



# What do test scores really capture? Evidence from a large-scale student assessment in Mexico



Rafael de Hoyos<sup>a,\*</sup>, Ricardo Estrada<sup>b</sup>, María José Vargas<sup>c</sup>

<sup>a</sup>XABER, ITAM, and the World Bank, 20433 Washington DC, NW, United States

<sup>b</sup>CAF—Development Bank of Latin America, Av. Eduardo Madero 900, Edificio Catalinas Plaza, piso 15, Ciudad de Buenos Aires, Argentina

<sup>c</sup>The World Bank, 1818 H Street, 20433 Washington DC, NW, United States

## ARTICLE INFO

### Article history:

Accepted 22 April 2021

### JEL Codes:

I20

J24

### Keywords:

Standardized testing

Student learning

Education policy

## ABSTRACT

This paper studies the relationship between test scores and cognitive skills using two longitudinal data sets that track student performance in a national standardized exam in grades 6, 9, and 12 and post-secondary school outcomes in Mexico. Exploiting a large sample of twins to control for all between-family differences in school, household, and neighborhood inputs, we find that primary school test scores are a strong predictor of secondary education outcomes. Using a data set that links results in the national standardized test to later outcomes, we find that secondary school test scores predict university enrollment and hourly wages. These results indicate that, despite their limitations, large-scale student assessments can capture the skills they are meant to measure and can therefore be used to monitor student learning in developing countries.

© 2021 Elsevier Ltd. All rights reserved.

## 1. Introduction

There is increasing recognition that education brings about individual and society-wide benefits when students acquire a set of relevant skills during their formative years. Literacy and numeracy stand out among these skills, as they are seen as the foundation for the acquisition of other skills and a direct determinant of critical factors for personal and social well-being, such as labor market outcomes, health conditions, and democratic participation (World Bank, 2018; CAF, 2016). There is, however, no consensus about whether education systems can measure these foundational skills and track their progress at scale.

In recent years, many developing countries have implemented large-scale student assessments via standardized tests to monitor cognitive skills such as literacy and numeracy.<sup>1</sup> However, critics of standardized tests argue that this type of testing promotes a reductionist approach to education that emphasizes literacy and numeracy to the detriment of other equally important subject areas. Moreover, critics point out that the reliability of standardized tests is compromised by either the absence or the presence of incentives. On

the one hand, students might not put enough effort into a test if no consequences are attached to the results (see, for example, Akyol, Krishna, & Wang (2018) and Gneezy et al. (2019), and a review in Finn (2015)). On the other hand, high-stakes tests can create perverse incentives that lead to game the system or teach to the test, which may raise test scores but not truly improve students' learning.<sup>2</sup> The effects of these incentives on the reliability of standardized tests may be larger in developing countries, where weak implementation capacity is more prevalent. Hence, this paper seeks to answer the following question: Do large-scale student assessments based on standardized testing capture the cognitive skills they are designed to measure?

One way to address this question is to estimate the relationship between test scores and future education and labor market outcomes.<sup>3</sup> A positive correlation between test scores and future outcomes is indicative but does not prove that standardized tests capture cognitive skills, as other unobserved factors could drive this correlation. A recent strand of papers has documented the direct effects of school, household, and neighborhood inputs on test scores

\* Corresponding author.

E-mail addresses: [rdehoyos@worldbank.org](mailto:rdehoyos@worldbank.org) (R. de Hoyos), [restrada@caf.com](mailto:restrada@caf.com) (R. Estrada), [mvargasmanquera@worldbank.org](mailto:mvargasmanquera@worldbank.org) (M.J. Vargas).

<sup>1</sup> Cheng and Gale (2014) survey 125 developing countries and find that 105 of them have implemented a national student assessment based on standardized testing.

<sup>2</sup> For a discussion on incentives and strategic behavior in testing, see Figlio and Loeb (2011), Koretz and Barron (1998), Koretz (2017) and Neal (2011).

<sup>3</sup> The other standard alternatives are to correlate test scores with other measures of student learning—notably GPA—or contemporaneous measures of labor market outcomes. However, the first strategy is subject to circularity and the second to reverse causality, as discussed in detail in Heckman and Kautz (2012).

(Salardi & Michaelsen, 2019; Chang & Padilla-Romo, 2019; Lavy, Ebenstein, & Roth, 2014). In addition, an important body of literature in psychology and economics has shown the significant influence of noncognitive skills on both test scores and labor market outcomes (Duckworth & Seligman, 2005; Heckman & Kautz, 2012; Heckman, Stixrud, & Urzua, 2006). We use a twin fixed-effects specification to control for all between-family differences in school, household, and neighborhood inputs. Then, in a second, more exploratory, specification we use within-individual variation in scores from different subject areas to control for the direct effect of noncognitive skills on test scores.

To estimate these specifications, we construct two longitudinal data sets to track students along their education trajectories and initial labor market outcomes in Mexico. Both data sets use information from the National Assessment of Academic Achievement in Schools (ENLACE, from its acronym in Spanish), a census-based standardized test applied to primary and secondary school students between 2006 to 2014. The first data set tracks students who took the test in 2007 at the end of primary school (grade 6), in 2010 at the end of lower secondary school (grade 9), and in 2013 at the end of upper secondary school (grade 12). Given the large size of this cohort (close to 2 million individuals), we are able to identify about 10,000 pairs of twins in this data set, using their last names, date of birth and the school that they attended in grade 6. The second data set merges the grade 12 test scores with a special module of the Mexican labor force survey (ENOE, from its acronym in Spanish) of 2010 that was applied to secondary school graduates between the ages of 18 and 20.

The results of this paper show that, despite their limitations, large-scale standardized tests can meaningfully capture skills. We find a positive and significant relationship between test scores and future education outcomes that remains even after controlling for observed and unobserved family heterogeneity. In the twin fixed-effects specification, a 1-standard-deviation (SD) higher score in grade 6 is correlated with a 3.3-percentage-point higher probability of on-time graduation from grade 9 and a 5.7-percentage point increase for grade 12, and with 0.49 and 0.53 SD higher test scores, respectively. These estimated coefficients can be interpreted as the composite effect of cognitive and noncognitive skills, measured by grade 6 test scores, on future education outcomes and the within-twin variation in factors with a direct effect both in grade 6 and later education outcomes (e.g., motivation to perform on the test). In an attempt to control for the effects of noncognitive skills on future education outcomes, we regress language (math) test scores at grades 9 and 12 on language (math) test scores at grades 6 controlling for grade 6 test scores in math (language) test scores at grade 6 and twin fixed effects. The rationale for this econometric specification is that the subject-specific relationship between grade 6 test scores and future test scores is driven mainly by cognitive skills.<sup>4</sup> The results from this specification show that grade 6 test scores have a strong relationship with test scores in the same subject in grades 9 and 12, controlling for twin fixed effects and grade 6 test scores in the other subject area. These results suggest that, although test scores from large-scale student assessments are a composite of different effects, they do capture the skills that they are designed to measure.

Finally, we analyze the relationship between grade 12 test scores and post-secondary school outcomes using the second panel data set. As this panel is less rich than the first one, the analysis follows a conventional strategy and substitutes the twin fixed effects with a set of covariates to control for family background. The

results show that grade 12 test scores are a good predictor of university enrollment and hourly wages. A 1-SD increase in grade 12 test scores is correlated with a 10-percentage point increase in the likelihood of being enrolled in university, and, conditional on being employed, individuals with 1-SD-higher test scores in grade 12 have wages that are 6% higher 1 or 2 years after graduating from high school.

This paper joins the papers that study the relationship between test scores and future outcomes in developed countries (Chetty et al., 2011; Lin, Lutter, & Ruhm, 2018; Murnane, Willett, & Levy, 1995; Rose, 2006 for the United States; Currie & Thomas, 2001 for the United Kingdom; and Lindqvist & Vestman, 2011 for Sweden) as the first paper to perform this analysis in a developing country. Conducting this study in a developing country is important because the relationship between test scores and future outcomes may differ from the one in developed countries, given the former's weaker implementation capacity and less mature labor markets. A second contribution of this paper is that it uses a census-based student assessment as opposed to the results of a controlled assessment conducted in the context of a research project—which is likely not subject to the same incentives and implementation challenges as large-scale student assessments. The third and most important contribution is that it presents the first analysis for a developing country exploiting a large sample of twins to study the relationship between test scores and future education outcomes among individuals with identical family background.<sup>5</sup> Overall, the evidence presented in this paper is of particular relevance to the large applied economics literature that uses large-scale student assessments to study the process of human capital formation or evaluate the effects of certain education policies.<sup>6</sup>

The results presented here have two important policy implications. First, although appropriate design and implementation matter for the quality and credibility of large-scale student assessments, the results presented here support the use of large-scale student assessments to measure and monitor the evolution of learning outcomes in complex education systems with relatively weak institutions and low implementation capacity. Second, the findings of this paper show that student test scores are a good predictor of future outcomes related to well-being. This is relevant for the design of effective policies to identify and address education disparities and, therefore, promote social mobility. Results from large-scale student assessments, especially those of early grades, could be used to target resources toward schools and students with the lowest learning levels—who also tend to be the poorest.

