L$	$p^M$	$p^R$	$p^D$	1
Control	$p^L$	$p^M$	$p^R$	$p^D$	1

As the first year end, some teachers change schools according to their preferences. Remainers stay in their school regardless of the treatment status. Likers that started in control schools will all move to APM schools. Likers that started in APM schools will remain in that type of school, although only a fraction of them will stay in the same school (we defined that proportion as  $\sigma$ ), and the rest moving to a different treated school. Dislikers that started in treated schools will all move to control schools, while dislikers that started in control schools will remain in that treatment arm, with a fraction remaining in their original school (defined as  $v$ ) and the rest moving to a different control school. All movers will change schools independently of their original placement. Irrespectively of where they started, a fraction will move to treated schools (defined as  $\mu$ ) and the rest will go to control schools.

**Table C1: Teacher Relocation between Years One and Two**

	Movement decision	Likers	Movers	Remainers	Dislikers	Row Sum
Assigned to treatment	Moves to treated	$p^L(1-\sigma)$	$p^M\mu$	0	0	$p^L(1-\sigma) + p^M\mu$
	Moves to control	0	$p^M(1-\mu)$	0	$p^D$	$p^M(1-\mu) + p^D$
	Stays in same school	$p^L\sigma$	0	$p^R$	0	$p^L\sigma + p^R$
Assigned to control	Moves to treated	$p^L$	$p^M\mu$	0	0	$p^L + p^M\mu$
	Moves to control	0	$p^M(1-\mu)$	0	$p^D(1-v)$	$p^M(1-\mu) + p^D(1-v)$
	Stays in same school	0	0	$p^R$	$p^Dv$	$p^R + p^Dv$