The rest of the paper is organized as follows. Section 2 presents the Mexican education system and the ENLACE test, Section 3 describes the panel data sets, Section 4 describes the analytical framework and empirical strategy, Section 5 presents and discusses the results, and Section 6 concludes.

<sup>5</sup> The findings of this paper are related in a complementary way to the rich psychometrics literature that studies the validity of large-scale skills assessments, either by validating the internal consistency of the constructs, testing structure and sampling design, or by studying the use of results coming from such assessments by policy makers or researchers outside psychology (see, for example, Braun & von Davier (2017) and Lin, Bumgarner, & Chatterji (2014)). This paper is also connected to the economics literature that uses structural models to identify the contributions of cognitive and noncognitive skills to test scores (Cunha, Heckman, & Schennach, 2010; Heckman et al., 2006) and recent work by Laajaj and Macours (2019), who validate measures of human capital that are widely used in applied economics (survey-based measures of skills for rural contexts, in their case).

<sup>6</sup> For examples using ENLACE see: Avitabile and de Hoyos (2018), Dustan, de Janvry, and Sadoulet (2017), Estrada (2019), Estrada and Gignoux (2017) and Salardi and Michaelsen (2019).

<sup>4</sup> Under the assumption that effort to perform well is constant across subject areas.

## 2. Context

### 2.1. The Mexican Education System

The basic education system in Mexico includes 3 years of preprimary, 6 years of primary, and 3 years of lower secondary (grades 7 to 9) school, and upper secondary education is for 3 years (grades 10 to 12). Both the basic and the upper secondary levels are mandatory. More than 30 million students are enrolled in mandatory education across 243,000 schools, that employ close to 1.5 million teachers. In the basic education system, public schools, which are decentralized at the state level, account for 90% of total school enrollment. In the upper secondary system, public schools are managed by the federal government, state governments, and public universities. Private sector education is a relatively small share of the education system, and it is heavily regulated by the Secretariat of Public Education (SEP, from the Spanish).

Most children ages 6 to 12 are enrolled in the education system and graduate from primary school. However, of every 100 students enrolled in lower secondary school, only 85 graduate on time, and this number falls to 65 in upper secondary. The weak performance of Mexican students in international assessments, notably the Program for International Student Assessment (PISA), gave the learning crisis a prominent place on the public agenda—Mexico ranked last among OECD (Organisation for Economic Co-operation and Development) countries in PISA 2000, the first wave of the assessment (OECD, 2001). These concerns contributed to the establishment of a national standardized assessment in 2006: ENLACE.<sup>7</sup>

### 2.2. The ENLACE Standardized Test

From 2006 to 2014, the SEP administered ENLACE, a census-based standardized test that gathered information on student achievement in math, literacy, and a rotating subject. Initially, ENLACE was given to students in grades 3 to 9 (primary and lower secondary), but starting in 2008, ENLACE was also given in grade 12 (the final year of secondary school).

ENLACE was designed as a low-stakes assessment and had no bearing for students on GPA, graduation, or admission to the next schooling level. The purpose of the assessment, as stated by SEP, was to increase parental and student participation in the learning process, improve lesson preparation, improve teacher and principal training, strengthen policy planning, and increase transparency and accountability.

School participation was mandatory for grades 3 to 9 and optional for grade 12. Nonetheless, the take-up of the test was consistently above 85%, even among 12th graders. The states of Michoacán and Oaxaca, which have a heavy presence of a radical teachers union (CNTE, from the Spanish), have consistently recorded the lowest participation rates (Table 1). A total of 15.1 million students in 136,000 schools took the examination in 2013, the last year ENLACE was given in most grades.

By design, ENLACE had a national mean score of 500 and a SD of 100 for every subject area and grade in its first year of implementation. ENLACE's methodology followed item response theory, with comparability of the results over time. The test had 50 to 75 questions per subject and was applied in eight 45-min sessions over 2 days. Several mechanisms were in place to prevent exam manipulation. An external coordinator was assigned to every school to oversee the test's administration, alongside the school principal. Teachers were not allowed to monitor their own classes, and parents were invited to act as exam monitors. Exams were centrally

**Table 1**

ENLACE: Student Take-Up (percentage).

	(1) Primary 2007	(2) Lower secondary 2010	(3) Upper secondary 2013
<i>National</i>	89.6	86.5	89.4
Michoacán	48.7	33.9	26.0
Oaxaca	65.1	0.7	92.9
National without Michoacán and Oaxaca	92.6	93.3	93.0

Notes: Primary school take-up includes grades 3 to 6, lower secondary grades 7 to 9, and upper secondary grade 12. Source: SEP.

marked by SEP, and computer software identified abnormal response patterns to detect cheating.

SEP produced school report cards that were distributed among school principals and an online website where parents and students could check their individual results. Yet, few schools viewed ENLACE as a diagnostic tool, which limited its effectiveness as an improvement tool (de Hoyos, García-Moreno, & Patrinos, 2017). Results from ENLACE received wide attention from the Mexican public. Every year, the results made the front page of most newspapers. NGOs produced and disseminated state and school rankings that made ENLACE a medium-stakes test from the point of view of the school directors and the school community. Despite the original objectives of the assessment, SEP used ENLACE scores to deliver monetary bonuses to teachers and principals participating in two different incentive programs—one starting in 2008 and a second one in 2009.<sup>8</sup> This decision made ENLACE a de facto high-stakes test for teachers and principals, encouraging strategic behavior and resulting in multiple concerns about teaching to the test and grade inflation.

Complaints about ENLACE—mainly related to grade inflation—and the creation of the national autonomous evaluation institute (INEE, from the Spanish), which was given the responsibility of regulating national student assessments, contributed to the cancellation of ENLACE. The test was administered for the last time in 2013 for grades 3–9 and in 2014 for grade 12. The INEE launched a new, survey-based test called *Plan Nacional para la Evaluación de los Aprendizajes* in 2015.

## 3. Data

Two longitudinal data sets are constructed for this paper. The first panel is formed from ENLACE test scores of students who completed primary school (grade 6) in 2007, lower secondary school (grade 9) in 2010, and upper secondary school (grade 12) in 2013. The second panel merges a special module of the Mexican labor survey (ENOE, from the Spanish) applied to individuals ages 18, 19, and 20 years during the third quarter of 2010, with students that sat the ENLACE test in grade 12 in May 2008, 2009, and 2010.<sup>9</sup>

### 3.1. The ENLACE Panel

The ENLACE data set recorded the unique national population identifier (CURP, from the Spanish) of all test takers, enabling the

<sup>8</sup> In 2008, SEP linked ENLACE to *Carrera Magisterial*, a national teacher incentive program that offered bonuses to primary and lower secondary teachers participating in the program. Students' ENLACE scores were given a weight of 20% in the program's total score; this weight was increased to 50% in 2011. In 2009, SEP launched the Program of Incentives for Teaching Quality, which delivered monetary bonuses to teachers and principals of classrooms and schools that performed highly on in specific categories on ENLACE.

<sup>9</sup> The data sets described in this paper can be requested from [www.xaber.org.mx](http://www.xaber.org.mx).

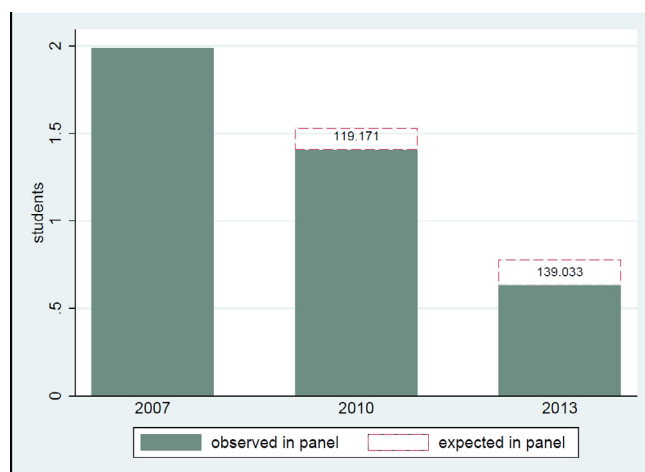
<sup>7</sup> For more information on ENLACE see: <http://www.enlace.sep.gob.mx/>

construction of a panel of students with learning outcomes at different points in their education trajectory. In addition to these outcomes, the ENLACE data set includes a school identifier, and the full name, birthday, state of birth, and sex of the student.<sup>10</sup> Using the CURP, we merged the information from all grade 6 students who took the exam in 2007 with their exam results from 2010 (grade 9) and 2013 (grade 12). We begin the panel in 2007 because of the relatively low take-up in 2006 (the first year ENLACE was administered) and the lack of use of CURPs in some states during this first year of application. Among the 1,881,470 students who sat the test in 2007 and had complete information on their CURP (94.5% of the total), we were able to identify 72.9% three years later in grade 9 and 34.5% six years later in grade 12 (Fig. 1).<sup>11</sup>

The large attrition observed in the panel is caused by (1) grade repetition, (2) school dropout, (3) exam take-up rates of less than 100%, and (4) imperfect matching due to administrative data errors. If a large share of the attrition in the panel is caused by low test take-up or imperfect matching, we might not be able to identify accurately the effects of grade 6 test scores on lower and upper secondary on-time graduation rates. To quantify the magnitude of the causes behind attrition, we use administrative data from the annual school census (known as *Formato 911*) to estimate the expected survival rates given state-level repetition and dropout rates in lower and upper secondary school. Given the school trajectories implied by administrative data, 77% of the student population that completed grade 6 in 2007 was expected to finish grade 9 in 2010 and 39% to graduate from upper secondary in 2013 (Fig. 1). Therefore, the ENLACE panel has a survival rate that is 4 and 5 percentage points lower vis-à-vis the survival rate implied by administrative data, in lower and upper secondary, respectively. This difference is the result of both less than 100% test take-up and imperfect matching. The ENLACE take-up rate is high, around 93% of the student population (excluding the states of Michoacán and Oaxaca). The difference in graduation or survival rates between the ENLACE panel and administrative data is not trivial, but it is reassuring that most of the attrition observed in the panel is explained by grade repetition and school dropouts, which are dimensions we want to examine in this paper.<sup>12</sup>

ENLACE take-up rates in the years under analysis were close to 93%, with the exception of the states of Oaxaca and Michoacán, which we exclude from the panel (see Table 1).<sup>13</sup> ENLACE included a context questionnaire which was applied to a random sample of students and parents. We link the panel data set we constructed to results from the context questionnaire to retrieve information on parental education and occupation. The sample size for Grade 6 was 7,557 observations in 2007, with 5,677 individuals with full responses to the variables of interest.

We use student's personal information to identify 20,252 twins in the ENLACE panel. We define twins as students with identical last names (the surname in Mexico is composed of father's and mother's last names) and birthday who attend grade 6 at the same school in 2007. After the matching, the twins identified account for



**Fig. 1.** ENLACE Panel Matching. Notes: (1) The graph presents the number of students in the ENLACE panel in 2007, 2010 and 2013 and the expected number of observations given the estimated school trajectories in secondary school using the Formato 911 data. (2) Sample: Students who took the ENLACE exam in grade 6 in 2007. (3) Data: ENLACE panel and Formato 911 data.

1.08% of the ENLACE panel in 2007, a level close to the prevalence of multiple pregnancies in Mexico. Table A.1 in the Appendix presents summary statistics for the ENLACE panel data set. The main estimations in the table include the sample of twins, but the results are similar if we use only the main sample.<sup>14</sup>

### 3.2. The ENILEMS-ENLACE Panel

Every quarter of the year, the Mexican statistics office, INEGI (from the Spanish), collects labor market information through the national labor force survey (ENOE), a rotating household survey.<sup>15</sup> In some quarters, ENOE's core survey is complemented by a thematic module that is usually requested and financed by different secretariats and government agencies. During the 3rd quarter of 2010 (July to September), ENOE's special module was the *Encuesta Nacional de Inserción Laboral de los Egresados de Educación Media Superior* (ENILEMS), a survey targeting upper secondary school graduates ages 18, 19, and 20. The objective of ENILEMS was to provide information on the transition between the end of mandatory schooling (grade 12) and higher education or the labor market.<sup>16</sup>

The ENILEMS-ENLACE panel merges information from the respondents of the ENILEMS 2010 survey with their results on ENLACE grade 12 from May 2008, 2009, or 2010. Although ENILEMS 2010 did not capture the CURP, it included all the necessary information to create a pseudo-CURP formed by combining letters and numbers coming from the full name, sex, birth date, and state of birth. The difference between the pseudo-CURP and the CURP is that the former does not have the last three digits of the population identification code that the Mexican government generates. We created the pseudo-CURP for 7,105 observations included in ENILEMS 2010 using the official algorithm for generating the CURP.<sup>17</sup>

<sup>10</sup> The CURPs as well as the students' names and dates of birth are confidential data protected by Mexico's personal information laws. For this study, we were able to use this information by closely collaborating with SEP.

<sup>11</sup> Some upper secondary schools—concentrated in the states of Nuevo León, Coahuila, and San Luis Potosí—follow a 2-year curriculum instead of the regular 3-year curriculum. Hence, we also merge the information from grade 6 students in 2007 with the results of those who took ENLACE at the end of their upper secondary education in 2012. From the 648,301 students observed at the end of upper secondary education, we find 4% in 2012 (2-year upper secondary students) and 96% in 2013 (3-year upper secondary students).

<sup>12</sup> Avitabile and de Hoyos (2018) use the same dataset as the one used in this paper, but for a previous year, and find that 88% of the attrition between grades 9 and 12 in ENLACE is explained by dropout and 12% by grade repetition.

<sup>13</sup> A fierce opposition to standardized testing by a local teachers union (CNTE, from the Spanish) explains the low ENLACE take-up rates in Oaxaca and Michoacán.

<sup>14</sup> Table A.1 in the Online Appendix reports differences in means between the sample of twins and non-twins. Although statistically significant at conventional levels, the differences are very small and support the idea that multiple (versus single) pregnancies are mostly a random event.

<sup>15</sup> For more information on ENOE, see <http://www.beta.inegi.org.mx/proyectos/enchogares/regulares/enoe/>.

<sup>16</sup> For more information on ENILEMS, see <http://www.beta.inegi.org.mx/proyectos/enchogares/modulos/enilems/>.

<sup>17</sup> The CURP is an 18-digit unique personal identifier formed by a combination of letters and numbers taken from the individual's full name, date of birth, sex, and state of birth plus a three-digit code assigned by the Mexican population council. For more information, see <https://renapo.gob.mx/swb/>.

A simple merge of ENILEMS and ENLACE grade 12 using the pseudo-CURP and CURP, respectively, was able to match 2,820 observations (40% of the ENILEMS sample). This relatively low matching rate is almost entirely explained by the lack of the last three digits in the pseudo-CURP and measurement error in the variables used to produce the merge. An additional 18% of the sample was recovered manually by identifying coding or registration errors in the CURP generation process (i.e., errors in birth date or misspelled names). Overall, 58% of the individuals in the ENILEMS sample were matched to their ENLACE grade 12 test scores. After eliminating missing observations in the ENLACE score, the panel reaches a total of 3,781 observations. The ENILEMS-ENLACE panel also includes the information from the ENOE regular questionnaire for 3,718 matched observations. Table A.2 in the Online Appendix reports differences in means between the observations of ENILEMS that were matched with ENLACE scores and the ones that were not; the differences are small. Table A.2 in the Appendix presents summary statistics for the ENILEMS-ENLACE panel data set.

#### 4. Analytical framework and empirical strategy

##### 4.1. Analytical framework

The main challenge to study cognitive skills is that these are not directly observable. As many countries have implemented large-scale student assessments via standardized tests to measure cognitive skills, an important question is to what degree such test scores really capture cognitive skills. One way to validate skills assessments is to look at the predictive power of test scores over future education and labor market outcomes (Heckman & Kautz, 2012). To better understand this strategy consider the following structural equation:

$$test_{it} = \psi(C_{it}, P_{it}, I_{it}) \tag{1}$$

Eq. (1) defines test scores of individual  $i$  at time  $t$  ( $test_{it}$ ) as a function of contemporaneous cognitive skills ( $C_{it}$ ), noncognitive skills ( $P_{it}$ ) and other inputs ( $I_{it}$ ). A test enjoys validity if higher cognitive skills produce higher test scores,  $\partial test / \partial C > 0$ . The higher the magnitude of  $\partial test / \partial C$ , the higher the validity of the test. But, test scores might also depend on noncognitive skills (e.g., perseverance or self-control) and other elements such as school and family inputs. Eq. (1) makes explicit that the same test given to individuals of the same cognitive ability will produce different scores if individuals exhibit differences in, say, intrinsic or extrinsic motivation to perform well on the test. Standardized tests are designed to measure cognitive skills in different subject areas. Hence, it is more precise to define  $C_{it}$  (and  $P_{it}$ ) as a vector of multiple sub-skills, which includes subject-specific skills ( $C_{it}^s$ )—like math ( $C_{it}^{math}$ ) and literacy ( $C_{it}^{literacy}$ )—and general skills ( $C_{it}^g$ ). This definition recognizes that test scores are produced by both subject-specific ( $s_{it}$ ) and general skills ( $g_{it}$ )—both of which can be mapped into cognitive and noncognitive skills,  $s_{it} = \{C_{it}^s, P_{it}^s\}$  and  $g_{it} = \{C_{it}^g, P_{it}^g\}$ .<sup>18</sup>

Let future education or labor market outcomes of individual  $i$  in period  $t + 1$ ,  $y_{it+1}$ , be defined by the following structural equation:

$$y_{it+1} = \phi(C_{it+1}, P_{it+1}, I_{it+1}) \tag{2}$$

where  $C_{it+1}$  is the level of cognitive skills of individual  $i$  in  $t + 1$ ,  $P_{it+1}$  is a vector of noncognitive skills, and  $I_{it+1}$  is a vector of other inputs. In Eq. (2),  $\partial y_{it+1} / \partial C_{it+1} > 0$  and  $\partial y_{it+1} / \partial P_{it+1} > 0$  show that cognitive and noncognitive skills have a positive impact on future education

<sup>18</sup> This definition follows the developmental psychology literature which decomposes cognitive skills in general (e.g. intelligence and working memory) and subject-specific skills (e.g. knowledge of arithmetical operations)—see for example Geary, Nicholas, Li, and Sun (2017) and Ritchie, Bates, and Deary (2015).

outcomes and a return in the labor market. Eq. (2) recognizes that family, school, and neighborhood inputs,  $I_{it+1}$ , can have a direct effect on outcomes—net of their effect on skills formation.

As the following equations recognize, skills formation is a dynamic process:

$$C_{it} = \chi(C_{it-1}, P_{it-1}, I_{it}) \tag{3}$$

$$P_{it} = \chi(C_{it-1}, P_{it-1}, I_{it}) \tag{4}$$

$$I_{it} = \chi(C_{it-1}, P_{it-1}, W_i) \tag{5}$$

Eqs. (3) and (4) state that skills are persistent over time but malleable. Cognitive and noncognitive skills follow a cumulative process whereby skills today depend on skills in the previous period plus educational investments (including school, family, and other inputs). Ample evidence recognizes the cross-dependence of cognitive and noncognitive skills (Heckman & Kautz, 2012).<sup>19</sup> Eq. (5) shows that inputs at period  $t$  depend on skills in the previous period and parental endowment ( $W_i$ ). Eqs. (3)–(5) also imply that skills and inputs are correlated across time. The framework distinguishes between the direct and indirect effects of skills on future outcomes. For instance, perseverance at time  $t-1$  can have a direct effect on test scores (at time  $t$ ) through its impact on perseverance at  $t$  (Eq. (2)), and an indirect effect through its influence on cognitive skills formation between  $t-1$  and  $t$  (Eq. (3)). Notice that only the direct effects are relevant for studying the relationship between test scores and cognitive skills at a given moment of time, as the indirect effects are part of the overall process of skill formation.

The framework presented here details the causal path linking cognitive skills, test scores and future outcomes: time-persistent cognitive skills that produce both test scores and future outcomes. It also outlines the existence of other causal paths—inputs and noncognitive skills—that can potentially link test scores to future outcomes without necessarily being mediated by cognitive skills.

##### 4.2. Empirical strategy

The starting point of the empirical strategy is an equation linking test scores at  $t$  with education and labor market outcomes at  $t + 1$ :

$$y_{it+1} = \beta_0 + \beta_1 test_{it} + \beta_k X_{ik} + \epsilon_i \tag{6}$$

where  $y_{it+1}$  is the education or labor market outcome of individual  $i$  in time period  $t + 1$ ,  $test_{it}$  is the individual's simple average ENLACE score in math and language in period  $t$ ,  $X_{ik}$  is a vector of controls for individual and family characteristics at  $t$ , and  $\epsilon_i$  is a disturbance term. A  $\hat{\beta}_1 > 0$  (or rather  $\hat{\beta}_1 \neq 0$ ) is enough to conclude that test scores—from ENLACE in this study—predict future outcomes, but it is not enough evidence to conclude that test scores capture cognitive skills. As outlined in the previous section, components such as school or family inputs and noncognitive skills can have a direct effect on both test scores and education and labor market outcomes and, therefore, could induce a correlation between test scores and future outcomes.

Ideally, one would like to include  $I_{it}$  from Eq. 1, the full vector of inputs at time  $t$ , in the estimation of Eq. (6). However, this requires access to longitudinal data with information on present education or labor market outcomes, and past test scores, and family, school, and neighborhood inputs, which is rarely available in most data sets. In the absence of this information, researchers usually include a vector  $X_{ik}$  of controls for family background as a proxy for the flow of inputs—justified by the dependence of inputs on the avail-

<sup>19</sup> Notice that given Eqs. (3) and (4), subject-specific and general skills can be written as a function of past subject-specific, general skills, and inputs, e.g.  $g_{it} = \{C_{it-1}^s, C_{it-1}^g, P_{it-1}^s, P_{it-1}^g, I_{it}\}$ .

ability of family resources. This strategy has strong limitations because of unobserved family heterogeneity that affects, in turn, school and household inputs. To improve on this strategy, we substitute  $X_{ik}$  for a vector of twin fixed effects ( $\tau_f$ ):

$$y_{it+1} = \beta_0 + \beta_1 test_{it} + \tau_f + \epsilon_i \tag{7}$$

This specification restricts the estimation to the sample of twins identified in the ENLACE panel data set. Estimation of Eq. (7) exploits only the within-twin variation in grade 6 ENLACE test scores, identifying, therefore, the relationship between test scores and education outcomes net of all observed and unobserved family characteristics. As twins are individuals who share the same family and were born on the same day, the specification controls for differences in family circumstances or birth rank that could lead to differences in the availability of family resources devoted to children from the same family. Of course, there might still be differences in the actual inputs that each twin receives, but these are likely much smaller than those between individuals from different families.  $\hat{\beta}_1$  can be interpreted as an estimate of the relationship between skills at grade 6, as captured by ENLACE, and future outcomes net of all between-family differences in household, school, and neighborhood inputs.

As discussed in Section 4.1, test scores could capture both cognitive and noncognitive skills—among the later particularly those leading to higher student effort at the time of sitting the test, such as perseverance, motivation, and grit. Unfortunately, ENLACE does not include information on noncognitive skills to control for, say, effort at the time of answering the test. As an exploratory exercise, we estimate a specification aimed at reducing the effects of noncognitive skills on future outcome by separating subject-specific and general skills. The rationale for this econometric specification is that the subject-specific relationship between grade 6 test scores and future test scores is driven mainly by cognitive skills, under the assumption that, at the individual level, there are no differences in effort across subject areas to perform well on the test. To fix ideas, assume that the production function of test scores can be modeled in the following way:

$$test_{it}^s = \alpha_0 + \alpha_1 s_{it} + \alpha_2 g_{it} + \tau_f + e_{it} \tag{8}$$

where  $test_{it}^s$  is the score of student  $i$  in subject area  $s$ ,  $s_{it}$  is the subject-specific skills of student  $i$ ,  $g_{it}$  are general skills of student  $i$ ,  $\tau_f$  is a vector of twin fixed effects, and  $e_{it}$  is a random disturbance term. This model makes explicit that test scores are produced by both subject-specific and general skills, as discussed in Section 4.1. We use the cross-subject correlation in Eq. (8) as a proxy for the effect of general skills on future test scores. Since ENLACE measures skills in two different subjects, math and literacy, we estimate the following specification:

$$enlace_{it+1}^{math} = \alpha_0 + \alpha_1 enlace_{it}^{math} + \alpha_2 enlace_{it}^{literacy} + \tau_f + \epsilon_{it} \tag{9}$$

In this specification, the score of student  $i$  in a particular subject area at grade 9 or 12 is regressed on the score for the same subject area at grade 6 controlling for grade 6 score in the other subject area plus a vector of twin fixed-effects. This strategy controls for the direct effects of general skills (including perseverance, motivation, and grit) on test scores. Notice that Eq. (9) allows for a differentiated impact of general skills on math and literacy—as we run separate specifications using math and literacy test scores as outcomes, and  $\hat{\alpha}_2$  is a subject-specific parameter. Hence,  $\hat{\alpha}_1$  in Eq. (9) can be interpreted as an estimate of the relationship between subject-specific skills in grade 6 and future scores net of family inputs and the direct impact of general skills.  $\hat{\alpha}_1$  in Eq. (9) can still include noncognitive skills that are subject-specific, such as the motivation to perform well in one but not the other subject tested.

We speculate though that while the presence of subject-specific noncognitive skills are possible, empirically they will likely have a small impact on the estimation of  $\hat{\alpha}_1$ .<sup>20</sup>

## 5. Results

### 5.1. Grade 6 test scores and secondary school outcomes

Fig. 2 plots local means of on-time graduation and test scores in grades 9 and 12 by percentiles of ENLACE grade 6 test scores. Higher test scores in grade 6 are associated with a higher probability of on-time graduation in grades 9 and 12, and, conditional on this, with higher test scores in grades 9 and 12. The differences in outcomes between students in the top and bottom of the grade 6 test score distribution are startling. For example, less than 20% of students in the bottom decile are enrolled in grade 12 six years later, compared to more than 50% of students in the top decile.

Table 2 (Columns 1–4) quantifies the relationships depicted in Fig. 2 by regressing graduation and test scores in grades 9 and 12 on grade 6 test scores and on a dummy variable for whether the student is female, using the individual-level data.<sup>21</sup> The results confirm that ENLACE test scores at grade 6 are a strong predictor of lower and upper secondary education outcomes. A 1-SD increase in grade 6 test scores is associated with 8.1-percentage point and 11.6-percentage point increases in the probability of on-time graduation from lower and upper secondary school, respectively. A similar story goes for future test scores. A 1-SD increase in grade 6 test scores is correlated with a 0.63–0.61 SD increase in test scores in grades 9 and 12, conditional on taking the ENLACE exam. All results are statistically significant at the one-percent level.

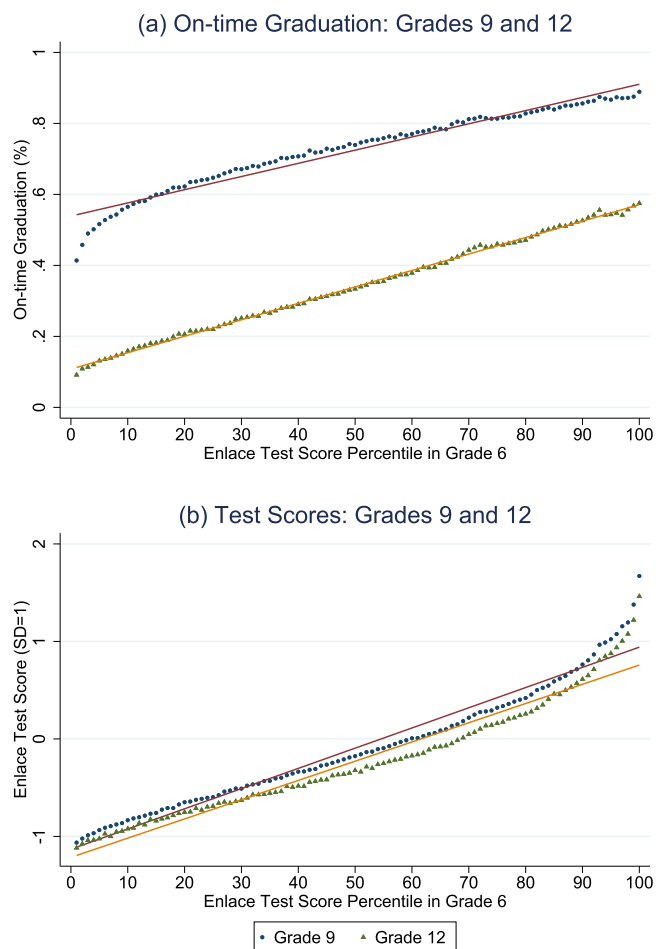
### 5.2. Controlling for family background

Columns 5–8 in Table 2 report the results of estimating Eq. (7), which includes a vector of twin fixed effects. Qualitatively, the findings are the same as the ones from the previous specification: Test scores at grade 6 predict on-time graduation from grades 9 and 12, and conditional on this, they also predict test scores in grades 9 and 12. In the four cases, the results have a high statistical and economical significance. However, as expected, the coefficients from the twin fixed-effects specification are smaller than those estimated without taking into account differences in family background (Columns 1–4). A 1-SD higher test score in grade 6 is correlated with a 3.4- and 5.6-percentage point, respectively, higher probability of on-time graduation from grades 9 and 12 (Columns 5–6), and with 0.49- and 0.53-SD, respectively, higher test scores in grades 9 and 12. In other words, taking family background differences into account reduces the magnitude of the estimated correlation between grade 6 test scores and secondary school outcomes by 59% and 48% in the case of on-time graduation from grades 9 and 12, respectively, and by 22% and 14%, respectively, in the case of test scores in the same grades. The larger reduction in the coefficient of interest in the dropout regressions suggests that family background plays a larger role in explaining on-time graduation from secondary school than future test scores.

The results show that there are significant gender gaps in outcomes. Conditional on the grade 6 score, girls are about 4

<sup>20</sup> We assume that the direct effect of subject-specific noncognitive skills on test scores is negligible—notably the motivation to perform well on the test in one but not the other subject tested. In contrast, we are agnostic about the indirect effect of such skills on test scores, i.e. the effect through the accumulation of subject specific cognitive skills—for example, due to the motivation to learn more about a subject.

<sup>21</sup> For ease of comparison with results from other specifications, the results presented in Columns 1 to 4 in Table 2 are estimated in the sample of twins. In the robustness section, we show that these results are very similar to the ones obtained with the full sample of the ENLACE panel.



**Fig. 2.** Grade 6 Test Scores and Secondary School Outcomes. Notes: (1) The graph plots local means of secondary school outcomes by ENLACE test score percentile in grade 6 in 2007. The solid line shows a linear fit estimated using the grouped data. Panel (a) reports the probability of on-time graduation from grades 9 and 12, proxied by sitting in the Enlace exam in those grades. Panel (b) reports Enlace test scores in grades 9 and 12 (normalised with mean 0 and SD 1) conditional on taking the Enlace exam in 2010 and 2013. (2) Sample: Students who took the ENLACE exam in grade 6 in 2007. (3) Data: ENLACE panel.

percentage points more likely than boys to follow an education trajectory that is free of age-grade distortions, a difference statistically significant at the 1% level (see Row 2 in Columns 5–6). When one looks at test scores, a surprising pattern emerges: The gender gap favors girls until it is reversed by grade 12. Conditional

**Table 2**  
OLS\_Grade 6 test scores and secondary school outcomes.

VARIABLES	(1) Graduation		(3) Score		(5) Graduation		(7) Score	
	Grade 9	Grade 12	Grade 9	Grade 12	Grade 9	Grade 12	Grade 9	Grade 12
Grade 6 score	0.0811*** (0.00298)	0.116*** (0.00329)	0.631*** (0.00687)	0.611*** (0.0101)	0.0335*** (0.00515)	0.0567*** (0.00620)	0.489*** (0.0148)	0.530*** (0.0236)
Girl	0.0273*** (0.00588)	0.0388*** (0.00673)	0.0554*** (0.0128)	-0.144*** (0.0184)	0.0396*** (0.00796)	0.0456*** (0.00950)	0.113*** (0.0208)	-0.119*** (0.0361)
Observations	20,252	20,252	15,494	8,108	20,252	20,252	15,494	8,108
R-squared	0.038	0.059	0.381	0.331	0.805	0.810	0.861	0.885
Twins FE	No	No	No	No	Yes	Yes	Yes	Yes
Mean Dep. Var.	0.764	0.400	0.0101	-0.0380	0.764	0.400	0.0101	-0.0380

Notes: (1) The table displays results of the estimation of Eqs. (6) and (7). (2) Dependent variable is on-time graduation and ENLACE test scores in grades 9 and 12. Graduation is measured as ENLACE take-up. ENLACE test score is the mean of the mathematics and literacy test scores. (3) Sample: Twins (students in the same school in grade 6, with identical last names and birth date) who took the ENLACE exam in grade 6 in 2007. (4) Data: ENLACE panel. (4) Robust standard errors are reported in parentheses. \*\*\*p<0.01, \*\*p<0.05, \*p<0.1.

on initial test scores (and staying in school), girls do better on average than boys in grade 9 (by 0.11 SD) but worse by grade 12 (by 0.12 SD). The switch in the sign of the gender coefficient is explained by girls' lower performance in mathematics in grade 12. See results by ENLACE subject in Table A.4 in the Online Appendix and Avitabile and de Hoyos (2018) for a discussion on this issue.

5.3. Accounting for general skills

Table 3 reports the results of estimating Eq. (9), the specification designed to control for general skills, in the sample of twins. Both math and literacy scores in grade 6 predict scores in those subject areas in grades 9 and 12, even when controlling for grade 6 scores in the other subject area. The coefficients of interest have a large statistical and economical significance. The magnitude of these point estimates is smaller than the magnitude of those presented in Table 2, as one would expect if test scores were driven by both subject-specific and general skills. Once the direct effects of general skills are removed, the association between 1-SD higher test scores in grade 6 is about 0.25–0.34 SD on test scores in grades 9 and 12 (as opposed to the 0.49–0.53 SD reported in Table 2). The implied estimates for general skills at grade 6 on future test scores are about 0.16–0.21 SD and are highly significant—see Row 2 in Columns 1 to 4 in Table 3. These results suggest that test scores at grade 6 are a strong predictor of future education outcomes, and that their predictive power is driven to a large extent by the link between subject-specific skills over time.

5.4. Grade 12 test scores and post-secondary school outcomes

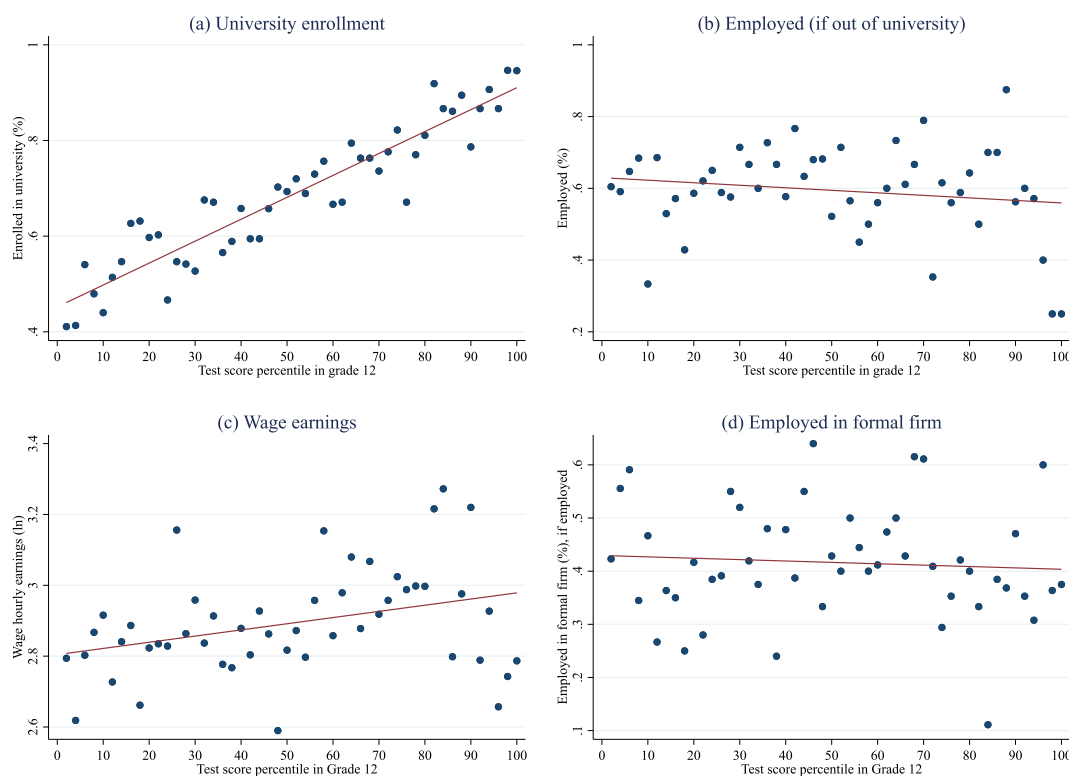
So far the analysis has focused on the effects of grade 6 test scores on secondary school outcomes. The information in the ENLACE panel allows us to implement a robust strategy to control for between-family differences in inputs and for general skills. However, the outcomes studied are limited to the education domain. In this section, we use the ENILEMS-ENLACE panel linking grade 12 test scores with post-secondary school outcomes to present evidence on the relationship test scores and access to university, employment, and wages. A limitation of the ENILEMS-ENLACE panel is that it does not allow for the identification of twins, so we have to rely on conventional and imperfect controls for differences in family background.

Fig. 3 shows visual evidence of the simple correlation between post-secondary school outcomes and grade 12 test scores. Table 4 reports the results of regressing post-secondary school outcomes on grade 12 test scores, gender, status of being a graduate of a private high school (a proxy for high family income), and a dummy

**Table 3**  
OLS\_Grade 6 test scores and secondary school test scores by subject.

VARIABLES	(1)	(2)	(3)	(4)
	Mathematics		Spanish	
	Grade 9	Grade 12	Grade 9	Grade 12
Grade 6 same subject score	0.254*** (0.0175)	0.339*** (0.0241)	0.309*** (0.0155)	0.310*** (0.0271)
Grade 6 other subject score	0.199*** (0.0166)	0.157*** (0.0256)	0.196*** (0.0154)	0.208*** (0.0261)
Girl	0.0208 (0.0235)	-0.307*** (0.0390)	0.187*** (0.0217)	0.152*** (0.0392)
Observations	15,494	8,108	15,494	8,110
R-squared	0.819	0.873	0.849	0.850
Twins FE	Yes	Yes	Yes	Yes
Mean Dep. Var.	-0.00162	-0.0558	0.0208	-0.00699

Notes: (1) The table displays the estimation of Eq. (9). (2) Dependent variable is mathematics and literacy ENLACE test scores in grades 9 and 12. (3) Sample: Twins (students in the same school in grade 6, with identical last names and birth date) who took the ENLACE exam in grade 6 in 2007. (3) Data: ENLACE panel. (4) Robust standard errors are reported in parentheses. \*\*\*p<0.01, \*\*p<0.05, \*p<0.1.



**Fig. 3.** Grade 12 Test Scores and Post-Secondary School Outcomes. Notes: (1) The graph plots local means of post-secondary school outcomes by ENLACE test score ventile in grade 12. The solid line shows a linear fit estimated using the grouped data. Panel (a) reports the probability of university enrollment; (b) reports the probability of being employed conditional on not being enrolled in college; and (c) and (d) report, respectively, the logarithm of the hourly wage and the probability of working in a formal firm conditional in both cases on being employed. (2) Outcomes are measured in the ENILEMS survey at ages 18 to 20 in the third quarter of 2010. ENLACE test scores come from the years 2008, 2009 and 2010. (3) Data: ENILEMS-ENLACE panel.

variable capturing rural/urban resident plus a vector of state and birth year fixed-effects. The results show a strong relationship between test scores and university enrollment. A 1-SD increase in the ENLACE score is associated with an 11-percentage point increase in the probability of university enrollment (statistically significant at the 1% level). Notably, conditional on ENLACE test scores, females are 4 percentage points less likely to be enrolled in university (statistical significance at the 10% level), and graduates of private secondary schools are 10 percentage points more likely to be enrolled in university (statistical significance at the 1% level). In other words, holding end-of-high school test scores constant, there is a gender and family background gap in

university enrollment against women and public secondary school graduates.

Among individuals who are not enrolled in higher education, there is no observed association between test scores and employment status or being employed in the formal sector. Among those who are employed, grade 12 test scores have a positive and significant relationship with future hourly wages, although not with the probability of being employed in a formal firm. A 1-SD increase in test scores is associated with an increase of about 6% in hourly wages (statistical significance at the 10% level). This hourly wage effect of a 1-SD increase in test scores is marginally smaller than those reported in by Lindqvist and Vestman (2011) using data for



**Table 4**  
OLS\_Grade 12 test scores and post-secondary school outcomes.

VARIABLES	(1) University Student	(2) Employed	(3) ln hourly wage	(4) Formal firm
Grade 12 score	0.111*** (0.00887)	-0.0132 (0.0192)	0.0582** (0.0244)	-0.00764 (0.0197)
Girl	-0.0372*** (0.0144)	-0.266*** (0.0283)	-0.0633 (0.0394)	0.00402 (0.0321)
Private secondary school	0.101*** (0.0181)	-0.103** (0.0425)	0.139** (0.0688)	-0.0697 (0.0439)
Urban resident	0.265*** (0.0270)	0.0122 (0.0342)	0.108* (0.0560)	0.112** (0.0468)
Observations	3,705	1,162	1,020	1,020
R-squared	0.173	0.122	0.106	0.072
Sample	All	Out of school	Employed	Employed
Birth Year Dummies	Yes	Yes	Yes	Yes
Birth State Dummies	Yes	Yes	Yes	Yes
Clusters	1778	822	706	706
Mean Dep. Var.	0.630	0.578	2.821	0.396

Notes: (1) The table displays the results of the estimation of Eq. (6). (2) The dependent variables are post-secondary school outcomes: a dummy indicating enrollment in university (column 1), a dummy indicating if employed (column 2), ln of hourly wage (column 3), a dummy for being employed in a formal firm (column 4). (3) Outcomes are measured in the ENILEMS survey at ages 18 to 20 in the third quarter of 2010. ENLACE test scores come from the years 2008, 2009 and 2010. (4) Data: ENILEMS-ENLACE panel. (5) Robust standard errors are reported in parentheses. \*\*\*p<0.01, \*\*p<0.05, \*p<0.1.

Sweden (between 6% and 15%) and significantly smaller than the effects found by Lin et al. (2018) using data for the United States (around 17%). These differences can be explained by a more imperfect labor market in Mexico—compared to the one in Sweden or the United States—, differences in the population subgroups analyzed and differences in the control variables included in the specifications.<sup>22</sup>

5.5. Robustness Checks

The results shown so far could be compromised by decisions we have made in various dimensions: (i) the sample used for analysis, (ii) the fact that in some specifications we are taking a simple average of math and language test scores, (iii) the estimation method, and (iv) the attrition process observed between the 6th and 12th grades. In this section, we address these issues and show that none of them are a threat to the results:

- Sample:** For ease of comparison, all regressions presented in Table 2 are estimated in the sample of twins of the ENLACE panel. That explains why the sample size in Columns 1 to 4 in Table 2 is the same as the sample size in Columns 5 to 8, although the latter includes the twin fixed effects. Using the full ENLACE data set, the sample size in Table A.3 in the Online Appendix varies between 600,000 and 1.8 million observations. However, despite the huge differences in sample size, the results are similar qualitatively and in magnitude.
- Average test scores:** For simplicity's sake, Tables 2 and 4 present results using aggregated ENLACE test scores as the simple average of math and literacy scores. To investigate if the reported results are driven by the scores in one of the two subject areas, Tables A.4 and A.5 in the Online Appendix present the results from math and literacy test scores in separate regressions. Although there are small differences between the estimates of each subject, the general findings are similar to the ones using the aggregated ENLACE test scores.

- Binary outcomes:** The estimation of a linear probability model on binary outcomes (e.g., the probability of on-time graduation) could lead to potential biases in  $\beta_1$  because of the linear projection of test scores. To address this concern, a probit model is estimated for the binary outcomes included in Tables 2 and 4, using the same vector of independent variables as in the ordinary least squares estimation (except for the fixed effects). The results are available in Tables A.6 and A.7 in the Online Appendix and show similar finding to the ones presented in the previous sections.
- Selective attrition:** The estimates of grade 6 test scores on scores in grades 9 and 12 are subject to selective attrition. As is shown in Fig. 1, there is a considerable proportion of the population that finished primary school (grade 6) but did not finish grades 9 or 12 on time, and therefore, we cannot observe their test scores. Is this driving part of the results? As we show in the results, attrition is not random, since students with low initial (grade 6) test scores have a lower probability of graduating on time. Despite this, as shown in Panel (b) of Graph 2, there are observations on grade 9 and grade 12 test scores along the full support of the grade 6 test scores distribution. In other words, although it is more likely for students with low test scores at grade 6 to drop out or repeat a grade, it is not a deterministic process, since many students with low scores in grade 6 still graduate on time from grades 9 and 12. Furthermore, conditional on grade 6 test scores, a priori, it is more likely for students in the lower part of the test score distribution—which also corresponds to poorer households—to have unobservable characteristics that make them more prone to drop out or repeat a grade. If this is the case, then the relationship between test scores from grade 6 and grades 9 and 12 becomes steeper than the one presented in Panel (b) of Graph 2. Therefore, the results presented above can be seen as a conservative estimation of the true relationship between scores from grade 6 to grades 9 or 12. A second argument supporting the robustness of the results to selective attrition is that the relationship between grade 6 test scores and scores at grade 9 is very similar to the effects on grade 12 scores, despite the fact that attrition is much larger at grade 12. Finally, we also estimated the relationship between grade 9 and grade 12—where there is considerably less attrition vis-à-vis the level observed between grades 6 and 12—and found very similar results (see results in Table A.8 in the Online Appendix)

<sup>22</sup> The ENILEMS - ENLACE data set includes only adults from ages 18 to 20 with completed upper secondary education, while Lindqvist and Vestman, 2011 uses data for men aged 30 to 40 of all education levels; Lin et al. (2018) include men and women ages 20 to 50 of all education levels. In addition, the estimations in Lindqvist and Vestman (2011) include controls for non-cognitive skills, household income, education and marital status of the household head; Lin et al. (2018) include non-cognitive characteristics and family background as controls.

## 6. Conclusions

The purpose of this paper was to test if a census-based student assessment implemented in a large developing country captures cognitive skills. To do so, we use the Mexican census-based standardized test ENLACE to construct a longitudinal data set tracking students' education trajectories along with test scores through grades 6, 9, and 12. The analysis shows that higher test scores in grade 6 have a large and significant relationship with the student's likelihood of finishing lower and upper secondary school on time; among those who finish, grade 6 test scores are a strong predictor of secondary school test scores. Using a sample of twins to deal with differences by family background, we find that a reduction of 1 SD in test scores in sixth grade reduces by 5.5 percentage points the probability of graduating from secondary school, and, among those who graduated, it reduces their grade 12 test scores by 0.53 SD. Using variation in test scores between the two subject areas included in the test (math and literacy) to control for general (including noncognitive) skills, we present suggestive evidence that ENLACE captures the cognitive skills that it is designed to measure.

The paper also performs a second test of the validity of ENLACE by identifying the short-term relationship between grade 12 test scores and post-secondary outcomes such as university enrollment and labor market outcomes. The results show that grade 12 test scores are a strong predictor of university enrollment and hourly wages. A positive change of 1 SD in test scores at the end of upper secondary is associated with a 11-percentage point increase in the likelihood of enrolling in university, and, conditional on being employed, with a 6% increase in hourly wages.

These findings indicate that, despite their limitations, large-scale standardized tests like ENLACE capture relevant life skills. That said, appropriate design and implementation matter for the quality of large-scale student assessments. Also, as Neal (2011) warns, assigning two competing goals to one test can backfire. Notably, the objective of assessing the evolution of learning in a school system over time is in conflict with the objective of producing student learning measures to reward teachers and principals performance. The ability to compare test scores over time requires keeping similar exam questions on the test, but doing so facilitates teaching to the test, a practice that is more likely to emerge when test scores are linked to financial incentives. Setting and keeping non-competing policy objectives and coupling it with careful monitoring of test implementation are necessary conditions for reaping the benefits of standardized testing in education systems.

The results of this study also shed light on how disadvantages at early stages of education can have important and persistent implications for future education and labor market outcomes. The con-

cept that learning begets learning is corroborated by evidence presented here that suggests that a low learning outcome in sixth grade can have a negative consequence in labor incomes 10 years later. The lower performance in sixth grade works its way forward in time, signaling lower chances of completing upper secondary school. If the student graduates, the individual's reduced learning outcomes diminish the probability of starting university and, among those who work, reduce the individual's wages. This should make a failing mark at the end of primary school—or earlier—a trigger to provide strong attention and support to the student. By identifying students who need more support early in their education trajectories, large-scale student assessments can help create an education system that promotes equal opportunities and social mobility rather than one that simply replicates or even exacerbates existing inequalities.

## Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Acknowledgments

The authors would not have been able to write this paper without the support of the Mexican Secretariat of Public Education (SEP), and we are particularly grateful to Ana María Aceves and Proceso Silva for granting us access to the data without compromising confidential information. The paper benefited from comments made by the editors, Arun Agrawal and Jampel Dell'Angelo; three anonymous referees; Lucila Berniell; José Luis Gaviria; Jérémie Gignoux; José Felipe Martínez, Liliانا Sousa; Juan Vargas; Paula Villaseñor; and participants at seminars at CAF, the World Bank/IZA Conference at Universidad Javeriana, the RISE Conference at Oxford University, the RIDGE Poverty and Inequality Workshop at EAFIT Medellín, and the LACEA Meeting at BUAP Puebla. We are grateful to Elizabeth Monroy for putting together the data sets used in this paper. The findings, interpretations, and conclusions expressed in this paper are entirely those of the authors. Previous versions of this manuscript were circulated under the titles "Do Large-Scale Student Assessments Really Capture Cognitive Skills?" and "Predicting Individual Wellbeing Through Test Scores: Evidence from a National Assessment in Mexico."

## Appendix A

See [Tables A.1 and A.2](#).

**Table A.1**  
ENLACE Panel\_Means and Standard Deviations.

VARIABLES	(1) Grade 6	(2) Grade 9	(3) Grade 12	(4) Twins Grade 6	(5) Survey Grade 6
Grade 6 score	0.0209 (0.991)	0.162 (0.975)	0.398 (0.942)	0.0407 (0.994)	0.213 (1.038)
Grade 9 enrollment	0.730 (0.444)	1 (0)	1 (0)	0.765 (0.424)	0.779 (0.415)
Grade 12 enrollment	0.345 (0.475)	0.473 (0.499)	1 (0)	0.400 (0.490)	0.372 (0.483)
Girl	0.494 (0.500)	0.515 (0.500)	0.545 (0.498)	0.536 (0.499)	0.492 (0.500)
Private primary school	0.0848 (0.279)	0.0988 (0.298)	0.124 (0.330)	0.0959 (0.295)	0.109 (0.311)
Mother has lower secondary					0.587 (0.492)
Father has lower secondary					0.628 (0.483)
Mother is white collar					0.135 (0.341)
Father is white collar					0.245 (0.430)
Observations	1,878,931	1,371,437	648,018	20,252	5,677

Notes: (1) The table shows the mean and standard deviations of all students in the ENLACE panel observed in 2007 (Column 1), 2010 (Column 2) and 2013 (Column 3). Column 4 reports statistics for the sample of students identified as twins in grade 6. Column 5 displays additional variables from student and parents surveys that were applied to a sample of ENLACE takers. (2) Sample: Students who took the ENLACE exam in grade 6 in 2007. (3) Data: ENLACE panel.

**Table A.2**  
ENILEMS-ENLACE panel\_Means and Standard Deviations.

VARIABLES	(1) All	(2) University	(3) Out of university	(4) Employed
Grade 12 score	0.212 (0.854)	0.374 (0.856)	-0.0634 (0.778)	0.113 (0.808)
University student	0.630 (0.483)	1 (0)	0 (0)	0.436 (0.496)
Employed	0.379 (0.485)	0.262 (0.440)	0.578 (0.494)	1 (0)
Girl	0.565 (0.496)	0.546 (0.498)	0.596 (0.491)	0.503 (0.500)
Private secondary school	0.175 (0.380)	0.203 (0.402)	0.128 (0.335)	0.133 (0.340)
Urban resident	0.848 (0.359)	0.912 (0.283)	0.739 (0.439)	0.814 (0.390)
Age	19.18 (0.701)	19.16 (0.688)	19.21 (0.723)	19.23 (0.697)
Observations	3,718	2,551	1,167	1,385

Notes: (1) The table displays the mean and standard deviations of several characteristics of all students matched in the ENILEMS-ENLACE (Column 1), students that reported to be in college (Column 2), out of college (Column 3) and employed (Column 4). (2) Sample: Respondents to the ENILEMS survey who were matched to their ENLACE results in 2008, 2009, or 2010. (3) Data: ENILEMS-ENLACE panel.

**Appendix B. Supplementary data**

Supplementary data associated with this article can be found, in the online version, at <https://doi.org/10.1016/j.worlddev.2021.105524>.

**References**

Akyol, P., Krishna, K., & Wang, J. (2018). Taking pisa seriously: How accurate are low stakes exams? NBER Working Papers 24930, National Bureau of Economic Research Inc.

Avitabile, C., & de Hoyos, R. (2018). The heterogeneous effect of information on student performance: Evidence from a randomized control trial in Mexico. *Journal of Development Economics*, 135, 318–348.

Braun, H., & von Davier, M. (2017). The use of test scores from large-scale assessment surveys: Psychometric and statistical considerations. *Journal of Human Resources*.

CAF (2016). RED 2016. Más habilidades para el trabajo y la vida: los aportes de la familia, la escuela, el entorno y el mundo laboral. CAF, Bogotá.

Chang, E., & Padilla-Romo, M. (2019). The effects of local violent crime on high-stakes tests. Working Papers 2019-03, University of Tennessee, Department of Economics.

Cheng, X., & Gale, C. (2014). List of national learning assessments. retrieved at <https://www.epdc.org/education-data-research/list-national-learning-assessments>.

Chetty, R., Friedman, J. N., Hilger, N., Saez, E., Schanzenbach, D. W., & Yagan, D. (2011). How does your kindergarten classroom affect your earnings? Evidence from project star. *The Quarterly Journal of Economics*, 126(4), 1593–1660.

Cunha, F., Heckman, J. J., & Schennach, S. M. (2010). Estimating the technology of cognitive and noncognitive skill formation. *Econometrica*, 78(3), 883–931.

Currie, J., & Thomas, D. (2001). *Early test scores, school quality and ses: Longrun effects on wage and employment outcomes*. Research in Labor Economics.

de Hoyos, R., García-Moreno, V. A., & Patrinos, H. A. (2017). The impact of an accountability intervention with diagnostic feedback: Evidence from Mexico. *Economics of Education Review*, 58, 123–140.

Duckworth, A. L., & Seligman, M. E. (2005). Self-discipline outdoes iq in predicting academic performance of adolescents. *Psychological Science*, 16(12), 939–944. PMID:16313657.

Dustan, A., de Janvry, A., & Sadoulet, E. (2017). Flourish or fail?: The risky reward of elite high school admission in Mexico city. *Journal of Human Resources*, 52(3), 756–799.

- Estrada, R. (2019). Rules versus discretion in public service: Teacher hiring in Mexico. *Journal of Labor Economics*, 37(2), 545–579.
- Estrada, R., & Gignoux, J. (2017). Benefits to elite schools and the expected returns to education: Evidence from Mexico City. *European Economic Review*, 95, 168–194.
- Figlio, D., & Loeb, S. (2011). School accountability. In E. Hanushek, S. Machin, L. Woessmann (Eds.), *Handbook of the economics of education* (Vol. 3, chapter 08, pp. 383–421). Elsevier, 1 edition.
- Finn, B. (2015). Measuring motivation in low-stakes assessments. *ETS Research Report Series*, 2015(2), 1–17.
- Geary, D., Nicholas, A., Li, Y., & Sun, J. (2017). Developmental change in the influence of domain-general abilities and domain-specific knowledge on mathematics achievement: An eight-year longitudinal study. *Journal of Educational Psychology*, 109(5), 680–693.
- Gneezy, U., List, J. A., Livingston, J. A., Qin, X., Sadoff, S., & Xu, Y. (2019). Measuring success in education: The role of effort on the test itself. *American Economic Review: Insights*, 1(3), 291–308.
- Heckman, J. J., & Kautz, T. (2012). Hard evidence on soft skills. *Labour Economics*, 19(4), 451–464. European Association of Labour Economists 23rd annual conference, Paphos, Cyprus, 22–24th September 2011.
- Heckman, J. J., Stixrud, J., & Urzua, S. (2006). The effects of cognitive and noncognitive abilities on labor market outcomes and social behavior. *Journal of Labor Economics*, 24(3), 411–482.
- Koretz, D. (2017). *The Testing Charade: Pretending to Make Schools Better*. University of Chicago Press.
- Koretz, D. M., & Barron, S. I. (1998). The Validity of Gains in Scores on the Kentucky Instructional Results Information System (KIRIS). ERIC.
- Laajaj, R., & Macours, K. (2019). Measuring skills in developing countries. *Journal of Human Resources*.
- Lavy, V., Ebenstein, A., & Roth, S. (2014). *The impact of short term exposure to ambient air pollution on cognitive performance and human capital formation*. NBER Working Papers 20648. National Bureau of Economic Research Inc.
- Lin, M., Bumgarner, E., & Chatterji, M. (2014). Understanding validity issues in international large scale assessments. *Quality Assurance in Education*.
- Lindqvist, E., & Vestman, R. (2011). The labor market returns to cognitive and noncognitive ability: Evidence from the Swedish Enlistment. *American Economic Journal: Applied Economics*, 3(1), 101–128.
- Lin, D., Lutter, R., & Ruhm, C. J. (2018). Cognitive performance and labour market outcomes. *Labour Economics*, 51, 121–135.
- Murnane, R. J., Willett, J. B., & Levy, F. (1995). The growing importance of cognitive skills in wage determination. *The Review of Economics and Statistics*, 77(2), 251–266.
- Neal, D. (2011). The design of performance pay in education. In E. Hanushek, S. Machin, L. Woessmann (Eds.), *Handbook of the economics of education* (Vol. 4, chapter 06, pp. 495–550). Elsevier, 1 edition.
- OECD (2001). *Knowledge and Skills for Life: First Results from PISA 2000*. OECD Publishing, Paris.
- Ritchie, S., Bates, T., & Deary, I. (2015). Is education associated with improvements in general cognitive ability, or in specific skills? *Developmental Psychology*, 51, 573–582.
- Rose, H. (2006). Do gains in test scores explain labor market outcomes?. *Economics of Education Review*, 25(4), 430–446.
- Salardi, P., & Michaelsen, M. (2019). Violence, psychological stress and educational performance during the war on drugs in Mexico. *Journal of Development Economics*, 102387.
- World Bank (2018). *World Development Report 2018: Learning to Realize Education's Promise*. Washington, DC: World Bank.